

Springer Optimization and Its Applications 173

Themistocles M. Rassias *Editor*

# Nonlinear Analysis, Differential Equations, and Applications



Springer

# Springer Optimization and Its Applications

Volume 173

## Series Editors

Panos M. Pardalos , University of Florida

My T. Thai , University of Florida

## Honorary Editor

Ding-Zhu Du, University of Texas at Dallas

## Advisory Editors

Roman V. Belavkin, Middlesex University

John R. Birge, University of Chicago

Sergiy Butenko, Texas A&M University

Vipin Kumar, University of Minnesota

Anna Nagurney, University of Massachusetts Amherst

Jun Pei, Hefei University of Technology

Oleg Prokopyev, University of Pittsburgh

Steffen Rebennack, Karlsruhe Institute of Technology

Mauricio Resende, Amazon

Tamás Terlaky, Lehigh University

Van Vu, Yale University

Michael N. Vrahatis, University of Patras

Guoliang Xue, Arizona State University

Yinyu Ye, Stanford University

## **Aims and Scope**

Optimization has continued to expand in all directions at an astonishing rate. New algorithmic and theoretical techniques are continually developing and the diffusion into other disciplines is proceeding at a rapid pace, with a spot light on machine learning, artificial intelligence, and quantum computing. Our knowledge of all aspects of the field has grown even more profound. At the same time, one of the most striking trends in optimization is the constantly increasing emphasis on the interdisciplinary nature of the field. Optimization has been a basic tool in areas not limited to applied mathematics, engineering, medicine, economics, computer science, operations research, and other sciences.

The series **Springer Optimization and Its Applications (SOIA)** aims to publish state-of-the-art expository works (monographs, contributed volumes, textbooks, handbooks) that focus on theory, methods, and applications of optimization. Topics covered include, but are not limited to, nonlinear optimization, combinatorial optimization, continuous optimization, stochastic optimization, Bayesian optimization, optimal control, discrete optimization, multi-objective optimization, and more. New to the series portfolio include Works at the intersection of optimization and machine learning, artificial intelligence, and quantum computing.

*Volumes from this series are indexed by Web of Science, zbMATH, Mathematical Reviews, and SCOPUS.*

More information about this series at <http://www.springer.com/series/7393>

Themistocles M. Rassias  
Editor

# Nonlinear Analysis, Differential Equations, and Applications

 Springer



*Editor*

Themistocles M. Rassias  
Department of Mathematics, Zografou  
Campus  
National Technical University of Athens  
Athens, Greece

ISSN 1931-6828                      ISSN 1931-6836 (electronic)  
Springer Optimization and Its Applications  
ISBN 978-3-030-72562-4              ISBN 978-3-030-72563-1 (eBook)  
<https://doi.org/10.1007/978-3-030-72563-1>

Mathematics Subject Classification: 26-XX, 30-XX, 34-XX, 35-XX, 47-XX

© Springer Nature Switzerland AG 2021

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors, and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Switzerland AG  
The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

# Preface

*Nonlinear Analysis, Differential Equations, and Applications* publishes research and research-survey papers devoted to a broad variety of topics on functional equations, ordinary differential equations, partial differential equations, stochastic differential equations, optimization theory, network games, generalized Nash equilibria, critical point theory, calculus of variations, nonlinear functional analysis, convex analysis, variational inequalities, topology, global differential geometry, curvature flows, perturbation theory, numerical analysis, mathematical finance, and a variety of applications in interdisciplinary topics. More specifically, the book chapters of this volume investigate compound superquadratic functions, the Hyers–Ulam stability of functional equations, edge degenerate pseudo-hyperbolic equations, Kirchhoff wave equation, BMO norms of operators on differential forms, equilibrium points of the perturbed R3BP, complex zeros of solutions to second order differential equations, a higher-order Ginzburg–Landau-type equation, multi-symplectic numerical schemes for differential equations, the Erdős–Rényi network model, strongly  $m$ -convex functions, higher order strongly generalized convex functions, factorization and solution of second order differential equations, generalized topologically open sets in relator spaces, graphical mean curvature flow, critical point theory in infinite dimensional spaces using the Leray–Schauder index, non-radial solutions of a supercritical equation in expanding domains, the semi-discrete method for the approximation of the solution of stochastic differential equations, homotopic metric-interval L-contractions in gauge spaces, Rhoades contractions theory, network centrality measures, the Radon transform in three space dimensions via plane integration and applications in positron emission tomography boundary perturbations on medical monitoring and imaging techniques, the KdV-B equation and biomedical applications.

We would like to express our warmest thanks to all the authors who contributed their valuable works for publication in this volume. We would also like to express our sincere appreciation to the staff of Springer for their help throughout the preparation of this book.

Athens, Greece

Themistocles M. Rassias

# Contents

|  |     |
|--|-----|
| <b>On Compound Superquadratic Functions</b> .....  | 1   |
| Shoshana Abramovich  |     |
| <b>Best Hyers–Ulam Stability Constants on a Time Scale with Discrete Core and Continuous Periphery</b> .....   | 17  |
| Douglas R. Anderson and Masakazu Onitsuka  |     |
| <b>Invariance Solutions and Blow-Up Property for Edge Degenerate Pseudo-Hyperbolic Equations in Edge Sobolev Spaces</b> .....                          | 39  |
| Carlo Cattani and Morteza Koozehgar Kalleji  |     |
| <b><math>\phi^4</math> Solitons in Kirchhoff Wave Equation</b> .....   | 71  |
| Y. Contoyiannis, P. Papadopoulos, M. Kampitakis, S. M. Potirakis, and N. L. Matiadou   |     |
| <b>Estimates for Lipschitz and BMO Norms of Operators on Differential Forms</b> .....  | 81  |
| Shusen Ding, Guannan Shi, and Yuming Xing  |     |
| <b>Application of Boundary Perturbations on Medical Monitoring and Imaging Techniques</b> .....  | 101 |
| M. Doschoris, A. Papargiri, V. S. Kalantonis, and P. Vafeas  |     |
| <b>Poynting–Robertson and Oblateness Effects on the Equilibrium Points of the Perturbed R3BP: Application on Cen X-4 Binary System</b> ...             | 131 |
| Aguda Ekele Vincent and Angela E. Perdiou  |     |
| <b>Localization and Perturbation of Complex Zeros of Solutions to Second Order Differential Equations with Polynomial Coefficients. A Survey</b> ..... | 149 |
| Michael Gil’   |     |
| <b>Dynamics of a Higher-Order Ginzburg–Landau-Type Equation</b> .....  | 187 |
| Theodoros P. Horikis, Nikos I. Karachalios, and Dimitrios J. Frantzeskakis   |     |

|   |     |
|---|-----|
| <b>The Role of Differential Equations in Applied Statistics</b> .....   | 209 |
| Christos P. Kitsos and C. S. A. Nisiotis  |     |
| <b>Geometric Derivation and Analysis of Multi-Symplectic Numerical Schemes for Differential Equations</b> .....   | 231 |
| Odysseas Kosmas, Dimitrios Papadopoulos, and Dimitrios Vlachos  |     |
| <b>Non-radial Solutions of a Supercritical Equation in Expanding Domains: The Limit Case</b> .....  | 253 |
| Nikos Labropoulos   |     |
| <b>Financial Contagion in Interbank Networks: The Case of Erdős–Rényi Network Model</b> .....   | 277 |
| K. Loukaki, P. Boufounou, and J. Leventides   |     |
| <b>Higher Order Strongly <math>m</math>-convex Functions</b> .....  | 319 |
| Muhammad Aslam Noor and Khalida Inayat Noor   |     |
| <b>Characterizations of Higher Order Strongly Generalized Convex Functions</b> .....  | 341 |
| Muhammad Aslam Noor, Khalida Inayat Noor, and Michael Th. Rassias   |     |
| <b>A Note on Generalized Nash Games Played on Networks</b> .....  | 365 |
| Mauro Passacantando and Fabio Raciti  |     |
| <b>Piecewise Polynomial Inversion of the Radon Transform in Three Space Dimensions via Plane Integration and Applications in Positron Emission Tomography</b> ..... | 381 |
| Nicholas E. Protonotarios, George A. Kastis, Nikolaos Dikaios, and Athanassios S. Fokas   |     |
| <b>Factorization and Solution of Linear and Nonlinear Second Order Differential Equations with Variable Coefficients and Mixed Conditions</b> .....                 | 397 |
| E. Providas   |     |
| <b>A General Framework for Studying Certain Generalized Topologically Open Sets in Relator Spaces</b> .....   | 415 |
| Themistocles M. Rassias and Árpád Szász   |     |
| <b>Graphical Mean Curvature Flow</b> .....  | 493 |
| Andreas Savas-Halilaj   |     |
| <b>Critical Point Theory in Infinite Dimensional Spaces Using the Leray–Schauder Index</b> .....  | 579 |
| Martin Schechter  |     |
| <b>Canonical Systems of Partial Differential Equations</b> .....  | 609 |
| Martin Schechter  |     |

**The Semi-discrete Method for the Approximation of the Solution of Stochastic Differential Equations** ..... 625  
Ioannis S. Stamatiou

**Homotopic Metric-Interval L-Contractions in Gauge Spaces** ..... 639  
Mihai Turinici

**Analytic Methods in Rhoades Contractions Theory** ..... 705  
Mihai Turinici

**Nonlinear Dynamics of the KdV-B Equation and Its Biomedical Applications**..... 765  
Michail A. Xenos and Anastasios C. Felias

# On Compound Superquadratic Functions



Shoshana Abramovich

**Abstract** By using compound superquadratic functions, the inequalities presented in this paper refine and extend Jensen and Jensen-Steffensen inequalities.

**2010 Mathematics Subject Classification:** 26D15, 26A51, 47A63, 47A64.

## 1 Introduction

By using compound superquadratic functions, the inequalities presented in this paper refine and extend Jensen and Jensen-Steffensen inequalities.

We start with some definitions, notations and lemmas that are used in this paper.

**Definition 1** [3, Definition 2.1] A function  $f : [0, \infty) \rightarrow \mathbb{R}$  is superquadratic provided that for all  $x \in [0, \infty)$  there exists a constant  $C_f(x) \in \mathbb{R}$  such that the inequality

$$f(y) - f(x) - C_f(x)(y - x) - f(|y - x|) \geq 0 \quad (1)$$

holds for all  $y \in [0, \infty)$ .

$f$  is called subquadratic if  $-f$  is superquadratic.

**Corollary 1** *The functions  $f(x) = x^p$ ,  $x \geq 0$  are superquadratic for  $p \geq 2$ , subquadratic for  $0 < p \leq 2$ , for which  $C_f(x) = px^{p-1} = f'(x)$ . When  $p = 2$  (1) is an equality.*

**Lemma 1** ([3, Lemma 2.1, Lemma 3.2]) *Let  $f$  be a superquadratic function with  $C_f(x)$  as in Definition 1.*

---

S. Abramovich (✉)

Department of Mathematics, University of Haifa, Haifa, Israel

e-mail: [abramos@math.haifa.ac.il](mailto:abramos@math.haifa.ac.il)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,

Springer Optimization and Its Applications 173,

[https://doi.org/10.1007/978-3-030-72563-1\\_1](https://doi.org/10.1007/978-3-030-72563-1_1)

Then:

- (i)  $f(0) \leq 0$ ,
- (ii) if  $f(0) = f'(0) = 0$  then  $C_f(x) = f'(x)$  whenever  $f$  is differentiable at  $x > 0$ , and  $\frac{f(x)}{x^2}$  is non-decreasing.
- (iii) if  $f \geq 0$ , then  $f$  is convex and  $f(0) = f'(0) = 0$ .

**Lemma 2 ([3, Lemma 3.1])** Suppose that  $\varphi : [0, b) \rightarrow \mathbb{R}$  is continuously differentiable and  $\varphi(0) \leq 0$ . If  $\varphi'$  is superadditive or  $\frac{\varphi'(x)}{x}$  is non-decreasing, then  $\varphi$  is superquadratic.

**Corollary 2 ([3])** Suppose that  $f$  is superquadratic. Let  $\xi_i \geq 0$ ,  $i = 1, \dots, m$ , and let  $\bar{\xi} = \sum_{i=1}^m p_i \xi_i$  where  $p_i \geq 0$ ,  $i = 1, \dots, m$ , and  $\sum_{i=1}^m p_i = 1$ . Then

$$\sum_{i=1}^m p_i f(\xi_i) - f(\bar{\xi}) \geq \sum_{i=1}^m p_i f(|\xi_i - \bar{\xi}|),$$

and in the special case that  $m = 2$ ,  $0 \leq t \leq 1$  and  $0 \leq a < b < \infty$

$$\begin{aligned} & (1-t)f(a) + tf(b) \\ & \geq f((1-t)a + tb) + (1-t)f(t(b-a)) + tf((1-t)(b-a)) \end{aligned}$$

hold.

**Definition 2 ([4])** The real  $n$ -tuple that satisfies

$$\begin{aligned} 0 \leq P_j \leq P_n, \quad j = 1, \dots, n, \quad P_n = 1, \quad (2) \\ P_j = \sum_{i=1}^j \rho_i, \quad \bar{P}_j = \sum_{i=j}^n \rho_i, \quad j = 1, \dots, n. \end{aligned}$$

are called Steffensen's coefficients.

On Jensen-Steffensen's inequality which extends Jensen's inequality for not necessarily positive coefficients see for example [7, page 57]. It states that:

**Theorem 1** Let  $\varphi : I \rightarrow \mathbb{R}$  be convex, then

$$\varphi \left( \sum_{i=1}^n \rho_i \zeta_i \right) \leq \sum_{i=1}^n \rho_i \varphi(\zeta_i)$$

holds, where  $I$  is an interval in  $\mathbb{R}$ ,  $\boldsymbol{\zeta} = (\zeta_1, \dots, \zeta_n)$  is any monotonic  $n$ -tuple in  $I^n$  and  $\boldsymbol{\rho} = (\rho_1, \dots, \rho_n)$  is a real  $n$ -tuple that satisfies (2).

After the introduction we obtain in Section 2 refinements of Jensen and Jensen-Steffensen's inequalities via compound superquadratic functions.

On compound functions related to superquadracity see [1] and [2] where the properties of

$$\varphi^{-1} \left( \varphi \left( \sum_{i=1}^n t_i a_i \right) + \sum_{i=1}^n t_i \varphi \left( \left| a_i - \sum_{j=1}^n t_j a_j \right| \right) \right)$$

when  $\sum_{i=1}^n t_i = 1$ ,  $t_i \geq 0$ ,  $a_i \geq 0$ ,  $i = 1, \dots, n$ , are investigated.

In Section 3 similar to the discussion in [1] and [2] we deal briefly with the behavior of

$$\varphi^{-1} \left( \sum_{i=1}^n t_i \varphi(a_i) - \sum_{i=1}^n t_i \varphi \left( \left| a_i - \sum_{j=1}^n t_j a_j \right| \right) \right)$$

where  $\sum_{i=1}^n t_i = 1$ ,  $t_i \geq 0$ ,  $a_i \geq 0$ ,  $i = 1, \dots, n$ .

On compound functions related to convexity see for instance [6, Theorems 3 and 4], [5], and [7] and their references.

## 2 Refinements of Jensen's Inequality via Compound Functions

In this section we deal with compound superquadratic functions. These functions lead to new inequalities that refine Jensen's inequality and Jensen-Steffensen's inequality.

We start with a theorem which guaranties that the compound function  $f = g \circ \varphi$  is non-negative superquadratic.

**Theorem 2** *Let  $g : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  be twice differentiable, non-decreasing, convex function and  $g(0) = 0$ . Let  $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  be twice differentiable, superquadratic function, and let  $f$  be the compound function defined as  $f = g \circ \varphi$ . Then, the function  $f : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  is convex, superquadratic and satisfies  $\left( \frac{f'(x)}{x} \right)' \geq 0$  when  $\left( \frac{\varphi'(x)}{x} \right)' \geq 0$ , or when  $g(x) = x^m$ ,  $m \in [2, \infty)$ .*

**Proof** The function  $f$  satisfies:

$$A : f'(x) = \varphi'(x) g'(\varphi(x)) \geq 0, \text{ because } \varphi' \text{ and } g' \geq 0.$$

$$B : f''(x) = \varphi''(x) g'(\varphi(x)) + \left( \varphi'(x) \right)^2 g''(\varphi(x)) \geq 0,$$

$$C : \left( \frac{f'(x)}{x} \right)' = g'(\varphi(x)) \left( \frac{\varphi'(x)}{x} \right)' + \frac{1}{x} \left( \varphi'(x) \right)^2 g''(\varphi(x)) \geq 0, \text{ and when } g(x) = x^m, m \geq 2$$



$$\begin{aligned}
D &: \left( \frac{f'(x)}{x} \right)' \\
&= \frac{mx\varphi^{m-1}(x)\varphi''(x)}{x^2} + \frac{m\varphi^{m-2}(x)\varphi'(x)\left((m-1)x\varphi'(x) - \varphi(x)\right)}{x^2} \\
&\geq \frac{m\varphi^{m-2}(x)\varphi'(x)\left((m-1)x\varphi'(x) - \varphi(x)\right)}{x^2} \geq 0
\end{aligned}$$

Indeed, when  $\left(\frac{\varphi'(x)}{x}\right)' \geq 0$ , and  $\varphi$  is non-negative superquadratic, then by Lemma 1(iii) it is also convex and  $\varphi(0) = \varphi'(0) = 0$ . The function  $g$  is convex, increasing, and non-negative, therefore according to the computations of A and B,  $f$  is convex, increasing,  $f(0) = f'(0) = 0$ , and by C,  $\left(\frac{f'(x)}{x}\right)' \geq 0$  holds. Therefore according to Lemma 2  $f$  is also superquadratic.

When  $m \geq 2$  the last inequality in D follows from  $\left(x\varphi'(x) - \varphi(x)\right)' = x\varphi''(x) \geq 0$ , and  $\varphi(0) = \varphi'(0) = 0$ . Hence  $\left(\frac{f'(x)}{x}\right)' \geq 0$ , and together with  $f(0) = f'(0) = 0$ , by Lemma 2,  $f$  is superquadratic too. The proof is complete.

We emphasize that there are non-negative superquadratic functions  $\varphi$  which do not satisfy the condition  $\left(\frac{\varphi'(x)}{x}\right)' \geq 0$  as shown in [3, Example 3.3].

The proofs in the sequel do not need this condition, but only the conditions that  $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  is strictly increasing and superquadratic, that  $g : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  is convex and  $g(0) = 0$  and  $f = g \circ \varphi$  is non-negative superquadratic or in some cases only convex.

Obviously, all the results in this section hold when the functions  $\varphi$ ,  $f$  and  $g$  satisfy the conditions of Theorem 2.

In the following theorem we present refinements of Jensen's inequality for compound functions  $f$  where  $f \circ \varphi^{-1}$  are convex. In [5] these functions are called **composite  $\varphi^{-1}$  convex functions**.

**Theorem 3** *Let  $g : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  be twice differentiable convex function and  $g(0) = 0$ . Let  $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  be twice differentiable, strictly increasing superquadratic function, and  $f = g \circ \varphi$ . Then, the following refinements of Jensen's inequality hold when  $t_i \geq 0$ ,  $a_i \geq 0$ ,  $i = 1, \dots, n$ ,  $\sum_{i=1}^n t_i = 1$  and  $\bar{a} = \sum_{i=1}^n t_i a_i$ :*

$$f\left(\sum_{i=1}^n t_i a_i\right) \leq f \circ \varphi^{-1}\left(\sum_{i=1}^n t_i \varphi(a_i) - \sum_{i=1}^n t_i \varphi(|a_i - \bar{a}|)\right) \quad (3)$$

$$\begin{aligned}
 &\leq f \circ \varphi^{-1} \left( \sum_{i=1}^n t_i \varphi(a_i) \right) - f \circ \varphi^{-1} \left( \sum_{i=1}^n t_i \varphi(|a_i - \bar{a}|) \right) \\
 &\leq \sum_{i=1}^n t_i f(a_i) - f \circ \varphi^{-1} \left( \sum_{i=1}^n t_i \varphi(|a_i - \bar{a}|) \right) \\
 &\leq \sum_{i=1}^n t_i f(a_i) - f \left( \sum_{i=1}^n t_i (|a_i - \bar{a}|) \right) \leq \sum_{i=1}^n t_i f(a_i).
 \end{aligned}$$

If  $f$  is also superquadratic then

$$\begin{aligned}
 f \left( \sum_{i=1}^n t_i a_i \right) &\leq \sum_{i=1}^n t_i f(a_i) - \sum_{i=1}^n t_i f(|a_i - \bar{a}|) \\
 &\leq \sum_{i=1}^n t_i f(a_i).
 \end{aligned} \tag{4}$$

In particular, when  $\frac{\varphi'(x)}{x}$  is increasing or when  $f \circ \varphi^{-1}(x) = x^m$ ,  $m \geq 2$  the inequalities (3) and (4) hold.

**Proof** The function  $\varphi$  is strictly increasing and therefore  $\varphi^{-1}$  exists.

To prove (3) we rewrite  $f$  as  $f = g \circ \varphi$  and get

$$\begin{aligned}
 f \left( \sum_{i=1}^n t_i a_i \right) &= g \circ \varphi \left( \sum_{i=1}^n t_i a_i \right) \\
 &\leq g \left( \sum_{i=1}^n t_i \varphi(a_i) - \sum_{i=1}^n t_i \varphi(|a_i - \bar{a}|) \right) \\
 &\leq g \left( \sum_{i=1}^n t_i \varphi(a_i) \right) - g \left( \sum_{i=1}^n t_i \varphi(|a_i - \bar{a}|) \right) \\
 &= f \circ \varphi^{-1} \left( \sum_{i=1}^n t_i \varphi(a_i) \right) - g \left( \sum_{i=1}^n t_i \varphi(|a_i - \bar{a}|) \right) \\
 &\leq \sum_{i=1}^n t_i f(a_i) - g \left( \sum_{i=1}^n t_i \varphi(|a_i - \bar{a}|) \right) \\
 &\leq \sum_{i=1}^n t_i f(a_i) - f \left( \sum_{i=1}^n t_i (|a_i - \bar{a}|) \right)
 \end{aligned} \tag{5}$$

$$\leq \sum_{i=1}^n t_i f(a_i).$$

Indeed, the first inequality in (5) is due to the monotonicity of  $g$  and the superquadracity of  $\varphi$ . The second inequality follows from the convexity of  $g$  and from  $g(0) = 0$ . The third inequality is because the inequality  $f \circ \varphi^{-1} \left( \sum_{i=1}^n t_i \varphi(a_i) \right) \leq \sum_{i=1}^n t_i f(a_i)$  which is the same as  $g \left( \sum_{i=1}^n t_i \varphi(a_i) \right) \leq \sum_{i=1}^n t_i g \circ \varphi(a_i)$  is deduced from the convexity of  $g$ . The fourth inequality follows from the convexity of  $\varphi$ , the monotonicity of  $g$  and because  $f = g \circ \varphi$ . The last inequality in (5) is due to positivity of  $f$ . Inequality (3) is proved because it is just a variation of (5). Inequalities (4) follow directly from the superquadracity of  $f$ , Corollary 2 and because of the positivity and convexity of  $f = g \circ \varphi$  which is derived directly from the conditions on  $\varphi$  and  $g$  (see A and B in Theorem 2). Under the conditions of Theorem 2  $f$  is superquadratic, hence (3) and (4) hold. The proof is complete.

**Corollary 3** *It is evident that we can extend Theorem 3 as follows: Let  $g_k : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ , be twice differentiable non-decreasing convex function and  $g_k(0) = 0$ ,  $k = 1, \dots, m$ . Let  $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  be twice differentiable, strictly increasing superquadratic function. Let  $f$  be the compound function  $f = g_m \circ g_{m-1} \circ \dots \circ g_1 \circ \varphi$ . Then, the following refinements of Jensen's inequality hold when  $k = 1, \dots, m$ ,  $t_i \geq 0$ ,  $a_i \geq 0$ ,  $i = 1, \dots, n$ ,  $\sum_{i=1}^n t_i = 1$  and  $\bar{a} = \sum_{i=1}^n t_i a_i$  :*

$$\begin{aligned} & f \left( \sum_{i=1}^n t_i a_i \right) \\ &= g_m \circ g_{m-1} \circ \dots \circ g_1 \circ \varphi \left( \sum_{i=1}^n t_i a_i \right) \\ &\leq g_m \circ g_{m-1} \circ \dots \circ g_1 \left( \sum_{i=1}^n t_i \varphi(a_i) - \sum_{i=1}^n t_i \varphi(|a_i - \bar{a}|) \right) \\ &\leq g_m \circ g_{m-1} \circ \dots \circ g_k \left( \sum_{i=1}^n t_i g_{k-1} \circ g_{k-2} \circ \dots \circ g_1 \varphi(a_i) \right. \\ &\quad \left. - g_{k-1} \circ g_{k-2} \circ \dots \circ g_1 \left( \sum_{i=1}^n t_i \varphi(|a_i - \bar{a}|) \right) \right) \\ &\leq \sum_{i=1}^n t_i g_m \circ g_{m-1} \circ \dots \circ g_1 \circ \varphi(a_i) \end{aligned}$$

$$\begin{aligned}
 & -g_m \circ g_{m-1} \circ \dots \circ g_1 \circ \varphi \left( \sum_{i=1}^n t_i (|a_i - \bar{a}|) \right) \\
 &= \sum_{i=1}^n t_i f(a_i) - f \left( \sum_{i=1}^n t_i (|a_i - \bar{a}|) \right) \\
 &\leq \sum_{i=1}^n t_i f(a_i).
 \end{aligned}$$

**Corollary 4** *Under the same conditions as in Theorem 3, if  $n = 2$ ,  $t_1 = t_2 = \frac{1}{2}$ , and  $a_1 = a < a_n = b$ , the inequalities in (3) and (4) translate into*

$$\begin{aligned}
 & f \left( \frac{a+b}{2} \right) \leq g \left( \frac{\varphi(a) + \varphi(b)}{2} - \varphi \left( \left| \frac{b-a}{2} \right| \right) \right) \tag{6} \\
 & \leq g \left( \frac{\varphi(a) + \varphi(b)}{2} \right) - g \left( \varphi \left( \left| \frac{b-a}{2} \right| \right) \right) \\
 & = f \left( \varphi^{-1} \left( \frac{\varphi(a) + \varphi(b)}{2} \right) \right) - f \left( \left| \frac{b-a}{2} \right| \right) \\
 & \leq \frac{f(a) + f(b)}{2} - f \left( \left| \frac{b-a}{2} \right| \right) \leq \frac{f(a) + f(b)}{2}.
 \end{aligned}$$

*Therefore, if  $\varphi$  is strictly convex or if  $g$  is strictly increasing, the inequalities in (6) are strict too. Hence, for  $t$  close enough to  $\frac{1}{2}$  the inequalities in*

$$\begin{aligned}
 & g \circ \varphi((1-t)a + tb) \\
 & \leq g((1-t)\varphi(a) + t\varphi(b) - t\varphi((1-t)|b-a|) - (1-t)\varphi(t|b-a|)) \\
 & \leq (1-t)g \circ \varphi(a) + tg \circ \varphi(b) \\
 & \quad -tg \circ \varphi((1-t)|b-a|) - (1-t)g \circ \varphi(t|b-a|) \\
 & \leq (1-t)g \circ \varphi(a) + tg \circ \varphi(b)
 \end{aligned}$$

*hold too.*

Corollary 4 is the motivation for proving Theorem 4, in which we deal with superquadratic functions  $\varphi$  and functions  $f$  such that the functions  $f \circ \varphi^{-1}$  are convex on  $x \geq 0$  and satisfy

$$\varphi^{-1} \left( \sum_{i=1}^2 t_i \varphi(a_i) - \sum_{i=1}^2 t_i \varphi \left( \left| a_i - \sum_{j=1}^2 t_j a_j \right| \right) \right)$$

$$\leq f^{-1} \left( \sum_{i=1}^2 t_i f(a_i) - \sum_{i=1}^2 t_i f \left( \left| a_i - \sum_{j=1}^2 t_j a_j \right| \right) \right)$$

when  $\sum_{i=1}^2 t_i = 1$ ,  $t_i \geq 0$ ,  $a_i \geq 0$ ,  $i = 1, 2$ .

**Theorem 4** Let  $g : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  be twice differentiable non-decreasing convex function and  $g(0) = 0$ . Let  $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  be twice differentiable, strictly increasing superquadratic function. Let  $f$  defined as  $f = g \circ \varphi$  be a superquadratic function. If

$$\begin{aligned} & g((1-t)\varphi(a) + t\varphi(b)) \\ & \leq (1-t)g \circ \varphi(a) + tg \circ \varphi(b) - t(1-t)g \circ \varphi(b-a) \end{aligned} \quad (7)$$

when  $0 \leq t \leq 1$ ,  $0 \leq a < b < \infty$  is satisfied, then the inequalities in

$$\begin{aligned} & f((1-t)a + tb) = g \circ \varphi((1-t)a + tb) \\ & \leq g((1-t)\varphi(a) + t\varphi(b) - t\varphi((1-t)(b-a)) - (1-t)\varphi(t(b-a))) \\ & \leq (1-t)g \circ \varphi(a) + tg \circ \varphi(b) \\ & \quad - t g \circ \varphi((1-t)(b-a)) - (1-t)g \circ \varphi(t(b-a)) \\ & \leq (1-t)g \circ \varphi(a) + tg \circ \varphi(b) = (1-t)f(a) + tf(b) \end{aligned} \quad (8)$$

hold. In particular, if also  $\frac{\varphi'(x)}{x}$  is increasing, or if  $g(x) = x^m$ ,  $m \geq 2$  and (7) holds then (8) holds.

**Proof** The function  $\varphi$  is non-negative and superquadratic, therefore it is convex and the inequalities in

$$\begin{aligned} & 0 \leq \varphi((1-t)a + tb) \\ & \leq (1-t)\varphi(a) + t\varphi(b) - t\varphi((1-t)|b-a|) - (1-t)\varphi(t|b-a|) \\ & \leq (1-t)\varphi(a) + t\varphi(b) \end{aligned} \quad (9)$$

hold. The function  $g$  is non-negative and increasing too and therefore from (9) we get

$$\begin{aligned} & 0 \leq g((1-t)\varphi(a) + t\varphi(b) - t\varphi((1-t)|b-a|) - (1-t)\varphi(t|b-a|)) \\ & \leq g((1-t)\varphi(a) + t\varphi(b)). \end{aligned} \quad (10)$$

The function  $f = g \circ \varphi$  is also superquadratic and non-negative and  $f(0) = f'(0) = 0$ . According to Lemma 1(ii)  $\frac{f(kx)}{x^2}$  is increasing when  $k > 0$ ,  $x \geq 0$  therefore the inequalities in

$$\begin{aligned}
0 &\leq t f((1-t)(b-a)) + (1-t) f(t(b-a)) \\
&= t(1-t)^2 \frac{f((1-t)(b-a))}{(1-t)^2} + (1-t)t^2 \frac{f(t(b-a))}{t^2} \\
&\leq t(1-t)^2 f(b-a) + (1-t)t^2 f(b-a) \\
&= t(1-t) f(b-a)
\end{aligned} \tag{11}$$

hold. From (9), (10) and (11), and (7), the inequalities

$$\begin{aligned}
&(1-t)g \circ \varphi(a) + tg \circ \varphi(b) \\
&-tg \circ \varphi((1-t)(b-a)) - (1-t)g \circ \varphi(t(b-a)) \\
&-g((1-t)\varphi(a) + t\varphi(b) - t\varphi((1-t)(b-a)) - (1-t)\varphi(t(b-a))) \\
&\geq (1-t)g \circ \varphi(a) + tg \circ \varphi(b) \\
&-t(1-t)g \circ \varphi(b-a) - g((1-t)\varphi(a) + t\varphi(b)) \\
&\geq 0
\end{aligned}$$

hold from which (8) follows. In the special case that also  $\frac{\varphi'(x)}{x}$  is increasing, or  $g(x) = x^m$ ,  $m \geq 2$ , then according to Theorem 2 the function  $f$  is superquadratic and therefore (8) holds. The proof of the theorem is complete.

**Lemma 3** *Let  $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  be an increasing function and let  $0 \leq a < b < \infty$ . Then for  $m = 2, 3, \dots$ , and  $0 \leq t \leq 1$  we get that*

$$\begin{aligned}
&(1-t)\varphi^m(a) + t\varphi^m(b) - ((1-t)\varphi(a) + t\varphi(b))^m \\
&= (\varphi(b) - \varphi(a))^m t \left(1 - t^{m-1}\right) + \sum_{k=2}^{m-1} \varphi^{m-k}(a) (\varphi(b) - \varphi(a))^k \left(t - t^k\right) \\
&= t(1-t) (\varphi(b) - \varphi(a))^m \sum_{k=0}^{m-2} t^k + \sum_{k=2}^{m-1} \varphi^{m-k}(a) (\varphi(b) - \varphi(a))^k \sum_{j=0}^{k-2} t^j \\
&\geq t(1-t) (\varphi(b) - \varphi(a))^m.
\end{aligned} \tag{12}$$

In the special cases that  $m = 2$  and  $m = 3$

$$\begin{aligned}
&(1-t)\varphi^2(a) + t\varphi^2(b) - ((1-t)\varphi(a) + t\varphi(b))^2 \\
&= t(1-t) (\varphi(b) - \varphi(a))^2
\end{aligned}$$

and

$$(1-t)\varphi^3(a) + t\varphi^3(b) - ((1-t)\varphi(a) + t\varphi(b))^3$$

$$\begin{aligned}
&= t(1-t) \left( (\varphi(b) - \varphi(a))^2 (\varphi(b) + 2\varphi(a)) + t(\varphi(b) - \varphi(a))^3 \right) \\
&\geq t(1-t) (\varphi(b) - \varphi(a))^3
\end{aligned}$$

respectively hold.

**Proof** We denote  $\varphi(a) = A$ ,  $\varphi(b) - \varphi(a) = x$  and define

$$F(t, x) = (1-t)A^m + t(A+x)^m - (A+tx)^m$$

and using the Newton Binomial Expansion of  $(A+x)^m$  and  $(A+tx)^m$ , we get that

$$\begin{aligned}
F(t, x) &= (1-t)A^m + \sum_{k=0}^m A^{m-k}x^k (t-t^k) \\
&= x^m t (1-t^{m-1}) + \sum_{k=2}^{m-1} A^{m-k}x^k (t-t^k) \\
&= t(1-t) \left( x^m \sum_{k=0}^{m-2} t^k + \sum_{k=2}^{m-1} Ax^k \sum_{j=0}^{k-2} t^j \right) \\
&\geq t(1-t)x^m \sum_{j=0}^{m-2} t^j \geq t(1-t)x^m,
\end{aligned}$$

which means that (12) holds.

From D in Theorem 2, Theorem 4 and Lemma 3 we get that:

**Theorem 5** Let  $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  be twice differentiable, superquadratic function and let  $f(x) = \varphi^m(x)$ ,  $m \geq 2$ . Then,  $f$  is superquadratic and when  $m = 2, 3, \dots$ , the following refinements of Jensen's inequality hold for  $0 \leq t \leq 1$

$$\begin{aligned}
&\varphi^m((1-t)a + tb) \tag{13} \\
&\leq ((1-t)\varphi(a) + t\varphi(b) - t\varphi((1-t)|b-a|) - (1-t)\varphi(t|b-a|))^m \\
&\leq (1-t)\varphi^m(a) + t\varphi^m(b) - t\varphi^m((1-t)|b-a|) - (1-t)\varphi^m(t|b-a|) \\
&\leq (1-t)\varphi^m(a) + t\varphi^m(b).
\end{aligned}$$

**Proof** When  $f(x) = g \circ \varphi(x) = (\varphi(x))^m$ ,  $m \geq 2$  then by D in Theorem 2,  $f$  is superquadratic. According to (7) in Theorem 4 and (12) in Lemma 3, in order to satisfy (8) it is enough to prove that the inequality

$$(1-t)\varphi^m(a) + t\varphi^m(b) - ((1-t)\varphi(a) + t\varphi(b))^m \tag{14}$$

$$\geq t(1-t)\varphi^m(b-a)$$

holds. This means that according to (12) and Theorem 4, it is enough to prove when  $m = 2, 3, \dots$  that

$$t(1-t)(\varphi(b) - \varphi(a))^m \geq t(1-t)\varphi^m(b-a). \tag{15}$$

holds. Inequality (15) holds because  $m > 0$  and  $\varphi$  is a convex function that satisfies  $\varphi(0) = 0$ , hence  $\varphi(b) - \varphi(a) \geq \varphi(b-a) - \varphi(0) = \varphi(b-a)$  and because  $0 \leq t \leq 1$ . Hence from Theorem 4 we get that for  $g(x) = x^m$ ,  $m = 2, 3, \dots$ , inequality (13) holds. The proof is complete.

Next we deal with two refinements of Theorem 1—the Jensen-Steffensen’s inequality for convex functions. First we quote part of the theorem on this subject from [4, Theorem 1] when the functions involved are non-negative superquadratic.

**Theorem 6** *Let  $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  be differentiable and superquadratic, let  $\mathbf{a} = (a_1, a_2, \dots, a_n)$  be a non-negative increasing  $n$ -tuple in  $\mathbb{R}^n$ , and  $\boldsymbol{\rho}$  be a real  $n$ -tuple satisfying Steffensen’s coefficients, that is  $\rho_i, i = 1, \dots, n$  satisfy the condition in (2). Let  $\bar{a}$  be defined by  $\bar{a} = \sum_{i=1}^n \rho_i a_i$ . Then*

$$\begin{aligned} & \sum_{i=1}^n \rho_i \varphi(a_i) - \varphi(\bar{a}) \\ & \geq \left( \sum_{i=1}^k P_i + \sum_{i=k+1}^n \bar{P}_i \right) \varphi \left( \frac{\sum_{i=1}^n \rho_i (|a_i - \bar{a}|)}{\sum_{i=1}^k P_i + \sum_{i=k+1}^n \bar{P}_i} \right) \\ & \geq (n-1) \varphi \left( \frac{\sum_{i=1}^n \rho_i (|a_i - \bar{a}|)}{n-1} \right), \end{aligned} \tag{16}$$

where  $k \in \{1, \dots, n-1\}$  satisfies

$$a_1 \leq a_2 \leq \dots \leq a_k < \bar{a} < a_{k+1} \leq \dots \leq a_n.$$

The theorem below is a refinement of Theorem 1 and is a Jensen-Steffensen’s inequality for non-negative compound superquadratic functions. Steffensen’s coefficients in the case  $n = 2$ , leads always to  $\rho_i > 0$ . Therefore we deal in this theorem with  $n \geq 3$ .

Using the compound function  $f = g \circ \varphi$  we get a superquadratic Jensen-Steffensen’s type inequality, and therefore we refine Jensen-Steffensen’s inequality for convex functions. The proof is similar to the proof of Theorem 3.

**Theorem 7** *Let  $g : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  be twice differentiable non-decreasing convex function and  $g(0) = 0$ . Let  $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  be twice differentiable, strictly increasing superquadratic function, and  $f$  defined as  $f = g \circ \varphi$  be superquadratic. If  $a_i$  and  $\rho_i, i = 1, \dots, n$ , satisfy the same conditions as in Theorem 6, then the inequalities*



$$\begin{aligned}
f\left(\sum_{i=1}^n \rho_i a_i\right) &\leq \sum_{i=1}^n \rho_i f(a_i) - (n-1) f\left(\frac{\sum_{i=1}^n \rho_i (|a_i - \bar{a}|)}{n-1}\right) \\
&\leq \sum_{i=1}^n \rho_i f(a_i) - f\left(\frac{\sum_{i=1}^n \rho_i (|a_i - \bar{a}|)}{n-1}\right) \leq \sum_{i=1}^n \rho_i f(a_i)
\end{aligned} \tag{17}$$

and

$$\begin{aligned}
&f\left(\sum_{i=1}^n \rho_i a_i\right) \\
&\leq f \circ \varphi^{-1}\left(\sum_{i=1}^n \rho_i \varphi(a_i) - (n-1) \varphi\left(\frac{\sum_{i=1}^n \rho_i (|a_i - \bar{a}|)}{n-1}\right)\right) \\
&\leq f \circ \varphi^{-1}\left(\sum_{i=1}^n \rho_i \varphi(a_i)\right) - f \circ \varphi^{-1}\left((n-1) \varphi\left(\frac{\sum_{i=1}^n \rho_i (|a_i - \bar{a}|)}{n-1}\right)\right) \\
&\leq f \circ \varphi^{-1}\left(\sum_{i=1}^n \rho_i \varphi(a_i)\right) - f \circ \varphi^{-1}\left(\varphi\left(\frac{\sum_{i=1}^n \rho_i (|a_i - \bar{a}|)}{n-1}\right)\right) \\
&\leq \sum_{i=1}^n \rho_i f(a_i) - f\left(\frac{\sum_{i=1}^n \rho_i (|a_i - \bar{a}|)}{n-1}\right) \leq \sum_{i=1}^n \rho_i f(a_i)
\end{aligned} \tag{18}$$

hold. In particular, if also  $\frac{\varphi'(x)}{x}$  is increasing, or if  $f \circ \varphi^{-1}(x) = x^m$ ,  $m \geq 2$ , then (17) and (18) hold.

**Proof** Inequality (17) follows from (16) and the superquadracity and positivity of  $f$ .

We use the fact that  $g$  is convex increasing and  $f = g \circ \varphi$  to prove (18). We follow step by step the proof of Theorem 3 and get that:

$$\begin{aligned}
f\left(\sum_{i=1}^n \rho_i a_i\right) &\leq g\left(\sum_{i=1}^n \rho_i \varphi(a_i) - (n-1) \varphi\left(\frac{\sum_{i=1}^n \rho_i (|a_i - \bar{a}|)}{n-1}\right)\right) \\
&\leq g\left(\sum_{i=1}^n \rho_i \varphi(a_i)\right) - g\left((n-1) \varphi\left(\frac{\sum_{i=1}^n \rho_i (|a_i - \bar{a}|)}{n-1}\right)\right) \\
&\leq g\left(\sum_{i=1}^n \rho_i \varphi(a_i)\right) - g \circ \varphi\left(\frac{\sum_{i=1}^n \rho_i (|a_i - \bar{a}|)}{n-1}\right) \\
&\leq \sum_{i=1}^n \rho_i g \circ \varphi(a_i) - g \circ \varphi\left(\frac{\sum_{i=1}^n \rho_i (|a_i - \bar{a}|)}{n-1}\right)
\end{aligned}$$

$$\leq \sum_{i=1}^n \rho_i f(a_i) - f\left(\frac{\sum_{i=1}^n \rho_i (|a_i - \bar{a}|)}{n-1}\right) \leq \sum_{i=1}^n \rho_i f(a_i).$$

Therefore (18) is proved. If also  $\frac{\varphi'(x)}{x}$  is increasing or  $f \circ \varphi^{-1}(x) = x^m, m \geq 2$ , then according to Theorem 2  $f = g \circ \varphi$  is superquadratic and therefore (17) and (18) hold. The proof of Theorem 7 is complete.

### 3 The Quasi-Mean $F_f(\mathbf{x}, \boldsymbol{\lambda})$

In [1] and [2] the following quasi-mean

$$W_f(\mathbf{x}, \boldsymbol{\lambda}) = f^{-1}\left(f\left(\sum_{r=1}^n \lambda_r x_r\right) + \sum_{r=1}^n \lambda_r f\left(\left|x_r - \sum_{i=1}^n \lambda_i x_i\right|\right)\right)$$

$$\sum_{r=1}^n \lambda_r = 1, \quad \lambda_r \geq 0, x_r \geq 0, r = 1, \dots, n,$$

is discussed. Similarly, the quasi-mean

$$F_f(\mathbf{x}, \boldsymbol{\lambda}) = f^{-1}\left(\sum_{i=1}^n \lambda_i f(x_i) - \sum_{i=1}^n \lambda_i f\left(\left|x_i - \sum_{j=1}^n \lambda_j x_j\right|\right)\right)$$

$$\sum_{r=1}^n \lambda_r = 1, \quad \lambda_r \geq 0, x_r \geq 0, r = 1, \dots, n,$$

is discussed here.

Both  $F_f(\mathbf{x}, \boldsymbol{\lambda})$  and  $W_f(\mathbf{x}, \boldsymbol{\lambda})$  consist the building blocks of superquadratic functions (see Corollary 2).

First we quote from [1] about  $W_f(x_1, x_2)$ :

**Definition 3 ([1, Definition 1])** Let a strictly increasing convex function  $f$  be defined on  $[0, b), 0 < b \leq \infty$ , and let  $f(0) = 0$ . For such  $f$  we define the quasi-mean  $W_f(x_1, x_2)$  as

$$W_f(x_1, x_2) = f^{-1}\left(f\left(\frac{x_1 + x_2}{2}\right) + f\left(\left|\frac{x_1 - x_2}{2}\right|\right)\right).$$

In the special case that  $f(x) = x^p, W_p(x_1, x_2)$  is defined as

$$W_p(x_1, x_2) = \left(\left(\frac{x_1 + x_2}{2}\right)^p + \left(\left|\frac{x_1 - x_2}{2}\right|\right)^p\right)^{\frac{1}{p}}$$

for  $x_1, x_2 \in [0, b)$ ,  $p \geq 1$ .

In Lemma 4 it is proved that  $W_f(x_1, x_2)$  is a quasi-mean:

**Lemma 4 ([1, Lemma 1])** *Under the conditions of Definition 3 on  $f$ ,  $W_f(x_1, x_2)$  is symmetric and satisfies:*

- (a)  $W_f(x, x) = x$ ,  $x \in [0, b)$ ,
- (b)  $x_1 \leq W_f(x_1, x_2) \leq x_2$ ,  $0 \leq x_1 \leq x_2 \leq b$ .  
When  $f(x) = x^p$ ,  $p \geq 1$ ,  $W_p(x_1, x_2)$  satisfies
- (c)  $W_p(\lambda x_1, \lambda x_2) = \lambda W_p(x_1, x_2)$ ,  $\lambda \geq 0$ ,  $0 \leq x_1, x_2 \leq b$ .

Similar to Definition 3 and Lemma 4, we discuss  $F_f(x_1, x_2)$  which is related to the theorems in this paper.

**Definition 4** Let a strictly increasing convex function  $f$  be defined on  $[0, b)$ ,  $0 < b \leq \infty$ , and let  $f(0) = 0$ . For such  $f$  we define the quasi-mean  $F_f(x_1, x_2)$  as

$$F_f(x_1, x_2) = f^{-1} \left( \frac{f(x_1) + f(x_2)}{2} - f \left( \left| \frac{x_1 - x_2}{2} \right| \right) \right).$$

In the special case that  $f(x) = x^p$ ,  $F_p(x_1, x_2)$  is defined as

$$F_p(x_1, x_2) = \left( \frac{x_1^p + x_2^p}{2} - \left( \left| \frac{x_1 - x_2}{2} \right| \right)^p \right)^{\frac{1}{p}}$$

for  $x_1, x_2 \in [0, b)$ ,  $p \geq 1$ .

**Lemma 5** *Under the conditions of Definition 4 on  $f$ ,  $F_f(x_1, x_2)$  is symmetric and satisfies:*

- (a)  $F_f(x, x) = x$ ,  $x \in [0, b)$ ,
- (b)  $x_1 \leq F_f(x_1, x_2) \leq x_2$ ,  $0 \leq x_1 \leq x_2 \leq b$ .  
When  $f(x) = x^p$ ,  $p \geq 1$ ,  $F_p(x_1, x_2)$  satisfies
- (c)  $F_p(\lambda x_1, \lambda x_2) = \lambda F_p(x_1, x_2)$ ,  $\lambda \geq 0$ ,  $0 \leq x_1, x_2 \leq b$ .

**Proof** Properties (a) and (c) are obvious. To prove that property (b) holds we have to show that

$$x_1 \leq f^{-1} \left( \frac{f(x_1) + f(x_2)}{2} - f \left( \left| \frac{x_1 - x_2}{2} \right| \right) \right) \leq x_2. \quad (19)$$

As  $f$  is strictly increasing (19) is equivalent to

$$f(x_1) \leq \frac{f(x_1) + f(x_2)}{2} - f \left( \left| \frac{x_1 - x_2}{2} \right| \right) \leq f(x_2) \quad (20)$$

and because  $f$  is non-negative and increasing the right hand-side of (20) follows from

$$f(x_2) \geq f(x_1) - 2f\left(\frac{x_2 - x_1}{2}\right),$$

and therefore the right hand-side of (19) holds.

To prove the left hand-side of (19) we deal with the left hand-side of (20) which is the same as to prove the inequality

$$f(x_2) - f(x_1) \geq 2f\left(\frac{x_2 - x_1}{2}\right). \quad (21)$$

We show that

$$f(x_2) - f(x_1) \geq f(x_2 - x_1) \geq 2f\left(\frac{x_2 - x_1}{2}\right). \quad (22)$$

Indeed the right hand-side inequality of (22) follows from the fact that when  $f$  is twice differentiable, convex and  $f(0) = 0$ , then  $\frac{f(kx)}{x}$  is increasing for  $x > 0$ . From (22) inequality (21) follows, and therefore the left hand-side inequality of (19) holds. The proof of the lemma is complete.

## References

1. S. Abramovich, Quasi-arithmetic means and superquadracity. *J. Math. Inequal.* **9**, 1157–1168 (2015)
2. S. Abramovich, The behavior of difference of two means, in *Developments in Functional Equations and Related Topics*, ed. by J. Brzdęk, K. Ciepliński, Th.M. Rassias (Springer, Cham, 2017), pp. 1–15
3. S. Abramovich, G. Jameson, G. Sinnamon, Refining Jensen's inequality. *Bull. Math. Soc. Math. Roum (Novel Series)* **47**(95), 3–14 (2004)
4. S. Abramovich, S. Banić, M. Matić, J. Pečarić, Jensen Steffensen's and related inequalities, for superquadratic functions. *Math. Inequal. Appl.* **11**, 23–41 (2008)
5. S.S. Dragomir, Inequalities of Hermite-Hadamard type for composite convex functions. *RGMIA. Res. Rap. Coll.* **23**, Art. 39, 21 pp. (2018)
6. J. Matkowski, A functional inequality characterizing convex functions, conjugacy and generalization of Hölder's and Minkowski's inequalities. *Aequat. Math.* **40**, 168–180 (1990)
7. J.E. Pečarić, F. Proschan, Y.L. Tong, *Convex Functions, Partial Ordering and Statistical Applications* (Academic Press, New York, 1992)

# Best Hyers–Ulam Stability Constants on a Time Scale with Discrete Core and Continuous Periphery



Douglas R. Anderson and Masakazu Onitsuka

**Abstract** Consider a time scale consisting of a discrete core with uniform step size, augmented with a continuous-interval periphery. On this time scale, we determine the best constants for the Hyers–Ulam stability of a first-order dynamic equation with complex constant coefficient, based on the placement of the complex coefficient in the complex plane, with respect to the imaginary axis and the Hilger circle. These best constants are then related to known results for the special cases of completely continuous and uniformly discrete time scales.

## 1 Introduction

In this paper we explore the Hyers–Ulam stability of a certain dynamic equation on a new time scale with a discrete, uniform core and continuous periphery. Ulam inaugurated this type of stability [37], followed by Hyers [22] and Rassias [34]. Since then, there has been wide-spread interest in this type of stability, including for difference equations, recurrence relations,  $h$ -difference equations, quantum equations, and dynamic equations on time scales. For early papers on difference equations, see Popa [31, 32]; more current works include Anderson and Onitsuka [5, 6], Baias and Popa [13], Brzdęk and Wójcik [16], Onitsuka [29, 30], Rasouli, Abbaszadeh, and Eshaghi [33], Xu and Brzdęk [38]. A related monograph is Brzdęk, Popa, Raşa, and Xu [17]. Quantum equations and Hyers–Ulam stability are investigated in Anderson and Onitsuka [7, 8]. For some work on matrix and nonlinear difference equations, see Jung and Nam [24, 25], and Nam [26–28]. For early papers on time scales, see András and Mészáros [12], Hua, Li, and Feng [21];

---

D. R. Anderson (✉)

Concordia College, Department of Mathematics, Moorhead, MN, USA

e-mail: [andersod@cord.edu](mailto:andersod@cord.edu)

M. Onitsuka

Okayama University of Science, Department of Applied Mathematics, Okayama, Japan

e-mail: [onitsuka@xmath.ous.ac.jp](mailto:onitsuka@xmath.ous.ac.jp)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,

Springer Optimization and Its Applications 173,

[https://doi.org/10.1007/978-3-030-72563-1\\_2](https://doi.org/10.1007/978-3-030-72563-1_2)

contemporary results include Anderson [2], Anderson, Jennissen, and Montplaisir [10], Anderson and Onitsuka [3, 4, 9], Shen [35], Shen and Li [36]. For recent papers with non-constant or periodic coefficients, see Anderson [1], Anderson, Onitsuka, and Rassias [11], Baias, Blaga, and Popa [14], Buşe, Lupulescu, and O'Regan [19], Buşe, O'Regan, and Sailerli [18].

This work will proceed as follows. In Section 2, we will define the time scale with discrete core and continuous periphery, introduce the basic derivative and exponential function for this time scale, and define Hyers–Ulam stability for the dynamic equation with a complex constant coefficient. In Section 3, we establish the best Hyers–Ulam stability constants in Theorem 5, based on the location of the complex coefficient with respect to the imaginary axis, and for negative real part, with respect to the left Hilger circle. If we expand the discrete core to all of  $h\mathbb{Z}$ , or shrink it to recover the continuum  $\mathbb{R}$ , we are able to relate our new results with the current literature in the field. As we do this, an interesting case arises when the real part of the complex coefficient is negative but it lies outside the Hilger circle; this case is explored in Section 4. After that, we provide a brief conclusion and future direction.

## 2 Time Scale with Discrete Core and Continuous Periphery

Let  $\mathbb{N}_0$  denote the non-negative integers  $\{0, 1, 2, \dots\}$ , let  $m \in \mathbb{N}_0$ , and let  $h > 0$ . Define the time scale with discrete core and continuous periphery via

$$\mathbb{T}_{hm} := (-\infty, -hm) \cup \{-hm, \dots, -h, 0, h, \dots, hm\} \cup (hm, \infty).$$

Here,  $h > 0$  is the uniform step size in the discrete core, with discrete spread  $m \in \mathbb{N}_0$  out to the continuous periphery. Define the graininess function  $\mu : \mathbb{T}_{hm} \rightarrow \mathbb{R}$  via

$$\mu(t) = \begin{cases} 0 & : t \in (-\infty, -hm) \cup [hm, \infty), \\ h & : t \in \{-hm, \dots, -h, 0, h, \dots, h(m-1)\}. \end{cases}$$

As  $h \rightarrow 0$ , or if  $m = 0$ , we have  $\mathbb{T}_{0,m} = \mathbb{T}_{h,0} = \mathbb{R}$ , and we recover results for classical differential equations; as  $m \rightarrow \infty$  for fixed  $h > 0$ , we have  $\mathbb{T}_{h,\infty} = h\mathbb{Z}$  and we recover results for standard  $h$ -difference equations.

In this section we introduce the first-order linear homogeneous equation with constant complex-valued coefficient

$$x^\Delta(t) - \lambda x(t) = 0, \quad \lambda \in \mathbb{C} \setminus \left\{ \frac{-1}{h} \right\}, \quad t \in \mathbb{T}_{hm}, \quad (1)$$

where

$$x^\Delta(t) := \begin{cases} \frac{d}{dt}x(t) & : t \in (-\infty, -hm) \cup [hm, \infty) \\ \frac{x(t+h)-x(t)}{h} & : t \in \{-hm, \dots, -h, 0, h, \dots, h(m-1)\}. \end{cases}$$

**Lemma 1 (Exponential Function)** Fix  $h > 0$ . For  $t \in \mathbb{T}_{hm}$ , define the function

$$e_\lambda(t, 0) := \begin{cases} (1 + h\lambda)^{-m} e^{\lambda(t+hm)} & : t \in (-\infty, -hm) \\ (1 + h\lambda)^{\frac{t}{h}} & : t \in \{-hm, \dots, -h, 0, h, \dots, hm\} \\ (1 + h\lambda)^m e^{\lambda(t-hm)} & : t \in (hm, \infty). \end{cases} \quad (2)$$

Then,  $x(t) = x_0 e_\lambda(t, 0)$  for  $e_\lambda(t, 0)$  given in (2) is the unique solution of (1) satisfying  $x(0) = x_0 \in \mathbb{C}$ .

### 3 Best Constants for First-Order Equations with Constant Complex Coefficient

In this section, we consider on  $\mathbb{T}_{hm}$  the Hyers–Ulam stability of (1), defined as follows.

**Definition 1 (HUS)** Let  $\varepsilon > 0$  be arbitrary. Equation (1) has Hyers–Ulam stability (HUS) if and only if given  $\phi : \mathbb{T}_{hm} \rightarrow \mathbb{C}$  satisfying  $|\phi^\Delta(t) - \lambda\phi(t)| \leq \varepsilon$  for all  $t \in \mathbb{T}_{hm}$ , there exists a solution  $x : \mathbb{T}_{hm} \rightarrow \mathbb{C}$  of (1) and a constant  $K > 0$  such that  $|\phi(t) - x(t)| \leq K\varepsilon$  for all  $t \in \mathbb{T}_{hm}$ . Such a constant  $K$  is called an HUS constant for (1) on  $\mathbb{T}_{hm}$ .

**Theorem 1** Let  $\lambda \in \mathbb{C} \setminus \left\{ \frac{-1}{h} \right\}$  with  $\operatorname{Re}(\lambda) > 0$ . Let  $\varepsilon > 0$  be a fixed arbitrary constant, and let  $\phi$  be a function on  $\mathbb{T}_{hm}$  satisfying the inequality

$$|\phi^\Delta(t) - \lambda\phi(t)| \leq \varepsilon, \quad t \in \mathbb{T}_{hm}.$$

Then,  $\lim_{t \rightarrow \infty} \frac{\phi(t)}{e_\lambda(t, 0)}$  exists, and the function  $x$  given by

$$x(t) := \left( \lim_{t \rightarrow \infty} \frac{\phi(t)}{e_\lambda(t, 0)} \right) e_\lambda(t, 0)$$

is the unique solution of (1) with

$$|\phi(t) - x(t)| \leq \varepsilon \left( \frac{1}{\operatorname{Re}(\lambda)} \right)$$

for all  $t \in \mathbb{T}_{hm}$ .

**Proof** Let  $\lambda \in \mathbb{C} \setminus \{\frac{-1}{h}\}$  with  $\operatorname{Re}(\lambda) > 0$ . Throughout this proof, as  $|\phi^\Delta(t) - \lambda\phi(t)| \leq \varepsilon$  for all  $t \in \mathbb{T}_{hm}$ , there exists a function  $q : \mathbb{T}_{hm} \rightarrow \mathbb{C}$  such that

$$\phi^\Delta(t) - \lambda\phi(t) = q(t), \quad |q(t)| \leq \varepsilon$$

for all  $t \in \mathbb{T}_{hm}$ . The variation of constants formula then yields

$$\phi(t) = \phi_0 e_\lambda(t, 0) + e_\lambda(t, 0) \int_0^t \frac{q(\tau)}{e_\lambda(\sigma(\tau), 0)} \Delta\tau.$$

Since  $\operatorname{Re}(\lambda) > 0$  and  $|q(t)| \leq \varepsilon$ , we can rewrite  $\phi$  as

$$\phi(t) = \left( \phi_0 + \int_0^\infty \frac{q(\tau)}{e_\lambda(\sigma(\tau), 0)} \Delta\tau \right) e_\lambda(t, 0) - e_\lambda(t, 0) \int_t^\infty \frac{q(\tau)}{e_\lambda(\sigma(\tau), 0)} \Delta\tau, \quad (3)$$

where

$$x_0 := \phi_0 + \int_0^\infty \frac{q(\tau)}{e_\lambda(\sigma(\tau), 0)} \Delta\tau \in \mathbb{C}$$

exists and is finite. Clearly

$$x(t) := x_0 e_\lambda(t, 0), \quad t \in \mathbb{T}_{hm}$$

is a solution of (1), and

$$\lim_{t \rightarrow \infty} \frac{\phi(t)}{e_\lambda(t, 0)} = \lim_{t \rightarrow \infty} \left( \phi_0 + \int_0^t \frac{q(\tau)}{e_\lambda(\sigma(\tau), 0)} \Delta\tau \right) = x_0$$

exists, so

$$x(t) = \left( \lim_{t \rightarrow \infty} \frac{\phi(t)}{e_\lambda(t, 0)} \right) e_\lambda(t, 0).$$

We take into account three cases based on the three branches of the exponential function in (2).

(a). For  $\operatorname{Re}(\lambda) > 0$  and  $t \in (hm, \infty)$ , using (3) we have that

$$\begin{aligned} |\phi(t) - x(t)| &= \left| -e_\lambda(t, 0) \int_t^\infty \frac{q(\tau)}{e_\lambda(\sigma(\tau), 0)} \Delta\tau \right| \\ &\leq \varepsilon |e_\lambda(t, 0)| \int_t^\infty \frac{\Delta\tau}{|e_\lambda(\sigma(\tau), 0)|} \\ &= \varepsilon |1 + h\lambda|^m e^{\operatorname{Re}(\lambda)(t-hm)} \int_t^\infty \frac{d\tau}{|1 + h\lambda|^m e^{\operatorname{Re}(\lambda)(\tau-hm)}} \\ &= \frac{\varepsilon}{\operatorname{Re}(\lambda)} \end{aligned}$$

holds for all  $t \in (hm, \infty)$ .



(b). For  $\operatorname{Re}(\lambda) > 0$  and  $t \in \{-hm, \dots, -h, 0, h, \dots, hm\}$ , using (3) we have

$$\begin{aligned}
 |\phi(t) - x(t)| &\leq \varepsilon |e_\lambda(t, 0)| \int_t^\infty \frac{\Delta\tau}{|e_\lambda(\sigma(\tau), 0)|} \\
 &= \varepsilon |1 + h\lambda|^{\frac{t}{h}} \left( \int_t^{hm} + \int_{hm}^\infty \right) \frac{\Delta\tau}{|e_\lambda(\sigma(\tau), 0)|} \\
 &= \varepsilon |1 + h\lambda|^{\frac{t}{h}} \left( \sum_{j=\frac{t}{h}}^{m-1} \frac{h}{|1 + h\lambda|^{j+1}} + \int_{hm}^\infty \frac{d\tau}{|1 + h\lambda|^m e^{\operatorname{Re}(\lambda)(\tau-hm)}} \right) \\
 &= \varepsilon |1 + h\lambda|^{\frac{t}{h}} \left( \frac{h \left( |1 + h\lambda|^{-\frac{t}{h}} - |1 + h\lambda|^{-m} \right)}{|1 + h\lambda| - 1} + \frac{1}{|1 + h\lambda|^m \operatorname{Re}(\lambda)} \right) \\
 &= \varepsilon \left( \frac{h}{|1 + h\lambda| - 1} + |1 + h\lambda|^{\frac{t}{h}-m} \left( \frac{1}{\operatorname{Re}(\lambda)} - \frac{h}{|1 + h\lambda| - 1} \right) \right) \\
 &\leq \frac{\varepsilon}{\operatorname{Re}(\lambda)}
 \end{aligned}$$

for all  $t \in \{-hm, \dots, -h, 0, h, \dots, hm\}$ , as  $\frac{t}{h} \leq m$  and  $\frac{1}{\operatorname{Re}(\lambda)} \geq \frac{h}{|1+h\lambda|-1}$  for  $\operatorname{Re}(\lambda) > 0$  and  $h > 0$ .

(c). For  $\operatorname{Re}(\lambda) > 0$  and  $t \in (-\infty, -hm)$ , using (3) we have

$$\begin{aligned}
 |\phi(t) - x(t)| &\leq \varepsilon |e_\lambda(t, 0)| \int_t^\infty \frac{\Delta\tau}{|e_\lambda(\sigma(\tau), 0)|} \\
 &= \frac{\varepsilon e^{\operatorname{Re}(\lambda)(t+hm)}}{|1 + h\lambda|^m} \left( \int_t^{-hm} + \int_{-hm}^{hm} + \int_{hm}^\infty \right) \frac{\Delta\tau}{|e_\lambda(\sigma(\tau), 0)|} \\
 &= \frac{\varepsilon e^{\operatorname{Re}(\lambda)(t+hm)}}{|1 + h\lambda|^{2m}} \left( \frac{1}{\operatorname{Re}(\lambda)} + \frac{(e^{-\operatorname{Re}(\lambda)(t+hm)} - 1)}{|1 + h\lambda|^{-2m} \operatorname{Re}(\lambda)} \right. \\
 &\quad \left. + \frac{h \left( |1 + h\lambda|^{2m} - 1 \right)}{|1 + h\lambda| - 1} \right) \\
 &= \varepsilon \left\{ \frac{1}{\operatorname{Re}(\lambda)} + e^{\operatorname{Re}(\lambda)(t+hm)} \left( \frac{1}{|1 + h\lambda|^{2m} \operatorname{Re}(\lambda)} \right. \right. \\
 &\quad \left. \left. + \left( \frac{h}{|1 + h\lambda| - 1} \right) \left( 1 - \frac{1}{|1 + h\lambda|^{2m}} \right) - \frac{1}{\operatorname{Re}(\lambda)} \right) \right\} \\
 &\leq \frac{\varepsilon}{\operatorname{Re}(\lambda)}
 \end{aligned}$$

for all  $t \in (-\infty, -hm)$ , as  $t < -hm$ , and the expression inside the square brackets is negative.

We next show that  $x$  is the unique solution of (1) such that  $|\phi(t) - x(t)| \leq K\varepsilon := \frac{1}{\operatorname{Re}(\lambda)}\varepsilon$  for all  $t \in \mathbb{T}_{hm}$ . Suppose  $\phi : \mathbb{T}_{hm} \rightarrow \mathbb{C}$  is an approximate solution of (1) such that

$$|\phi^\Delta(t) - \lambda\phi(t)| \leq \varepsilon \text{ for all } t \in \mathbb{T}_{hm}$$

for some  $\varepsilon > 0$ . Suppose further that  $x_1, x_2 : \mathbb{T}_{hm} \rightarrow \mathbb{C}$  are two different solutions of (1) such that  $|\phi(t) - x_j(t)| \leq K\varepsilon$  for all  $t \in \mathbb{T}_{hm}$ , for  $j = 1, 2$ . Then, we have for constants  $c_j \in \mathbb{C}$  that

$$x_j(t) = c_j e_\lambda(t, 0), \quad c_1 \neq c_2,$$

and

$$|c_1 - c_2| \cdot |e_\lambda(t, 0)| = |x_1(t) - x_2(t)| \leq |x_1(t) - \phi(t)| + |\phi(t) - x_2(t)| \leq 2K\varepsilon;$$

letting  $t \rightarrow \infty$  yields  $\infty < 2K\varepsilon$ , a contradiction. Consequently,  $x$  is the unique solution of (1) such that  $|\phi(t) - x(t)| \leq \frac{\varepsilon}{\operatorname{Re}(\lambda)}$  for all  $t \in \mathbb{T}_{hm}$ . This completes the proof.  $\square$

**Theorem 2** *Let  $\lambda \in \mathbb{C} \setminus \left\{ \frac{-1}{h} \right\}$  with  $\operatorname{Re}(\lambda) < 0$ . Let  $\varepsilon > 0$  be a fixed arbitrary constant, and let  $\phi$  be a function on  $\mathbb{T}_{hm}$  satisfying the inequality*

$$|\phi^\Delta(t) - \lambda\phi(t)| \leq \varepsilon, \quad t \in \mathbb{T}_{hm}.$$

Then,  $\lim_{t \rightarrow -\infty} \frac{\phi(t)}{e_\lambda(t, 0)}$  exists, and the function  $x$  given by

$$x(t) := \left( \lim_{t \rightarrow -\infty} \frac{\phi(t)}{e_\lambda(t, 0)} \right) e_\lambda(t, 0)$$

is the unique solution of (1) with  $|\phi(t) - x(t)| \leq K\varepsilon$  for all  $t \in \mathbb{T}_{hm}$ , where

$$K := \begin{cases} \frac{-1}{\operatorname{Re}(\lambda)} + 2hm & : |1 + h\lambda| = 1 \\ \max \left\{ \frac{-1}{\operatorname{Re}(\lambda)}, \frac{h}{1 - |1 + h\lambda|} + |1 + h\lambda|^{2m} \left( \frac{h}{|1 + h\lambda| - 1} - \frac{1}{\operatorname{Re}(\lambda)} \right) \right\} & : |1 + h\lambda| \neq 1. \end{cases} \quad (4)$$

In particular, the following holds.

(i) If  $t \in (-\infty, -hm)$ , then

$$|\phi(t) - x(t)| \leq \varepsilon \left( \frac{-1}{\operatorname{Re}(\lambda)} \right)$$

for all  $t \in (-\infty, -hm)$ .

(ii) If  $|1 + h\lambda| = 1$  and  $t \in \{-hm, \dots, -h, 0, h, \dots, hm\} \cup (hm, \infty)$ , then

$$|\phi(t) - x(t)| \leq \varepsilon \left( \frac{-1}{\operatorname{Re}(\lambda)} + 2hm \right)$$

for all  $t \in \{-hm, \dots, -h, 0, h, \dots, hm\} \cup (hm, \infty)$ .

(iii) If  $0 < |1 + h\lambda| < 1$  and  $t \in \{-hm, \dots, -h, 0, h, \dots, hm\} \cup (hm, \infty)$ , then

$$|\phi(t) - x(t)| \leq \varepsilon \left( \frac{h}{1 - |1 + h\lambda|} + |1 + h\lambda|^{2m} \left( \frac{h}{|1 + h\lambda| - 1} - \frac{1}{\operatorname{Re}(\lambda)} \right) \right)$$

for all  $t \in \{-hm, \dots, -h, 0, h, \dots, hm\} \cup (hm, \infty)$ .

(iv) If  $|1 + h\lambda| > 1$  and  $t \in \{-hm, \dots, -h, 0, h, \dots, hm\} \cup (hm, \infty)$ , then

$$|\phi(t) - x(t)| \leq \varepsilon \left( \frac{h}{1 - |1 + h\lambda|} + |1 + h\lambda|^{2m} \left( \frac{h}{|1 + h\lambda| - 1} - \frac{1}{\operatorname{Re}(\lambda)} \right) \right)$$

for all  $t \in \{-hm, \dots, -h, 0, h, \dots, hm\} \cup (hm, \infty)$ .

**Proof** Let  $\lambda \in \mathbb{C} \setminus \{\frac{-1}{h}\}$  with  $\operatorname{Re}(\lambda) < 0$ . Supposing  $|\phi^\Delta(t) - \lambda\phi(t)| \leq \varepsilon$  for all  $t \in \mathbb{T}_{hm}$ , there exists a function  $q : \mathbb{T}_{hm} \rightarrow \mathbb{C}$  such that

$$\phi^\Delta(t) - \lambda\phi(t) = q(t), \quad |q(t)| \leq \varepsilon$$

for all  $t \in \mathbb{T}_{hm}$ . Then, we have

$$\phi(t) = \phi_0 e_\lambda(t, 0) + e_\lambda(t, 0) \int_0^t \frac{q(\tau)}{e_\lambda(\sigma(\tau), 0)} \Delta\tau.$$

Since  $\operatorname{Re}(\lambda) < 0$  and  $|q(t)| \leq \varepsilon$ , we can rewrite  $\phi$  as

$$\phi(t) = \left( \phi_0 - \int_{-\infty}^0 \frac{q(\tau)}{e_\lambda(\sigma(\tau), 0)} \Delta\tau \right) e_\lambda(t, 0) + e_\lambda(t, 0) \int_{-\infty}^t \frac{q(\tau)}{e_\lambda(\sigma(\tau), 0)} \Delta\tau, \quad (5)$$

where

$$x_0 := \phi_0 - \int_{-\infty}^0 \frac{q(\tau)}{e_\lambda(\sigma(\tau), 0)} \Delta\tau \in \mathbb{C}$$

exists and is finite. As in the previous case,

$$x(t) := x_0 e_\lambda(t, 0), \quad t \in \mathbb{T}_{hm}$$

is a solution of (1), and

$$\lim_{t \rightarrow -\infty} \frac{\phi(t)}{e_\lambda(t, 0)} = \lim_{t \rightarrow -\infty} \left( \phi_0 - \int_t^0 \frac{q(\tau)}{e_\lambda(\sigma(\tau), 0)} \Delta\tau \right) = x_0$$

exists, so

$$x(t) = \left( \lim_{t \rightarrow -\infty} \frac{\phi(t)}{e_\lambda(t, 0)} \right) e_\lambda(t, 0).$$

We again work our way through the three cases based on the three branches of the exponential function in (2).

(i). For  $\operatorname{Re}(\lambda) < 0$  and  $t \in (-\infty, -hm)$ , using (5) we have

$$\begin{aligned} |\phi(t) - x(t)| &= \left| e_\lambda(t, 0) \int_{-\infty}^t \frac{q(\tau)}{e_\lambda(\sigma(\tau), 0)} \Delta\tau \right| \\ &\leq \varepsilon |e_\lambda(t, 0)| \int_{-\infty}^t \frac{\Delta\tau}{|e_\lambda(\sigma(\tau), 0)|} \\ &= \varepsilon |1 + h\lambda|^{-m} e^{\operatorname{Re}(\lambda)(t+hm)} \int_{-\infty}^t \frac{|1 + h\lambda|^m d\tau}{e^{\operatorname{Re}(\lambda)(\tau+hm)}} \\ &= -\frac{\varepsilon}{\operatorname{Re}(\lambda)} \end{aligned}$$

holds for  $\operatorname{Re}(\lambda) < 0$  and for all  $t \in (-\infty, -hm)$ .

(ii) (a). For  $\operatorname{Re}(\lambda) < 0$  with  $|1 + h\lambda| = 1$  and  $t \in \{-hm, \dots, -h, 0, h, \dots, hm\}$ , using (5) we have

$$\begin{aligned} |\phi(t) - x(t)| &\leq \varepsilon |e_\lambda(t, 0)| \int_{-\infty}^t \frac{\Delta\tau}{|e_\lambda(\sigma(\tau), 0)|} \\ &= \varepsilon \left( \int_{-\infty}^{-hm} + \int_{-hm}^t \right) \frac{\Delta\tau}{|e_\lambda(\sigma(\tau), 0)|} \\ &= \varepsilon \left( \int_{-\infty}^{-hm} \frac{d\tau}{e^{\operatorname{Re}(\lambda)(\tau+hm)}} + \sum_{j=-m}^{\frac{t-h}{h}} h \right) \\ &= \varepsilon \left( \frac{-1}{\operatorname{Re}(\lambda)} + hm + t \right) \\ &\leq \varepsilon \left( \frac{-1}{\operatorname{Re}(\lambda)} + 2hm \right), \end{aligned}$$

for  $\operatorname{Re}(\lambda) < 0$  with  $|1 + h\lambda| = 1$  and  $t \in \{-hm, \dots, -h, 0, h, \dots, hm\}$ .

(ii) (b). For  $\operatorname{Re}(\lambda) < 0$  with  $|1 + h\lambda| = 1$  and  $t \in (hm, \infty)$ , using (5) we have

$$|\phi(t) - x(t)| \leq \varepsilon e^{\operatorname{Re}(\lambda)(t-hm)} \left( \int_{-\infty}^{-hm} + \int_{-hm}^{hm} + \int_{hm}^t \right) \frac{\Delta\tau}{|e_\lambda(\sigma(\tau), 0)|}$$

$$\begin{aligned}
&= \varepsilon e^{\operatorname{Re}(\lambda)(t-hm)} \left( 2hm - \frac{e^{-\operatorname{Re}(\lambda)(t-hm)}}{\operatorname{Re}(\lambda)} \right) \\
&\leq \varepsilon \left( \frac{-1}{\operatorname{Re}(\lambda)} + 2hm \right),
\end{aligned}$$

as  $t > hm$  and  $\operatorname{Re}(\lambda) < 0$ .

(iii) (a). For  $\operatorname{Re}(\lambda) < 0$  with  $0 < |1 + h\lambda| < 1$  and  $t \in \{-hm, \dots, -h, 0, h, \dots, hm\}$ , using (5) we have

$$\begin{aligned}
|\phi(t) - x(t)| &\leq \varepsilon |e_\lambda(t, 0)| \int_{-\infty}^t \frac{\Delta\tau}{|e_\lambda(\sigma(\tau), 0)|} \\
&= \varepsilon |1 + h\lambda|^{\frac{t}{h}} \left( \int_{-\infty}^{-hm} + \int_{-hm}^t \right) \frac{\Delta\tau}{|e_\lambda(\sigma(\tau), 0)|} \\
&= \varepsilon |1 + h\lambda|^{\frac{t}{h}} \left( \int_{-\infty}^{-hm} \frac{|1 + h\lambda|^m d\tau}{e^{\operatorname{Re}(\lambda)(\tau+hm)}} + \sum_{j=-m}^{\frac{t-h}{h}} \frac{h}{|1 + h\lambda|^{j+1}} \right) \\
&= \varepsilon |1 + h\lambda|^{\frac{t}{h}} \left( \frac{h \left( |1 + h\lambda|^m - |1 + h\lambda|^{-\frac{t}{h}} \right)}{|1 + h\lambda| - 1} - \frac{|1 + h\lambda|^m}{\operatorname{Re}(\lambda)} \right) \\
&= \varepsilon \left( \frac{h}{1 - |1 + h\lambda|} + |1 + h\lambda|^{\frac{t}{h}+m} \left( \frac{h}{|1 + h\lambda| - 1} - \frac{1}{\operatorname{Re}(\lambda)} \right) \right) \\
&\leq \varepsilon \left( \frac{h}{1 - |1 + h\lambda|} + |1 + h\lambda|^{2m} \left( \frac{h}{|1 + h\lambda| - 1} - \frac{1}{\operatorname{Re}(\lambda)} \right) \right),
\end{aligned}$$

as  $\frac{t}{h} \leq m$  and  $\frac{h}{|1+h\lambda|-1} - \frac{1}{\operatorname{Re}(\lambda)} \leq 0$  for  $\operatorname{Re}(\lambda) < 0$  with  $0 < |1 + h\lambda| < 1$  and  $h > 0$ .

(iii) (b). For  $\operatorname{Re}(\lambda) < 0$  with  $0 < |1 + h\lambda| < 1$  and  $t \in (hm, \infty)$ , using (5) we have

$$\begin{aligned}
|\phi(t) - x(t)| &\leq \frac{\varepsilon e^{\operatorname{Re}(\lambda)(t-hm)}}{|1 + h\lambda|^{-m}} \left( \int_{-\infty}^{-hm} + \int_{-hm}^{hm} + \int_{hm}^t \right) \frac{\Delta\tau}{|e_\lambda(\sigma(\tau), 0)|} \\
&= \varepsilon e^{\operatorname{Re}(\lambda)(t-hm)} \left( \frac{1 - e^{-\operatorname{Re}(\lambda)(t-hm)} - |1 + h\lambda|^{2m}}{\operatorname{Re}(\lambda)} \right. \\
&\quad \left. + \frac{h \left( |1 + h\lambda|^{2m} - 1 \right)}{|1 + h\lambda| - 1} \right)
\end{aligned}$$

$$\begin{aligned}
&= \varepsilon \left( \frac{-1}{\operatorname{Re}(\lambda)} + e^{\operatorname{Re}(\lambda)(t-hm)} \left( \frac{|1+h\lambda|^{2m}-1}{-\operatorname{Re}(\lambda)} + \frac{h(|1+h\lambda|^{2m}-1)}{|1+h\lambda|-1} \right) \right) \\
&\leq \varepsilon \left( \frac{h}{1-|1+h\lambda|} + |1+h\lambda|^{2m} \left( \frac{h}{|1+h\lambda|-1} - \frac{1}{\operatorname{Re}(\lambda)} \right) \right),
\end{aligned}$$

as  $t > hm$ ,  $\operatorname{Re}(\lambda) < 0$  with  $0 < |1+h\lambda| < 1$ , and the expression inside the square brackets is non-negative.

(iv) (a). For  $\operatorname{Re}(\lambda) < 0$  with  $|1+h\lambda| > 1$  and  $t \in \{-hm, \dots, -h, 0, h, \dots, hm\}$ , using the same calculation as in case (iii)(a), we get

$$\begin{aligned}
|\phi(t) - x(t)| &\leq \varepsilon \left( \frac{h}{1-|1+h\lambda|} + |1+h\lambda|^{\frac{t}{h}+m} \left( \frac{h}{|1+h\lambda|-1} - \frac{1}{\operatorname{Re}(\lambda)} \right) \right) \\
&\leq \varepsilon \left( \frac{h}{1-|1+h\lambda|} + |1+h\lambda|^{2m} \left( \frac{h}{|1+h\lambda|-1} - \frac{1}{\operatorname{Re}(\lambda)} \right) \right),
\end{aligned}$$

as  $\frac{t}{h} \leq m$  and  $\frac{h}{|1+h\lambda|-1} - \frac{1}{\operatorname{Re}(\lambda)} > 0$  for  $\operatorname{Re}(\lambda) < 0$  with  $|1+h\lambda| > 1$  and  $h > 0$ .

(iv) (b). For  $\operatorname{Re}(\lambda) < 0$  with  $|1+h\lambda| > 1$  and  $t \in (hm, \infty)$ , using the same calculation as in case (iii)(b), we get

$$\begin{aligned}
|\phi(t) - x(t)| &\leq \varepsilon \left( \frac{-1}{\operatorname{Re}(\lambda)} + e^{\operatorname{Re}(\lambda)(t-hm)} \left( \frac{|1+h\lambda|^{2m}-1}{-\operatorname{Re}(\lambda)} + \frac{h(|1+h\lambda|^{2m}-1)}{|1+h\lambda|-1} \right) \right) \\
&\leq \varepsilon \left( \frac{h}{1-|1+h\lambda|} + |1+h\lambda|^{2m} \left( \frac{h}{|1+h\lambda|-1} - \frac{1}{\operatorname{Re}(\lambda)} \right) \right),
\end{aligned}$$

as  $t > hm$ ,  $\operatorname{Re}(\lambda) < 0$  with  $|1+h\lambda| > 1$ , and the expression inside the square brackets is positive.

We next show that  $x$  is the unique solution of (1) such that  $|\phi(t) - x(t)| \leq K\varepsilon$  for all  $t \in \mathbb{T}_{hm}$ , where  $K$  is given by (4). Suppose  $\phi : \mathbb{T}_{hm} \rightarrow \mathbb{C}$  is an approximate solution of (1) such that

$$|\phi^\Delta(t) - \lambda\phi(t)| \leq \varepsilon \text{ for all } t \in \mathbb{T}_{hm}$$

for some  $\varepsilon > 0$ . Suppose further that  $x_1, x_2 : \mathbb{T}_{hm} \rightarrow \mathbb{C}$  are two different solutions of (1) such that  $|\phi(t) - x_j(t)| \leq K\varepsilon$  for all  $t \in \mathbb{T}_{hm}$ , for  $j = 1, 2$ . Then, we have for constants  $c_j \in \mathbb{C}$  that

$$x_j(t) = c_j e_\lambda(t, 0), \quad c_1 \neq c_2,$$

and

$$|c_1 - c_2| \cdot |e_\lambda(t, 0)| = |x_1(t) - x_2(t)| \leq |x_1(t) - \phi(t)| + |\phi(t) - x_2(t)| \leq 2K\varepsilon;$$

letting  $t \rightarrow -\infty$  yields  $\infty < 2K\varepsilon$ , a contradiction. Consequently,  $x$  is the unique solution of (1) such that  $|\phi(t) - x(t)| \leq \varepsilon K$  for all  $t \in \mathbb{T}_{hm}$ . This completes the proof.  $\square$

**Theorem 3** Let  $\lambda \in \mathbb{C} \setminus \left\{ \frac{-1}{h} \right\}$  with  $\operatorname{Re}(\lambda) = 0$ . Then, (1) is not Hyers–Ulam stable on  $\mathbb{T}_{hm}$ .

**Proof** Assume  $\operatorname{Re}(\lambda) = 0$  for  $\lambda \in \mathbb{C}$ . Let arbitrary  $\varepsilon > 0$  be given, and let  $\lambda = i\beta$  for some  $\beta \in \mathbb{R}$ . Then,

$$\phi(t) := \frac{\varepsilon t e_{i\beta}(t, 0)}{(1 + h^2 \beta^2)^{\frac{m+1}{2}}}, \quad t \in \mathbb{T}_{hm}$$

satisfies the inequality

$$|\phi^\Delta(t) - i\beta\phi(t)| = \frac{\varepsilon |(1 + i\beta\mu(t))e_{i\beta}(t, 0)|}{(1 + h^2 \beta^2)^{\frac{m+1}{2}}} \leq \frac{\varepsilon |e_{i\beta}(t, 0)|}{(1 + h^2 \beta^2)^{\frac{m}{2}}} \leq \varepsilon$$

for all  $t \in \mathbb{T}_{hm}$ . Since  $x(t) = x_0 e_{i\beta}(t, 0)$  is the general solution of (1) when  $\lambda = i\beta$ , then

$$|\phi(t) - x(t)| = \frac{|e_{i\beta}(t, 0)|}{(1 + h^2 \beta^2)^{\frac{m+1}{2}}} \left| \varepsilon t - x_0 (1 + h^2 \beta^2)^{\frac{m+1}{2}} \right| \rightarrow \infty$$

as  $t \rightarrow \pm\infty$  for  $t \in \mathbb{T}_{hm}$  and for any  $x_0 \in \mathbb{C}$ ,  $\beta \in \mathbb{R}$ ,  $h > 0$ . So, (1) lacks HUS on  $\mathbb{T}_{hm}$  if  $\lambda = i\beta$ .  $\square$

Using the previous theorems, we can establish the following results.

**Theorem 4** Let  $\lambda \in \mathbb{C} \setminus \left\{ \frac{-1}{h} \right\}$ . Equation (1) has HUS on  $\mathbb{T}_{hm}$  if and only if  $\operatorname{Re}(\lambda) \neq 0$ .

**Proof** By Theorems 1, 2 and 3, we obtain the result, immediately.  $\square$

**Lemma 2** Let  $\lambda \in \mathbb{C} \setminus \left\{ \frac{-1}{h} \right\}$  with  $\operatorname{Re}(\lambda) \neq 0$ .

(i) If  $\operatorname{Re}(\lambda) > 0$ , then the HUS constant  $K$  for (1) satisfies

$$K \geq \frac{1}{\operatorname{Re}(\lambda)}.$$

(ii) If  $\operatorname{Re}(\lambda) < 0$ , then the HUS constant  $K$  for (1) satisfies

$$K \geq \begin{cases} \frac{-1}{\operatorname{Re}(\lambda)} + 2hm & : |1 + h\lambda| = 1 \\ \max \left\{ \frac{-1}{\operatorname{Re}(\lambda)}, \frac{h}{1 - |1 + h\lambda|} + |1 + h\lambda|^{2m} \left( \frac{h}{|1 + h\lambda| - 1} - \frac{1}{\operatorname{Re}(\lambda)} \right) \right\} & : |1 + h\lambda| \neq 1. \end{cases}$$

**Proof** Since  $\lambda \in \mathbb{C} \setminus \left\{ \frac{-1}{h} \right\}$  with  $\operatorname{Re}(\lambda) \neq 0$ , Equation (1) has HUS by Theorem 4. We will proceed by cases.

(i). Let  $\lambda = \alpha + i\beta \in \mathbb{C} \setminus \left\{ \frac{-1}{h} \right\}$ , and assume  $\operatorname{Re}(\lambda) = \alpha > 0$ ; set

$$\phi(t) := \frac{-\varepsilon e_{i\beta}(t, 0)}{\alpha (1 + h^2 \beta^2)^{\frac{m}{2}}} + \frac{\varepsilon}{\alpha} e_{\lambda}(t, 0).$$

It follows that

$$|\phi^{\Delta}(t) - \lambda\phi(t)| = \frac{|\varepsilon\alpha e_{i\beta}(t, 0)|}{\alpha (1 + h^2 \beta^2)^{\frac{m}{2}}} = \frac{\varepsilon |e_{i\beta}(t, 0)|}{(1 + h^2 \beta^2)^{\frac{m}{2}}} \leq \varepsilon.$$

Since  $x(t) = \frac{\varepsilon}{\alpha} e_{\lambda}(t, 0)$  is a solution of (1),

$$|\phi(t) - x(t)| = \frac{\varepsilon |e_{i\beta}(t, 0)|}{\alpha (1 + h^2 \beta^2)^{\frac{m}{2}}} \leq \frac{\varepsilon}{\alpha},$$

with equality at  $t = hm$ , so the minimal HUS constant  $K$  for (1) satisfies

$$K \geq \frac{1}{\alpha} = \frac{1}{\operatorname{Re}(\lambda)}.$$

This ends the proof of case (i).

(ii) (a). Assume  $\operatorname{Re}(\lambda) < 0$  with  $|1 + h\lambda| = 1$ . Let

$$\phi(t) = e_{\lambda}(t, 0) \int_0^t \frac{q(\tau)}{e_{\lambda}(\sigma(\tau), 0)} \Delta\tau, \quad q(\tau) = \frac{\varepsilon e_{\lambda}(\sigma(\tau), 0)}{|e_{\lambda}(\sigma(\tau), 0)|}, \quad (6)$$

for all  $t \in \mathbb{T}_{hm}$ . Then,

$$\phi^{\Delta}(t) - \lambda\phi(t) = q(t), \quad |q(t)| = \varepsilon,$$

and, employing (6), we see that  $\phi$  takes the form

$$\phi(t) = \varepsilon \begin{cases} \left( \frac{1}{\operatorname{Re}(\lambda)} - hm \right) e_{\lambda}(t, 0) - \frac{e^{i \operatorname{Im}(\lambda)(t+hm)}}{(1+h\lambda)^m \operatorname{Re}(\lambda)} & : t \in (-\infty, -hm) \\ t e_{\lambda}(t, 0) & : t \in \{-hm, \dots, hm\} \\ \left( \frac{1}{\operatorname{Re}(\lambda)} + hm \right) e_{\lambda}(t, 0) - \frac{e^{i \operatorname{Im}(\lambda)(t-hm)}}{(1+h\lambda)^{-m} \operatorname{Re}(\lambda)} & : t \in (hm, \infty). \end{cases}$$

If we take

$$x(t) := \varepsilon \left( \frac{1}{\operatorname{Re}(\lambda)} - hm \right) e_{\lambda}(t, 0),$$



then  $x$  is a solution of (1), and

$$|\phi(t) - x(t)| = \begin{cases} \varepsilon \left| \frac{-e^{i \operatorname{Im}(\lambda)(t+hm)}}{(1+h\lambda)^m \operatorname{Re}(\lambda)} \right| = \frac{-\varepsilon}{\operatorname{Re}(\lambda)} & : t \in (-\infty, -hm) \\ \varepsilon \left| -\frac{1}{\operatorname{Re}(\lambda)} + t + hm \right| \leq \varepsilon \left( \frac{-1}{\operatorname{Re}(\lambda)} + 2hm \right) & : t \in \{-hm, \dots, hm\} \\ \varepsilon \left| -\frac{1}{\operatorname{Re}(\lambda)} + 2hme^{\operatorname{Re}(\lambda)(t-hm)} \right| \leq \varepsilon \left( \frac{-1}{\operatorname{Re}(\lambda)} + 2hm \right) & : t \in (hm, \infty), \end{cases}$$

where we have equality at  $t = hm$ . This shows that the HUS constant  $K$  must satisfy

$$K \geq \left( \frac{-1}{\operatorname{Re}(\lambda)} + 2hm \right)$$

for  $\operatorname{Re}(\lambda) < 0$  with  $|1 + h\lambda| = 1$ . Here ends the proof of case (ii)(a).

(ii)(b). Assume  $\operatorname{Re}(\lambda) < 0$  with  $|1 + h\lambda| \neq 1$ . Again, let  $\phi$  be given by (6) for all  $t \in \mathbb{T}_{hm}$ . Then,

$$\phi^\Delta(t) - \lambda\phi(t) = q(t), \quad |q(t)| = \varepsilon,$$

and in this case  $\phi$  takes the form

$$\phi(t) = \varepsilon \begin{cases} \left( \frac{|1+h\lambda|^m}{\operatorname{Re}(\lambda)} - \frac{h(|1+h\lambda|^m-1)}{|1+h\lambda|-1} \right) e_\lambda(t, 0) - \frac{|1+h\lambda|^m e^{i \operatorname{Im}(\lambda)(t+hm)}}{(1+h\lambda)^m \operatorname{Re}(\lambda)} & : t \in (-\infty, -hm) \\ \frac{h \left( |1+h\lambda|^{\frac{t}{h}} - 1 \right)}{|1+h\lambda|^{\frac{t}{h}} (|1+h\lambda|-1)} e_\lambda(t, 0) & : t \in \{-hm, \dots, hm\} \\ \left( \frac{|1+h\lambda|^{-m}}{\operatorname{Re}(\lambda)} + \frac{h(|1+h\lambda|^m-1)}{|1+h\lambda|^m(|1+h\lambda|-1)} \right) e_\lambda(t, 0) - \frac{(1+h\lambda)^m e^{i \operatorname{Im}(\lambda)(t-hm)}}{|1+h\lambda|^m \operatorname{Re}(\lambda)} & : t \in (hm, \infty). \end{cases}$$

If we take

$$x(t) := \varepsilon \left( \frac{|1 + h\lambda|^m}{\operatorname{Re}(\lambda)} - \frac{h(|1 + h\lambda|^m - 1)}{|1 + h\lambda| - 1} \right) e_\lambda(t, 0), \tag{7}$$

then  $x$  is a solution of (1), and

$$|\phi(t) - x(t)| = \frac{-\varepsilon}{\operatorname{Re}(\lambda)}, \quad t \in (-\infty, -hm).$$

For  $t \in \{-hm, \dots, -h, 0, h, \dots, hm\}$ ,

$$|\phi(t) - x(t)| = \varepsilon \left( \frac{h \left( |1 + h\lambda|^{m+\frac{t}{h}} - 1 \right)}{|1 + h\lambda| - 1} - \frac{|1 + h\lambda|^{m+\frac{t}{h}}}{\operatorname{Re}(\lambda)} \right).$$

If  $0 < |1 + h\lambda| < 1$ , then as in the proof of Theorem 2 (iii)(a), we have

$$\begin{aligned} |\phi(t) - x(t)| &= \varepsilon \left( \frac{h}{1 - |1 + h\lambda|} + |1 + h\lambda|^{\frac{t}{h}+m} \left( \frac{h}{|1 + h\lambda| - 1} - \frac{1}{\operatorname{Re}(\lambda)} \right) \right) \\ &\leq \varepsilon \left( \frac{h}{1 - |1 + h\lambda|} + |1 + h\lambda|^{2m} \left( \frac{h}{|1 + h\lambda| - 1} - \frac{1}{\operatorname{Re}(\lambda)} \right) \right), \end{aligned}$$

as  $\frac{t}{h} \leq m$  and  $\frac{h}{|1+h\lambda|-1} - \frac{1}{\operatorname{Re}(\lambda)} \leq 0$  for  $\operatorname{Re}(\lambda) < 0$  with  $0 < |1 + h\lambda| < 1$  and  $h > 0$ , with equality at  $t = hm$ . If  $|1 + h\lambda| > 1$ , then as in the proof of Theorem 2 (iv)(a),

$$\begin{aligned} |\phi(t) - x(t)| &= \varepsilon \left( \frac{h}{1 - |1 + h\lambda|} + |1 + h\lambda|^{\frac{t}{h}+m} \left( \frac{h}{|1 + h\lambda| - 1} - \frac{1}{\operatorname{Re}(\lambda)} \right) \right) \\ &\leq \varepsilon \left( \frac{h}{1 - |1 + h\lambda|} + |1 + h\lambda|^{2m} \left( \frac{h}{|1 + h\lambda| - 1} - \frac{1}{\operatorname{Re}(\lambda)} \right) \right), \end{aligned}$$

as  $\frac{t}{h} \leq m$  and  $\frac{h}{|1+h\lambda|-1} - \frac{1}{\operatorname{Re}(\lambda)} > 0$  for  $\operatorname{Re}(\lambda) < 0$  with  $|1 + h\lambda| > 1$  and  $h > 0$ . For  $t \in (hm, \infty)$  and  $0 < |1 + h\lambda| < 1$ , then as in the proof of Theorem 2 (iii)(b),

$$\begin{aligned} |\phi(t) - x(t)| &= \varepsilon \left( \frac{-1}{\operatorname{Re}(\lambda)} + e^{\operatorname{Re}(\lambda)(t-hm)} \left( \frac{|1 + h\lambda|^{2m} - 1}{-\operatorname{Re}(\lambda)} + \frac{h(|1 + h\lambda|^{2m} - 1)}{|1 + h\lambda| - 1} \right) \right) \\ &\leq \varepsilon \left( \frac{h}{1 - |1 + h\lambda|} + |1 + h\lambda|^{2m} \left( \frac{h}{|1 + h\lambda| - 1} - \frac{1}{\operatorname{Re}(\lambda)} \right) \right), \end{aligned}$$

as  $t > hm$ ,  $\operatorname{Re}(\lambda) < 0$  with  $0 < |1 + h\lambda| < 1$ , and the expression inside the square brackets is non-negative. For  $t \in (hm, \infty)$  and  $|1 + h\lambda| > 1$ , then as in the proof of Theorem 2 (iv)(b), we have

$$\begin{aligned} |\phi(t) - x(t)| &= \varepsilon \left( \frac{-1}{\operatorname{Re}(\lambda)} + e^{\operatorname{Re}(\lambda)(t-hm)} \left( \frac{|1 + h\lambda|^{2m} - 1}{-\operatorname{Re}(\lambda)} + \frac{h(|1 + h\lambda|^{2m} - 1)}{|1 + h\lambda| - 1} \right) \right) \\ &\leq \varepsilon \left( \frac{h}{1 - |1 + h\lambda|} + |1 + h\lambda|^{2m} \left( \frac{h}{|1 + h\lambda| - 1} - \frac{1}{\operatorname{Re}(\lambda)} \right) \right), \end{aligned}$$

as  $t > hm$ ,  $\operatorname{Re}(\lambda) < 0$  with  $|1 + h\lambda| > 1$ , and the expression inside the square brackets is positive. This ends the proof of case (ii)(b), and thus the overall result holds.  $\square$

**Theorem 5** Let  $\lambda \in \mathbb{C} \setminus \{\frac{-1}{h}\}$ . If  $\operatorname{Re}(\lambda) \neq 0$ , then (1) has HUS on  $\mathbb{T}_{hm}$ .

(i) If  $\operatorname{Re}(\lambda) > 0$ , then

$$K = \frac{1}{\operatorname{Re}(\lambda)}$$

is the best (minimal) HUS constant.

(ii) If  $\operatorname{Re}(\lambda) < 0$ , then

$$K = \begin{cases} \frac{-1}{\operatorname{Re}(\lambda)} + 2hm & : |1 + h\lambda| = 1 \\ \max \left\{ \frac{-1}{\operatorname{Re}(\lambda)}, \frac{h}{1 - |1 + h\lambda|} + |1 + h\lambda|^{2m} \left( \frac{h}{|1 + h\lambda| - 1} - \frac{1}{\operatorname{Re}(\lambda)} \right) \right\} & : |1 + h\lambda| \neq 1 \end{cases}$$

is the best (minimal) HUS constant.

**Proof** This result follows immediately from the definitions of HUS and HUS constant, Theorems 1–4, and Lemma 2.  $\square$

*Remark 1* If  $m = 0$ , then  $\mathbb{T}_{h,0} = \mathbb{R}$ , and the results in Theorems 1 and 2 (i) – (iv) match exactly the known results for  $\mathbb{T} = \mathbb{R}$ , namely that  $x'(t) - \lambda x(t) = 0$  has HUS on  $\mathbb{R}$ , and

$$K = \frac{1}{|\operatorname{Re}(\lambda)|}$$

is the best possible HUS constant. If  $h \rightarrow 0$ , then  $\mathbb{T}_{0,m} = \mathbb{R}$ , and the results in Theorems 1 and 2 (i) – (iv) also recover the known results for  $\mathbb{T} = \mathbb{R}$ , because

$$\lim_{h \rightarrow 0^+} \left( \frac{-1}{\operatorname{Re}(\lambda)} + 2hm \right) = \frac{-1}{\operatorname{Re}(\lambda)}, \quad \lim_{h \rightarrow 0^+} \frac{h}{|1 + h\lambda| - 1} = \lim_{h \rightarrow 0^+} \frac{1}{\operatorname{Re}_h(\lambda)} = \frac{1}{\operatorname{Re}(\lambda)}$$

hold, where  $\operatorname{Re}_h(\lambda)$  represents the Hilger real part [20] for  $h$ -difference equations.

For fixed  $h > 0$ , if  $m \rightarrow \infty$ , then  $\mathbb{T}_{h,\infty} = h\mathbb{Z}$ , and the results in Theorem 1 and Theorem 2 (i) and (iii) match exactly the known results for  $\mathbb{T} = h\mathbb{Z}$ , namely that  $\Delta_h x(t) - \lambda x(t) = 0$  has HUS on  $h\mathbb{Z}$ , and

$$K = \frac{h}{|1 - |1 + h\lambda||} = \frac{1}{|\operatorname{Re}_h(\lambda)|}$$

is the best possible HUS constant. Theorem 2 (ii) shows an interesting connection; as  $m \rightarrow \infty$ , the HUS constant in (ii) goes to infinity as well. This is accurate, as  $\operatorname{Re}(\lambda) < 0$  with  $|1 + h\lambda| = 1$  makes the  $h$ -difference equation version of (1) Hyers–Ulam unstable on  $h\mathbb{Z}$ , as  $\lambda \in \mathbb{C}$  is then on the left Hilger circle [5]; see [15, Chapter 2.1], [20], and [23] for more on the Hilger complex plane, and [2, 5, 10] for more on the Hilger circle and HUS. On the other hand, in case (iv)  $\operatorname{Re}(\lambda) < 0$  with  $|1 + h\lambda| > 1$ , a result that does not match is obtained, that is,

$$\lim_{m \rightarrow \infty} \left( \frac{h}{1 - |1 + h\lambda|} + |1 + h\lambda|^{2m} \left( \frac{h}{|1 + h\lambda| - 1} - \frac{1}{\operatorname{Re}(\lambda)} \right) \right) = \infty$$

when  $\operatorname{Re}(\lambda) < 0$  with  $|1 + h\lambda| > 1$ , but, we know that  $\Delta_h x(t) - \lambda x(t) = 0$  has HUS when  $|1 + h\lambda| > 1$  (see [5]). Why does this logical gap occur? According to the information of Theorem 2.5 (ii) in [5], in this case, the unique solution  $x$

is determined when  $t \rightarrow \infty$ . As you can see from the claim of Theorem 2, even in this case, the unique solution  $x$  is determined by the information of  $t \rightarrow \infty$ . Therefore, we can say that the case  $\operatorname{Re}(\lambda) < 0$  with  $|1 + h\lambda| > 1$  is a distinguishing characteristic of Hyers–Ulam stability on this time scale with discrete core and continuous periphery. We explore this anomaly in the next section.

#### 4 Connection with $h$ -Difference Equations in the Case $|1 + h\lambda| > 1$

The following result is effective for clarifying the connection with the  $h$ -difference equation  $\Delta_h x(t) - \lambda x(t) = 0$  with  $|1 + h\lambda| > 1$ .

**Theorem 6** *Let  $\lambda \in \mathbb{C} \setminus \left\{ \frac{-1}{h} \right\}$  with  $\operatorname{Re}(\lambda) < 0$  and  $|1 + h\lambda| > 1$ . Let  $\varepsilon > 0$  be a fixed arbitrary constant, and let  $\phi$  be a function on  $\mathbb{T}_{hm}$  satisfying the inequality*

$$|\phi^\Delta(t) - \lambda\phi(t)| \leq \varepsilon, \quad t \in \mathbb{T}_{hm}.$$

*Then, the function  $x$  given by*

$$x(t) := \left( \frac{\phi(hm)}{e_\lambda(hm, 0)} \right) e_\lambda(t, 0)$$

*is a solution of (1) with*

$$|\phi(t) - x(t)| \leq \varepsilon \max \left\{ \frac{h(1 - |1 + h\lambda|^{-2m})}{|1 + h\lambda| - 1}, \frac{-1}{\operatorname{Re}(\lambda)} \right\}$$

*for all  $t \in \{-hm, \dots, -h, 0, h, \dots, hm\} \cup (hm, \infty)$ . In particular, the following holds.*

(i) *If  $t \in \{-hm, \dots, -h, 0, h, \dots, hm\}$ , then*

$$|\phi(t) - x(t)| \leq \varepsilon \left( \frac{h(1 - |1 + h\lambda|^{-2m})}{|1 + h\lambda| - 1} \right)$$

*for all  $t \in \{-hm, \dots, -h, 0, h, \dots, hm\}$ .*

(ii) *If  $t \in (hm, \infty)$ , then*

$$|\phi(t) - x(t)| \leq \varepsilon \left( \frac{-1}{\operatorname{Re}(\lambda)} \right)$$

*for all  $t \in (hm, \infty)$ .*

**Proof** Let  $\lambda \in \mathbb{C} \setminus \left\{ \frac{-1}{h} \right\}$  with  $\operatorname{Re}(\lambda) < 0$  and  $|1 + h\lambda| > 1$ . Suppose that  $|\phi^\Delta(t) - \lambda\phi(t)| \leq \varepsilon$  for all  $t \in \mathbb{T}_{hm}$ , there exists a function  $q : \mathbb{T}_{hm} \rightarrow \mathbb{C}$  such that

$$\phi^\Delta(t) - \lambda\phi(t) = q(t), \quad |q(t)| \leq \varepsilon$$

for all  $t \in \mathbb{T}_{hm}$ . Then, we have

$$\phi(t) = \phi_0 e_\lambda(t, 0) + e_\lambda(t, 0) \int_0^t \frac{q(\tau)}{e_\lambda(\sigma(\tau), 0)} \Delta\tau.$$

Let  $x(t) = x_0 e_\lambda(t, 0)$  be the solution of (1) with

$$x_0 := \phi_0 + \int_0^{hm} \frac{q(\tau)}{e_\lambda(\sigma(\tau), 0)} \Delta\tau \in \mathbb{C}.$$

It follows that

$$\phi(t) - x(t) = -e_\lambda(t, 0) \left( \int_t^{hm} \frac{q(\tau)}{e_\lambda(\sigma(\tau), 0)} \Delta\tau \right). \quad (8)$$

(a) For  $\operatorname{Re}(\lambda) < 0$  with  $|1 + h\lambda| > 1$  and  $t \in \{-hm, \dots, -h, 0, h, \dots, hm\}$ , using (8) we have

$$\begin{aligned} |\phi(t) - x(t)| &\leq \varepsilon |e_\lambda(t, 0)| \int_t^{hm} \frac{\Delta\tau}{|e_\lambda(\sigma(\tau), 0)|} = \varepsilon |1 + h\lambda|^{\frac{t}{h}} \sum_{j=\frac{t}{h}}^{m-1} \frac{h}{|1 + h\lambda|^{j+1}} \\ &= \varepsilon \left( \frac{h \left( 1 - |1 + h\lambda|^{-m + \frac{t}{h}} \right)}{|1 + h\lambda| - 1} \right) \leq \varepsilon \left( \frac{h \left( 1 - |1 + h\lambda|^{-2m} \right)}{|1 + h\lambda| - 1} \right), \end{aligned}$$

as  $\frac{t}{h} \leq m$ ,  $h > 0$ , and  $\operatorname{Re}(\lambda) < 0$  with  $|1 + h\lambda| > 1$ .

(b) For  $\operatorname{Re}(\lambda) < 0$  with  $|1 + h\lambda| > 1$  and  $t \in [hm, \infty)$ , using (8) we have

$$\begin{aligned} |\phi(t) - x(t)| &= \left| e_\lambda(t, 0) \left( \int_{hm}^t \frac{q(\tau)}{e_\lambda(\sigma(\tau), 0)} \Delta\tau \right) \right| \\ &\leq \varepsilon |1 + h\lambda|^m e^{\operatorname{Re}(\lambda)(t-hm)} \int_{hm}^t \frac{\Delta\tau}{|e_\lambda(\sigma(\tau), 0)|} \\ &= \varepsilon \left( \frac{-1 + e^{\operatorname{Re}(\lambda)(t-hm)}}{\operatorname{Re}(\lambda)} \right) \leq \varepsilon \left( \frac{-1}{\operatorname{Re}(\lambda)} \right), \end{aligned}$$

as  $t \geq hm$  and  $\operatorname{Re}(\lambda) < 0$  with  $|1 + h\lambda| > 1$ . This completes the proof.  $\square$

*Remark 2* For fixed  $h > 0$ , if  $m \rightarrow \infty$ , then  $\mathbb{T}_{h,\infty} = h\mathbb{Z}$ , and the result in Theorem 6 (i) reproduces exactly the known result for  $\mathbb{T} = h\mathbb{Z}$  (see, Theorem 2.5 (ii) in [5]). Actually, we will explain this fact. Since

$$\left| \int_0^{hm} \frac{q(\tau)}{e_\lambda(\sigma(\tau), 0)} \Delta\tau \right| \leq \varepsilon \int_0^{hm} \frac{\Delta\tau}{|e_\lambda(\sigma(\tau), 0)|} \leq \varepsilon \sum_{j=0}^{m-1} \frac{h}{|1+h\lambda|^{j+1}} < \frac{h}{|1+h\lambda| - 1}$$

holds, we see that

$$\lim_{m \rightarrow \infty} \frac{\phi(hm)}{e_\lambda(hm, 0)} = \lim_{t \rightarrow \infty} \frac{\phi(t)}{e_\lambda(t, 0)}$$

exists. In addition, an HUS constant is

$$\lim_{m \rightarrow \infty} \frac{h(1 - |1 + h\lambda|^{-2m})}{|1 + h\lambda| - 1} = \frac{h}{|1 + h\lambda| - 1}.$$

Theorem 6 (i) says that the function  $x$  given by

$$x(t) := \left( \lim_{t \rightarrow \infty} \frac{\phi(t)}{e_\lambda(t, 0)} \right) e_\lambda(t, 0)$$

is a solution of (1) with

$$|\phi(t) - x(t)| \leq \frac{h\varepsilon}{|1 + h\lambda| - 1}$$

for all  $t \in h\mathbb{Z}$ . As  $m \rightarrow \infty$ , our exponential function  $e_\lambda(t, 0)$  corresponds to  $(1 + h\lambda)^{\frac{t}{h}}$  for all  $t \in h\mathbb{Z}$ . In this case, we can prove the uniqueness of the solution. Let  $\varepsilon > 0$ , and let  $\phi : \{-hm, \dots, -h, 0, h, \dots, hm\} \rightarrow \mathbb{C}$  satisfy

$$|\phi^\Delta(t) - \lambda\phi(t)| \leq \varepsilon \text{ for all } t \in \{-hm, \dots, -h, 0, h, \dots, hm\}.$$

Suppose that  $x_1, x_2 : \{-hm, \dots, -h, 0, h, \dots, hm\} \rightarrow \mathbb{C}$  are two different solutions of (1) such that  $|\phi(t) - x_j(t)| \leq K\varepsilon := \frac{h}{|1+h\lambda|-1}\varepsilon$  for all  $t \in \{-hm, \dots, -h, 0, h, \dots, hm\}$ , for  $j = 1, 2$ . Then, we have for constants  $c_j \in \mathbb{C}$  that

$$x_j(t) = c_j(1 + h\lambda)^{\frac{t}{h}}, \quad c_1 \neq c_2,$$

and

$$|c_1 - c_2||1 + h\lambda|^{\frac{t}{h}} = |x_1(t) - x_2(t)| \leq |x_1(t) - \phi(t)| + |\phi(t) - x_2(t)| \leq 2K\varepsilon;$$

letting  $m \rightarrow \infty$  and  $t \rightarrow \infty$  yields  $\infty < 2K\varepsilon$ , a contradiction. Consequently,  $x$  is the unique solution of  $\Delta_h x(t) - \lambda x(t) = 0$  such that  $|\phi(t) - x(t)| \leq \varepsilon K$  for all  $t \in h\mathbb{Z}$ .

## 5 Conclusion and Future Directions

In this paper we determined the best Hyers–Ulam stability constants for a first-order complex constant coefficient dynamic equation on a time scale with a discrete core and continuous periphery. In the future, we will study a time scale with a discrete periphery and continuous core, whose exponential function for  $\lambda \in \mathbb{C} \setminus \{-\frac{1}{h}\}$  is

$$e_\lambda(t, 0) := \begin{cases} (1 + h\lambda)^{\frac{t}{h}+m} e^{-hm\lambda} & : t \in \{\dots, -h(m+2), -h(m+1)\} \\ e^{\lambda t} & : t \in [-hm, hm] \\ (1 + h\lambda)^{\frac{t}{h}-m} e^{hm\lambda} & : t \in \{h(m+1), h(m+2), \dots\}. \end{cases}$$

on  $\mathbb{T}_{hm} = \{\dots, -h(m+2), -h(m+1)\} \cup [-hm, hm] \cup \{h(m+1), h(m+2), \dots\}$ , where  $h > 0$  is the discrete step size and  $m$  is a non-negative integer.

**Acknowledgments** The second author was supported by JSPS KAKENHI Grant Number JP20K03668.

## References

1. D.R. Anderson, Hyers–Ulam stability for a first-order linear proportional nabla difference operator, in *Frontiers in Functional Equations and Analytic Inequalities*, ed. by G. Anastassiou, J.M. Rassias (Springer Nature Switzerland AG, Cham, 2019)
2. D.R. Anderson, The discrete diamond-alpha imaginary ellipse and Hyers–Ulam stability. *Int. J. Difference Equations* **14**(1), 25–38 (2019)
3. D.R. Anderson, M. Onitsuka, Hyers–Ulam stability of first-order homogeneous linear dynamic equations on time scales. *Demonstratio Math.* **51**, 198–210 (2018)
4. D.R. Anderson, M. Onitsuka, Hyers–Ulam stability for a discrete time scale with two step sizes. *Appl. Math. Comput.* **344–345**, 128–140 (2019)
5. D.R. Anderson, M. Onitsuka, Best constant for Hyers–Ulam stability of second-order  $h$ -difference equations with constant coefficients. *Results Math* **74**, 151 (2019). <https://doi.org/10.1007/s00025-019-1077-9>
6. D.R. Anderson, M. Onitsuka, Hyers–Ulam stability and best constant for Cayley  $h$ -difference equations. *Bull. Malaysian Math. Sci. Soc.* (2020). <https://doi.org/10.1007/s40840-020-00920-z>
7. D.R. Anderson, M. Onitsuka, Hyers–Ulam stability for quantum equations of Euler type. *Discrete Dyn. Nature Soc.* **2020**, Article ID 5626481, 10 pp. (2020). <https://doi.org/10.1155/2020/5626481>
8. D.R. Anderson, M. Onitsuka, Hyers–Ulam stability for quantum equations. *Aequationes Math.* (2020). <https://doi.org/10.1007/s00010-020-00734-1>

9. D.R. Anderson, M. Onitsuka, Best constant for Hyers–Ulam stability of two step sizes linear difference equations. *J. Math. Anal. Appl.* (2021). <https://doi.org/10.1016/j.jmaa.2020.124807>
10. D.R. Anderson, A.J. Jennissen, C.J. Montplaisir, Hyers–Ulam stability for a continuous time scale with discrete uniform jumps. *Int. J. Difference Equations* **15**(2), 1–21 (2020)
11. D.R. Anderson, M. Onitsuka, J.M. Rassias, Best constant for Ulam stability of first-order  $h$ -difference equations with periodic coefficient. *J. Math. Anal. Appl.* **491**, 124363 (2020). <https://doi.org/10.1016/j.jmaa.2020.124363>
12. S. András, A.R. Mészáros, Ulam–Hyers stability of dynamic equations on time scales via Picard operators. *Appl. Math. Computation*. **219**, 4853–4864 (2013)
13. A.R. Baias, D. Popa, On Ulam stability of a linear difference equation in Banach spaces. *Bull. Malays. Math. Sci. Soc.* (2019). <https://doi.org/10.1007/s40840-019-00744-6>
14. A.R. Baias, F. Blaga, D. Popa, On the best Ulam constant of a first order linear difference equation in Banach spaces. *Acta Math. Hungar.* (2020). <https://doi.org/10.1007/s10474-020-01098-3>
15. M. Bohner, A. Peterson, *Dynamic Equations on Time Scales, An Introduction with Applications* (Birkhäuser, Boston, 2001)
16. J. Brzdęk, P. Wójcik, On approximate solutions of some difference equations. *Bull. Australian Math. Soc.* **95**(3), 76–481 (2017)
17. J. Brzdęk, D. Popa, I. Raşa, B. Xu, *Ulam Stability of Operators*, a volume in *Mathematical Analysis and Its Applications* (Academic Press, New York, 2018)
18. C. Buşe, D. O'Regan, O. Sailerli, Hyers–Ulam stability for linear differences with time dependent and periodic coefficients. *Symmetry* **11**, 512 (2019). <https://doi.org/10.3390/sym11040512>
19. C. Buşe, V. Lupulescu, D. O'Regan, Hyers–Ulam stability for equations with differences and differential equations with time-dependent and periodic coefficients. *Proc. R. Soc. Edinburgh A Math.* **150**(5), 2175–2188 (2020). <https://doi.org/10.1017/prm.2019.12>
20. S. Hilger, Special functions, Laplace and Fourier transform on measure chains. *Dynamic Syst. Appl.* **8**(3–4), 471–488 (1999)
21. L. Hua, Y. Li, J. Feng, On Hyers–Ulam stability of dynamic integral equation on time scales. *Mathematica Aeterna* **4**(6), 559–571 (2014)
22. D.H. Hyers, On the stability of the linear functional equation. *Proc. Nat. Acad. Sci. U. S. A.* **27**, 222–224 (1941)
23. B.J. Jackson, J.M. Davis, D. Poulsen, A polar representation of the Hilger complex plane. *Int. J. Difference Equations* **15**(2), 419–427 (2020)
24. S.-M. Jung, Y.W. Nam, Hyers–Ulam stability of Pielou logistic difference equation. *J. Nonlinear Sci. Appl.* **10**, 3115–3122 (2017)
25. S.-M. Jung, Y.W. Nam, Hyers–Ulam stability of the first order inhomogeneous matrix difference equation. *J. Comput. Anal. Appl.* **23**(8), 1368–1383 (2017)
26. Y.W. Nam, Hyers–Ulam stability of hyperbolic Möbius difference equation. *Filomat* **32**(13), 4555–4575 (2018). <https://doi.org/10.2298/FIL1813555N>
27. Y.W. Nam, Hyers–Ulam stability of elliptic Möbius difference equation. *Cogent Math. Stat.* **5**(1), 1–9 (2018)
28. Y.W. Nam, Hyers–Ulam stability of loxodromic Möbius difference equation. *Appl. Math. Comput.* **356**(1), 119–136 (2019). <https://doi.org/10.1016/j.amc.2019.03.033>
29. M. Onitsuka, Influence of the step size on Hyers–Ulam stability of first-order homogeneous linear difference equations. *Int. J. Difference Equations* **12**(2), 281–302 (2017)
30. M. Onitsuka, Hyers–Ulam stability of second-order nonhomogeneous linear difference equations with a constant step size. *J. Comput. Anal. Appl.* **28**(1), 152–165 (2020)
31. D. Popa, Hyers–Ulam stability of the linear recurrence with constant coefficients. *Adv. Differential Equations* **2005**, 407076 (2005)
32. D. Popa, Hyers–Ulam–Rassias stability of a linear recurrence. *J. Math. Anal. Appl.* **309**, 591–597 (2005)
33. H. Rasouli, S. Abbaszadeh, M. Eshaghi, Approximately linear recurrences. *J. Appl. Anal.* **24**(1), 81–85 (2018)



34. Th.M. Rassias, On the stability of linear mapping in Banach spaces. *Proc. Amer. Math. Soc.* **72**, 297–300 (1978)
35. Y.H. Shen, The Ulam stability of first order linear dynamic equations on time scales. *Results Math.* **72**(4), 1881–1895 (2017)
36. Y.H. Shen, Y.J. Li, Hyers–Ulam stability of first order nonhomogeneous linear dynamic equations on time scales. *Commun. Math. Res.* **35**(2), 139–148 (2019). <https://doi.org/10.13447/j.1674-5647.2019.02.05>
37. S.M. Ulam, *A Collection of the Mathematical Problems* (Interscience, New York, 1960)
38. B. Xu, J. Brzdęk, Hyers–Ulam stability of a system of first order linear recurrences with constant coefficients. *Discrete Dyn. Nat. Soc.* **2015**, Article ID 269356, 5 pp. (2015)

# Invariance Solutions and Blow-Up Property for Edge Degenerate Pseudo-Hyperbolic Equations in Edge Sobolev Spaces



Carlo Cattani and Morteza Koozehgar Kalleji

**Abstract** This article is dedicated to study of the initial-boundary value problem of edge pseudo-hyperbolic system with damping term on the manifold with edge singularity. First, we will discuss about the invariance of solution set of a class of edge degenerate pseudo-hyperbolic equations on the edge Sobolev spaces. Then, by using a family of modified potential wells and concavity methods, it is obtained existence and nonexistence results of global solutions with exponential decay and is shown the blow-up in finite time of solutions on the manifold with edge singularities.

## 1 Introduction

Initial-boundary value problems written for hyperbolic semilinear partial differential equations emerged in several applications to physics, mechanics and engineering sciences [9, 24, 25]. Interesting phenomena are often connected with geometric singularities, for instance, in mechanics or cracks in a medium are described by hypersurfaces with a boundary. In this cases, configurations of that kind belong to the category of spaces (manifolds) with geometric singularities, here with edges. Also, when one asks physics to calculate the self-energy of an electron, or the structure of space time at the center of a black hole, one encounter with mathematical bad behaviour, that is the singularities from the point view of mathematics. In recent years, from a mathematical point of view, the analysis on such (in general, stratified) spaces has become a mathematical structure theory with many deep relations with geometry, topology, and mathematical physics [10, 15, 23, 25]. In [21], Melrose, Vasy and Wunsch investigated the geometric propagation and

---

C. Cattani (✉)  
Engineering School, University of Tuscia, Viterbo, Italy  
e-mail: [cattani@unitus.it](mailto:cattani@unitus.it)

M. K. Kalleji  
Department of Mathematics, Faculty of Sciences, Arak University, Arak, Iran  
e-mail: [m-koozehgarkalleji@araku.ac.ir](mailto:m-koozehgarkalleji@araku.ac.ir)

diffraction of singularities of solutions to the wave equation on manifolds with edge singularities. Let  $X$  be an  $n$ -dimensional manifold with boundary, where the boundary  $\partial X$  is endowed with a fibration  $Z \rightarrow \partial X$  and  $\partial X \rightarrow Y$  where  $Y, Z$  are without boundary. By an edge metric  $g$  on  $X$ , we mean a metric  $g$  on the interior of  $X$  which is a smooth 2-cotensor up to the boundary but which degenerates there in a way compatible with the fibration. A manifold with boundary equipped with such an edge metric also is called an edge manifold or a manifold with edge structure. If  $Z$  is point, then an edge metric on  $X$  is simply a metric in the usual sense, smooth up to the boundary, while if  $Y$  is a point,  $X$  is conic manifold [4]. A simple example of a more general edge metric is obtained by performing a real blowup on a submanifold  $B$  of a smooth, boundaryless manifold  $A$ . The blowup operation simply introduces polar coordinates near  $B$ , i.e., it replaces  $B$  by its spherical normal bundle, thus yielding a manifold  $X$  with boundary. The pullback of a smooth metric on  $A$  to  $X$  is then an edge metric [21].

Up to now, elliptic boundary value problems in domains with point singularities have been thoroughly investigated [1–4, 7, 8, 14]. The natures of the solutions to these equations have been investigated by several means. For instance, problems with the Dirichlet boundary conditions were investigated in [1, 2, 7, 10, 14] in which the unique existence, the multiplicity, the regularity and the asymptotic behaviour near the conical points of the solutions are established. Finite time blow-up of solutions of generalized hyperbolic equations have been studied by many authors [1, 2, 5, 7, 18, 28]. In these references, the authors consider problems either for negative energy or for weaker conditions than a condition of negative initial energy. Other authors have assumed a condition of positive energy under other two conditions on the initial functions. However, the mentioned authors have not studied the compatibility of these conditions, which is come times hard to understand. These authors have used the classic concavity Levine’s method [17]. In this article, we use the edge Sobolev inequality and Poincaré inequality and modified method in [7, 8] to prove on the global well-posedness of solutions to initial-boundary value problems for semilinear degenerate pseudo-hyperbolic equations with dissipative term on manifolds with edge singularities. More precisely, we study the following initial-boundary value problem for semilinear hyperbolic equation

$$\begin{cases} \partial_t^2 u - \Delta_{\mathbb{E}} u + V(z)u + \gamma \Delta_{\mathbb{E}} \partial_t u = g_t(z)|u|^{p-1}u, & z \in \text{int}\mathbb{E}, t > 0, \\ u(z, 0) = u_0(z), \quad \partial_t u(z, 0) = u_1(z), & z \in \text{int}\mathbb{E} \\ u(z, t) = 0, & z \in \partial\mathbb{E}, t \geq 0, \end{cases} \quad (1)$$

where,  $2 < p + 1 < \frac{2n}{n-2} = 2^*$  is the critical cone Sobolev exponents,  $z = (r, x, y)$ ,  $u = u(z, t)$  is unknown function and  $\gamma$  is a non-negative parameter. Also,  $u_0 \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E})$ ,  $u_1 \in L_2^{\frac{n+1}{2}}(\mathbb{E})$ ,  $N = 1 + n + q \geq 3$  is a dimension of  $\mathbb{E}$  and coordinates  $z := (r, x, y) = (r, x_1, \dots, x_n, y_1, \dots, y_q) \in \mathbb{E}$ . Here the domain  $\mathbb{E}$  is  $[0, 1) \times X \times Y$ ,  $X$  is an  $(n - 1)$ -dimensional closed compact manifold,  $Y \subset \mathbb{R}^q$  is a bounded domain, which is regarded as the local model near the edge points on manifolds with edge singularities, and  $\partial\mathbb{E} = \{0\} \times X \times Y$ . Moreover, the operator  $\Delta_{\mathbb{E}}$  in 1 is

defined by  $(r\partial_r)^2 + \partial_{x_1}^2 + \dots + \partial_{x_n}^2 + (r\partial_{y_1})^2 + \dots + (r\partial_{y_q})^2$ , which is an elliptic operator with totally characteristic degeneracy on the boundary  $r = 0$ , we also call it Fuchsian type edge-Laplace operator, and the corresponding gradient operator by  $\nabla_{\mathbb{E}} := (r\partial_r, \partial_{x_1}, \dots, \partial_{x_n}, r\partial_{y_1}, \dots, r\partial_{y_q})$ . In the Equation 1, we assume that  $V(z) \in L^{\frac{n+1}{4}}(\text{int}\mathbb{E}) \cap C(\text{int}\mathbb{E})$  is a positive potential function such that  $\inf_{z \in \mathbb{E}} V(z) > 0$ . For every  $t \geq 0$ , we suppose that  $g_t : \mathbb{E} \rightarrow \mathbb{R}$  is a non-negative function which  $g_t(z) := g(z, t)$  for every  $z \in \text{int}\mathbb{E}$  and  $g(z, t) \in L^\infty(\text{int}\mathbb{E}) \cap C^1(\text{int}\mathbb{E})$ . The through of this paper we consider the following constants:

$$C_* = \inf \left\{ \frac{\|\sqrt{V(z)}u(z)\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}}{\|\nabla_{\mathbb{E}}u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}} ; u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}) \right\},$$

$$C_{**} = \sup \left\{ \frac{\|g_t(z)^{\frac{1}{p+1}}u\|_{L_{p+1}^{\frac{n+1}{2}}(\mathbb{E})}}{\|\nabla_{\mathbb{E}}u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}} ; u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}) \right\}.$$

Our investigation is in fact provoked by the study of [20] and we shall apply a potential method which was established by Sattinger [26]. So based on edge Sobolev spaces [10, 27], we study the existence and non-existence global weak solutions for semilinear pseudo-hyperbolic differential equations with respect to variable time with a positive potential function and a non-negative weighted function. The well-known operator  $(\Delta_{\mathbb{E}} + V(x) + \Delta_{\mathbb{E}}\partial_t)u$  and other special types of it (see [11]) appears naturally in the nonlinear heat and wave equations [25], nonlinear Schrödinger equation with potential function [12] and the references therein for a complete description of the model. In the sitting of parabolic type system, the authors [6, 18] studied global existence, exponential decay and finite time blow-up of solutions for a class of semilinear pseudo-parabolic equations with conical degeneration. Also, our problem can be seen as a class of degenerate hyperbolic type equations in case that  $V(z) = 0$  and  $g_t(z) \equiv 1$  then the problem 1 is reduced to problem 1.1 in [13] and in the classical sense our problem include the classical problem

$$\begin{cases} \partial_t^2 u - \Delta u + \gamma \partial_t u = f(u), & x \in \Omega, t > 0, \\ u(x, 0) = u_0(x), \quad \partial_t u(x, 0) = u_1(x), & x \in \Omega \\ u(t, x) = 0, & x \in \partial\Omega, t \geq 0, \end{cases} \quad (2)$$

where  $\Omega$  is bounded domain of  $\mathbb{R}^n$  with smooth boundary  $\partial\Omega$  and  $\Delta$  is the standard Laplace operator and  $f$  is a suitable function [13, 17, 19]. It is well-known that problem 2 has been studied by many authors, for example [19, 20] and the references therein.

In Section 2, we recall the definition of the edge Sobolev space and the corresponding properties. In Section 3, we will give some properties of potential

wells for problem 1 on the manifold with edge singularity, which is very useful in the process of our main results. In Section 4, we give the proofs of the results of global existence and non-existence, exponential decay and finite time blowing-up of problem 1.

## 2 Edge Sobolev Spaces

Consider  $X$  as a closed compact  $C^\infty$ -manifold of dimension  $n$  of the unit sphere in  $\mathbb{R}^{n+1}$ . We define an infinite cone in  $\mathbb{R}^{n+1}$  as a quotient space  $X^\Delta = \frac{\bar{\mathbb{R}}_+ \times X}{\{0\} \times X}$ , with base  $X$ . The cylindrical coordinates  $(r, \theta) \in X^\Delta - \{0\}$  in  $\mathbb{R}^{n+1} - \{0\}$  are the standard coordinates. This gives us the description of  $X^\Delta - \{0\}$  in the form  $\mathbb{R}_+ \times X$ . Then the stretched cone can be defined as  $\bar{\mathbb{R}}_+ \times X = X^\wedge$ . Now, consider  $B = X^\Delta$  with a conical point, then by the similar way in [8, 10, 27], one can define the stretched manifold  $\mathbb{B}$  with respect to  $B$  as a  $C^\infty$ -manifold with smooth boundary  $\partial\mathbb{B} \cong X(0)$ , where  $X(0)$  is the cross section of singular point zero such that there is a diffeomorphism  $B - \{0\} \cong \mathbb{B} - \partial\mathbb{B}$ , the restriction of which to  $U - \{0\} \cong V - \partial\mathbb{B}$  for an open neighborhood  $U \subset B$  near the conic point zero and a collar neighborhood  $V \subset \mathbb{B}$  with  $V \cong [0, 1) \times X(0)$ . Therefore, we can take  $\mathbb{B} = [0, 1) \times X \subset \bar{\mathbb{R}}_+ \times X = X^\wedge$ . In order to consider another type of a manifold with singularity of order one so-called wedge manifold, we consider a bounded domain  $Y$  in  $\mathbb{R}^q$ . Set  $W = X^\Delta \times Y = B \times Y$ . Then  $W$  is a corresponding wedge in  $\mathbb{R}^{1+n+q}$ . Therefore, the stretched wedge manifold  $\mathbb{W}$  to  $W$  is  $X^\wedge \times Y$  which is a manifold with smooth boundary  $\{0\} \times X \times Y$ . Set  $(r, x) \in X^\wedge$ . In order to define a finite wedge, it sufficient to consider the case  $r \in [0, 1)$ . Thus, we define a finite wedge as

$$E = \frac{[0, 1) \times X}{\{0\} \times X} \times Y \subset X^\Delta \times Y = W.$$

The stretched wedge manifold with respect to  $E$  is

$$\mathbb{E} = [0, 1) \times X \times Y = \mathbb{B} \times Y \subset X^\wedge \times Y = W^\wedge,$$

with smooth boundary  $\partial\mathbb{E} = \{0\} \times X \times Y$ .

**Definition 1** For  $(r, x, y) \in \mathbb{R}_+^N$  with  $N = 1 + n + q$ , assume that  $u(r, x, y) \in \mathcal{D}'(\mathbb{R}_+^N)$ . We say that  $u(r, x, y) \in L_p(\mathbb{R}_+^N; d\mu)$  if

$$\|u\|_{L_p} = \left( \int_{\mathbb{R}_+^N} r^N |u(r, x, y)|^p d\mu \right)^{\frac{1}{p}} < +\infty,$$

where  $d\mu = \frac{dr}{r} dx_1 \dots dx_n \frac{dy_1}{r} \dots \frac{dy_q}{r}$  and for  $1 \leq p < \infty$ .

Moreover, the weighted  $L_p$  spaces with wight  $\gamma \in \mathbb{R}$  is denoted by  $L_p^\gamma(\mathbb{R}_+^N; d\mu)$ , which consists of function  $u(r, x, y)$  such that

$$\|u\|_{L_p^\gamma} = \left( \int_{\mathbb{R}_+^N} r^N |r^{-\gamma} u(r, x, y)|^p d\mu \right)^{\frac{1}{p}} < +\infty.$$

Now, we can define the weighted  $p$ -Sobolev spaces with natural scale for all  $1 \leq p < \infty$  on  $\mathbb{R}_+^{N=1+n+q}$ .

**Definition 2** For  $m \in \mathbb{N}$ ,  $\gamma \in \mathbb{R}$  and  $N = 1 + n + q$ , the spaces

$$\mathcal{H}_p^{m,\gamma}(\mathbb{R}_+^N) = \left\{ u \in \mathcal{D}'(\mathbb{R}_+^N) \mid r^{\frac{N}{p}-\gamma} (r\partial_r)^k \partial_x^\alpha (r\partial_y)^\beta u \in L_p(\mathbb{R}_+^N; d\mu) \right\}$$

for  $k \in \mathbb{N}$ , multi-indices  $\alpha \in \mathbb{N}^n$  and  $\beta \in \mathbb{N}^q$  with  $k + |\alpha| + |\beta| \leq m$ . In other words, if  $u(r, x, y) \in \mathcal{H}_p^{m,\gamma}(\mathbb{R}_+^N)$  then  $(r\partial_r)^k \partial_x^\alpha (r\partial_y)^\beta u \in L_p^\gamma(\mathbb{R}_+^N; d\mu)$ . Therefore,  $\mathcal{H}_p^{m,\gamma}(\mathbb{R}_+^N)$  is a Banach space with the following norm

$$\|u\|_{\mathcal{H}_p^{m,\gamma}(\mathbb{R}_+^N)} = \sum_{k+|\alpha|+|\beta|\leq m} \left( \int_{\mathbb{R}_+^N} r^N |r^{-\gamma} (r\partial_r)^k \partial_x^\alpha (r\partial_y)^\beta u|^p d\mu \right)^{\frac{1}{p}}.$$

Moreover, the subspace  $\mathcal{H}_{p,0}^{m,\gamma}(\mathbb{R}_+^N)$  of  $\mathcal{H}_p^{m,\gamma}(\mathbb{R}_+^N)$  denotes the closure of  $C_0^\infty(\mathbb{R}_+^N)$  in  $\mathcal{H}_p^{m,\gamma}(\mathbb{R}_+^N)$ . Now, similarly to the definitions above, we can introduce the following weighted  $p$ -Sobolev spaces on  $X^\wedge \times Y$ , where  $X^\wedge = \mathbb{R}_+ \times X$  and  $X^\wedge \times Y$  is an open stretched wedge.

$$\mathcal{H}_p^{m,\gamma}(X^\wedge \times Y) := \left\{ u \in \mathcal{D}'(X^\wedge \times Y) \mid r^{\frac{N}{p}-\gamma} (r\partial_r)^k \partial_x^\alpha (r\partial_y)^\beta u \in L_p(X^\wedge \times Y; d\mu) \right\}$$

for  $k \in \mathbb{N}$ , multi-indices  $\alpha \in \mathbb{N}^n$  and  $\beta \in \mathbb{N}^q$  with  $k + |\alpha| + |\beta| \leq m$ .

Then  $\mathcal{H}_p^{m,\gamma}(X^\wedge \times Y)$  is a Banach space with the following norm

$$\|u\|_{\mathcal{H}_p^{m,\gamma}(X^\wedge \times Y)} = \sum_{k+|\alpha|+|\beta|\leq m} \left( \int_{X^\wedge \times Y} r^N |r^{-\gamma} (r\partial_r)^k \partial_x^\alpha (r\partial_y)^\beta u|^p d\mu \right)^{\frac{1}{p}}.$$

The subspace  $\mathcal{H}_{p,0}^{m,\gamma}(X^\wedge \times Y)$  of  $\mathcal{H}_p^{m,\gamma}(X^\wedge \times Y)$  is defined as the closure of  $C_0^\infty(X^\wedge \times Y)$ .

**Definition 3** Let  $\mathbb{E}$  be the stretched wedge to the finite wedge  $E$ , then  $\mathcal{H}_p^{m,\gamma}(\mathbb{E})$  for  $m \in \mathbb{N}$ ,  $\gamma \in \mathbb{R}$  denotes the subset of all  $u \in W_{loc}^{m,p}(int\mathbb{E})$  such that  $\omega u \in$

$\mathcal{H}_p^{m,\gamma}(X^\wedge \times Y)$  for any cut-off function  $\omega$ , supported by a collar neighborhood of  $(0, 1) \times \partial\mathbb{E}$ . Moreover, the subspace  $\mathcal{H}_{p,0}^{m,\gamma}(\mathbb{E})$  of  $\mathcal{H}_p^{m,\gamma}(\mathbb{E})$  is defined as follows

$$\mathcal{H}_{p,0}^{m,\gamma}(\mathbb{E}) := [\omega]\mathcal{H}_{p,0}^{m,\gamma}(X^\wedge \times Y) + [1 - \omega]W_0^{m,p}(int\mathbb{E})$$

where the classical Sobolev space  $W_0^{m,p}(int\mathbb{E})$  denotes the closure of  $C_0^\infty(int\mathbb{E})$  in  $W^{m,p}(\tilde{\mathbb{E}})$  for  $\tilde{\mathbb{E}}$  that is a closed compact  $C^\infty$  manifold with boundary.

If  $u \in L_{p'}^{\frac{n+1}{p}}(\mathbb{E})$  and  $v \in L_{p'}^{\frac{n+1}{p}}(\mathbb{E})$  with  $p, p' \in (1, \infty)$  and  $\frac{1}{p} + \frac{1}{p'} = 1$ , then one can obtain the following edge type Hölder inequality

$$\int_{\mathbb{E}} r^q |uv| d\mu \leq \left( \int_{\mathbb{E}} r^q |u|^p d\mu \right)^{\frac{1}{p}} \left( \int_{\mathbb{E}} r^q |v|^{p'} d\mu \right)^{\frac{1}{p'}}.$$

In the case  $p = 2$ , we have the corresponding edge type Schwartz inequality

$$\int_{\mathbb{E}} r^q |uv| d\mu \leq \left( \int_{\mathbb{E}} r^q |u|^2 d\mu \right)^{\frac{1}{2}} \left( \int_{\mathbb{E}} r^q |v|^2 d\mu \right)^{\frac{1}{2}}.$$

In the sequel, for convenience we denote

$$(u, v)_2 = \int_{\mathbb{E}} r^q uv d\mu, \quad \|u\|_{L_{p'}^{\frac{n+1}{p}}(\mathbb{E})} = \left( \int_{\mathbb{E}} r^q |u|^p d\mu \right)^{\frac{1}{p}}.$$

**Proposition 1 (Poincaré Inequality [7])** *Let  $\mathbb{E} = [0, 1) \times X \times Y$  be a stretched edge manifold,  $\gamma \in \mathbb{R}$  and  $p \in (1, \infty)$ . If  $u \in \mathcal{H}_p^{1,\gamma}(\mathbb{E})$  then*

$$\|u(z)\|_{L_p^\gamma(\mathbb{E})} \leq c \|\nabla_{\mathbb{E}} u(z)\|_{L_p^\gamma(\mathbb{E})} \quad (3)$$

where  $\nabla_{\mathbb{E}} := (r\partial_r, \partial_{x_1}, \dots, \partial_{x_n}, r\partial_{y_1}, \dots, r\partial_{y_q})$  and the constant  $c$  depending only on  $\mathbb{E}$ .

**Proposition 2 ([7])** *For  $1 < p < 2^*$  the embedding  $\mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}) \hookrightarrow \mathcal{H}_{p,0}^{0, \frac{n+1}{p}}(\mathbb{E})$  is continuous.*

**Proposition 3 ([7])** *There exist  $0 < \lambda_1 \leq \lambda_2 \leq \lambda_3 \leq \dots \leq \lambda_j \leq \dots$ , and  $\lambda_j \rightarrow \infty$  such that for all  $j \geq 1$ , the following Dirichlet problem*

$$\begin{cases} -\Delta_{\mathbb{E}} \phi_j = \lambda_j \phi_j, & x \in int\mathbb{E}, \\ \phi_j = 0, & x \in \partial\mathbb{E}, \end{cases} \quad (4)$$

admits non-trivial solution in  $\mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E})$ . Moreover, we can choose positive  $\{\phi_j\}_{j \geq 1}$  which constitute an orthonormal basis of Hilbert space  $\mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E})$ , and the inequality

$$\lambda_1^{\frac{1}{2}} \|u(z)\|_{L_2^{\frac{n}{2}}(\mathbb{E})} \leq \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})},$$

holds.

### 3 Some Auxiliary Results

In this section we give some results about the potential wells for problem 1 and we obtain some properties of energy functional that we will use to prove the main results in Section 4.

Similar to the classical case, we introduce the following functionals on the cone Sobolev space  $\mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E})$ :

$$J(u) = \frac{1}{2} \int_{\mathbb{E}} r^q |\nabla_{\mathbb{E}} u|^2 d\mu + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |u|^2 d\mu - \frac{1}{p+1} \int_{\mathbb{E}} r^q g_t(z) |u|^{p+1} d\mu,$$

$$K(u) = \int_{\mathbb{E}} r^q |\nabla_{\mathbb{E}} u|^2 d\mu + \int_{\mathbb{E}} r^q V(x) |u|^2 d\mu - \int_{\mathbb{E}} r^q g_t(z) |u|^{p+1} d\mu.$$

Then  $J(u)$  and  $K(u)$  are well-defined and belong to space  $C^1(\mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}), \mathbb{R})$ . Now we define

$$\mathcal{N} = \left\{ u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}) \ ; \ K(u) = 0, \ \int_{\mathbb{E}} r^q |\nabla_{\mathbb{E}} u|^2 d\mu \neq 0 \right\},$$

$$d = \inf \left\{ \sup_{\lambda \geq 0} J(\lambda u) \ ; \ u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}), \ \int_{\mathbb{E}} r^q |\nabla_{\mathbb{E}} u|^2 d\mu \neq 0 \right\}.$$

Thus, similar to the results in [20] we obtain that  $0 < d = \inf_{u \in \mathcal{N}} J(u)$ . For  $0 < \delta$  we define

$$K_\delta(u) = \delta \left[ \int_{\mathbb{E}} r^q |\nabla_{\mathbb{E}} u|^2 d\mu + \int_{\mathbb{E}} r^q V(z) |u|^2 d\mu \right] - \int_{\mathbb{E}} r^q g_t(z) |u|^{p+1} d\mu,$$



$$\mathcal{N}_\delta = \left\{ u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{B}) \quad ; \quad K_\delta(u) = 0, \quad \int_{\mathbb{B}} r^q |\nabla_{\mathbb{E}}|^2 d\mu \neq 0 \right\},$$

$$d(\delta) = \inf_{u \in \mathcal{N}_\delta} J(u).$$

**Proposition 4** *If  $0 < \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} < r(\delta)$  where  $r(\delta) = \left(\frac{(C_*^2 + 1)\delta}{C^{p+1}}\right)^{\frac{1}{p-1}}$ , then  $K_\delta(u) > 0$ . In particular, if*

$$0 < \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} < r(1)$$

then  $K(u) > 0$ .

**Proof** We conclude the following

$$\begin{aligned} \|g_t(z)^{\frac{1}{p+1}} u\|_{L_{p+1}^{\frac{n+1}{2}}(\mathbb{E})}^{p+1} &= \int_{\mathbb{E}} r^q |g_t(z)^{\frac{1}{p+1}} u(z)|^{p+1} d\mu = \int_{\mathbb{E}} r^q |g_t(z)| |u(z)|^{p+1} d\mu \leq \\ &\|g_t\|_{L^\infty} \int_{\mathbb{E}} r^q |u|^{p+1} d\mu \quad \Rightarrow \\ \|g_t(z)^{\frac{1}{p+1}} u\|_{L_{p+1}^{\frac{n+1}{2}}(\mathbb{E})}^{p+1} &\leq C_g \|u\|_{L_{p+1}^{\frac{n+1}{2}}(\mathbb{E})}^{p+1}. \end{aligned} \quad (5)$$

Also from definition of  $C_*$ :

$$\|V(z)^{\frac{1}{2}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \geq C_*^2 \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2. \quad (6)$$

Then by definition of  $K_\delta$  and using the assumption we get that

$$\begin{aligned} K_\delta(u) &= \delta \left[ \int_{\mathbb{E}} r^q |\nabla_{\mathbb{E}} u|^2 d\mu + \int_{\mathbb{E}} r^q V(z) |u|^2 d\mu \right] - \int_{\mathbb{E}} r^q g_t(z) |u|^{p+1} d\mu \\ &\geq \delta(1 + C_*^2) \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - C_*^{p+1} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^{p+1} \\ &= \left( \delta(1 + C_*^2) - C_*^{p+1} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^{p-1} \right) \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 > 0. \end{aligned}$$

In case that  $\delta = 1$  then by definition of functional  $K$  we obtain that  $K(u) > 0$ .

**Proposition 5** *If  $K_\delta(u) < 0$ , then  $\|\nabla_{\mathbb{E}}u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} > r(\delta)$ . In particular, if  $K(u) < 0$ , then  $\|\nabla_{\mathbb{E}}u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} > r(1)$ .*

**Proof** Since  $K_\delta(u) < 0$ , then by definition of  $K_\delta(u)$ , we get that  $\|\nabla_{\mathbb{E}}u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \neq 0$ . Now, we have

$$\begin{aligned} \delta \|\nabla_{\mathbb{E}}u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 &< \int_{\mathbb{E}} r^q g_t(z) |u(z)|^{p+1} d\mu - \delta \int_{\mathbb{E}} r^q V(z) |u(z)|^2 d\mu \\ &\leq \|g_t(x)^{\frac{1}{p+1}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^{p+1} - \delta \|V(z)^{\frac{1}{2}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \\ &< C_{**}^{p+1} \|\nabla_{\mathbb{E}}u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^{p-1} \|\nabla_{\mathbb{E}}u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - \delta C_*^2 \|\nabla_{\mathbb{E}}u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2. \end{aligned}$$

Therefore,

$$\|\nabla_{\mathbb{E}}u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^{p-1} > \left( \frac{\delta(1 + C_*^2)}{C_{**}^{p+1}} \right) = r^{p-1}(\delta).$$

**Corollary 1** *Let  $u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E})$ ,  $K_\delta(u) = 0$  and  $\|\nabla_{\mathbb{E}}u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} \neq 0$ . Then  $\|\nabla_{\mathbb{E}}u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} \geq r(\delta)$ . In particular, if  $K(u) = 0$  and  $\|\nabla_{\mathbb{E}}u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} \neq 0$ , then  $\|\nabla_{\mathbb{E}}u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} \geq r(1)$ .*

**Lemma 1** (i) *The functional  $J(\lambda u)$  admits its maximum for  $\lambda = \lambda_*$  where*

$$\lambda_* = \left( \frac{\|\nabla_{\mathbb{E}}u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \int_{\mathbb{E}} r^q V(z) |u(z)|^2 d\mu}{\int_{\mathbb{E}} r^q g_t(z) |u(z)|^{p+1} d\mu} \right)^{\frac{1}{p-1}}.$$

Also for  $0 \leq \lambda < \lambda_*$ ,  $J(\lambda u)$  is strictly increasing and for  $\lambda_* < \lambda$ , it is strictly decreasing.

- (ii)  $K(\lambda_* u) = 0$  and  $K(\lambda u) > 0$  if  $0 < \lambda < \lambda_*$ . Also if  $\lambda_* < \lambda$  then  $K(\lambda u) < 0$ .  
 (iii) By results in i and ii we obtain that

$$\begin{aligned} d &= \inf \left\{ \sup_{\lambda \geq 0} J(\lambda u) \quad ; \quad u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}), \quad \int_{\mathbb{E}} |\nabla_{\mathbb{E}}u|^2 d\mu \neq 0 \right\} \\ &= \frac{p-1}{2(p+1)} (1 + C_*^2)^{\frac{p+1}{p-1}} C_{**}^{-2\frac{p+1}{1-p}} \|\nabla_{\mathbb{E}}u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2. \end{aligned}$$

**Proof** For proof of *i* and *ii* we obtain the following conclusions. Let  $u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E})$  and  $\int_{\mathbb{E}} |\nabla_{\mathbb{E}} u|^2 d\mu \neq 0$ . Then by definition of  $J$  we obtain that

$$\begin{aligned}
\lim_{\lambda \rightarrow +\infty} J(\lambda u) &= \lim_{\lambda \rightarrow +\infty} \left[ \frac{1}{2} \int_{\mathbb{E}} r^q |\nabla_{\mathbb{E}} \lambda u|^2 d\mu + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |\lambda u(z)|^2 d\mu \right. \\
&\quad \left. - \frac{1}{p+1} \int_{\mathbb{E}} r^q g_t(z) |\lambda u(z)|^{p+1} d\mu \right] \\
&= \lim_{\lambda \rightarrow +\infty} \left[ \frac{1}{2} \|\nabla_{\mathbb{E}} \lambda u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \|V(z)^{\frac{1}{2}} \lambda u(z)\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - \frac{1}{p+1} \|g_t(z)^{\frac{1}{p+1}} \lambda u(z)\|_{L_{p+1}^{\frac{n+1}{2}}(\mathbb{E})}^{p+1} \right] \\
&= \lim_{\lambda \rightarrow +\infty} \left[ \frac{\lambda^2}{2} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{\lambda^2}{2} \|V(z)^{\frac{1}{2}} u(z)\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - \frac{\lambda^{p+1}}{p+1} \|g_t(z)^{\frac{1}{p+1}} u(z)\|_{L_{p+1}^{\frac{n+1}{2}}(\mathbb{E})}^{p+1} \right] \\
&\geq \lim_{\lambda \rightarrow +\infty} \left[ \frac{\lambda^2}{2} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{\lambda^2}{2} C_*^2 \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - \frac{\lambda^{p+1}}{p+1} C_{**}^{p+1} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^{p+1} \right] \\
&= \lim_{\lambda \rightarrow +\infty} \left[ \frac{\lambda^2}{2} + \frac{\lambda^2}{2} C_*^2 - \frac{\lambda^{p+1}}{p+1} C_{**}^{p+1} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^{p-1} \right] \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 = -\infty.
\end{aligned}$$

Also we have

$$\begin{aligned}
J(\lambda u) &= \frac{1}{2} \int_{\mathbb{E}} |\nabla_{\mathbb{E}} \lambda u|^2 d\mu + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |\lambda u(z)|^2 d\mu - \frac{1}{p+1} \int_{\mathbb{E}} r^q g_t(z) |\lambda u(z)|^{p+1} d\mu \\
&= \frac{\lambda^2}{2} \int_{\mathbb{E}} |\nabla_{\mathbb{E}} u|^2 d\mu + \frac{\lambda^2}{2} \int_{\mathbb{E}} V(z) |u(z)|^2 d\mu - \frac{\lambda^{p+1}}{p+1} \int_{\mathbb{E}} g_t(z) |u(z)|^{p+1} d\mu.
\end{aligned}$$

Then

$$\begin{aligned}
\frac{\partial J(\lambda u)}{\partial \lambda} &= \lambda \int_{\mathbb{E}} |\nabla_{\mathbb{E}} u|^2 d\mu + \lambda \int_{\mathbb{E}} r^q V(z) |u(z)|^2 d\mu - \lambda^p \int_{\mathbb{E}} r^q g_t(z) |u(z)|^{p+1} d\mu \\
&= \lambda \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \lambda \|V(z)^{\frac{1}{2}} u(z)\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - \lambda^p \|g_t(z)^{\frac{1}{p+1}} u(z)\|_{L_{p+1}^{\frac{n+1}{2}}(\mathbb{E})}^{p+1}.
\end{aligned}$$

Now,  $\frac{\partial J(\lambda u)}{\partial \lambda} = 0$ , it follows that

$$\lambda_* := \left( \frac{\|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \int_{\mathbb{E}} r^q V(z) |u(z)|^2 d\mu}{\int_{\mathbb{E}} r^q g_t(z) |u(z)|^{p+1} d\mu} \right)^{\frac{1}{p-1}}$$

is a maximum of  $J(\lambda u)$  since  $\frac{\partial^2 (J(\lambda u))}{\partial \lambda^2} \Big|_{\lambda=\lambda_*} < 0$ .

(iii) Using of *i* and *ii*  $\sup_{\lambda \geq 0} J(\lambda u) = J(\lambda_* u)$ . Thus,

$$\begin{aligned}
J(\lambda_* u) &= \frac{1}{2} \int_{\mathbb{E}} r^q |\nabla_{\mathbb{E}} \lambda_* u|^2 d\mu + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |\lambda_* u(z)|^2 d\mu \\
&- \frac{1}{p+1} \int_{\mathbb{E}} r^q g_t(z) |\lambda_* u|^{p+1} d\mu \\
&= \lambda_*^2 \left[ \frac{1}{2} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |u(z)|^2 d\mu - \frac{\lambda_*^{p-1}}{p+1} \int_{\mathbb{E}} r^q g_t(z) |u|^{p+1} d\mu \right] \\
&= \lambda_*^2 \left[ \frac{1}{2} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |u(z)|^2 d\mu \right. \\
&- \frac{1}{p+1} \left( \frac{\|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \int_{\mathbb{E}} r^q V(z) |u(z)|^2 d\mu}{\int_{\mathbb{E}} r^q g_t(z) |u(z)|^{p+1} d\mu} \right)^{\frac{p-1}{p-1}} \left. \left( \int_{\mathbb{E}} r^q g_t(z) |u|^{p+1} d\mu \right) \right] \\
&= \lambda_*^2 \left[ \frac{1}{2} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |u(z)|^2 d\mu - \frac{1}{p+1} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \right. \\
&- \left. \frac{1}{p+1} \int_{\mathbb{E}} r^q V(z) |u(z)|^2 d\mu \right] \\
&= \lambda_*^2 \left[ \left( \frac{1}{2} - \frac{1}{p+1} \right) \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \left( \frac{1}{2} - \frac{1}{p+1} \right) \int_{\mathbb{E}} r^q V(z) |u(z)|^2 d\mu \right] \\
&= \left( \frac{\|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \int_{\mathbb{E}} r^q V(z) |u(z)|^2 d\mu}{\int_{\mathbb{E}} r^q g_t(z) |u(z)|^{p+1} d\mu} \right)^{\frac{2}{p-1}} \\
&\times \left( \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \int_{\mathbb{E}} r^q V(z) |u(z)|^2 d\mu \right) \times \frac{p-1}{2(p+1)} \\
&\geq \left( \frac{\|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + C_*^2 \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2}{\int_{\mathbb{E}} r^q g_t |u|^{p+1} d\mu} \right)^{\frac{2}{p-1}} \left[ \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + C_*^2 \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \right] \frac{p-1}{2(p+1)} \\
&\geq \frac{p-1}{2(p+1)} (1 + C_*^2)^{\frac{p+1}{p-1}} C_{**}^{\frac{-2(p+1)}{p-1}} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2.
\end{aligned}$$

Therefore,

$$\begin{aligned} d &= \inf_{\substack{u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}), \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \neq 0}} J(\lambda_* u) \\ &= \frac{p-1}{2(p+1)} (1 + C_*^2)^{\frac{p+1}{p-1}} C_{**}^{\frac{-2(p+1)}{p-1}} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2. \end{aligned}$$

**Proposition 6** *Let  $0 < \delta < \frac{p+1}{2}$ , then  $d(\delta) \geq a(\delta)r^2(\delta)$  where  $a(\delta) = \left(\frac{1}{2} - \frac{\delta}{p+1}\right)(1 + C_*^2)$ . Moreover, we have*

$$d(\delta) = \inf_{u \in \mathcal{N}_\delta} J(u) = d \lambda(\delta)^2 a(\delta) [1 + c_*^2]^{-1} \frac{2(p+1)}{p-1}.$$

**Proof** Let  $u \in \mathcal{N}_\delta$ , so by Proposition 5 we get that  $\|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} > r(\delta)$ . Then by definition of  $J$  and  $K_\delta$  we obtain that

$$\begin{aligned} J(u) &= \frac{1}{2} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |u(z)|^2 d\mu - \frac{1}{p+1} \int_{\mathbb{E}} r^q g_t(z) |u(z)|^{p+1} d\mu \\ &= \frac{1}{2} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |u(z)|^2 r^q \\ &\quad - \frac{1}{p+1} \left( \delta \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - K_\delta(u) + \delta \int_{\mathbb{E}} r^q V(z) |u(z)|^2 d\mu \right). \end{aligned}$$

Since  $K_\delta(u) = 0$ ,

$$\begin{aligned} J(u) &\geq \left(\frac{1}{2} - \frac{\delta}{p+1}\right) \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{\delta(p-1)}{2(p+1)} \|V(x)^{\frac{1}{2}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \\ &\geq \left(\frac{1}{2} - \frac{\delta}{p+1}\right) \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{\delta(p-1)}{2(p+1)} C_*^2 \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \\ &= \left(\frac{1}{2} - \frac{\delta}{p+1}\right) (1 + C_*^2) \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2. \end{aligned}$$

Since  $\|\nabla u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \geq r^2(\delta)$  then,

$$d(\delta) \geq a(\delta)r^2(\delta).$$

Now, we prove the second part of the assertion. By definition of  $\mathcal{N}_\delta$  and  $\mathcal{N}$ , for  $\bar{u} \in \mathcal{N}_\delta$  and  $\lambda\bar{u} \in \mathcal{N}$ , we obtain

$$\lambda^2 \|\nabla_{\mathbb{E}} \bar{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \lambda^2 \int_{\mathbb{E}} r^q V(z) |\bar{u}|^2 d\mu = \lambda^{p+1} \int_{\mathbb{E}} r^q g_t(z) |\bar{u}|^{p+1} d\mu, \quad (7)$$

and

$$\delta \|\nabla_{\mathbb{E}} \bar{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \delta \int_{\mathbb{E}} r^q V(z) |\bar{u}|^2 d\mu = \int_{\mathbb{E}} r^q g_t(z) |\bar{u}|^{p+1} d\mu. \quad (8)$$

Then 7 gives

$$\lambda = \left( \frac{\|\nabla_{\mathbb{E}} \bar{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \int_{\mathbb{E}} r^q V(z) |\bar{u}|^2 d\mu}{\int_{\mathbb{E}} r^q g_t(z) |\bar{u}|^{p+1} d\mu} \right)^{\frac{1}{p-1}}, \quad (9)$$

and 8 gives that

$$\delta = \frac{\int_{\mathbb{E}} r^q g_t(z) |\bar{u}|^{p+1} d\mu}{\|\nabla_{\mathbb{E}} \bar{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \int_{\mathbb{E}} r^q V(z) |\bar{u}|^2 d\mu}. \quad (10)$$

By 10 and 9, we define

$$\lambda = \lambda(\delta) = \left( \frac{1}{\delta} \right)^{\frac{1}{p-1}}. \quad (11)$$

Moreover, for such  $\lambda$ ,  $\lambda\bar{u} \in \mathcal{N}$ , so by definition of  $d$  we get that

$$\begin{aligned} d \leq J(\lambda\bar{u}) &= \frac{1}{2} \|\nabla_{\mathbb{E}} \lambda\bar{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |\lambda\bar{u}|^2 d\mu \\ &\quad - \frac{1}{p+1} \int_{\mathbb{E}} r^q g_t(z) |\lambda\bar{u}|^{p+1} d\mu \\ &= \frac{\lambda^2}{2} \|\nabla_{\mathbb{E}} \bar{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{\lambda^2}{2} \int_{\mathbb{E}} r^q V(z) |\bar{u}|^2 d\mu \\ &\quad - \frac{1}{p+1} \left[ \|\nabla_{\mathbb{E}} \lambda\bar{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \int_{\mathbb{E}} r^q V(z) |\lambda\bar{u}|^2 d\mu - K(\lambda\bar{u}) \right] \end{aligned}$$

$$\begin{aligned}
&= \lambda^2 \left[ \frac{1}{2} \|\nabla_{\mathbb{E}} \bar{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |\bar{u}|^2 d\mu - \frac{1}{p+1} \|\nabla_{\mathbb{E}} \bar{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \right. \\
&\quad \left. - \frac{1}{p+1} \int_{\mathbb{E}} r^q V(z) |\bar{u}|^2 d\mu \right] \\
&\leq \left(\frac{1}{\delta}\right)^{\frac{2}{p-1}} \left[ \frac{p-1}{2(p+1)} \|\nabla_{\mathbb{E}} \bar{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - \frac{(1-p)C_*^2}{2(p+1)} \|\nabla_{\mathbb{E}} \bar{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \right].
\end{aligned}$$

On the other hand,

$$\begin{aligned}
d(\delta) &= J(\bar{u}) = \frac{1}{2} \|\nabla_{\mathbb{E}} \bar{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |\bar{u}|^2 d\mu - \frac{1}{p+1} \int_{\mathbb{E}} r^q g_t(z) |\bar{u}|^{p+1} d\mu \\
&= \frac{1}{2} \|\nabla_{\mathbb{E}} \bar{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |\bar{u}|^2 d\mu - \frac{1}{p+1} \left( \delta \|\nabla_{\mathbb{E}} \bar{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \right. \\
&\quad \left. + \delta \int_{\mathbb{E}} r^q V(z) |\bar{u}|^2 d\mu - K_{\delta}(\bar{u}) \right) \\
&= \left(\frac{1}{2} - \frac{\delta}{p+1}\right) \|\nabla_{\mathbb{E}} \bar{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \left(\frac{1}{2} - \frac{\delta}{p+1}\right) \int_{\mathbb{E}} r^q V(z) |\bar{u}|^2 d\mu \\
&\geq \left(\frac{1}{2} - \frac{\delta}{p+1}\right) (1 + C_*^2) \|\nabla_{\mathbb{E}} \bar{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 = a(\delta) \|\nabla_{\mathbb{E}} \bar{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2.
\end{aligned}$$

Indeed,

$$d \leq J(\lambda \bar{u}) \leq \left(\frac{1}{\delta}\right)^{\frac{2}{p-1}} \left[ \frac{p-1}{2(p+1)} (1 + C_*^2) \right] \frac{d(\delta)}{a(\delta)}.$$

Hence,

$$d(\delta) \geq a(\delta) \left(\frac{1}{\delta}\right)^{-\frac{2}{p-1}} [1 + C_*^2]^{-1} \left[\frac{2(p+1)}{p-1}\right] d.$$

Now, we let  $0 < \delta$  and  $\tilde{u} \in \mathcal{N}$  is minimizer of  $d$  that is

$$d = J(\tilde{u}) = \frac{1}{2} \|\nabla_{\mathbb{E}} \tilde{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |\tilde{u}|^2 d\mu - \frac{1}{p+1} \int_{\mathbb{E}} r^q g_t(z) |\tilde{u}|^{p+1} d\mu.$$

we define  $\lambda = \lambda(\delta)$  by

$$\delta \|\nabla_{\mathbb{E}} \lambda \tilde{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \delta \int_{\mathbb{E}} r^q V(z) |\lambda \tilde{u}|^2 d\mu = \int_{\mathbb{E}} r^q g_t(z) |\lambda \tilde{u}|^{p+1} d\mu.$$

Then for any  $0 < \delta$ , there exists a unique  $\lambda$  which satisfies

$$\lambda = \delta^{\frac{1}{p-1}}.$$

Hence, for such  $\lambda$ ,  $\lambda \tilde{u} \in \mathcal{N}_{\delta}$  by definition of  $d(\delta)$  we get that

$$\begin{aligned} d &= \frac{1}{2} \|\nabla_{\mathbb{E}} \lambda \tilde{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |\lambda \tilde{u}|^2 d\mu \\ &\quad - \frac{1}{p+1} \left( \delta \|\nabla_{\mathbb{E}} \lambda \tilde{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \delta \int_{\mathbb{E}} r^q V(z) |\lambda \tilde{u}|^2 d\mu - K_{\delta}(\lambda \tilde{u}) \right) \\ &= \left( \frac{1}{2} - \frac{\delta}{p+1} \right) \|\nabla_{\mathbb{E}} \lambda \tilde{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \left( \frac{1}{2} - \frac{\delta}{p+1} \right) \int_{\mathbb{E}} r^q V(z) |\lambda \tilde{u}|^2 d\mu \\ &\geq \left[ \frac{1}{2} - \frac{\delta}{p+1} + C_*^2 \left( \frac{1}{2} - \frac{\delta}{p+1} \right) \right] \|\nabla_{\mathbb{E}} \lambda \tilde{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2. \end{aligned}$$

On the other hand,

$$\begin{aligned} d(\delta) &\leq J(\lambda \tilde{u}) = \frac{1}{2} \|\nabla_{\mathbb{E}} \lambda \tilde{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |\lambda \tilde{u}|^2 d\mu - \frac{1}{p+1} \int_{\mathbb{E}} r^q g_t(z) |\lambda \tilde{u}|^{p+1} d\mu \\ &= \frac{\lambda^2}{2} \|\nabla_{\mathbb{E}} \tilde{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{\lambda^2}{2} \int_{\mathbb{E}} r^q V(z) |\tilde{u}|^2 d\mu \\ &\quad - \frac{1}{p+1} \left( \delta \|\nabla_{\mathbb{E}} \lambda \tilde{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \delta \int_{\mathbb{E}} r^q V(z) |\lambda \tilde{u}|^2 d\mu - K_{\delta}(\lambda \tilde{u}) \right) \\ &= \lambda^2 \left[ \left( \frac{1}{2} - \frac{\delta}{p+1} \right) \|\nabla_{\mathbb{E}} \tilde{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - \left( \frac{\delta}{p+1} - \frac{1}{2} \right) \int_{\mathbb{E}} r^q V(z) |\tilde{u}|^2 d\mu \right] \\ &\leq \delta^{\frac{2}{p-1}} \left[ \left( \frac{1}{2} - \frac{\delta}{p+1} \right) \|\nabla_{\mathbb{E}} \tilde{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - \left( \frac{\delta}{p+1} - \frac{1}{2} \right) C_*^2 \|\nabla_{\mathbb{E}} \tilde{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \right] \\ &\leq \delta^{\frac{2}{p+1}} \left[ \frac{1}{2} - \frac{\delta}{p+1} + \left( \frac{1}{2} - \frac{\delta}{p+1} \right) C_*^2 \right] \|\nabla_{\mathbb{E}} \tilde{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 = \delta^{\frac{2}{p-1}} a(\delta) \|\nabla_{\mathbb{E}} \tilde{u}\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2. \end{aligned}$$



Then,

$$d(\delta) \leq \delta^{\frac{2}{p-1}} a(\delta) d [1 + C_*^2]^{-1} \frac{2(p+1)}{p-1}.$$

Therefore,

$$d(\delta) = \inf_{u \in \mathcal{N}_\delta} J(u) = \delta^{\frac{2}{p-1}} a(\delta) d [1 + C_*^2]^{-1} \frac{2(p+1)}{p-1}.$$

*Remark 1* According to  $d(\delta)$  in Proposition 6, we obtain that

- (i)  $\lim_{\delta \rightarrow 0} d(\delta) = 0$ .  
(ii)  $d(\delta) = d^{\frac{2(p+1)}{p-1}} \left[ \frac{1}{2} \delta^{\frac{2}{p-1}} - \frac{1}{p+1} \delta^{\frac{p+1}{p-1}} \right]$ . Then

$$d'(\delta) = \frac{d2(p+1)}{(p-1)^2} \delta^{\frac{2}{p-1}} [\delta^{-1} - 1] = 0 \Rightarrow \delta = 1.$$

Hence, if  $0 < \delta < 1$  then  $d(\delta)$  is strictly increasing function and if  $\delta > 1$  then  $d(\delta)$  is strictly decreasing function.

## 4 Invariance of the Solutions

Now, we introduce the following potential wells

$$W = \left\{ u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}) \quad ; \quad K(u) > 0, \quad J(u) < d \right\} \cup \{0\},$$

$$W_\delta = \left\{ u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}) \quad ; \quad K_\delta(u) > 0, \quad J(u) < d(\delta) \right\} \cup \{0\},$$

for  $0 < \delta$ , and corresponding potentials outside of the set that defined as above

$$E = \left\{ u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}) \quad ; \quad K(u) < 0, \quad J(u) < d \right\},$$

$$E_\delta = \left\{ u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}) \quad ; \quad K_\delta(u) < 0, \quad J(u) < d(\delta) \right\}$$

for any  $0 < \delta$ . According to the definition of potential wells  $W_\delta$  and potential outside  $E_\delta$  one can get the following inclusions:

- (I)  $W_{\delta_1} \subset W_{\delta_2}$  whenever  $0 < \delta_1 < \delta_2 \leq 1$ ,

(II)  $E_{\delta_1} \subset E_{\delta_2}$  whenever  $1 \leq \delta_2 < \delta_1 < \frac{p+1}{2}$ . Furthermore, from the above results on can define the following sets

$$V_\delta = \left\{ u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}) : \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} < r(\delta) \right\}$$

$$\bar{V}_\delta = V_\delta \cup \partial V_\delta = \left\{ u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}) : \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} \leq r(\delta) \right\}$$

$$V_\delta^c = \left\{ u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}) : \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} > r(\delta) \right\}.$$

Then for every  $0 < \delta < \frac{p+1}{2}$  one gets that

$$V_{t(\delta)} \subset W_\delta \subset V_{s(\delta)}, \quad E_\delta \subset V_\delta^c$$

where

$$V_{t(\delta)} = \left\{ u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}) : \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} < \min\{r^2(\delta), r_0^2(\delta)\} \right\}$$

$$V_{s(\delta)} = \left\{ u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}) : \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} < \frac{d(\delta)}{a(\delta)} \right\}$$

where  $r_0(\delta)$  is the unique real root of equation  $\frac{r^2}{2} = d(\delta)$ .

**Definition 4** Suppose that  $u(t)$  is a weak solution of problem 1.  $T_{\max}$  is called maximal existence time of solution  $u(t)$  if one the following conditions hold:

- (1) If  $u(t)$  exists for every  $0 \leq t < +\infty$  then  $T_{\max} = +\infty$ . In this case, we say that the solution is global.
- (2) If there exists a  $t_0 \in (0, \infty)$  such that  $u(t)$  exists for every  $0 \leq t < t_0$ , but does not exist at  $t = t_0$ , then  $T_{\max} = t_0$ .

**Definition 5**  $u = u(z, t) \in L^\infty(0, T_{\max}; \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}))$  with  $\partial_t u \in L^\infty(0, T_{\max}; L_2^{\frac{n+1}{2}}(\mathbb{E}))$  is called a weak solution of the problem 1 on  $int \mathbb{E} \times [0, T_{\max})$  if

$$\begin{aligned} (u_t, v)_2 + \gamma(\nabla_{\mathbb{E}} u, \nabla_{\mathbb{E}} v)_2 + \int_0^t (\nabla_{\mathbb{E}} u, \nabla_{\mathbb{E}} v)_2 d\tau + \int_0^t (V(x)u, v)_2 d\tau \\ = \int_0^t (g_t(z)|u|^{p-1}u, v)_2 d\tau \\ + (\gamma u_0, v)_2 + (u_1, v)_2 \quad \forall v \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}), \end{aligned}$$

$u(z, 0) = u_0$  in  $\mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E})$  and hold the following energy inequality

$$I(t) + \gamma \int_0^t \|\nabla_{\mathbb{E}}(\partial_{\tau} u)\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 d\tau \leq I(0), \quad \forall t \in (0, T_{\max}),$$

where  $0 \leq T_{\max} \leq \infty$  and

$$I(t) = \frac{1}{2} \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{p+1} \int_0^t \left\| \left( \frac{d}{d\tau} g_{\tau}(z) \right)^{\frac{1}{p+1}} u \right\|_{L_p^{\frac{n+1}{2}}(\mathbb{E})}^{p+1} d\tau + J(u).$$

We note, since  $u \in L^{\infty}\left(0, T_{\max}; \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E})\right)$  and  $\partial_t u \in L^{\infty}\left(0, T_{\max}; L_2^{\frac{n+1}{2}}(\mathbb{E})\right)$  from the first equation of the problem 1 as similar in [13], one can obtain that  $\partial_t^2 u \in L^{\infty}\left(0, T_{\max}; \mathcal{H}_{2,0}^{-1, \frac{n+1}{2}}(\mathbb{E})\right)$ .

Now we discuss the invariance of some sets corresponding to the problem 1.

**Proposition 7** *Let  $0 < J(u) < d$  for  $u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E})$ . Suppose that  $\delta_1 < 1 < \delta_2$  be roots of equation  $d(\delta) = J(u)$ . Then  $K_{\delta}(u)$  has no change in its sign for  $\delta \in (\delta_1, \delta_2)$ .*

**Proof** Since  $0 < J(u) < d$  then by Propositions 1 and 2 we can assume that  $\|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \neq 0$ . We assume that there exists a  $\delta_0 \in (\delta_1, \delta_2)$  for which  $K_{\delta_0}(u) = 0$ . Hence, by definition of  $d(\delta)$  we have  $J(u) \geq d(\delta)$ . But, we have two cases the following for  $\delta_0$

$$\begin{cases} \delta_1 < \delta_0 < 1 < \delta_2 \\ \delta_1 < 1 < \delta_0 < \delta_2 \end{cases}$$

Now, by Remark 1 We get that  $d(\delta_1) < d(\delta_0)$  or  $d(\delta_2) < d(\delta_0)$  then we obtain that  $d(\delta_1) = d(\delta_2) = J(u) < d(\delta_0)$  that this is contradiction .

**Theorem 1** *Let  $u_0 \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E})$ ,  $0 < e < d$ . Suppose that  $\delta_1 < \delta_2$  are roots of equations  $d(\delta) = e$  then*

- (i) *all solutions of problem 1 with  $0 < J(u_0) \leq e$  belong to set  $W_{\delta}$  for  $\delta_1 < \delta < \delta_2$  provided  $K(u_0) > 0$  or  $\|\nabla_{\mathbb{E}} u_0\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 = 0$ .*
- (ii) *all solutions of problem 1 with  $0 < J(u_0) \leq e$  belong to  $E_{\delta}$  for  $\delta \in (\delta_1, \delta_2)$  provided  $K(u_0) < 0$ .*

**Proof**

- (i) Let  $u(t)$  be a solution of the problem 1 with initial value  $u_0$  for which satisfies in conditions  $0 < J(u_0) \leq e < d$ ,  $K(u_0) > 0$  or  $\|\nabla_{\mathbb{E}} u_0\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 = 0$ . Let  $T$

be existence time for solution  $u(t)$ . If  $\|\nabla_{\mathbb{E}}u_0\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 = 0$ , then since  $u_0$  has compact support  $u_0 = 0$ , so by definition of  $W_\delta$  we obtain that  $u_0 \in W_\delta$ . If  $K(u_0) > 0$  then by assumption we have

$$0 < J(u_0) \leq e = d(\delta_1) = d(\delta_2) < d(\delta) \leq d$$

for  $\delta_1 < \delta < \delta_2$ . Hence,  $K_\delta(u_0(t)) > 0$  for  $\delta_1 < \delta < \delta_2$ , by Proposition 7. Therefore, by definition of  $W_\delta$ ,  $u_0 \in W_\delta$  for  $\delta_1 < \delta < \delta_2$ . Now, we have to show that for  $\delta_1 < \delta < \delta_2$  and  $0 < t < T$ ,  $u(t) \in W_\delta$ . Suppose that, there exist  $t_0 \in (0, T)$  such that for  $\delta_1 < \delta < \delta_2$ ,  $u(t_0) \in \partial W_\delta$ . Then we can imply that,  $K_\delta(u(t_0)) = 0$  and  $\|\nabla_{\mathbb{E}}u_0\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \neq 0$ , or by definition of  $W_\delta$ ,  $J(u(t_0)) = d(\delta)$ . Since  $u(t_0)$  is a solution of problem 1, so it satisfies in energy inequality i.e.

$$\begin{aligned} \frac{1}{2}\|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{p+1} \int_0^t \left\| \left( \frac{d}{d\tau} g_\tau(z) \right)^{\frac{1}{p+1}} u \right\|_{L_{p+1}^{\frac{n+1}{2}}(\mathbb{Z})}^{p+1} d\tau + J(u(t)) \\ + \gamma \int_0^t \|\nabla_{\mathbb{E}}(\partial_\tau u)\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 d\tau \\ \leq I(0) = J(u_0) \leq e < d(\delta), \end{aligned}$$

for any  $\delta \in (\delta_1, \delta_2)$  and  $t \in (0, T)$ . Therefore, the equality  $J(u(t_0)) = d(\delta)$  for any  $\delta \in (\delta_1, \delta_2)$  and  $t \in (0, T)$  is not possible. If  $K_\delta(u(t_0)) = 0$  and  $\|\nabla_{\mathbb{E}}u_0\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \neq 0$ , then by definition of  $d(\delta)$  we get that  $d(\delta) \leq J(u_0(t))$ , that is in contradiction with energy inequality. Therefore,  $u(t) \in W_\delta$  for any  $\delta \in (\delta_1, \delta_2)$  and  $t \in (0, T)$ .

- (ii) similar to first case it can be prove that  $u_0 \in E_\delta$  for  $\delta \in (\delta_1, \delta_2)$  provided  $K_\delta(u_0) < 0$ . Now, we should prove  $u(t) \in E_\delta$  for any  $\delta \in (\delta_1, \delta_2)$  and  $t \in (0, T)$ . Suppose that there exist  $t_0 \in (0, T)$ , such that for  $t \in [0, t_0)$ ,  $u(t) \in E_\delta$  and  $u(t_0) \in \partial E_\delta$ , that is,  $K_\delta(u_0) = 0$  or  $J(u(t_0)) = d(\delta)$  for  $\delta \in (\delta_1, \delta_2)$ . According to energy inequality the equality  $J(u(t_0)) = d(\delta)$  is not possible similar to first case. Hence, we assume that  $K_\delta(u(t_0)) = 0$ , then  $K_\delta(u(t)) < 0$  for  $t \in (0, t_0)$ , since for  $t \in [0, t_0)$ ,  $u(t) \in E_\delta$ , then by definition of  $E_\delta$ ,  $K_\delta(u(t)) < 0$ . Now, using the Proposition 5 we obtain that  $\|\nabla_{\mathbb{E}}u(t)\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} > r(\delta)$  and  $\|\nabla_{\mathbb{E}}u(t_0)\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} > r(\delta) \neq 0$ . Hence by definition of  $d(\delta)$ ,  $J(u(t_0)) \geq d(\delta)$  which is in contradiction with energy inequality.

*Remark 2* suppose that all assumptions in Theorem 1 hold. Then for any  $\delta \in (\delta_1, \delta_2)$  both seta  $W_\delta$  and  $E_\delta$  are invariant. Moreover, both sets

$$W_{\delta_1\delta_2} = \bigsqcup_{\delta_1 < \delta < \delta_2} W_\delta, \quad E_{\delta_1\delta_2} = \bigsqcup_{\delta_1 < \delta < \delta_2} E_\delta$$

are invariant respectively under flow of the problem 1. Hence, we can get for all weak solutions of the problem 1

$$u(t) \notin \mathcal{N}_{\delta_1\delta_2} = \bigsqcup_{\delta_1 < \delta < \delta_2} \mathcal{N}_\delta.$$

To discuss about the invariant of the solutions with negative level energy, we introduce the following results.

**Proposition 8** *Let  $u_0 \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E})$  and  $u_1 \in L_2^{\frac{n+1}{2}}(\mathbb{E})$ . Suppose that  $I(0) = 0$  and  $\|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} \neq 0$ . Then all weak solutions of the problem 1 satisfy*

$$\|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^{p-1} \geq M = \frac{(p+1)(1+C_*^2)}{2C_{**}^{p+1}}.$$

**Proof** Let us consider  $u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E})$  as a weak solution of the problem 1. According to the Definition 5

$$I(t) + \gamma \int_0^t \|\nabla_{\mathbb{E}}(\partial_\tau u)\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 d\tau \leq I(0) = 0.$$

Therefore, by definition of constants  $C_*$  and  $C_{**}$

$$\begin{aligned} \frac{1}{2} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{C_*^2}{2} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 &\leq \frac{1}{2} \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \\ &\quad + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |u|^2 d\mu \\ &\leq \frac{C_{**}^{p+1}}{p+1} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^{p+1}. \end{aligned}$$

Hence,

$$\|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^{p-1} \geq \frac{(p+1)(1+C_*^2)}{2C_{**}^{p+1}} = M.$$

**Theorem 2** Let  $u_0 \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E})$  and  $u_1 \in L_2^{\frac{n+1}{2}}(\mathbb{E})$ . Suppose that either  $I(0) < 0$  or  $I(0) = 0$  and  $\|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} \neq 0$ . Then all weak solutions of the problem 1 belong to  $E_\delta$  for any  $\delta > 0$ .

**Proof** Let  $u(t)$  be an arbitrary weak solution of the problem 1 with expressed assumptions in face of the Theorem and  $T$  be the existence time of  $u(t)$ . From Definition 5, for every  $\delta > 0$  and  $t \in [0, T)$ , we can obtain

$$\begin{aligned}
& \frac{1}{2} \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + a(\delta) \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{p+1} K_\delta(u) \\
&= \frac{1}{2} \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \left( \frac{1}{2} - \frac{\delta}{p+1} \right) (1 + C_*^2) \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \\
&+ \frac{1}{p+1} \left( \delta \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \delta \int_{\mathbb{E}} r^q V(z) |u|^2 d\mu \right. \\
&\quad \left. - \int_{\mathbb{E}} r^q g_t(z) |u|^{p+1} d\mu \right) \\
&= \frac{1}{2} \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - \frac{\delta C_*^2}{p+1} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \\
&+ \frac{C_*^2}{2} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{\delta}{p+1} \int_{\mathbb{E}} r^q V(z) |u|^2 d\mu \\
&\quad - \frac{1}{p+1} \int_{\mathbb{E}} r^q g_t(z) |u|^{p+1} d\mu \\
&= \frac{1}{2} \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - \frac{\delta C_*^2}{p+1} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \\
&+ \frac{C_*^2}{2} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - \left( \frac{1}{2} - \frac{\delta}{p+1} \right) \int_{\mathbb{E}} r^q V(z) |u|^2 d\mu \\
&\quad + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |u|^2 d\mu \\
&- \frac{1}{p+1} \int_{\mathbb{E}} r^q g_t(z) |u|^{p+1} d\mu \leq \frac{1}{2} \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + J(u) \\
&+ \left( \frac{C_*^2}{2} - \frac{\delta C_*^2}{p+1} - C_*^2 \left( \frac{1}{2} - \frac{\delta}{p+1} \right) \right) \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \\
&+ \frac{1}{p+1} \int_0^t \left\| \left( \frac{d}{d\tau} g_\tau(x) \right)^{\frac{1}{p+1}} u \right\|_{L_{p+1}^{\frac{n+1}{2}}(\mathbb{E})}^{p+1} \\
&\quad + \gamma \int_0^t \|\nabla_{\mathbb{E}}(\partial_\tau u)\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \leq I(0). \tag{12}
\end{aligned}$$

If  $I(0) < 0$ , then 12 implies that  $K_\delta(u) < 0$  and  $J(u) < 0 < d(\delta)$  for every  $\delta > 0$  and  $t \in [0, T)$ . If  $I(0) = 0$  and  $\|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} \neq 0$ , then Proposition 8 gives  $\|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} \geq M$  for  $t \in [0, T)$ . Again by relation 12 we get  $K_\delta(u) < 0$  and  $J(u) < 0 < d(\delta)$  for  $\delta > 0$  and  $t \in [0, T)$ . Therefore, for two cases discussed above, for every  $\delta > 0$  and  $t \in [0, T)$ , we have  $u \in E_\delta$ .

## 5 Global Existence and Finite-Time of the Solutions

In this section, we prove the global existence and nonexistence of solutions and give a sharp condition for global existence of solutions for problem 1 with  $I(0) < d$ .

**Theorem 3** *Let  $\gamma \geq 0$ ,  $u_0 \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E})$  and  $u_1 \in L_2^{\frac{n+1}{2}}(\mathbb{E})$ . Suppose that  $I(0) < d$ ,  $K(u_0) > 0$  or  $\|\nabla_{\mathbb{E}} u_0\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} = 0$ . Then problem 1 admits a global weak solution  $u(t) \in L^\infty(0, \infty; \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}))$  with  $\partial_t u \in L_2^{\frac{n+1}{2}}(\mathbb{E})$  and  $u(t) \in W$  for  $t \in [0, \infty)$ .*

**Proof** By Proposition 3 we can choose  $\{w_j(z)\}$  as orthonormal basis of space  $\mathcal{H}_{2,0}^{1, \frac{n}{2}}(\mathbb{B})$ . Then we construct approximation solution  $u_m(z, t)$  similar to [20] as following:

$$u_m(z, t) = \sum_{j=1}^m h_{jm}(t) w_j(z),$$

for  $m = 1, 2, \dots$  that satisfies in problem 1 then,

$$\begin{aligned} (\partial_t^2 u_m, w_k)_2 + (\nabla_{\mathbb{E}} u_m, \nabla_{\mathbb{E}} w_k)_2 + (V(z) u_m, w_k)_2 + \gamma (\nabla_{\mathbb{E}}(\partial_t u_m), \nabla_{\mathbb{E}} w_k)_2 \\ = (g_t(z) u_m |u_m|^{p-1}, w_k)_2, \end{aligned} \quad (13)$$

$$u_m(z, 0) = \sum_{j=1}^m h_{jm}(0) w_j(z) \rightarrow u_0(z), \quad (14)$$

in  $\mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E})$  and

$$\partial_t u_m(z, 0) = \sum_{j=1}^m h'_{jm}(0) w_j(z) \rightarrow u_1(z), \quad (15)$$

in  $L_2^{\frac{n+1}{2}}(\mathbb{E})$ . Multiplying 13, 14 and 15 by  $h'_{km}(t)$  and forming the crossmarklogo sum on  $k = 1, 2, \dots$ ,

$$\begin{aligned} & \sum_{k=1}^m (\partial_t^2 u_m, w_k)_2 h'_{km}(t) + (\nabla_{\mathbb{E}} u_m, \nabla_{\mathbb{E}} w_k)_2 h'_{km}(t) + (V(z)u_m, w_k)_2 h'_{km}(t) \\ & \quad + \sum_{k=1}^m \gamma (\nabla_{\mathbb{E}}(\partial_t u_m), \nabla_{\mathbb{E}} w_k)_2 h'_{km}(t) \\ & \quad = \sum_{k=1}^m (g_t(z)u_m |u_m|^{p-1}, w_k)_2 h'_{km}(t), \end{aligned}$$

for  $m = 1, 2, 3, \dots$ . Therefore,

$$\begin{aligned} & \int_{\mathbb{E}} r^q \partial_t^2 u_m \partial_t u_m d\mu + \int_{\mathbb{E}} r^q \nabla_{\mathbb{E}} u_m \partial_t \nabla_{\mathbb{E}} u_m d\mu + \int_{\mathbb{E}} r^q V(z)u_m \partial_t u_m d\mu \\ & \quad + \gamma \int_{\mathbb{E}} r^q \nabla_{\mathbb{E}}(\partial_t u_m) \nabla_{\mathbb{E}}(\partial_t u_m) d\mu \\ & \quad = \int_{\mathbb{E}} r^q g_t(z)u_m |u_m|^{p-1} \partial_t u_m d\mu. \end{aligned} \quad (16)$$

Using The Leibniz rule one can get

$$\begin{aligned} & \frac{1}{2} \frac{d}{dt} \int_{\mathbb{E}} r^q |\partial_t^2 u_m|^2 d\mu + \frac{1}{2} \frac{d}{dt} \int_{\mathbb{E}} r^q |\nabla_{\mathbb{E}} u_m|^2 d\mu + \frac{1}{2} \frac{d}{dt} \int_{\mathbb{E}} r^q V(z)|u_m|^2 d\mu \\ & \quad + \gamma \int_{\mathbb{E}} r^q |\nabla_{\mathbb{E}}(\partial_t u_m)|^2 d\mu = \frac{1}{p+1} \frac{d}{dt} \int_{\mathbb{E}} r^q g_t(z)|u_m|^{p+1} d\mu \\ & \quad - \frac{1}{p+1} \int_{\mathbb{E}} r^q \left(\frac{d}{dt} g_t(z)\right) |u_m|^{p+1} d\mu. \end{aligned} \quad (17)$$

By integration of the relation 17 with respect to  $t$

$$\begin{aligned} & \frac{1}{2} \|\partial_t u_m\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \|\nabla_{\mathbb{E}} u_m\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \int_0^t \int_{\mathbb{E}} r^q V(z)|u_m|^2 d\mu + \gamma \int_0^t \|\nabla_{\mathbb{E}}(\partial_\tau u_m)\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 d\tau \\ & \quad - \frac{1}{p+1} \int_{\mathbb{E}} r^q g_t(z)|u_m|^{p+1} d\mu + \frac{1}{p+1} \int_0^t \left\| \left(\frac{d}{d\tau} g_\tau(z)\right)^{p+1} u_m \right\|_{\frac{1}{p+1}} d\tau \\ & \quad = I(t) + \gamma \int_0^t \|\nabla_{\mathbb{E}}(\partial_\tau u_m)\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 d\tau \leq I(0) < d, \end{aligned} \quad (18)$$

where the last equal is upon Definition 5. Hence, for sufficiently large  $m$  and  $0 \leq t < \infty$  we obtain that  $u_m \in W$  by Proposition 1. Using 18 and definition of functional  $K$ ,



$$\begin{aligned}
J(u_m) &= \frac{1}{2} \|\nabla_{\mathbb{E}} u_m\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |u_m|^2 d\mu - \frac{1}{p+1} \int_{\mathbb{E}} r^q g_t(z) |u_m|^{p+1} d\mu \\
&= \frac{1}{2} \|\nabla_{\mathbb{E}} u_m\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \int_{\mathbb{E}} r^q V(z) |u_m|^2 d\mu \\
&\quad - \frac{1}{p+1} \left( \|\nabla_{\mathbb{E}} u_m\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \int_{\mathbb{E}} r^q V(z) |u_m|^2 d\mu - K(u_m) \right) \\
&\geq \left( \frac{p-1}{2(p+1)} \right) \left[ \|\nabla_{\mathbb{E}} u_m\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \int_{\mathbb{E}} r^q V(z) |u_m|^2 d\mu \right] \\
&\geq \frac{p-1}{2(p+1)} (1 + C_*^2) \|\nabla_{\mathbb{E}} u_m\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2.
\end{aligned}$$

Then

$$\begin{aligned}
&\int_0^t \frac{1}{2} \|\partial_\tau u_m\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 d\tau + \frac{1}{p+1} \int_0^t \left\| \left( \frac{d}{d\tau} g_\tau \right)^{\frac{1}{p+1}} u_m \right\|_{L_{p+1}^{\frac{n+1}{2}}(\mathbb{E})}^{p+1} d\tau \\
&\quad + \frac{p-1}{2(p+1)} (1 + C_*^2) \|\nabla_{\mathbb{E}} u_m\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \\
&\leq I(t) + \gamma \int_0^t \|\nabla_{\mathbb{E}}(\partial_\tau u_m)\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 d\tau \leq I(0) < d. \quad (19)
\end{aligned}$$

for  $t \in [0, \infty)$  and sufficiently large  $m$ . Now, by relation 19 we can get that

$$\|\nabla_{\mathbb{E}} u_m\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 < \frac{2(p+1)}{p-1} (1 + C_*^2)^{-1} d, \quad (20)$$

for  $t \in [0, \infty)$  and

$$\frac{1}{2} \int_0^t \|\partial_\tau u_m\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 d\tau + \frac{1}{p+1} \int_0^t \left\| \left( \frac{d}{d\tau} g_\tau \right)^{\frac{1}{p+1}} u_m \right\|_{L_{p+1}^{\frac{n+1}{2}}(\mathbb{E})}^{p+1} d\tau < d, \quad (21)$$

for  $t \in [0, \infty)$ . Also we obtain that

$$\begin{aligned}
\int_{\mathbb{E}} r^q |g_t(z)|^{\frac{p}{p+1}} |u_m|^{p-1} |u_m|^{\frac{p+1}{p}} d\mu &= \int_{\mathbb{E}} r^q g_t(z) |u_m|^{p+1} d\mu \leq C_{**}^{p+1} \|\nabla_{\mathbb{E}} u_m\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^{p+1} \\
&< C_{**}^{p+1} \left( \frac{2(p+1)}{p-1} (1 + C_*^2)^{-1} d \right)^{\frac{p+1}{2}} \quad (22)
\end{aligned}$$

and

$$\begin{aligned}
\int_{\mathbb{E}} r^q |V(z)^{\frac{1}{2}} u_m|^2 d\mu &= \int_{\mathbb{E}} r^q V(z) |u_m|^2 d\mu \leq C_*^2 \|\nabla_{\mathbb{E}} u_m\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \\
&< C_*^2 \left( \frac{2(p+1)}{p-1} (1 + C_*^2)^{-1} d \right)^2. \tag{23}
\end{aligned}$$

From 20, 21, 22 and 23, it follows that there exists  $u$  and a subsequence still denotes  $\{u_m\}$  for which as  $m \rightarrow \infty$ ,  $u_m \rightarrow u$  in  $L^\infty(0, \infty; \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}))$  weakly star and a.e. in  $\text{int}\mathbb{E} \times [0, \infty)$ ,  $\partial_t u_m \rightarrow \partial_t u$  in  $L^2(0, \infty; L_2^{\frac{n+1}{2}}(\mathbb{E}))$ , weakly star. Also we have  $V(z)|u_m|^2 \rightarrow V(z)|u|^2$  in  $L^\infty(0, \infty; \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}))$  weakly star and a.e. in  $\text{int}\mathbb{E} \times [0, \infty)$ , and  $g_t(z)u_m|u_m|^{p-1} \rightarrow g_t(z)u|u|^{p-1}$  in  $L^\infty(0, \infty; \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}))$  weakly star and a.e. in  $\text{int}\mathbb{E} \times [0, \infty)$ . Therefore, in 13 for  $k$  fixed and  $m \rightarrow \infty$  we get that

$$\begin{aligned}
(\gamma u, w_k)_2 + (u_t, w_k)_2 + \int_0^t (\nabla_{\mathbb{E}} u, \nabla_{\mathbb{E}} w_k)_2 d\tau + \int_0^t (V(z)u, w_k)_2 d\tau \\
= \int_0^t (g_t(z)u|u|^{p-1}, w_k)_2 d\tau \\
+ (\gamma u_0, w_k)_2 + (u_1, w_k)_2.
\end{aligned}$$

On the other hand, from the relation 14,  $u(z, 0) = u_0(z)$  in  $\mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E})$  and from 15  $\partial_t u(z, 0) = u_1$  in  $L_2^{\frac{n+1}{2}}(\mathbb{E})$ . By density we obtain  $u \in L^\infty(0, \infty; \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}))$  with  $\partial_t u \in L^2(0, \infty; L_2^{\frac{n+1}{2}}(\mathbb{E}))$  is global weak solution of problem 1. Since  $u$  satisfies problem 1, so by definition of  $K$  we have  $K(u) = 0$ . Hence,  $u(t) \in W$  for  $0 \leq t < \infty$ .

**Corollary 2** *If we replace the assumption  $I(0) < d$ ,  $K(u_0) > 0$  by  $0 < I(0) < d$ ,  $K_{\delta_2}(u_0) > 0$  where  $(\delta_1, \delta_2)$  is the maximal interval including  $\delta = 1$ , (see Remark 1) such that  $I(0) < d(\delta)$  for  $\delta \in (\delta_1, \delta_2)$ . Then problem 1 admits a global weak solution  $u(t) \in L^\infty(0, \infty; \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}))$  with  $\partial_t u \in L^\infty(0, \infty, L_2^{\frac{n+1}{2}}(\mathbb{E}))$  and  $u(t) \in W_\delta$  for  $\delta \in (\delta_1, \delta_2)$ ,  $t \in [0, \infty)$ .*

**Proof** It is immediately implied from Theorems 1 and 3.

**Corollary 3** *If we replace the assumption  $K_{\delta_2}(u_0) > 0$  or  $\|\nabla_{\mathbb{E}} u_0\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} = 0$ , by  $\|\nabla_{\mathbb{E}} u_0\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} < r(\delta_2)$ , then problem 1 admits a global weak solution  $u(t) \in L^\infty(0, \infty; \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}))$  with  $\partial_t u(t) \in L^\infty(0, \infty; L_2^{\frac{n+1}{2}}(\mathbb{E}))$  satisfying*

$$\|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \leq \frac{I(0)}{a(\delta_1)}, \quad \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \leq 2I(0), \quad 0 \leq t \leq \infty \quad (24)$$

**Proof** From assumption  $\|\nabla_{\mathbb{E}} u_0\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} < r(\delta_2)$ , we can get that  $K_{\delta_2}(u_0) > 0$  or  $\|\nabla_{\mathbb{E}} u_0\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} = 0$ . Then it follows from Corollary 2 that problem 1 admits a global weak solution such that for any  $\delta_1 < \delta < \delta_2$ ,  $0 \leq t < \infty$ ,  $u(t) \in L^\infty(0, \infty; \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}))$  with  $\partial_t u \in L^\infty(0, \infty; L_2^{\frac{n+1}{2}}(\mathbb{E}))$  and  $u(t) \in W_\delta$ . Moreover, similar of the proof Theorem 2 for every  $\delta_1 < \delta < \delta_2$ ,  $0 \leq t < \infty$ ,

$$\frac{1}{2} \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + a(\delta) \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{p+1} K_\delta(u) \leq I(0).$$

If we tend  $\delta$  to  $\delta_1$  then we achieve 24.

Now we discuss the global non-existence of solutions of the problem 1.

**Theorem 4** Let  $0 \leq \gamma \leq (p-1)\sqrt{1+C_*^2\lambda_1^{\frac{1}{2}}}$ ,  $u_0 \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E})$ ,  $u_1 \in L_2^{\frac{n+1}{2}}(\mathbb{E})$ . Suppose that  $I(0) < d$  and  $K(u_0) < 0$ . Then the existence time of solution for problem 1 is finite, where  $\lambda_1$  is the first eigenvalue in Proposition 3 i.e.

$$\lambda_1^{\frac{1}{2}} = \inf_{u \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E}), \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} \neq 0} \frac{\|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}}{\|u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}}.$$

**Proof** Let  $u(t)$  be any weak solution of problem 1 with  $I(0) < d$  and  $K(u_0) < 0$ ,  $T$  be the maximal existence time of  $u(t)$ . We will prove  $T < \infty$  by contradiction. Let  $M(t) := \|u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2$ , then

$$\dot{M}(t) = \frac{d}{dt} \int_E r^q |u(z, t)|^2 d\mu = 2(\partial_t u, u)_2,$$

from definition of functional  $K$ ,

$$\ddot{M}(t) = 2\|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + 2(\partial_t^2 u, u)_2 = 2\|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - 2\gamma(\nabla_{\mathbb{E}}(\partial_t u), \nabla_{\mathbb{E}} u)_2 - 2K(u). \quad (25)$$

Using proof of Theorem 2 we can get,

$$\frac{1}{2} \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + a(1) \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{p+1} K(u)$$

$$\begin{aligned}
 &= \frac{1}{2} \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \left( \frac{1}{2} - \frac{1}{p+1} \right) (1 + C_*^2) \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \\
 &+ \frac{1}{p+1} \left( \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \int_{\mathbb{E}} r^q V(z) |u|^2 d\mu - \int_{\mathbb{E}} r^q g_r(z) |u|^{p+1} d\mu \right) \\
 &\leq \frac{1}{2} \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{1}{2} \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \\
 &\quad + \left[ \left( \frac{1}{2} - \frac{1}{p+1} \right) + \frac{1}{p+1} \right] \int_{\mathbb{E}} r^q V(x) |u|^2 d\mu \\
 &- \frac{1}{p+1} \int_{\mathbb{E}} r^q g_r(z) |u|^{p+1} d\mu \leq \frac{1}{2} \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \\
 &\quad + \frac{1}{p+1} \int_0^t \left\| \left( \frac{d}{d\tau} g_\tau(z) \right)^{\frac{1}{p+1}} u \right\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^{p+1} d\tau \\
 &+ J(u) + \gamma \int_0^t \|\nabla_{\mathbb{E}}(\partial_\tau u)\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 d\tau \\
 &= I(t) + \gamma \int_0^t \|\nabla_{\mathbb{E}}(\partial_\tau u)\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 d\tau \leq I(0). \tag{26}
 \end{aligned}$$

Thus inequality 26 implies that

$$\begin{aligned}
 \ddot{M}(t) &\geq 2 \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \\
 &- 2\gamma(\partial_t u, u)_2 - 2(p+1) \left[ I(0) - \frac{1}{2} \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \right. \\
 &\quad \left. - \frac{p-1}{2(p+1)} (1 + C_*^2) \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \right] \\
 &= (p+3) \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + (p-1)(1 + C_*^2) \|\nabla_{\mathbb{E}} u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \\
 &- 2\gamma(\nabla_{\mathbb{E}}(\partial_t u), \nabla_{\mathbb{E}} u)_2 - 2(p+1)I(0). \tag{27}
 \end{aligned}$$

In first, let us consider  $I(0) \leq 0$ . Then,

$$\ddot{M}(t) \geq (p+3) \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + (p-1)(1 + C_*^2) \lambda_1 \|u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - 2\gamma(\nabla_{\mathbb{E}}(\partial_t u), \nabla_{\mathbb{E}} u)_2.$$

condition  $\gamma < (p-1)(1+C_*^2)\lambda_1$  implies that, there exists a constant  $\epsilon \in \left(0, (p-1)(1+C_*^2)\right)$  such that

$$\gamma^2 < (p-1-\epsilon)(1+C_*^2)\lambda_1.$$

Therefore,

$$\begin{aligned} \ddot{M}(t) &\geq (4+\epsilon)\|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + (p-1-\epsilon)\|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - 2\gamma(\nabla_{\mathbb{E}}(\partial_t u), \nabla_{\mathbb{E}}u)_2 \\ &\quad + (p-1)(1+C_*^2)\lambda_1^2\|u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2. \end{aligned} \quad (28)$$

On the other hand,

$$\begin{aligned} 2\gamma(\nabla_{\mathbb{E}}(\partial_t u), \nabla_{\mathbb{E}}u)_2 &\leq (p-1-\epsilon)\|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{\gamma^2}{p-1-\epsilon}\|u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 \\ &\leq (p-1-\epsilon)\|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + (p-1)(1+C_*^2)\lambda_1^2\|u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2. \end{aligned} \quad (29)$$

From 28 and 29, we can get that

$$\ddot{M}(t) \geq (4+\epsilon)\|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2. \quad (30)$$

By Edge Hölder inequality we get

$$M(t)\ddot{M}(t) - \frac{4+\epsilon}{4}\dot{M}(t) \geq (4+\epsilon)\left(\|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2\|u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - (\nabla_{\mathbb{E}}(\partial_t u), \nabla_{\mathbb{E}}u)_2\right) \geq 0,$$

$$(M^{-\alpha})'' = \frac{-\alpha}{M^{\alpha+2}(t)}\left(M(t)\ddot{M}(t) - (\alpha+1)\dot{M}(t)^2\right) \leq 0,$$

for  $\alpha = \frac{\epsilon}{4}$  and  $0 \leq t < \infty$ . Hence, there exists a  $T_1 > 0$  such that

$$\lim_{t \rightarrow T_1} M^{-\alpha}(t) = 0$$

and  $\lim_{t \rightarrow T_1} M(t) = +\infty$ , which is contradicts  $T_{\max} = +\infty$ .

In second case, we consider  $0 < I(0) < d$ . In this case from Theorem 1 we have  $u \in E_\delta$  for  $0 \leq t < \infty$  and  $\delta \in (1, \delta_2)$  (see Remark 1) where  $(\delta_1, \delta_2)$  is the maximal interval including  $\delta = 1$  such that  $d(\delta) > I(0)$  for  $\delta \in (\delta_1, \delta_2)$ .

Therefore,  $K_\delta(u) < 0$  and  $\|\nabla_{\mathbb{E}}u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})} > r(\delta)$  for  $1 < \delta < \delta_2$ ,  $0 \leq t < \infty$ .  
 Consequent,  $K_\delta(u) \leq 0$  and  $\|\nabla_{\mathbb{E}}u\| \geq r(\delta)$  for  $0 \leq t < \infty$ . From 25,

$$\begin{aligned} \frac{d}{dt}(e^{\gamma t} \dot{M}(t)) &= e^{\gamma t} \left( \gamma \dot{M}(t) + \ddot{M}(t) \right) = 2e^{\gamma t} \left( \|\partial_t u\|_{L_2^{\frac{n}{2}}(\mathbb{E})}^2 - K(u) \right) \\ &= 2e^{\gamma t} \left( \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 + (\delta_2 - 1) \|\nabla_{\mathbb{E}}u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - K_{\delta_2}(u) \right) \\ &\geq 2e^{\gamma t} (\delta_2 - 1) r^2 \delta_2 = C \delta_2 e^{\gamma t}. \end{aligned}$$

Hence,

$$\begin{aligned} e^{\gamma t} \dot{M}(t) &\geq C \delta_2 \int_0^t e^{\gamma \tau} d\tau + \dot{M}(0) = \frac{C \delta_2}{\gamma} (e^{\gamma t} - 1) + \dot{M}(0), \\ \dot{M}(t) &\geq \frac{C \delta_2}{\gamma} (1 - e^{-\gamma t}) + e^{-\gamma t} \dot{M}(0). \end{aligned}$$

Hence there exists  $t_0 > 0$  for which

$$\dot{M}(t) \geq \frac{C \delta_2}{2\gamma} \quad \forall t \geq t_0$$

and

$$M(t) \geq \frac{C \delta_2}{2\gamma} (t - t_0) + M(t_0) \geq \frac{C \delta_2}{2\gamma} (t - t_0), \quad t \geq t_0. \tag{31}$$

From assumption  $\gamma < (p - 1)(1 + C_*^2)\lambda_1$ , it follows there exists a constant

$$\epsilon \in \left( 0, (p - 1)(1 + C_*^2) \right)$$

such that

$$\gamma^2 < (p - 1 - \epsilon) \left[ (p - 1)(1 + C_*^2)\lambda_1 - \epsilon \right].$$

From 27,

$$\begin{aligned} \ddot{M}(t) &\geq (p + 3) \|\partial_t u\|_{L_2^{\frac{n+1}{2}}(\mathbb{E})}^2 - 2\gamma (\nabla_{\mathbb{E}}(\partial_t u), \nabla_{\mathbb{E}}u)_2 + (p - 1)(1 + C_*^2)\lambda_1^2 \|u\|_{L_2^{\frac{n}{2}}(\mathbb{E})}^2 \\ &\quad - 2(p + 1)I(0) \end{aligned}$$

$$\begin{aligned}
 &= \|(4 + \epsilon)\|\partial_t u\|_{L^2_{\frac{n}{2}}(\mathbb{E})}^2 + (p - 1 - \epsilon)\|\partial_t u\|_{L^2_{\frac{n+1}{2}}(\mathbb{E})}^2 - 2\gamma(\nabla_{\mathbb{E}}(\partial_t u), \nabla_{\mathbb{E}}u)_2 \\
 &+ [(p - 1)(1 + C_*^2)\lambda_1^2 - \epsilon]\|\partial_t u\|_{L^2_{\frac{n}{2}}(\mathbb{E})}^2 + \epsilon M(t) - 2(p + 1)I(0). \tag{32}
 \end{aligned}$$

Also we can obtain

$$\begin{aligned}
 2\gamma(\nabla_{\mathbb{E}}(\partial_t u), \nabla_{\mathbb{E}}u)_2 &\leq (p - 1 - \epsilon)\|\partial_t u\|_{L^2_{\frac{n+1}{2}}(\mathbb{E})}^2 + \frac{\gamma^2}{p - 1 - \epsilon}\|u\|_{L^2_{\frac{n+1}{2}}(\mathbb{E})}^2 \\
 &\leq (p - 1 - \epsilon)\|\partial_t u\|_{L^2_{\frac{n+1}{2}}(\mathbb{E})}^2 \\
 &+ [(p - 1)(1 + C_*^2)\lambda_1^2 - \epsilon]\|u\|_{L^2_{\frac{n+1}{2}}(\mathbb{E})}^2. \tag{33}
 \end{aligned}$$

From 32 and 33 we get

$$\ddot{M}(t) \geq (4 + \epsilon)\|\partial_t u\|_{L^2_{\frac{n+1}{2}}(\mathbb{E})}^2 + \epsilon M(t) - 2(p + 1)I(0). \tag{34}$$

From 31, it follows that there exists a  $t_1 > 0$  such that

$$\epsilon M(t) > 2(p + 1)I(0) \quad \forall t > t_1,$$

and then

$$\ddot{M}(t) > (4 + \epsilon)\|\partial_t u\|_{L^2_{\frac{n+1}{2}}(\mathbb{E})}^2, \quad \forall t > t_1.$$

Now, similar to first case we can obtain a contradiction. Hence we always have  $T_{\max} < \infty$ .

From Theorems 13 and 4 we can obtain the following theorem for global existence and non-existence of solutions for problem 1.

**Theorem 5** *Let  $0 \leq \gamma \leq (p-1)\sqrt{1 + C_*^2}\lambda_1^{\frac{1}{2}}$ ,  $u_0 \in \mathcal{H}_{2,0}^{1, \frac{n+1}{2}}(\mathbb{E})$  and  $u_1 \in L^2_{\frac{n+1}{2}}(\mathbb{E})$ . Suppose that  $I(0) < 0$ . Then, when  $K(u_0) > 0$ , problem 1 admits a global weak solution and when  $K(u_0) < 0$ , problem 1 does not admits any global weak solution.*

## References

1. M. Alimohammady, M.K. Kalleji, Existence results for a class of semilinear totally characteristic hypoelliptic equations with conic degeneration. *J. Func. Anal.* **265**, 2331–2356 (2013)
2. M. Alimohammady, C. Carlo, M.K. Kalleji, Invariance and existence analysis for semilinear hyperbolic equations with damping and conical singularity. *J. Math. Anal. Appl.* **455**, 569–591 (2017)

3. M. Alimohammady, M.K. Kalleji, Gh. Karamali, Global results for semilinear hyperbolic equations with damping term on manifolds with conical singularity. *Math. Methods Appl. Sci.* **40**(11), 4160–4178 (2017)
4. G. AustinFord, J. Wunsch, The diffractive wave trace on manifolds with conic singularities. *Adv. Math.* **304**, 1330–1385 (2017)
5. G. Chen, F. Da, Blow-up of solution of Cauchy problem for three-dimensional damped nonlinear hyperbolic equation. *Nonlinear. Anal.* **71**(1–2), 358–372 (2009)
6. H. Chen, G. Liu, Global existence and nonexistence for semilinear parabolic equations with conical degeneration. *J. Pseudo-Differ. Oper. Appl.* **3**, 329–349 (2012)
7. H. Chen, X. Liu, Asymptotic stability and blow-up of solutions for semi-linear degenerate parabolic equations with singular potential. *Discrete Contin. Dyn. Syst.* **36**(2), 661–682 (2016)
8. H. Chen, X. Liu, Y. Wei, Dirichlet problem for semilinear edge-degenerate elliptic equations with singular potential term. *J. Differ. Equ.* **252**, 4289–4314 (2012)
9. N.J. Daras, Th.M. Rassias, *Computational Mathematics and Variational Analysis* (Springer, Cham, 2020). <https://doi.org/10.1007/978-3-030-44625-3>
10. Y.V. Egorov, B.W. Schulze, *Pseudo-differential Operators, Singularities, Applications* (Springer, Basel, 1997)
11. S. Ervedoza, Control and stabilization properties for a singular heat equation with an inverse-square potential. *Commun. Partial Differential Equations* **33**(10–12), 1996–2019 (2008)
12. V. Felli, E.M. Marchini, S. Terracini, On Schrödinger operators with multipolar inverse-square potentials. *J. Funct. Anal.* **250**(2), 265–316 (2007)
13. X. Jiang, R. Xu, Global well-posedness for semilinear hyperbolic equations with dissipative term. *J. Appl. Math. Comput.* **38**, 467–687 (2012)
14. M.K. Kalleji, M. Alimohammady, A.A. Jafari, Multiple solutions for class of nonhomogeneous semilinear equations with critical cone Sobolev exponent. *Proc. Amer. Math. Soc.* **147**, 597–608 (2019)
15. A. Kashuri, T.M. Rassias, R. Liko, Some new integral inequalities via general fractional operators, in *Computational Mathematics and Variational Analysis*, ed. by N. Daras, T. Rassias. Springer Optimization and Its Applications, vol. 159 (Springer, Cham, 2020). [https://doi.org/10.1007/978-3-030-44625-3\\_9](https://doi.org/10.1007/978-3-030-44625-3_9)
16. D.H. Lehmer, Euler constants for arithmetical progressions. *Acta Arith.* **27**, 125–142 (1975)
17. H.A. Levine, Instability and non-existence of global solutions to nonlinear wave equations of the form  $Pu_{tt} = -Au + F(u)$ . *Trans. Am. Math. Soc.* **192**, 1–21 (1974)
18. G. Li, J. Yu, W. Liu, Global existence, exponential decay and finite time blow-up of solutions for a class of semilinear pseudo-parabolic equations with conical degeneration. *J. Pseudo-Differ. Oper. Appl.* **8**(2–17), 629–660 (2017)
19. Y. Liu, On potential wells and vacuum isolating of solutions for semilinear wave equations. *J. Differential Equations* **192**, 155–169 (2003)
20. Y. Liu, J. Zhao, On potential wells and applications to semilinear hyperbolic and parabolic equations. *Nonlinear Anal.* **64**, 2665–2687 (2006)
21. R. Melrose, A. Vasy, J. Wunsch, Propagation of singularities for the wave equation on edge manifolds. *Duke Math. J.* **144**(1), 109–193 (2008)
22. M. Mikolás, On certain sums generating the Dedekind sums and their reciprocity laws. *Pacific J. Math.* **7**, 1167–1178 (1957)
23. M. Ramazannejad, M. Alimohammady, C. Cattani, On algorithms for difference of monotone operators, in *Computational Mathematics and Variational Analysis*, ed. by N. Daras, T. Rassias. Springer Optimization and Its Applications, vol. 159 (Springer, Cham, 2020). [https://doi.org/10.1007/978-3-030-44625-3\\_21](https://doi.org/10.1007/978-3-030-44625-3_21)
24. T.M. Rassias, M. Pardalos, *Mathematical Analysis and Applications* (Springer International Publishing, Cham, 2019)
25. M. Reed, B. Simon, *Methods of Modern Mathematical Physics* (Academic Press, New York, 1980)
26. H.D. Sattinger, On global solutions of nonlinear hyperbolic equations. *Arch. Ration. Mech. Anal.* **30**, 148–172 (1975)



27. B.W. Schulze, *Boundary Value Problems and Singular Pseudo-Differential Operators* (Wiley, Chichester 1998)
28. Z. Yang, Global existence, asymptotic behavior and blow-up solutions for a class of nonlinear wave equations with dissipative term. *J. Differential Equations* **187**(2), 520–540 (2003)

# $\phi^4$ Solitons in Kirchhoff Wave Equation



Y. Contoyiannis, P. Papadopoulos, M. Kampitakis, S. M. Potirakis,  
and N. L. Matiadou

**Abstract** We express the Kirchhoff wave equation in terms of classic field theory. This permits us to introduce the spontaneous symmetry breaking phenomenon in the study of linear structures, such as strings in order to investigate the existence of solitons solutions. We find  $\phi^4$  solitons in the space of spatial gradient of lateral displacement of a string. This helps us detect stable states in deformations of strings.

## 1 Introduction

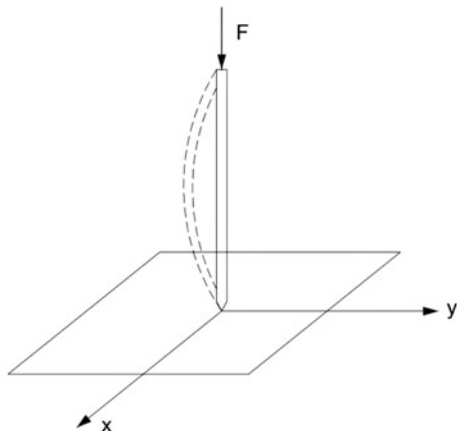
In the last 30 years important progress has been made in understanding of properties of certain non-linear differential equations which arise in many different areas of Physics, e.g., physics of plasma, solid state physics, biophysics, field theory etc. [1–5]. A common interesting feature is the occurrence of solitons, i.e., stable, non-dissipative and localized configurations behaving in many ways like particles. In the analysis of these equations many interesting mathematical structures have been discovered which surprisingly also appear in quantum mechanics and quantum field theory [6]. From a pragmatic point of view these completely soluble non-linear equations are a substantial extension of the ‘tool kit’ of a physicist which otherwise is mainly restricted to solving linear systems. They also serve as valuable source for intuition about the behavior of non-linear systems. In Mathematics and Physics, a soliton, or solitary wave, is a self-reinforcing wave-packet that maintains its shape while it propagates at a constant velocity. Solitons are caused by a cancellation of nonlinear and dispersive effects in the medium. Solitons are the solutions of

---

Y. Contoyiannis · P. Papadopoulos (✉) · S. M. Potirakis · N. L. Matiadou  
Department of Electrical and Electronics Engineering, University of West Attica, Athens, Greece  
e-mail: [yiaconto@uniwa.gr](mailto:yiaconto@uniwa.gr); [ppapadop@uniwa.gr](mailto:ppapadop@uniwa.gr); [spoti@uniwa.gr](mailto:spoti@uniwa.gr); [lmatiadou@uniwa.gr](mailto:lmatiadou@uniwa.gr)

M. Kampitakis  
Hellenic Electricity Distribution Network Operator SA, Network Major Installations Department,  
N.Faliro, Greece  
e-mail: [m.kampitakis@deddie.gr](mailto:m.kampitakis@deddie.gr)

**Fig. 1** The cylindrical symmetry of the system around  $z$  axis is broken by the buckling of the string



a widespread class of weakly nonlinear dispersive partial differential equations describing physical systems.

$\phi^4$  solitons are stable solutions which appear in spontaneous symmetry breaking (SSB) in scalar field theories [7]. A category of systems for which SSB might happen are linear structures such as rod, string, needle etc. Thus, if a string is compressed by the application of a force  $F$  along its axis the obvious solution is that it stays in the configuration  $x = y = 0$  (see Figure 1). However, if the force gets too large ( $F > F_{cr}$ ), the string will jump into a bent position. It does this because the energy in this state is lower than in meta-stable state, where it stays aligned along the  $z$  axis. The cylindrical symmetry of the system around  $z$  axis is broken by the buckling of the string [7].

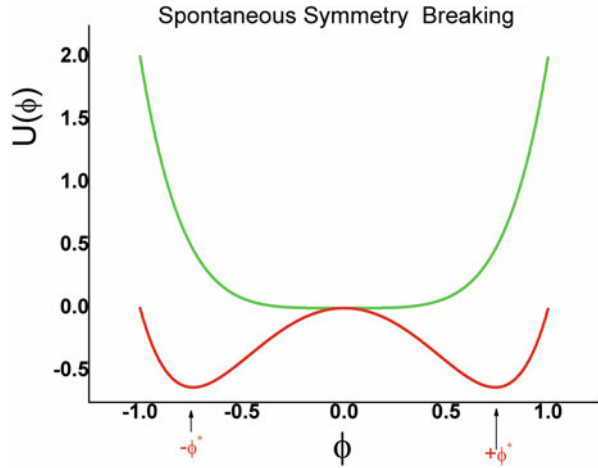
In a  $\phi^4$ —scalar field the SSB is sourced from a concrete type of potential density. This SSB produces solitary waves which ensure the stable behavior of  $\phi$  scalar field. The wave function of the lateral displacement  $u$  for a string such as the one in Figure 1 is the solution of the Kirchhoff wave equation [8]. In this work we attempt to reveal similarities between the potential density produced from Kirchhoff wave equation and  $\phi^4$ —scalar field potential density. This would permit us to reveal the existence of soliton in the Kirchhoff description. This would help us in order to find stable states when the string suffers lateral deformation under the action of axial tensions. This is the main motivation of the present work.

## 2 SSB in $\phi^4$ Scalar Field Theory: The Kink Solitons

The Lagrangian density of scalar field  $\phi(x)$  with a  $\phi^4$  interaction is given as [9]:

$$\mathcal{L} = \frac{1}{2}(\partial_\mu\phi)(\partial^\mu\phi) - \left\{ \frac{1}{2}\alpha\phi^2 + \frac{1}{4}\lambda\phi^4 \right\}, \quad (1)$$

**Fig. 2** The SSB in scalar field theory. The critical point (0,0) behaves as a saddle-point



where  $\mu = 0, 1, 2, 3$  with  $\partial_0\phi = \frac{\partial\phi}{\partial t}, \partial_1\phi = -\frac{\partial\phi}{\partial x}, \partial_2\phi = -\frac{\partial\phi}{\partial y}, \partial_3\phi = -\frac{\partial\phi}{\partial z}$

The term of kinetic density is  $\frac{1}{2}(\partial_\mu\phi)(\partial^\mu\phi)$  and the potential density is:

$$U(\phi) = \frac{1}{2}\alpha\phi^2 + \frac{1}{4}\lambda\phi^4 \tag{2}$$

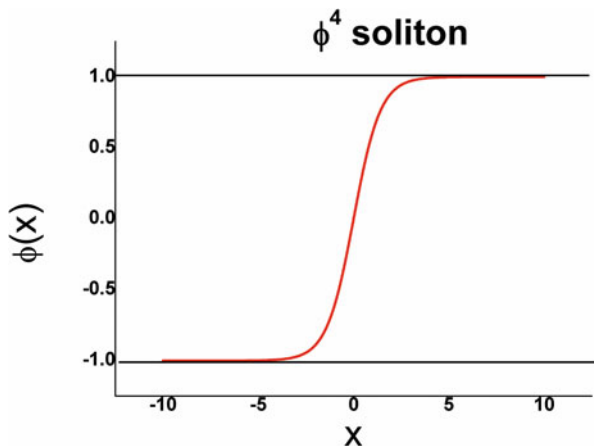
When  $\alpha, \lambda$  have positive values we have the symmetric phase (green line in Figure 2). When  $\alpha < 0, \lambda > 0$  we have the phase of symmetry breaking (SB) (red line in Figure 2). Thus, the ground state of energy has shifted from  $\phi = 0$  to  $\phi^* = \pm\sqrt{\frac{|\alpha|}{\lambda}}$ . This is the SSB phenomenon where the system should select a new vacuum. In a thermal system the parameter  $\alpha$  is a function of  $\frac{T-T_c}{T_c}$  where  $T$  is the temperature and  $T_c$  is the critical temperature. In a similar way, in a string SSB the parameter  $\alpha$  could be a function of  $\frac{F-F_c}{F_c}$ . A model which demonstrates SSB is the  $\phi^4$  theory [9] where the potential has the form:

$$U(\phi) = \frac{\lambda}{4}\left(\phi^2 - \frac{\alpha}{\lambda}\right)^2 = -\frac{1}{2}|\alpha|\phi^2 + \frac{1}{4}\lambda\phi^4 + \frac{\alpha^2}{4\lambda} \tag{3}$$

The above potential refers to the SB phase ( $\alpha < 0, \lambda > 0$ ) and a constant term has been added. So, the critical state (0, 0) is excited to  $(0, \frac{\alpha^2}{4\lambda})$ . This state describes the meta-stable state of the string, before the cylindrical symmetry of the system around z axis is broken by the buckling of the string. The potential of Equation (3) has the same minima as the potential of Equation (2) namely  $\phi^* = \pm\sqrt{\frac{|\alpha|}{\lambda}}$ . This means that solitons solutions, if they exist, must asymptotically tend toward these values as  $x \rightarrow \pm\infty$ , that is:

$$\phi(|x| = \infty) = \pm\sqrt{\frac{|\alpha|}{\lambda}} \tag{4}$$

**Fig. 3** The  $\phi^4$  soliton with  $\alpha = \lambda = 1$  and  $x_0 = 0$



We can integrate the  $\phi^4$  theory given by the Equation (3) to yield [9]:  $x - x_0 = \pm \int_{\phi(x_0)}^{\phi(x)} \frac{d\phi}{\sqrt{\frac{\lambda}{2} (\phi^2 - \frac{\alpha}{\lambda})}}$  and inverting, we find that [9]:

$$\phi(x) = \pm \left( \sqrt{\frac{|\alpha|}{\lambda}} \right) \tanh \left( \frac{\sqrt{\alpha}}{\sqrt{2}} (x - x_0) \right) \quad (5)$$

This is the kink soliton of  $\phi^4$  scalar field (Figure 3).

The most important property of solitons, as it has already been mentioned, is that they are stable structures which behave as particles. The energy density of these solitons is [9]:

$$\varepsilon(x) = \left( \frac{\alpha^2}{2\lambda} \right) \text{sech}^4 \left[ \frac{m(x - x_0)}{\sqrt{2}} \right]. \quad (6)$$

The mass of particle-soliton is given by the integral over the energy density:

$$M = \int_{-\infty}^{\infty} \varepsilon(x) dx = \frac{2\sqrt{2}}{3} \frac{|\alpha|^{\frac{3}{2}}}{\lambda} \quad (7)$$

### 3 Kirchhoff Wave Function, Energy and Potential

The Kirchhoff wave equation without damping term is defined [8, 10, 11], as:

$$\frac{\partial^2 u}{\partial t^2} - \left( 1 + \int_{\Omega} |\nabla u|^2 dx \right) \nabla^2 u = 0 \quad (8)$$

Where  $u(t, x)$  is the lateral displacement of a string at the space coordinate  $x$  and the time  $t$ , while  $\Omega$  is a bounded domain in  $\mathbb{R}^N$  with a smooth boundary  $\partial\Omega$ .

Here the energy is given as  $E(t) = \int_{\Omega} \varepsilon dx$  where the energy density  $\varepsilon$  is [12–14]:

$$\varepsilon(t) = \frac{1}{2} \left| \frac{\partial u}{\partial t} \right|^2 + \left\{ \frac{1}{2} \left| \frac{\partial u}{\partial x} \right|^2 + \frac{1}{4} \left| \frac{\partial u}{\partial x} \right|^2 \left| \frac{\partial u}{\partial x} \right|^2 \right\} \quad (9)$$

The first term in Equation (9) is the kinetic term and the term in the curly brackets is the Kirchhoff potential density (Kpd),  $U_{Kirchhoff}$ . So, we have that:

$$U_{Kirchhoff} = \frac{1}{2} \left| \frac{\partial u}{\partial x} \right|^2 + \frac{1}{4} \left| \frac{\partial u}{\partial x} \right|^2 \left| \frac{\partial u}{\partial x} \right|^2 = \frac{1}{2} \left| \frac{\partial u}{\partial x} \right|^2 + \frac{1}{4} \left| \frac{\partial u}{\partial x} \right|^4 \quad (10)$$

## 4 The Kpd Produced from Classical Field Theory

In this section we will try to produce the Kpd through the classical field theory. Thus, the investigation of solitons in the wave equation is not just the result of comparing potentials but it has a fundamental origin. Let's start from the classical wave equation:

$$\partial_{\mu}^2 \phi = \frac{\partial^2 \phi}{\partial^2 t} - \nabla^2 \phi = 0 \quad (11)$$

Note that in Equation (11) the wave speed constant factor  $c^2$  has been omitted. This is done here since both Equation (8) and Equation (10) that represent Kirchhoff wave equation without damping term and Kpd, respectively, appear in the cited references without constant factors. However, we will restore the specific factor later, during the derivation of the complete form of Kpd (Equation (22)).

If we substitute the  $\nabla^2 \phi$  as  $\nabla^2 \phi (1 + (\frac{\partial \phi}{\partial x})^2)$ , then Equation (8) is written as:

$$\frac{\partial^2 \phi}{\partial^2 t} - \nabla^2 \phi \left( 1 + \left( \frac{\partial \phi}{\partial x} \right)^2 \right) = 0 \quad (12)$$

$$\partial_{\mu}^2 \phi = \frac{\partial^2 \phi}{\partial^2 t} - \nabla^2 \phi = -\frac{\partial U}{\partial \phi} \quad (13)$$

For static solution  $(\frac{\partial^2 \phi}{\partial^2 t}) = 0$  the E-L equation of Equation (13) becomes:

$$-\nabla^2 \phi = -\frac{\partial U}{\partial \phi} \quad (14)$$

where  $U(\phi)$  the potential density of  $\phi$ -field. From Equation (14) we can estimate the potential density  $U$  as follows:

We multiply Equation (14) by  $\frac{\partial\phi}{\partial x}$  and we take the following:

$$-\nabla^2\phi \cdot \frac{\partial\phi}{\partial x} = -\frac{\partial U}{\partial\phi} \cdot \frac{\partial\phi}{\partial x}, \quad (15)$$

which can be integrated over  $x$ , yielding [9]:

$$U(\phi) = \frac{1}{2} \left( \frac{\partial\phi}{\partial x} \right)^2 \quad (16)$$

Following the above procedure, we estimate (see Appendix) the  $U'$  for the case of the wave described by Equation (12) as:

$$U'(\phi) = \frac{1}{2} \left( \frac{\partial\phi}{\partial x} \right)^2 + \frac{1}{4} \left( \frac{\partial\phi}{\partial x} \right)^4 \quad (17)$$

The potential density  $U'(\phi)$  has the same form with Kpd, as expressed in Equation (10).

## 5 SSB in the Kpd

Now we will investigate the SSB in the Kpd. If we compare the Kpd from Equation (17) and the potential density of SSB in the  $\phi^4$  theory from Equation (3), we find out that Equation (3) refers to a scalar field  $\phi$  while Equation (17) refers to gradient of  $\phi$ , that is  $\frac{\partial\phi}{\partial x}$  (or  $\nabla\phi$ ). We face this by defining a new field  $\xi$  as  $\xi \equiv \nabla\phi$ . Thus, we can research solitons in  $\xi$ -field. This means for the string case, that the solitons solutions exist not at the lateral displacement space but in its spatial gradient space. The next thing we have to do, is to introduce coefficients in the terms of Equation (17).

The original equation of Kirchhoff without damping is written as [8]:

$$\frac{\partial^2 u}{\partial t^2} = \nabla^2 u \cdot \left( \frac{p_0}{p_h} + \frac{Y}{p_2 L} \int_0^L \left( \frac{\partial u}{\partial x} \right)^2 dx \right), \quad (18)$$

where  $0 < x < L$ , with  $L$  the length of string,  $Y$  the Young modulus,  $p$  the mass density,  $h$  the cross-section area,  $p_0$  the initial external force. In classical wave equation Equation (11), we normally have a coefficient  $c^2 = \frac{p_0}{p_h}$  in front of the term  $\nabla^2\phi$ . This is in agreement with Equation (18). The substitution  $\nabla^2\phi \rightarrow \frac{p_0}{p_h} \nabla^2\phi$  in Equation (11) transports the coefficient  $\frac{p_0}{p_h}$  in Equation (17) in front of the first term. From Equation (18), we have that the second coefficient is the quantity  $\frac{Y}{p_2 L} > 0$ . Thus, the Kpd is written as:

$$U'(\phi) = \frac{1}{2}\alpha\left(\frac{\partial\phi}{\partial x}\right)^2 + \frac{1}{4}\lambda\left(\frac{\partial\phi}{\partial x}\right)^4, \quad (19)$$

where:

$$\alpha = \frac{p_0}{ph} \quad (20)$$

and

$$\lambda = \frac{Y}{p2L} \quad (21)$$

The next thing we have to do is to add the constant term in Equation (19). This quantity is  $\frac{\alpha^2}{4\lambda} = \frac{p_0^2}{2ph^2Y}L$ . We consider the  $p_0$  as the resultant of axial forces. The quantity  $\frac{p_0^2}{2ph^2Y}L > 0$  expresses the excited potential of string before the cylindrical symmetry of the system around z axis is broken by the buckling of the string. This meta-stable state due to the axial compression  $\Delta L$  from the external force. We have to give an explanation for the negative sign of the coefficient  $\alpha$ , which indicates the symmetry breaking whenever the string leaves the axis and goes to the lateral positions. The symmetry breaking is accomplished when the external axial force overcomes a critical value. Then the internal elasticity forces obtain measure greater than external forces and  $p_0$  obtains negative sign. From Equation (20) the coefficient  $\alpha$  becomes negative too. Therefore, the Kpd for the field  $\xi$  is written in the final form as:

$$U_{Kirchhoff}(\xi) = -\frac{1}{2}|\alpha|\xi^2 + \frac{1}{4}\lambda\xi^4 + \frac{\alpha^2}{4\lambda} \quad (22)$$

Now the Kpd has taken the form of the potential density of SSB (see Equation (3)), which provides the theoretical basis for the formation of kink solitons. The solitons solutions from Equation (5) and (Equations (20) and (21)) is written as:

$$\xi(x) = \pm\sqrt{\frac{p_02L}{Yh}} \tanh\left[\sqrt{\frac{p_0}{2ph}}(x - \Delta L_{cr})\right], \quad (23)$$

where  $\Delta L_{cr}$  is the axial compression when the force overcomes its critical value. The existence of solitons depends from the asymptotic behavior  $\xi(\pm\infty) = \pm\sqrt{\frac{p_02L}{Yh}}$ .

The mass of Kirchhoff soliton which is given in Equation (7) has the form:

$$M_{Kirchhoff} = \frac{4\sqrt{2}}{3} \frac{p_0^{3/2}}{Yh^{3/2}\sqrt{p}}L \quad (24)$$



The transmission length where the solitons survive is determined by their mass, according to the proportion which connects the mass of a particle and its transmission length in the field theory is:

$$R \sim \frac{1}{M} \quad (25)$$

Thus, for materials with high Young modulus  $Y$ , such as steel, from Equation (24) we obtain that the  $M_{Kirchhoff}$  is small. This means from Equation (25) that transmission length  $R$  where the solitons survive is long and the range of stable state is long too. So, the steel string could be found in lateral positions (Figure 1) with greater stability and without breaking. Nevertheless, from Equation (24) one can determine values of parameters which are possible to give stable states.

## 6 Conclusions

In this work we have produced the potential density of Kirchhoff wave equation, for static solution without damping term, through the classical field theory. Thus, we attempt to study a linear structure such as a string, which suffers lateral deformation under the action of axial tensions through the spontaneous symmetry breaking phenomenon. The result is that  $\phi^4$  solitons in the space of gradient of lateral displacement of string, emerge. The existence of these stable solutions permits us to determine the stability of string deformation, through the extension of spatial range of solitons propagation. This approximation we applied on the issue of elasticity, is a new way to face the limits of elasticity for linear structures.

### Appendix: Estimation of the Potential Density $U'(\phi)$

For the Kirchhoff wave equation  $\frac{\partial^2 \phi}{\partial t^2} - \nabla^2 \phi \left( 1 + \left( \frac{\partial \phi}{\partial x} \right)^2 \right) = 0$  (Equation (12)) we initially consider that exists a potential density  $U'(\phi)$  which satisfies a generalized ‘‘Euler-Lagrange’’ equation that could be written as:

$$-\nabla^2 \phi \left( 1 + \left( \frac{\partial \phi}{\partial x} \right)^2 \right) = -\frac{\partial U'}{\partial \phi}, \quad (26)$$

by proceeding to the replacement  $-\nabla^2 \phi \rightarrow -\nabla^2 \phi \left( 1 + \left( \frac{\partial \phi}{\partial x} \right)^2 \right)$  in the E-L equation for static solution as presented in Equation (14).

We set as potential density  $U'(\phi)$ :

$$U'(\phi) = U(\phi) + U_1(\phi) \tag{27}$$

So, Equation (26) is written as:

$$-\nabla^2\phi - \nabla^2\phi \cdot \left(\frac{\partial\phi}{\partial x}\right)^2 = -\frac{\partial U}{\partial\phi} - \frac{\partial U_1}{\partial\phi}, \tag{28}$$

and by using the E-L of Equation (14) we obtain:

$$\nabla^2\phi \cdot \left(\frac{\partial\phi}{\partial x}\right)^2 = \frac{\partial U_1}{\partial\phi} \tag{29}$$

This equation tells us that the correction term in E-L equation corresponds to a potential density  $U_1$  as the E-L equation corresponds to the potential density  $U$ .

We multiply Equation (29) by  $\frac{\partial\phi}{\partial x}$  to take:

$$\nabla^2\phi \cdot \frac{\partial\phi}{\partial x} \left(\frac{\partial\phi}{\partial x}\right)^2 = \frac{\partial U_1}{\partial\phi} \frac{\partial\phi}{\partial x} \tag{30}$$

Moreover, Equation (30) can be integrated over  $x$ , yielding  $\int \frac{\partial^2\phi}{\partial x^2} \cdot \frac{\partial\phi}{\partial x} \left(\frac{\partial\phi}{\partial x}\right)^2 dx = \int \frac{\partial U_1}{\partial x} \frac{\partial\phi}{\partial x} dx$ .

Thus, we have that:

$$\int \frac{\partial}{\partial x} \left(\frac{\partial\phi}{\partial x}\right) \left(\frac{\partial\phi}{\partial x}\right)^3 dx = \int \frac{\partial U_1}{\partial x} dx = U_1(\phi) \tag{31}$$

The first part of Equation (31) is estimated as:

$$\begin{aligned} \int \frac{\partial}{\partial x} \left(\frac{\partial\phi}{\partial x}\right) \left(\frac{\partial\phi}{\partial x}\right)^3 dx &= \left(\frac{\partial\phi}{\partial x}\right)^4 - \int \left(\frac{\partial\phi}{\partial x}\right) \frac{\partial}{\partial x} \left(\frac{\partial\phi}{\partial x}\right)^3 dx = \left(\frac{\partial\phi}{\partial x}\right)^4 \\ &\quad - \int \frac{\partial\phi}{\partial x} 3\left(\frac{\partial\phi}{\partial x}\right)^2 \frac{\partial}{\partial x} \left(\frac{\partial\phi}{\partial x}\right) dx \Rightarrow \\ 4 \int \frac{\partial}{\partial x} \left(\frac{\partial\phi}{\partial x}\right) \left(\frac{\partial\phi}{\partial x}\right)^3 dx &= \left(\frac{\partial\phi}{\partial x}\right)^4 \Rightarrow \\ \int \frac{\partial}{\partial x} \left(\frac{\partial\phi}{\partial x}\right) \left(\frac{\partial\phi}{\partial x}\right)^3 &= \frac{1}{4} \left(\frac{\partial\phi}{\partial x}\right)^4 \end{aligned} \tag{32}$$

Using Equations (16), (27), (31) and (32) we finally obtain:

$$U'(\phi) = \frac{1}{2} \left(\frac{\partial\phi}{\partial x}\right)^2 + \frac{1}{4} \left(\frac{\partial\phi}{\partial x}\right)^4 \tag{33}$$

## References

1. M. Remoissenet, *Waves Called Solitons: Concepts and Experiments* (Springer, Berlin-Heidelberg, 1999), p. 11
2. S.T. Cundiff, B.C. Collings, N.N. Akhmediev, J.M. Soto-Crespo, K. Bergman, W.H. Knox, Observation of polarization-locked vector solitons in an optical fiber. *Phys. Rev. Lett.* **82**, 3988 (1999)
3. A.S. Davydov, *Solitons in Molecular Systems. Mathematics and its Applications (Soviet Series)*, vol. 61, 2nd edn. (Kluwer Academic Publishers, Dordrecht, 1991)
4. T. Heimburg, A.D. Jackson, *On Soliton Propagation in Biomembranes and Nerves*. *Proc. Natl. Acad. Sci. U. S. A.* **102**(2), 9790–9795 (2005)
5. W. Craig, P. Guyenne, J. Hammack, D. Henderson, C. Sulem, Solitary water wave interactions. *Phys Fluids* **18**, 57106 (2006)
6. L.H. Ryder, *Quantum Field Theory* (Cambridge University Press, Cambridge, 1985)
7. F. Halzen, A.D. Martin *Quarks and Leptons: An Introductory Course in Modern Particle Physics* (Wiley, New York, 1983)
8. G. Kirchhoff, *Vorlesungen über mechanik* (B.G. Teubner, Leipzig, 1883)
9. M. Kaku, *Quantum Field Theory* (Oxford University Press, New York, 1993)
10. P. Papadopoulos, N. Stavrakakis, Central manifold theory for the generalized equation of Kirchhoff strings on  $\mathbb{R}^N$ . *Nonlinear Anal. TMA* **61**, 1343–1362 (2005)
11. P. Papadopoulos, N. Stavrakakis, Compact invariant sets for some quasilinear nonlocal Kirchhoff strings on  $\mathbb{R}^N$ . *Appl. Anal.* **87**, 133–148 (2008)
12. K. Ono, On global existence, asymptotic stability and blowing-up of solutions for some degenerate non-linear wave equations of Kirchhoff type with a strong dissipation. *Math Methods Appl. Sci.* **20**, 151–177 (1997)
13. P. Papadopoulos, N. Stavrakakis, Global existence and blow-up results for an equation of Kirchhoff type on  $\mathbb{R}^N$ . *Topol. Methods Nonlinear Anal.* **17**, 91–109 (2001)
14. P. Papadopoulos, M. Karamolengos, A. Pappas, Global existence and energy decay for mildly degenerate Kirchhoff's equations on  $\mathbb{R}^N$ . *J. Interdisciplinary Math.* **12**(6), 767–783 (2009)

# Estimates for Lipschitz and BMO Norms of Operators on Differential Forms



Shusen Ding, Guannan Shi, and Yuming Xing

**Abstract** In this paper, we introduce the generalized Lipschitz and BMO norms of differential forms and establish the upper bound estimates for the generalized Lipschitz and BMO norms of operators applied to differential forms. We also demonstrate applications of our main results using examples.

## 1 Introduction

The main purpose of this paper is to establish the upper bound estimates for the generalized Lipschitz and BMO norms of the iterated operators  $D^k G^k$  and  $D^{k+1} G^k$  applied to differential forms  $u$  defined in  $\mathbb{R}^n$  in terms of the  $L^p$  norms of  $u$ , where  $k$  is a positive integer;  $G$  is Green's operator and  $D = d + d^*$  is the Hodge-Dirac operator on differential forms. The Dirac operator  $D$  and Green's operator  $G$  are very well studied and widely used in many fields of mathematics and physics. They play a critical role in the study of the nonlinear problems in PDEs and nonlinear potential theory. For example, in the case  $k = 1$ , the composition  $D^2 G$  is used to define the well-known Poisson's equation  $D^2 G(u) = u - H(u)$  (or  $\Delta G(u) = u - H(u)$ ), where  $H$  is the harmonic projection operator. In the same sense as the  $L^p$  theory, the estimates for the BMO norms of differential forms and the related operators are also decisive on the investigation of the solution properties of PDEs, especially on the study of Harnack's inequality for solutions

---

S. Ding

Department of Mathematics, Seattle University, Seattle, WA, USA

e-mail: [sding@seattleu.edu](mailto:sding@seattleu.edu)

G. Shi (✉)

School of Mathematics and Statistics, Northeast Petroleum University, Daqing, China

e-mail: [sgncx@163.com](mailto:sgncx@163.com)

Y. Xing

Department of Mathematics, Harbin Institute of Technology, Harbin, China

e-mail: [xyuming@hit.edu.cn](mailto:xyuming@hit.edu.cn)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,

Springer Optimization and Its Applications 173,

[https://doi.org/10.1007/978-3-030-72563-1\\_5](https://doi.org/10.1007/978-3-030-72563-1_5)

to certain partial differential equations, see [1] for example. Some estimates for the BMO norm and local Lipschitz norm of differential forms or related operators can be found in [2–5]. We should notice that  $D^k G^k$  and  $D^{k+1} G^k$  are more general operators which include the composite operator  $DG$  as a special case, see [5] where  $DG$  has been investigated. However, there is no systematic study on the BMO norm and local Lipschitz norm of the iterated operators  $D^k G^k$  and  $D^{k+1} G^k$  for the case  $k > 1$  in the literature. Hence, we are motivated to establish the upper bound estimates for the generalized Lipschitz and BMO norms of the composite operator in this paper. We first extend the definitions of the classical  $\text{locLip}_\alpha$  and BMO norms into the generalized  $\text{locLip}_\alpha^s$  and  $\text{BMO}^s$  norms, respectively. Then, we study the relationship between these two norms and  $L^p$  norms. The estimates for norms and comparisons of norms are very important in the investigation of the corresponding spaces in analysis. For example, it is well known that the BMO space, the dual of Hardy space, is a substitute of  $L^\infty$  space and has been playing a very indispensable role in harmonic analysis and exterior differential analysis, as well as in the study of the characterization of singular integral operators since it was set forth by John and Nirenberg in 1961. We refer the readers to Chapter IV in [6] and [1, 7] for the function case of the BMO space, and Chapter 9 in [2] and [8–11] for the case to differential forms. Our main results are presented and proved in Section 3. These results will enrich the theory of operators on differential forms.

Unless stated otherwise, we keep using the traditional notation and symbols throughout this paper. Let  $\Omega$  be a smoothly bounded domain without the boundary in  $\mathbb{R}^n$ ,  $n \geq 2$ , and  $B = B(x, \rho)$  be the ball in  $\mathbb{R}^n$  with radius  $\rho$  centered at  $x$ , which satisfies  $\text{diam}(\sigma B) = \sigma \text{diam}(B)$ . Let the direct sum  $\Lambda = \Lambda(\mathbb{R}^n) = \bigoplus_{l=0}^n \Lambda^l(\mathbb{R}^n)$  be a graded algebra with respect to the exterior product, and  $\Lambda^l = \Lambda^l(\mathbb{R}^n)$  be the space of  $l$ -covectors in  $\mathbb{R}^n$ , which is spanned by the dual orthogonal basis  $dx_{i_1}, \dots, dx_{i_l}$ , where  $x_{i_1}, \dots, x_{i_l}$  are the coordinate functions on  $\mathbb{R}^n$ . For the set  $\Lambda$ , we denote the pointwise inner product by  $\langle \cdot, \cdot \rangle$  and the module by  $|\cdot|$ . Then, every differential form  $u(x) \in \Lambda^l(\mathbb{R}^n)$  can be uniquely written as

$$u(x) = \sum_I u_I(x) dx_I = \sum u_{i_1 i_2 \dots i_l}(x) dx_{i_1} \wedge dx_{i_2} \wedge \dots \wedge dx_{i_l},$$

where the coefficients  $u_{i_1 i_2 \dots i_l}(x)$  are differentiable functions and  $I = (i_1, i_2, \dots, i_l)$ ,  $1 \leq i_1 < i_2 < \dots < i_l \leq n$ . Actually, differential forms are the generalizations of the functions, which include functions as their special cases (functions are called 0-forms). The Hodge-star operator  $\star : \Lambda^l(\mathbb{R}^n) \rightarrow \Lambda^{n-l}(\mathbb{R}^n)$  is defined by the rule that  $\star 1 = dx_{i_1} \wedge dx_{i_2} \wedge \dots \wedge dx_{i_l}$  and  $\alpha \wedge \star \beta = \langle \alpha, \beta \rangle (\star 1)$  for every  $\alpha, \beta \in \Lambda^l$ ,  $l = 0, 1, \dots, n$ . By this definition, it induces that  $\star$  is an isometric isomorphism on  $\Lambda^l$ . The linear operator  $d : \mathbb{D}'(\Omega, \Lambda^l) \rightarrow \mathbb{D}'(\Omega, \Lambda^{l+1})$ ,  $0 \leq l \leq n - 1$ , is called the exterior differential and  $d^\star = (-1)^{nl+1} \star d \star : \mathbb{D}'(\Omega, \Lambda^{l+1}) \rightarrow \mathbb{D}'(\Omega, \Lambda^l)$ , the formal adjoint of  $d$ , is known as Hodge codifferential. The interested readers could see [10–13] for further introduction and appropriate properties. Also, we use  $L^p(\Omega, \Lambda)$  to denote the classical  $L^p$  space for

differential forms,  $1 < p < \infty$ , equipped with the norm  $\|u\|_{p,\Omega} = (\int_{\Omega} |u|^p dx)^{\frac{1}{p}} = (\int_{\Omega} (\sum_I |u_I|^2)^{\frac{p}{2}} dx)^{\frac{1}{p}}$ .  $W^{1,p}(\Omega, \Lambda)$  is the classical Sobolev space for differential forms with the norm  $\|u\|_{W^{1,p}(\Omega)} = (\text{diam}(\Omega))^{-1} \|u\|_{p,\Omega} + \|\nabla u\|_{p,\Omega}$ .  $W_d^p(\Omega, \Lambda^l)$  is the space of differential  $l$ -forms such that  $du \in L^p(\Omega, \Lambda^l)$ . Analogously,  $W_{d^*}^p(\Omega, \Lambda^l)$  is the space of differential  $l$ -forms such that  $d^*u \in L^p(\Omega, \Lambda^l)$ . Inspired by these classical spaces for differential forms, we generalize the BMO space and local Lipschitz space as follows.

**Definition 1.1** For every  $\omega \in L_{loc}^s(\Omega, \Lambda^l)$ ,  $s \geq 1$ , we say  $\omega \in \text{BMO}^s(\Omega, \Lambda^l)$  with the norm defined by

$$\|\omega\|_{*,s,\Omega} = \sup_{\sigma Q \subset \Omega} |Q|^{-1/s} \|\omega - \omega_Q\|_{s,Q}, \tag{1.1}$$

if  $\omega$  satisfies  $\sup_{\sigma Q \subset \Omega} |Q|^{-1/s} \|\omega - \omega_Q\|_{s,Q} < \infty$ , where  $l = 0, 1, \dots, n$  and  $\sigma > 1$  is some expansion factor.

**Definition 1.2** For every  $\omega \in L_{loc}^s(\Omega, \Lambda^l)$ ,  $s \geq 1$ ,  $l = 0, 1, \dots, n$  and  $0 < \alpha \leq 1$ , we call  $\omega \in \text{locLip}_{\alpha}^s(\Omega, \Lambda^l)$  with the norm denoted by

$$\|\omega\|_{\text{locLip}_{\alpha}^s(\Omega)} = \sup_{\sigma Q \subset \Omega} |Q|^{-(n+\alpha s)/sn} \|\omega - \omega_Q\|_{s,Q}, \tag{1.2}$$

if  $\omega$  satisfies  $\sup_{\sigma Q \subset \Omega} |Q|^{-(n+\alpha s)/sn} \|\omega - \omega_Q\|_{s,Q} < \infty$ , where  $\sigma > 1$  is some expansion factor.

Especially, for the case  $s = 1$ , the  $\text{BMO}^s$  norm and  $\text{locLip}_{\alpha}^s$  norm just reduce to the following classical BMO norm and  $\text{locLip}_{\alpha}$  norm given in [10] by C. Nolder, respectively.

$$\|\omega\|_{*,1,\Omega} = \|\omega\|_{*,\Omega} = \sup_{\sigma Q \subset \Omega} |Q|^{-1} \|\omega - \omega_Q\|_{1,Q} \tag{1.3}$$

and

$$\|\omega\|_{\text{locLip}_{\alpha}^1(\Omega)} = \|\omega\|_{\text{locLip}_{\alpha}(\Omega)} = \sup_{\sigma Q \subset \Omega} |Q|^{-(n+\alpha)/n} \|\omega - \omega_Q\|_{1,Q} \tag{1.4}$$

Furthermore, notice that  $|Q|^{\alpha/n} \leq |\Omega|^{1/n}$  since  $0 < \alpha \leq 1$  and  $n \geq 1$ , which results in that

$$|Q|^{-1/s} = |Q|^{\alpha/n} |Q|^{-1/s-\alpha/n} \leq |\Omega|^{1/n} |Q|^{-1/s-\alpha/n}$$

So, similarly as the result in [14], we have that there is a constant  $C > 0$ , independent of  $\omega$ , such that

$$\|\omega\|_{*,s,\Omega} \leq C \|\omega\|_{locLip_\alpha^s(\Omega)} \quad (1.5)$$

for every  $\omega \in W^{1,s}(\Omega, \Lambda^l)$ , which enables us to compare the  $locLip_\alpha^s$  norm and the  $BMO^s$  norm for  $D^k G^k$  and  $D^{k+1} G^k$  simply. In addition, from now on, we point out that the constants  $C$  and  $C_i$  employed in this paper,  $i = 1, 2, \dots$ , may differ from one line to the next.

## 2 Local Poincaré-Type Inequalities

In this section, as preparation for the principle assertion, we show the explicit formulas of  $D^k G^k$  and  $D^{k+1} G^k$  and the Poincaré-type inequalities of  $D^k G^k$  and  $D^{k+1} G^k$  by applying the explicit representation in Lemma 2.4 and Lemma 2.5. First, let us start with the brief review of Green's operator  $G$ . For any fixed integer  $l = 0, 1, \dots, n$ , let  $\mathbb{H}$  be the harmonic  $l$ -field denoted by

$$\mathbb{H} = \{u \in W(\Omega, \Lambda) : du = d^*u = 0, u \in L^p, \text{ for some } 1 < p < \infty\}.$$

In the meantime, we take the operator  $\delta : L^p(\Omega, \Lambda) \cap \mathbb{H}^\perp \rightarrow W^{1,p}(\Omega, \Lambda) \cap \mathbb{H}^\perp$  defined by Morrey in [15], which satisfies that for every  $u \in L^p(\Omega, \Lambda) \cap \mathbb{H}^\perp$ ,  $\delta(u)$  is the unique form in  $W^{1,p}(\Omega, \Lambda) \cap \mathbb{H}^\perp$  such that  $\Delta\delta(u) = u$ , where  $\Delta = D^2 = dd^* + d^*d$  is the Laplace operator, and  $\mathbb{H}^\perp$  is the complement space of harmonic field  $\mathbb{H}$ . Therefore, we are given the definition as follows.

**Definition 2.1 ([16])** Green's operator  $G : L^p(\Omega, \Lambda) \rightarrow W^{1,p}(\Omega, \Lambda) \cap \mathbb{H}^\perp$ ,  $1 < p < \infty$ , is defined by

$$G(u) = \delta(u - H(u))$$

for every  $u \in L^p(\Omega, \Lambda)$ , where  $H : L^p(\Omega, \Lambda) \rightarrow \mathbb{H}$  is the projection operator. Moreover, observe that  $\Delta\delta(u) = u$ , so we have that

$$\Delta G(u) = u - H(u). \quad (2.1)$$

By employing the classical dominated convergence theorem, C. Scott in [16] further gave the upper bound estimate of Green's operator  $G$ .

**Lemma 2.2** *Let  $u \in L^s(\Omega, \Lambda)$ ,  $1 < s < \infty$ , be a differential form defined in the domain  $\Omega$ . Then, there exists a positive constant  $C$ , independent of  $u$ , such that*

$$\|dd^*G(u)\|_{s,B} + \|d^*dG(u)\|_{s,B} + \|dG(u)\|_{s,B} + \|d^*G(u)\|_{s,B} + \|G(u)\|_{s,B} \leq C(s)\|u\|_{s,\sigma B} \quad (2.2)$$

for any ball  $B \subset \sigma B \subset \Omega$  with some constant  $\sigma > 1$ , where  $\Omega$  is a smoothly bounded domain without boundary.

*Remark 1* For any  $v \in L^p(\Omega, \Lambda) \cap \mathbb{H}^\perp$ , by the definition of the projection operator  $H$ , it is easy to obtain that  $H(v) = 0$ . Since  $G(u) \in W^{1,p}(\Omega, \Lambda) \cap \mathbb{H}^\perp$  for every  $u \in L^p(\Omega, \Lambda)$ , replacing  $v$  with  $G(u)$  yields that  $H(v) = HG(u) = 0$ . In other words, the harmonic projection of Green's operator  $G$  on  $L^p(\Omega, \Lambda)$  is always equal to zero.

*Remark 2* Also, applying Lemma 2.2 repeatedly, it is obvious to achieve that there is a constant  $C > 0$ , independent of  $u$ , such that

$$\|G^m(u)\|_{p,B} \leq C \|u\|_{p,\sigma B}. \tag{2.3}$$

In particular, if  $u \in W_d^p(\Omega, \Lambda)$  (or  $u \in W_{d^*}^p(\Omega, \Lambda)$ ), we know that Green's operator  $G$  can commute with  $d$  (or  $d^*$ ), which implies that

$$dG(u) = G(du) \text{ or } d^*G(u) = G(d^*u).$$

Similarly to the method employed in (2.3), we have that

$$\|dG^m(u)\|_{p,B} \leq C \|du\|_{p,\sigma B} \text{ or } \|d^*G^m(u)\|_{p,B} \leq C \|d^*u\|_{p,\sigma B} \tag{2.4}$$

for any integer  $m \geq 1$ , where  $\sigma > 1$  is some constant.

Meanwhile, to facilitate the upcoming argument about the Poincaré-type estimates in Theorem 2.6 and Theorem 2.7, we need the following results as well.

**Lemma 2.3 ([17])** *Let  $v \in L_{loc}^p(\Omega, \Lambda^l)$ ,  $1 < p < \infty$ , be a differential form defined in  $\Omega$  and  $T : L^p(\Omega, \Lambda^l) \rightarrow W^{1,p}(\Omega, \Lambda^{l-1})$  be the homotopy operator,  $l = 1, 2, \dots, n$ . Then, we have*

$$v = d(Tv) + T(dv), \tag{2.5}$$

$$\|\nabla(Tv)\|_{p,\Omega} \leq C|\Omega| \|v\|_{p,\Omega} \text{ and } \|Tv\|_{p,\Omega} \leq C|\Omega| \text{diam}(\Omega) \|v\|_{p,\Omega} \tag{2.6}$$

hold for any bounded and convex domain  $\Omega$ .

Before starting the primary argument in this section, it is worth to note that the explicit representations in Lemma 2.4 and Lemma 2.5 are the essential steps for the argument of the Poincaré-type inequalities. In precise, if our attention is only to estimate  $\|D^k G^k(u)\|_{locLip_\alpha^s}$  (or  $\|D^{k+1} G^k(u)\|_{locLip_\alpha^s(\Omega)}$ ) in terms of the  $L^p$  norm  $\|u\|_{p,\Omega}$ , we can prove it directly with the aid of the higher imbedding inequality given in [18]. Otherwise, while we are concerned on the upper boundedness of  $\|D^k G^k(u)\|_{locLip_\alpha^s(\Omega)}$  (or  $\|D^{k+1} G^k(u)\|_{locLip_\alpha^s(\Omega)}$ ) in terms of the BMO<sup>s</sup> norm  $\|u\|_{*,s,\Omega}$ , the higher imbedding result is not valid for this case any more. Thus, to overcome this difficulty, the key tools in our approach are Lemma 2.4 and Lemma 2.5, which are established by adapting the technique developed in [19] with the inductive method.



**Lemma 2.4** *Let  $u \in L^p_{loc}(\Omega, \Lambda)$ ,  $1 < p < \infty$ , be a differential form in the domain  $\Omega$ ,  $D$  be the Hodge-Dirac operator and  $G$  be Green's operator. Then, we have that*

$$D^k G^k(u) = G^m(u), \quad (2.7)$$

$$D^{k+1} G^k(u) = dG^m(u) + d^*G^m(u), \quad (2.8)$$

for every even integer  $k = 2m$  and  $m = 1, 2, \dots$ .

**Proof** First, since  $\Delta = D^2 = (d + d^*)^2 = dd^* + d^*d$ , we know that it holds

$$u = \Delta G(u) + H(u) = dd^*G(u) + d^*dG(u) + H(u) \quad (2.9)$$

for every  $u \in L^p_{loc}(\Omega, \Lambda)$ , which also implies that

$$dd^*G(u) + d^*dG(u) = u - H(u). \quad (2.10)$$

Due to the fact that  $HG(u) = 0$  always holds by Remark 1, replacing  $u$  with  $G^m(u)$  in (2.10) gives that

$$dd^*G(G^m(u)) + d^*dG(G^m(u)) = G^m(u) \quad (2.11)$$

whenever the positive integer  $m \geq 1$ .

Now, we will assert the representation (2.7) by using the inductive method. In the case of  $k = 2$  and  $m = 1$ , we have

$$D^2 G^2(u) = (d + d^*)^2 G^2(u) = dd^*G(G(u)) + d^*dG(G(u)). \quad (2.12)$$

Substituting 2.11 with  $m = 1$  into (2.12) yields that  $D^2 G^2(u) = G(u)$ . Assume that the desired result holds for any  $k = 2(m - 1)$ ,  $m = 2, 3, \dots$ , that is,

$$D^k G^k(u) = D^{2(m-1)} G^{2(m-1)}(u) = G^{m-1}(u). \quad (2.13)$$

Then, when  $k$  is taken as  $2m$ , it continues with (2.13) and (2.11) that

$$\begin{aligned} D^k G^k(u) &= D^2 D^{2(m-1)} G^{2(m-1)}(G^2(u)) = D^2 G^{m+1}(u) \\ &= dd^*G(G^m(u)) + d^*dG(G^m(u)) = G^m(u). \end{aligned} \quad (2.14)$$

So, the desired result (2.7) holds. Moreover, for the operator  $D^{k+1} G^k(u)$ , making use of (2.7) and the fact  $D = d + d^*$ , we obtain that

$$D^{k+1} G^k(u) = D(D^k G^k(u)) = D(G^m(u)) = dG^m(u) + d^*G^m(u).$$

Therefore, we finish the proof of Lemma 2.4.  $\square$

In analogue to the method developed in Lemma 2.4, we also derive the following results for the case  $k = 2m + 1$ .

**Lemma 2.5** *Let  $u \in L^p_{loc}(\Omega, \Lambda)$ ,  $1 < p < \infty$ , be a differential form defined in the domain  $\Omega$ ,  $D$  be the Hodge-Dirac operator and  $G$  be Green's operator. Then, we derive that*

$$D^k G^k(u) = dG^{m+1}(u) + d^*G^{m+1}(u), \quad (2.15)$$

$$D^{k+1} G^k(u) = G^m(u) \quad (2.16)$$

for every odd integer  $k = 2m + 1$  and  $m = 1, 2, \dots$ .

Now, we are ready to give the local Poincaré-type estimates of the iterated operator  $D^k G^k$  and  $D^{k+1} G^k$  in terms of the  $L^p$  norms of  $du$  and  $d^*u$ , respectively.

**Theorem 2.6** *Assume that the differential form  $u$  is of the Sobolev class  $W^{1,p}_{loc}(\Omega, \Lambda)$ ,  $1 < p < \infty$ ,  $D$  is the Hodge-Dirac operator and  $G$  is Green's operator. Then, for any even integer  $k = 2m$ ,  $m = 1, 2, \dots$ , there exists a constant  $C > 0$ , independent of  $u$ , such that*

$$\|D^k G^k(u) - (D^k G^k(u))_B\|_{p,B} \leq C|B|^{1+1/n} \|du\|_{p,\sigma B}, \quad (2.17)$$

$$\|D^{k+1} G^k(u) - (D^{k+1} G^k(u))_B\|_{p,B} \leq C|B|^{1+1/n} \|d^*u\|_{p,\sigma B} \quad (2.18)$$

for all balls  $B \subset \sigma B \subset \Omega$  with some constant  $\sigma > 1$ .

**Proof** Initially, to prove (2.17), applying the decomposition (2.5) to  $D^k G^k(u)$ , we have

$$D^k G^k(u) = dT(D^k G^k(u)) + Td(D^k G^k(u)). \quad (2.19)$$

Since  $dT(D^k G^k(u)) = (D^k G^k(u))_B$ , for every  $p > 1$ , using (2.19), (2.7) and (2.6), it follows that

$$\begin{aligned} \|D^k G^k(u) - (D^k G^k(u))_B\|_{p,B} &= \|Td(D^k G^k(u))\|_{p,B} \\ &\leq C_1|B|\text{diam}(B)\|d(D^k G^k(u))\|_{p,B} \\ &= C_1|B|\text{diam}(B)\|d(G^m(u))\|_{p,B} \\ &\leq C_2|B|^{1+1/n}\|d(G^m(u))\|_{p,B}. \end{aligned} \quad (2.20)$$

Due to the definition of the Sobolev space and the facts that  $\|du\|_{p,\Omega'} \leq \|\nabla u\|_{p,\Omega'} < \infty$  and  $\|d^*u\|_{p,\Omega'} \leq \|\nabla u\|_{p,\Omega'} < \infty$  for any  $\Omega' \subset\subset \Omega$ , one may readily see that Green's operator  $G$  can commute with  $d$  and  $d^*$ . Then, combining (2.20) with (2.4) follows that

$$\|D^k G^k(u) - (D^k G^k(u))_B\|_{p,B} \leq C_3 |B|^{1+1/n} \|du\|_{p,\sigma_1 B}$$

for any even integer  $k > 0$ . Thus, we have (2.17) always holds for all balls  $B \subset \sigma B \subset \Omega$  with some constant  $\sigma_1 > 1$ .

Now, we turn to the proof of the inequality (2.18). First, applying the commute property between  $G$  and  $d^*$  and (2.2), we have

$$\|dd^* G^k(u)\|_{p,B} = \|dG^k(d^*(u))\|_{p,B} \leq C_4 \|d^*u\|_{p,\sigma_2 B}. \quad (2.21)$$

Making use of the similar treatment as in the proof of  $D^k G^k$  with (2.8) and (2.21), we attain that

$$\begin{aligned} \|D^{k+1} G^k(u) - (D^{k+1} G^k(u))_B\|_{p,B} &= \|Td(D^{k+1} G^k(u))\|_{p,B} \\ &\leq C_5 |B| \text{diam}(B) \|d(D^{k+1} G^k(u))\|_{p,B} \\ &\leq C_6 |B|^{1+1/n} \|d(dG^m(u) + d^* G^m(u))\|_{p,B} \\ &= C_6 |B|^{1+1/n} \|dd^* G^m(u)\|_{p,B} \\ &\leq C_7 |B|^{1+1/n} \|d^*u\|_{p,\sigma_2 B} \end{aligned} \quad (2.22)$$

for every even integer  $k > 0$  and some constant  $\sigma_2 > 1$  with all balls  $B \subset \sigma_2 B \subset \Omega$ . Therefore, the proof of Theorem 2.6 is ended.  $\square$

Next, it is natural to take the case of the odd integer  $k > 1$  into account. Using the same process as the case  $k = 2m$  by Lemma 2.5 instead of Lemma 2.4, we derive the results for the odd integer  $k = 2m + 1$ . Considering the length of the paper, we only state the results of Theorem 2.7.

**Theorem 2.7** *Assume that the differential form  $u$  is of the Sobolev class  $W_{loc}^{1,p}(\Omega, \Lambda)$ ,  $1 < p < \infty$ ,  $D$  is the Hodge-Dirac operator and  $G$  is Green's operator. Then, for any odd integer  $k = 2m + 1$ ,  $m = 1, 2, \dots$ , there exists a constant  $C > 0$ , independent of  $u$ , such that*

$$\|D^k G^k(u) - (D^k G^k(u))_B\|_{p,B} \leq C |B|^{1+1/n} \|d^*u\|_{p,\sigma B}, \quad (2.23)$$

$$\|D^{k+1} G^k(u) - (D^{k+1} G^k(u))_B\|_{p,B} \leq C |B|^{1+1/n} \|du\|_{p,\sigma B} \quad (2.24)$$

for all balls  $B \subset \sigma B \subset \Omega$  with some constant  $\sigma > 1$ .

*Remark 3* It should be noticed that the results in Theorem 2.6 and Theorem 2.7 will play a significant role in latter discussion. Specifically, just because of the right terms  $du$  and  $d^*u$  in Theorem 2.6 and Theorem 2.7, it provides us an effective way to derive the upper boundedness of the iterated operators  $D^k G^k$  and  $D^{k+1} G^k$  in terms of the  $BMO^s$  norm for the conjugate  $A$ -harmonic tensors  $u$  and  $v$ .

### 3 Estimates for $BMO^S$ and $locLip_\alpha^S$ Norms

In this section, we present our principal results about the estimates for  $BMO^S$  norm and  $locLip_\alpha^S$  norm for  $D^k G^k$  and  $D^{k+1} G^k$  applied to differential forms  $u$  and  $v$  associated with some conjugate  $A$ -harmonic equation.

During the recent years, the study in the conjugate  $A$ -harmonic tensors is of growing interest and has made much progress, see [2, 10, 20, 21] for examples. Here, we consider the conjugate  $A$ -harmonic tensors of the form as follows.

**Definition 3.1 ([10])** Differential forms  $u \in W^{1,p}(\Omega, \Lambda)$  and  $v \in W^{1,q}(\Omega, \Lambda)$  are called the conjugate  $A$ -harmonic tensors if  $u$  and  $v$  satisfy the conjugate  $A$ -harmonic equation of the form

$$A(du) = d^*v, \tag{3.1}$$

where the operator  $A : \Lambda(\Omega) \rightarrow \Lambda(\Omega)$  is restricted by the following structural assumptions:

- (i) the mapping  $\xi \rightarrow A(\xi)$  is continuous;
- (ii)  $|A(\xi)| \leq a_1|\xi|^{p-1}, \langle A(\xi), \xi \rangle \geq b_1|\xi|^p$ ;
- (iii)  $A(\lambda\xi) = \lambda|\lambda|^{p-2}A(\xi)$  whenever  $\lambda \in \mathbb{R}, \lambda \neq 0$ ;
- (iv) the monotonicity inequality:  $|\langle A(\xi) - A(\eta), \xi - \eta \rangle| \geq L_1(|\xi|^2 + |\eta|^2)^{\frac{p-2}{2}}|\xi - \eta|^2$ .

for all  $\xi \in \Lambda(\mathbb{R}^n)$ . Here,  $a_1, b_1$  and  $L_1 > 0$  are the positive constants and  $1 < p, q < \infty$  are the conjugate exponents with  $1/p + 1/q = 1$ .

According to Definition 3.1, together with the facts  $dd = 0$  and  $d^*d^* = 0$ , it is obvious to see that such a differential form  $u$  in (3.1) is also a solution to the  $A$ -harmonic equation

$$d^*A(du) = 0. \tag{3.2}$$

Moreover, if the operator  $A$  is invertible, in view of the isometric property of the Hodge-star operator  $\star$ , there exists an operator  $B$  such that the differential form  $v$  in (3.1) meanwhile satisfies

$$d^*B(d(\star v)) = 0, \tag{3.3}$$

where the operator  $B : \Lambda(\Omega) \rightarrow \Lambda(\Omega)$  is given the similar conditions i)–iv) that

- (b-i) the mapping  $\xi \rightarrow B(\xi)$  is continuous on  $\Lambda(\mathbb{R}^n)$ ;
- (b-ii)  $|B(\xi)| \leq a_2|\xi|^{q-1}, \langle B(\xi), \xi \rangle \geq b_2|\xi|^q$ ;
- (b-iii)  $B(\kappa\xi) = \kappa|\kappa|^{q-2}B(\xi)$  whenever  $\kappa \in \mathbb{R}, \kappa \neq 0$ ;
- (b-iv) the monotonicity inequality:  $|\langle B(\xi) - B(\eta), \xi - \eta \rangle| \geq L_2(|\xi|^2 + |\eta|^2)^{\frac{q-2}{2}}|\xi - \eta|^2$ .

for almost every  $x \in \Omega$  and all  $\xi \in \wedge^l(\mathbb{R}^n)$ . Here,  $a_2, b_2$  and  $L_2$  are the positive constants and  $1 < q < \infty$  is associated with (3.3).

Observe that  $A$ -harmonic equation is a special case of Dirac-harmonic equation, so we derive the Caccioppoli inequality and the weak reverse Hölder inequality, respectively, by Corollary 2.3 and Theorem 4.3 in [22].

**Lemma 3.2** *Let  $u \in W^{1,p}(\Omega, \Lambda)$  and  $v \in W^{1,q}(\Omega, \Lambda)$  satisfy the conjugate  $A$ -harmonic equation (3.1), and the operator  $A$  be invertible, where  $1 < p, q < \infty$  are the given conjugate exponents with  $1/p + 1/q = 1$ . Then, there exists a constant  $C > 0$ , independent of  $u$  and  $v$ , such that*

$$\|du\|_{p,B} \leq C|B|^{-1/n}\|u - c\|_{p,\sigma B}, \quad (3.4)$$

$$\|d^*v\|_{q,B} \leq C|B|^{-1/n}\|*v - c^*\|_{q,\sigma B} \quad (3.5)$$

for some constant  $\sigma > 1$  and any ball  $B \subset \sigma B \subset \Omega$ , where  $c$  and  $c^*$  are both closed forms.

**Lemma 3.3** *Let  $\omega \in W^{1,p}(\Omega, \Lambda)$  be a solution to the homogenous  $A$ -harmonic equation,  $1 < p < \infty$ . Then, for every  $0 < s, t < \infty$ , there exists a constant  $C > 0$ , independent of  $\omega$ , such that*

$$\|\omega\|_{s,B} \leq C|B|^{1/s-1/t}\|\omega\|_{t,\sigma B}, \quad (3.6)$$

where all balls  $B \subset \sigma B \subset \Omega$  and  $\sigma > 1$  is some constant.

In addition, the local higher order inequality is also necessary for our latter argument.

**Lemma 3.4** *Let  $u \in L^p_{loc}(\Omega, \Lambda)$ ,  $1 < p < \infty$ , be a differential form,  $D$  be the Hodge-Dirac operator and  $G$  be Green's operator. Then, for any positive integer  $k \geq 1$ , we have that*

(i) *if  $1 < p < n$ , for any real number  $0 < s < np/(n - p)$ , there exists a constant  $C > 0$ , independent of  $u$ , such that*

$$\|D^k G^k(u) - (D^k G^k(u))_B\|_{s,B} \leq C|B|^{1+1/n+1/s-1/p}\|u\|_{p,\sigma B}, \quad (3.7)$$

$$\|D^{k+1} G^k(u) - (D^{k+1} G^k(u))_B\|_{p,B} \leq C|B|^{1+1/n+1/s-1/p}\|u\|_{p,\sigma B} \quad (3.8)$$

(ii) *if  $p \geq n$ , for any real number  $s > 0$ , there is a constant  $C > 0$ , independent of  $u$ , such that*

$$\|D^k G^k(u) - (D^k G^k(u))_B\|_{s,B} \leq C|B|^{1+1/n+1/s-1/p}\|u\|_{p,\sigma B}, \quad (3.9)$$

$$\|D^{k+1} G^k(u) - (D^{k+1} G^k(u))_B\|_{p,B} \leq C|B|^{1+1/n+1/s-1/p}\|u\|_{p,\sigma B} \quad (3.10)$$

for all balls  $B \subset \sigma B \subset \Omega$  with some constant  $\sigma > 1$ .

Now, with these facts in mind, let us first prove Theorem 3.5.

**Theorem 3.5** *Let  $u \in L^p(\Omega, \Lambda)$ ,  $1 < p < n$ , be a differential form defined on the smoothly bounded domain  $\Omega$ ,  $D$  be the Hodge-Dirac operator and  $G$  be Green's operator. Then, for any positive integer  $k > 1$  and any real number  $0 < s < np/(n-p)$ , there exist two constants  $C_1, C_2 > 0$ , independent of  $u$ , such that*

$$\|D^k G^k(u)\|_{*,s,\Omega} \leq C_1 \|D^k G^k(u)\|_{locLip_\alpha^s(\Omega)} \leq C_2 \|u\|_{p,\Omega}, \quad (3.11)$$

$$\|D^{k+1} G^k(u)\|_{*,s,\Omega} \leq C_1 \|D^{k+1} G^k(u)\|_{locLip_\alpha^s(\Omega)} \leq C_2 \|u\|_{p,\Omega}, \quad (3.12)$$

where  $0 < \alpha \leq 1$  is some constant.

**Proof** First, we notice that  $1 + \frac{1}{n} - \frac{1}{p} - \frac{\alpha}{n} = \left(1 - \frac{1}{p}\right) + \left(\frac{1}{n} - \frac{\alpha}{n}\right) > 0$  because  $0 < \alpha \leq 1$  and  $1 < p < \infty$ . Then, for any ball  $B \subset \Omega$ , we have

$$|B|^{1+1/n-1/p-\alpha/n} \leq |\Omega|^{1+1/n-1/p-\alpha/n}. \quad (3.13)$$

In the meantime, by replacing  $\omega$  with  $D^k G^k(u)$  and  $D^{k+1} G^k(u)$  in (1.5), respectively, it is immediate to achieve that

$$\|D^k G^k(u)\|_{*,s,\Omega} \leq C_1 \|D^k G^k(u)\|_{locLip_\alpha^s(\Omega)}, \quad (3.14)$$

$$\|D^{k+1} G^k(u)\|_{*,s,\Omega} \leq C_1 \|D^{k+1} G^k(u)\|_{locLip_\alpha^s(\Omega)}. \quad (3.15)$$

Thus, to estimate (3.11), applying (3.7) and (3.13) gives

$$\begin{aligned} \|D^k G^k(u)\|_{locLip_\alpha^s(\Omega)} &= \sup_{\sigma_2 B \subset \Omega} |B|^{-\frac{n+\alpha s}{3n}} \|D^k G^k(u) - (D^k G^k(u))_B\|_{s,B} \\ &\leq \sup_{\sigma_2 B \subset \Omega} |B|^{-1/s-\alpha/n} C_2 |B|^{1+1/s+1/n-1/p} \|u\|_{p,\sigma_1 B} \\ &= \sup_{\sigma_2 B \subset \Omega} C_2 |B|^{1+1/n-1/p-\alpha/n} \|u\|_{p,\sigma_1 B} \\ &\leq \sup_{\sigma_2 B \subset \Omega} C_2 |\Omega|^{1+1/n-1/p-\alpha/n} \|u\|_{p,\sigma_1 B} \\ &\leq C_3 \sup_{\sigma_2 B \subset \Omega} \|u\|_{p,\sigma_1 B} \\ &\leq C_4 \|u\|_{p,\Omega}, \end{aligned} \quad (3.16)$$

where the constants  $\sigma_2 > \sigma_1 > 1$  and all balls  $B \subset \sigma_1 B \subset \sigma_2 B \subset \Omega$ . So, according to (3.14) and (3.16), we have that (3.11) holds as desired. Moreover, using the same treatment to the operator  $D^{k+1} G^k(u)$  with (3.8) and (3.15), the inequality (3.12) holds as well. Therefore, the proof of Theorem 3.5 is finished.  $\square$

For the case  $p \geq n$ , repeating the process as in Theorem 3.5 with (3.9) and (3.10), we obtain the analogue results.

**Theorem 3.6** *Let  $u \in L^p(\Omega, \Lambda)$ ,  $p \geq n$ , be a differential form defined on the smoothly bounded domain  $\Omega$ ,  $D$  be the Hodge-Dirac operator and  $G$  be Green's operator. Then, for any positive integer  $k > 1$  and any real number  $s > 0$ , there exist two constants  $C_1, C_2 > 0$ , independent of  $u$ , such that*

$$\|D^k G^k(u)\|_{*,s,\Omega} \leq C_1 \|D^k G^k(u)\|_{locLip_\alpha^s(\Omega)} \leq C_2 \|u\|_{p,\Omega}, \quad (3.17)$$

$$\|D^{k+1} G^k(u)\|_{*,s,\Omega} \leq C_1 \|D^{k+1} G^k(u)\|_{locLip_\alpha^s(\Omega)} \leq C_2 \|u\|_{p,\Omega}, \quad (3.18)$$

where  $0 < \alpha \leq 1$  is some constant.

Next, we begin to establish our principle relationship between  $BMO_\alpha^s$  norm and  $locLip_\alpha^s$  norm of the iterated operators in terms of the norms of the conjugate harmonic tensors  $u$  and  $v$ . From Theorem 3.7 and Theorem 3.8 to Corollary 3.9 and Corollary 3.10 below, we always assume that  $\Omega \subset \mathbb{R}^n$  is smoothly bounded domain without boundary, the operator  $A$  in (3.1) is invertible.

**Theorem 3.7** *Let  $u \in W^{1,p}(\Omega, \Lambda)$  and  $v \in W^{1,q}(\Omega, \Lambda)$ ,  $1 < p, q < \infty$  with  $1/p + 1/q = 1$ , be the conjugate  $A$ -harmonic tensors satisfying the Equation (3.1),  $D$  be the Hodge-Dirac operator and  $G$  be Green's operator. Then, for every integer  $k = 2m$  and any real number  $s > 0$ ,  $m = 1, 2, \dots$ , there are two constants  $C_1, C_2 > 0$ , independent of  $u$  and  $v$ , such that*

$$\|D^k G^k(u)\|_{*,s,\Omega} \leq C_1 \|D^k G^k(u)\|_{locLip_\alpha^s(\Omega)} \leq C_2 \|u\|_{*,p,\Omega}, \quad (3.19)$$

$$\|D^{k+1} G^k(v)\|_{*,s,\Omega} \leq C_1 \|D^{k+1} G^k(v)\|_{locLip_\beta^s(\Omega)} \leq C_2 \|\star v\|_{*,p,\Omega}, \quad (3.20)$$

where  $0 < \alpha, \beta \leq 1$  are the expansion factors.

**Proof** First, without loss of generality, we assume that the conjugate  $A$ -harmonic tensor  $u$  is a solution to the  $A$ -harmonic equation (3.2). Then, it is natural to view the corresponding  $v$  as a solution to Equation (3.3). Next, we will divide our proof into two parts.

(i) For every  $1 < p < \infty$ , applying (2.17) into Definition 1.2, we have that

$$\begin{aligned} \|D^k G^k(u)\|_{locLip_\alpha^s(\Omega)} &= \sup_{\sigma_1 B \subset \Omega} |B|^{-1/s-\alpha/n} \|D^k G^k(u) - (D^k G^k(u))_B\|_{s,B} \\ &\leq \sup_{\sigma_1 B \subset \Omega} C_1 |B|^{-1/s-\alpha/n} |B|^{1+1/n} \|du\|_{s,\sigma_2 B} \\ &\leq C_1 \sup_{\sigma_1 B \subset \Omega} |B|^{1+1/n-1/s-\alpha/n} \|du\|_{s,\sigma_2 B} \end{aligned} \quad (3.21)$$

Observe that  $du$  is a solution for the  $A$ -harmonic equation since  $du$  is a closed form. Then, for any real number  $s > 0$ , using Lemma 3.3 yields that

$$\|du\|_{s,\sigma_2 B} \leq C_2 |B|^{1/s-1/p} \|du\|_{p,\sigma_3 B}. \quad (3.22)$$

where  $\sigma_3 > \sigma_2 > 1$ . Under the assumption, we know that  $u$  satisfies the Caccippoli inequality (3.4). Especially, choosing  $c = u_B$  in (3.4) follows

$$\|du\|_{p,\sigma_3 B} \leq C_3 |B|^{-1/n} \|u - u_B\|_{p,\sigma_4 B} \quad (3.23)$$

for some constant  $\sigma_4 > \sigma_3 > 1$  with any ball  $\sigma_3 B \subset \sigma_4 B \subset \Omega$ . Moreover, combining (3.22) and (3.23) gives

$$\|du\|_{s,\sigma_2 B} \leq C_4 |B|^{1/s-1/n-1/p} \|u - u_B\|_{p,\sigma_4 B} \quad (3.24)$$

So, substituting (3.24) into (3.21), together with Definition 1.1, yields that

$$\begin{aligned} \|D^k G^k(u)\|_{locLip_\alpha^s(\Omega)} &\leq C_5 \sup_{\sigma_1 B \subset \Omega} |B|^{1+1/n-1/s-\alpha/n} |B|^{1/s-1/n-1/p} \|u - u_B\|_{p,\sigma_4 B} \\ &\leq C_5 \sup_{\sigma_1 B \subset \Omega} |B|^{1-1/p-\alpha/n} \|u - u_B\|_{p,\sigma_4 B} \\ &\leq C_5 \sup_{\sigma_1 B \subset \Omega} |\Omega|^{1-\alpha/n} |B|^{-1/p} \|u - u_B\|_{p,\sigma_4 B} \\ &\leq C_6 \sup_{\sigma_1 B \subset \Omega} |B|^{-1/p} \|u - u_B\|_{p,\sigma_4 B} \\ &\leq C_6 \|u\|_{*,p,\Omega}, \end{aligned} \quad (3.25)$$

where the constants  $\sigma_1 > \sigma_4 > 1$ . Therefore, we have that (3.19) holds for any even integer  $k > 1$  and any real number  $s > 0$ .

The proof of (3.20) is similar to that of (3.19). Next, we only present the different steps.

- (ii) For every conjugate  $A$ -harmonic tensor  $v \in W^{1,q}(\Omega, \Lambda)$ , employing the same treatment used in the proof of (3.19), along with (2.18), we have that

$$\|D^k G^k(v)\|_{locLip_\beta^s(\Omega)} \leq C_7 \sup_{\eta_1 B \subset \Omega} |B|^{1+1/n-1/s-\beta/n} \|d^*v\|_{s,\eta_2 B}. \quad (3.26)$$

According to the isometric property of the Hodge-star operator  $\star$ , we know that  $|d^*v| = |d\star v|$ . Notice that  $d\star v$  is a closed form satisfying  $A$ -harmonic equation. So, for any real number  $s > 0$ , using Lemma 3.3 again, we derive that

$$\|d^*v\|_{s,\eta_2 B} = \|d\star v\|_{s,\eta_2 B} \leq C_8 |B|^{1/s-1/q} \|d\star v\|_{q,\eta_3 B} \quad (3.27)$$



Also, by the comments after Definition 3.1, it implies that  $\star v$  is a solution to the  $A$ -harmonic equation (3.3). Then, by Lemma 3.2, letting  $c^\star = (\star v)_B$  shows that

$$\|d \star v\|_{q, \eta_3 B} \leq C_9 |\eta_3 B|^{-1/n} \|\star v - (\star v)_B\|_{q, \eta_4 B}. \quad (3.28)$$

So, combining (3.27) with (3.28) and plugging it into (3.26), we have that

$$\begin{aligned} \|D^{k+1} G^k(v)\|_{loc Lip_\beta^s(\Omega)} &\leq \sup_{\eta_1 B \subset \Omega} |B|^{1+1/n-1/s-\beta/n} \|d^\star v\|_{s, \eta_2 B} \\ &\leq \sup_{\eta_1 B \subset \Omega} C_{10} |B|^{1+1/n-1/s-\beta/n} |B|^{1/s-1/q-1/n} \\ &\quad \|\star v - (\star v)_B\|_{q, \eta_4 B} \\ &\leq \sup_{\eta_1 B \subset \Omega} C_{10} |\Omega|^{1-\beta/n} |B|^{-1/q} \|\star v - (\star v)_B\|_{q, \eta_4 B} \\ &\leq C_{11} \sup_{\eta_1 B \subset \Omega} |B|^{-1/q} \|\star v - (\star v)_B\|_{q, \eta_4 B} \\ &= C_{11} \|\star v\|_{*, q, \Omega} \end{aligned} \quad (3.29)$$

as desired, where the constants  $\eta_1 > \eta_4 > \eta_3 > \eta_2 > 1$ .

□

Now, in the odd case  $k = 2m + 1$ , we have the similar estimates as follows. It should be pointed out that the proof of Theorem 3.8 is the analogue of Theorem 3.7, so we only state the results and leave the proof of the odd case  $k > 1$  to the readers.

**Theorem 3.8** *Let  $u \in W^{1,p}(\Omega, \Lambda)$  and  $v \in W^{1,q}(\Omega, \Lambda)$ ,  $1 < p, q < \infty$  with  $1/p + 1/q = 1$ , be the conjugate  $A$ -harmonic tensors satisfying the Equation (3.1),  $D$  be the Hodge-Dirac operator and  $G$  be Green's operator. Then, for every odd integer  $k = 2m + 1$  and any real number  $s > 0$ ,  $m = 1, 2, \dots$ , there are two constants  $C_1, C_2 > 0$ , independent of  $u$  and  $v$ , such that*

$$\|D^k G^k(v)\|_{*, s, \Omega} \leq C_1 \|D^k G^k(v)\|_{loc Lip_\beta^s(\Omega)} \leq C_2 \|\star v\|_{*, q, \Omega}, \quad (3.30)$$

$$\|D^{k+1} G^k(u)\|_{*, s, \Omega} \leq C_1 \|D^{k+1} G^k(u)\|_{loc Lip_\alpha^s(\Omega)} \leq C_2 \|u\|_{*, p, \Omega}, \quad (3.31)$$

where  $0 < \alpha, \beta \leq 1$  are the expansion factors.

In particular, if  $p(\alpha - 1) = q(\eta - 1)$ , as a consequence of Theorem 3.7 and 3.8, the following estimates are established simply by means of Theorem 6.6 in [10]. It is worth to notice that the treatment applied in Corollary 3.9 and Corollary 3.10 are very similar, so we only give the complete proof of Corollary 3.9 in details.

**Corollary 3.9** *Let  $u \in W^{1,p}(\Omega, \Lambda)$  and  $v \in W^{1,q}(\Omega, \Lambda)$ ,  $1 < p, q < \infty$  with  $1/p + 1/q = 1$ , be the conjugate  $A$ -harmonic tensors satisfying the Equation (3.1),*

$D$  be the Dirac operator and  $G$  be the Green's operator. If  $0 < \alpha, \beta \leq 1$  satisfy  $p(\alpha - 1) = q(\beta - 1)$ , for any real  $s > 0$ , then there exist two constants  $C_1, C_2 > 0$ , independent of  $u$  and  $v$ , such that

$$\|D^k G^k(u)\|_{*,s,\Omega} \leq C_1 \|D^k G^k(u)\|_{locLip_\alpha^s(\Omega)} \leq C_2 \|\star v\|_{locLip_\beta^q(\Omega)}^{q/p}, \quad (3.32)$$

$$\|D^{k+1} G^k(v)\|_{*,s,\Omega} \leq C_1 \|D^{k+1} G^k(v)\|_{locLip_\beta^s(\Omega)} \leq C_2 \|u\|_{locLip_\alpha^p(\Omega)}^{p/q}, \quad (3.33)$$

whenever  $k = 2m, m = 1, 2, \dots$ .

**Proof** First, combining (1.5) and Theorem 6.6 in [10], we have

$$\|u\|_{*,s,\Omega} \leq C_1 \|u\|_{locLip_\alpha^s(\Omega)} \leq C_2 \|\star v\|_{locLip_\beta^q(\Omega)}^{q/p}, \quad (3.34)$$

$$\|\star v\|_{BMO,\Omega} \leq C_3 \|\star v\|_{locLip_\beta(\Omega)} \leq C_4 \|u\|_{locLip_\alpha^p(\Omega)}^{p/q}. \quad (3.35)$$

Then, substituting (3.34) into (3.19) and (3.35) into (3.20), respectively, it yields that

$$\|D^k G^k(u)\|_{*,s,\Omega} \leq C_5 \|D^k G^k(u)\|_{locLip_\alpha^s(\Omega)} \leq C_6 \|\star v\|_{locLip_\beta^q(\Omega)}^{q/p},$$

$$\|D^{k+1} G^k(v)\|_{*,s,\Omega} \leq C_7 \|D^{k+1} G^k(v)\|_{locLip_\beta^s(\Omega)} \leq C_8 \|u\|_{locLip_\alpha^p(\Omega)}^{p/q}$$

as desired.  $\square$

**Corollary 3.10** Suppose that  $0 < \alpha, \beta \leq 1$  satisfy  $p(\alpha - 1) = q(\beta - 1)$ , for any real number  $s > 0$ , then there exist two constants  $C_1, C_2 > 0$ , independent of  $u$  and  $v$ , such that

$$\|D^k G^k(v)\|_{*,s,\Omega} \leq C_1 \|D^k G^k(v)\|_{locLip_\beta^s(\Omega)} \leq C_2 \|u\|_{locLip_\alpha^p(\Omega)}^{p/q}, \quad (3.36)$$

$$\|D^{k+1} G^k(u)\|_{*,s,\Omega} \leq C_1 \|D^{k+1} G^k(u)\|_{locLip_\alpha^s(\Omega)} \leq C_2 \|\star v\|_{locLip_\beta^q(\Omega)}^{q/p}, \quad (3.37)$$

whenever  $k = 2m + 1, m = 1, 2, \dots$ .

What is more, for each pair of conjugate  $A$ -harmonic tensors  $u$  and  $v$ , in accord to the facts that  $|du|^p \leq |d^\star v|^q \leq a_1^q |du|^p$  and  $|d^\star v| = |d \star v|$ , one may easily establish such a useful  $L^p$ -equivalence with respect to  $u$  and  $v$  as follows:

$$\|du\|_{p,\Omega'} \leq \|d \star v\|_{q,\Omega'}^{q/p} \leq a_1^{q/p} \|du\|_{p,\Omega'}, \quad (3.38)$$

whenever  $\Omega' \subset \Omega$ , where  $1 < p, q < \infty$  are the conjugate Hölder exponents. In view of the equivalence (3.38), if  $u$  and  $v$  are the conjugate  $A$ -harmonic tensors,

it further reveals the relations (3.39)–(3.42) below. Namely, when  $k$  is any positive even integer, there exist two constants  $C_1, C_2 > 0$ , independent of  $u$  and  $v$ , such that

$$\|D^k G^k(u)\|_{*,s,\Omega} \leq C_1 \|D^k G^k(u)\|_{locLip_\alpha^s(\Omega)} \leq C_2 \|d \star v\|_{q,\Omega}^{q/p}, \tag{3.39}$$

$$\|D^{k+1} G^k(v)\|_{*,s,\Omega} \leq C_1 \|D^{k+1} G^k(v)\|_{locLip_\beta^s(\Omega)} \leq C_2 \|du\|_{p,\Omega}^{p/q}. \tag{3.40}$$

for any real number  $s > 0$ . As such, when  $k > 1$  is any odd integer, there also exist two constants  $C_1, C_2 > 0$ , independent of  $u$  and  $v$ , such that

$$\|D^k G^k(v)\|_{*,s,\Omega} \leq C_1 \|D^k G^k(v)\|_{locLip_\beta^s(\Omega)} \leq C_2 \|du\|_{p,\Omega}^{p/q}, \tag{3.41}$$

$$\|D^{k+1} G^k(u)\|_{*,s,\Omega} \leq C_1 \|D^{k+1} G^k(u)\|_{locLip_\alpha^s(\Omega)} \leq C_2 \|d \star v\|_{q,\Omega}^{q/p}. \tag{3.42}$$

It should be pointed out that the proof of the above assertions are parallel to the those of Theorem 3.7. Therefore, we omit the details.

### 4 Applications

In this section, we use some concrete examples to illustrate the applications of the main results obtained in Section 3.

Let the mapping  $f : \Omega \rightarrow \mathbb{R}^n, f = (f^1, \dots, f^n)$ , be of Sobolev class  $W_{loc}^{1,p}(\Omega, \Lambda)$  and  $J(x, f) = \det(Df(x))$  be the Jacobian determinant of  $f$ . Then, we have that

$$u = J(x_{i_1}, x_{i_2}, \dots, x_{i_l}; f^{j_1}, f^{j_2}, \dots, f^{j_l}) dx_{i_1} \wedge dx_{i_2} \wedge \dots \wedge dx_{i_l}, \tag{4.1}$$

is a differential  $l$ -form, where  $J(x_{i_1}, x_{i_2}, \dots, x_{i_l}; f^{j_1}, f^{j_2}, \dots, f^{j_l})$  is the subdeterminant of  $J(x, f)$  of the form:

$$J(x_{i_1}, x_{i_2}, \dots, x_{i_l}; f^{j_1}, f^{j_2}, \dots, f^{j_l}) = \begin{vmatrix} f_{x_{i_1}}^{j_1} & f_{x_{i_2}}^{j_1} & \dots & f_{x_{i_l}}^{j_1} \\ f_{x_{i_1}}^{j_2} & f_{x_{i_2}}^{j_2} & \dots & f_{x_{i_l}}^{j_2} \\ \dots & \dots & \dots & \dots \\ f_{x_{i_1}}^{j_l} & f_{x_{i_2}}^{j_l} & \dots & f_{x_{i_l}}^{j_l} \end{vmatrix}$$

Referring to Chapter 1 in [2], we find that Theorem 3.5 and Theorem 3.6 are applicable to such sort of the differential form  $u$ . Here, take a special case of 2-dimensional Euclidean space for instance.

*Example 4.1* Assume that  $u = J(x, y; f^1, f^2)dx \wedge dy$  is the differential 2-form defined on the domain  $\Omega = \{(x, y) \in \mathbb{R}^2 : 0 < x^2 + y^2 < r^2\}$ , where the mapping  $f : \Omega \rightarrow \mathbb{R}^2$  is of the Sobolev class  $W_{loc}^{1,p}(\Omega, \Lambda)$  denoted by

$$f(x, y) = (f^1(x, y), f^2(x, y)) = \left( \frac{x}{(x^2 + y^2)^{1/8}}, \frac{y}{(x^2 + y^2)^{1/8}} \right) \quad (4.2)$$

for any  $r > 0$  and  $p > 1$ . After a simple calculation, one may derive that

$$u = J(x, y; f^1, f^2)dx \wedge dy = \frac{3}{4}(x^2 + y^2)^{-1/4}dx \wedge dy.$$

Thus, by the spherical coordinate transformation, it is easy to see that  $u \in L^p(\Omega, \Lambda^2)$  for any  $p < 4$ . For example, choosing  $p = 3/2$ , we know that  $u \in L^{3/2}(\Omega, \Lambda^2)$ . However, by the direct integral calculation with Definition 1.1 and Definition 1.2, it is quite hard to infer the higher order boundedness of  $\text{BMO}^s$  norm and  $\text{locLip}_\alpha^s$  norm with respect to  $D^k G^k(u)$  and  $D^{k+1} G^k(u)$ . Then, applying Theorem 3.5 to  $D^k G^k$  and  $D^{k+1} G^k$ , for any  $0 < s < np/(n-p) = \frac{2 \cdot 3/2}{2-3/2} = 6$ , we have that  $D^k G^k(u) \in \text{BMO}^s(\Omega, \Lambda^2)$  and  $D^{k+1} G^k(u) \in \text{BMO}^s(\Omega, \Lambda^2)$ . Moreover, there are two constants  $C_1, C_2 > 0$ , independent of  $u$ , such that

$$\|D^k G^k(u)\|_{*,s,\Omega} \leq C_1 \|D^k G^k(u)\|_{\text{locLip}_\alpha^s(\Omega)} \leq C_2 r^{5/4}, \quad (4.3)$$

$$\|D^{k+1} G^k(u)\|_{*,s,\Omega} \leq C_1 \|D^{k+1} G^k(u)\|_{\text{locLip}_\alpha^s(\Omega)} \leq C_2 r^{5/4} \quad (4.4)$$

for every  $0 < \alpha \leq 1$  and all positive integer  $k \geq 1$ .

Especially, if the homeomorphism  $f : \Omega \rightarrow \mathbb{R}^n$  of Sobolev class  $W_{loc}^{1,n}(\Omega, \mathbb{R}^n)$ , as mentioned above, is the  $K$ -quasiregular mapping,  $K \geq 1$ . From [23], we know that

$$u = f^l df^1 \wedge \dots \wedge df^{l-1} \quad \text{and} \quad v = \star f^{l+1} df^{l+2} \wedge \dots \wedge df^n$$

are the conjugate  $A$ -harmonic tensors, whenever  $l = 1, 2, \dots, n-1$ . Here, consider the 4-dimensional space as an example.

*Example 4.2* Let  $f = (f^1, f^2, f^3, f^4)$  be the  $K$ -quasiregular mapping defined on the domain  $\Omega = \{(x_1, x_2, x_3, x_4) : |x_i| < a, i = 1, 2, 3, 4\} \subset \mathbb{R}^4$ , and choose the conjugate  $A$ -harmonic tensors as follows:

$$u = f^2 df^1 \quad \text{and} \quad v = \star f^3 df^4.$$

where  $0 < a < \infty$  is some real number. If  $u \in \text{BMO}^p(\Omega, \Lambda)$  and  $\star v \in \text{BMO}^q(\Omega, \Lambda)$ , where  $p$  and  $q$  are conjugate exponents with  $1/p + 1/q = 1$ , by applying Theorem 3.7 and Theorem 3.8, respectively, we have that for any even

integer  $k = 2m$  and any real number  $s > 0$ ,  $m = 1, 2, \dots$ , there are two constants  $C_1, C_2 > 0$ , such that

$$\|D^k G^k(u)\|_{*,s,\Omega} \leq C_1 \|D^k G^k(u)\|_{locLip^s_\alpha(\Omega)} \leq C_2 \|f^2 df^1\|_{*,p,\Omega},$$

$$\|D^{k+1} G^k(v)\|_{*,s,\Omega} \leq C_1 \|D^{k+1} G^k(v)\|_{locLip^s_\beta(\Omega)} \leq C_2 \|f^3 df^4\|_{*,q,\Omega},$$

where  $0 < \alpha, \beta \leq 1$  are two factors. While the integer  $k = 2m + 1$ ,  $m = 1, 2, \dots$ , it holds that

$$\|D^k G^k(v)\|_{*,s,\Omega} \leq C_1 \|D^k G^k(v)\|_{locLip^s_\beta(\Omega)} \leq C_2 \|f^3 df^4\|_{*,q,\Omega},$$

$$\|D^{k+1} G^k(u)\|_{*,s,\Omega} \leq C_1 \|D^{k+1} G^k(u)\|_{locLip^s_\alpha(\Omega)} \leq C_2 \|f^2 df^1\|_{*,p,\Omega},$$

for any real number  $s > 0$ , where  $0 < \alpha, \beta \leq 1$  are two factors.

*Remark 4* In general, all results we establish here provide us an impressive description about the relation between  $BMO^s$  norm and  $locLip^s_\alpha$  norm for the iterated operators. Also, from the results, one may realize that  $locLip^s_\alpha$ -norm estimates for differential forms are fairly essential for the process to derive the  $BMO^s$  estimate with respect to  $D^k G^k$  and  $D^{k+1} G^k$  for differential forms.

**Acknowledgments** The author G. Shi are very grateful to the Chinese Scholarship Council (CSC) for its financial support on her visit in USA, and sincerely appreciates the Department of Mathematics at Seattle University for providing the precious opportunity for her academic visit.

## References

1. G. Lu, Embedding theorems into Lipschitz and BMO spaces and applications to quasilinear subelliptic differential equations. *Publ. Mat.* **40**, 301–329 (1996)
2. R.P. Agarwal, S. Ding, C.A. Nolder, *Inequalities for Differential Forms* (Springer, New York, 2009)
3. X. Li, Y. Wang, Y. Xing, Lipschitz and BMO norm inequalities for the composite operator on differential forms. *J. Inequalities Appl.* (2015). <https://doi.org/10.1186/s13660-015-0896-9>
4. Y. Xing, Y. Wang, BMO and Lipschitz norm estimates for composite operators. *Potential Anal.* **31**, 335–344 (2009)
5. D. Deng, X. Duong, A. Sikora, L. Yan, Comparison of the classical BMO spaces associated with operators and applications. *Rev. Mat. Iberoamericana* **24**, 267–296 (2008)
6. S. Stein, *Harmonic Analysis: Real-variable Methods, Orthogonality and Oscillatory Integrals* (Princeton University Press, New Jersey, 1993)
7. A. Bonami, T. Iwaniec, P. Jones, M. Zinsmeister, On the product of functions in BMO and  $H^1$ . *Annu. Inst. Fourier (Grenoble)* **57**, 1405–1439 (2007)
8. A. Carbonaro, A. McIntosh, A.J. Morris, Local Hardy spaces of differential forms on Riemannian manifolds. *J. Geom. Anal.* **23**, 106–169 (2013)
9. Y. Xing, S. Ding, Norm comparison inequalities for the composite operator. *J. Inequalities Appl.* (2009). <https://doi.org/10.1155/2009/212915>

10. C.A. Nolder, Hardy-Littlewood theorems for  $A$ -harmonic tensors. III. *J. Math.* **43**, 613–632 (1999)
11. G.F.D. Duff, D.C. Spencer, Harmonic tensors on Riemannian manifolds with boundary. *Ann. Math.* **56**, 128–156 (1952)
12. B. Stroffolini, On weakly  $A$ -harmonic tensors. *Stud. Math.* **114**, 289–301 (1995)
13. V. Gol'dshtein, M. Troyanov, Sobolev inequalities for differential forms and  $L_{q,p}$ -cohomology. *J. Geom. Anal.* **16**, 597–631 (2006)
14. S. Ding, Lipschitz and BMO norm inequalities for operators. *Nonlinear Anal. Theory Methods Appl.* **71**, 2350–2357 (2009)
15. C.B. Morrey, *Multiple Integrals in the Calculus of Variations* (Springer, Berlin, 1966)
16. C. Scott,  $L^p$  theory of differential forms on manifolds. *Trans. Am. Math. Soc.* **347**, 2075–2096 (1995)
17. T. Iwaniec, A. Lutoborski, Integral estimates for null Lagrangians. *Arch. Ration. Mech. Anal.* **125**, 25–79
18. S. Ding, G. Shi, Y. Xing, Higher integrability of iterated operators on differential forms. *Nonlinear Anal. Theory Methods Appl.* **145**, 83–96 (2016)
19. S. Ding, G. Shi, Y. Xing, Representation and regularity of higher order Dirac and Green's operators. Preprint
20. R.P. Agarwal, S. Ding, Advances in differential forms and the  $A$ -harmonic equation. *Math. Comput. Model.* **37**, 1393–1426 (2003)
21. S. Ding, Weighted Hardy-Littlewood inequality for  $A$ -harmonic tensors. *Proc. Am. Math. Soc.* **125**, 1727–1735 (1997)
22. S. Ding, B. Liu, Dirac-harmonic equations for differential forms. *Nonlinear Anal. Theory Methods Appl.* **122**, 43–57 (2015)
23. T. Iwaniec, G. Martin, Quasiregular mappings in even dimensions. *Acta Math.* **170**, 29–81 (1933)

# Application of Boundary Perturbations on Medical Monitoring and Imaging Techniques



M. Doschoris, A. Papargiri, V. S. Kalantonis, and P. Vafeas

**Abstract** We present an overview on the application of boundary perturbations for Electroencephalography and Magnetoencephalography predominantly for the spherical geometry. With the mathematical tools produced, both forward and inverse problems can be tackled providing explicit computationally efficient solutions. Utilizing perturbation analysis in the framework of medical monitoring and imaging techniques, possesses the advantage introducing geometric variations without limiting the installation of analytic, or at least semi-analytic solutions, in view of complicated surfaces. In our example, surfaces which do not allow an analytic mathematical treatment can be handled if considered as small deviations from the sphere. In that setting, irregularities in head shapes, e.g. craniofacial alterations can be investigated theoretically.

**MSC** 78A25, 35Q61, 92C05, 35Q92, 45K05, 35R30

---

M. Doschoris

Institute of Genetics and Biometry, Leibniz Institute for Farm Animal Biology, Dummerstorf, Germany

e-mail: [doschoris@fbn-dummerstorf.de](mailto:doschoris@fbn-dummerstorf.de)

A. Papargiri · V. S. Kalantonis (✉)

Department of Electrical and Computer Engineering, University of Patras, Patras, Greece

e-mail: [apapargiri@upatras.gr](mailto:apapargiri@upatras.gr); [kalantonis@upatras.gr](mailto:kalantonis@upatras.gr)

P. Vafeas

Department of Chemical Engineering, University of Patras, Patras, Greece

e-mail: [vafeas@chemeng.upatras.gr](mailto:vafeas@chemeng.upatras.gr)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,

Springer Optimization and Its Applications 173,

[https://doi.org/10.1007/978-3-030-72563-1\\_6](https://doi.org/10.1007/978-3-030-72563-1_6)

## 1 Introduction

Medical monitoring techniques such as Electroencephalography (EEG) and functional neuroimaging techniques such as Magnetoencephalography (MEG) provide the possibility of recognizing brain structures and map cortical activity at its core functional level. Both consist noninvasive diagnostic techniques, playing a crucial role in many aspects of today's clinical research either as the principal investigative tool or providing supportive evidence. As Ilmoniemi and Näätänen assert (for MEG) [1] “[. . .] tool to probe into the real-time operation of the human brain in experimental conditions that are suitable for studying sensory and cognitive brain functions as well as their disturbances.” As an illustration, EEG is practised in reporting schizophrenia [2], in detection and analysis of epileptic activity [3], in the study of language processing [4] and many others. On the other hand, MEG is used tracing the dynamics and connectivity of large-scale brain activity [5].

The core question of EEG and MEG is to gain insight of the working human brain by recognizing the precise position and strength of the underlying neuronal activity. In the case of EEG, measurements are based on the post-synaptic potentials, which originate in the pyramidal neurons recorded on the scalp [6], whereas MEG measures oscillatory magnetic fields [7]. The identification of which brain areas have been activated, based on either EEG or MEG measurements, is termed the inverse problem or source reconstruction, while the obverse is labeled the forward problem [8]. Historically, the forward and inverse EEG problems have been extensively scrutinized since the 1950s when Wilson and Bayley [9] attempted to quantify the interplay between neuronal activity and the potentials, generated at the scalp [10, 11]. Similar hold for MEG, since David Cohen first measured magnetic fields produced by Humans in the late 1960s [12].

From the mathematical point of view, the formulation for the Electroencephalography and Magnetoencephalography problems are derived from the quasi-static theory of electromagnetism [13, 14]. We recall that microscopic currents are largely approximated by equivalent dipole sources and that an exclusive source configuration for each measurement of either EEG or MEG does not exist, comprising the corresponding inverse problems ill-posed. Regarding non-uniqueness, Albanese and Monk [15] demonstrated that the reconstruction of a three-dimensional current based on EEG measurements is unrealizable and has been practically demonstrated by Dassios and Fokas [16]. Nevertheless, if the source representing neuronal activity has dimensionality less than three, the inverse problems at hand admit a unique solution. The classic example is the dipole approximation of zero dimensionality. The same is true stipulating that neuronal activity is characterized via continuously distributed currents in one [17] and two dimensions [18]. Evidently, discarding the dipole hypothesis, allows the investigation and analysis of complicated activation patterns in terms of distributed currents.

Another interesting aspect regards the complementarity of EEG and MEG, a consequence due to sensitivity of EEG and MEG to radial and tangential sources. The notion of complementarity has been mathematically defined spherical model



of the brain and a continuously distributed neuronal current, demonstrating that the information which is missing in EEG data is contained in MEG data and vice versa [19]. Importantly, it has been proved in [20] that, independent of the geometry, the inverse EEG problem is not capable of reconstructing more than 1/3 of any neuronal activity, whereas the inverse MEG problem cannot recover the 2/3 of any brain activity, including the 1/3 of EEG.

Scope of the paper is to provide an analytic overview on the application of boundary perturbations for EEG and MEG predominantly for the spherical geometry, allowing a systematic presentation of the mathematical complexity associated. Extensive use is made of special functions and especially spherical harmonics. Details and technicalities regarding manipulation of formulas involving spherical harmonics can be found in [21, 22]. Our work is organized as follows: In the next subsection, we provide a brief discussion on boundary perturbations. In Section 2, the mathematical formulation of medical monitoring techniques are presented and the influence of geometric variations on the forward as well as inverse problem to EEG and MEG are studied. Finally, in Section 3, two specific examples are provided for the obtained solutions of both EEG and MEG, while in Section 4 we discuss our results and conclude.

### 1.1 *Boundary Perturbations. A Short Introduction*

Consider the action of the operator  $\mathcal{L}(\boldsymbol{\tau})$  on a function  $w(\boldsymbol{\tau})$  within a suitable domain  $\Omega$ , such that

$$(\mathcal{L}w)(\boldsymbol{\tau}) = \phi(\boldsymbol{\tau}), \quad \boldsymbol{\tau} \in \Omega, \quad (1)$$

where  $\phi(\boldsymbol{\tau})$  is the outcome of the operation. In the setting of Boundary Value Problems (BVPs) the latter must be accompanied by appropriate boundary conditions, namely

$$(\mathcal{B}w)(\boldsymbol{\tau}) = \psi(\boldsymbol{\tau}), \quad \boldsymbol{\tau} \in \partial\Omega. \quad (2)$$

In above setting everything is known except  $w(\boldsymbol{\tau})$ . For a variety of reasons, closed form solutions to BVP (1) and (2) in its present form are not feasible. The operator  $\mathcal{L}(\boldsymbol{\tau})$  may be linear but of high dimensions or even worst, non-linear. On the other hand, for complex boundaries even simple linear operators do not possess eigenfunctions.

We approach BVP (1) and (2) by substituting  $\Omega$  and its boundary by an auxiliary domain  $D$  and boundary  $\partial D$  such that (1) and (2) can be solved in terms of the corresponding eigenfunctions of the operator in question, which we consider to be linear for simplicity. In order to proceed, a connection between boundaries has to be established. To this end, any point  $\boldsymbol{\tau}$  lying on  $\partial\Omega$  can be written as the sum of the position vector  $\mathbf{r}$  plus their difference, measuring the deviation between the so-

called perturbed boundary  $\partial\Omega$  and the corresponding unperturbed  $\partial D$ . Denoting by  $f(\hat{\mathbf{r}})$  the aforementioned difference and introducing a scaling factor  $\epsilon$ , we have that

$$\hat{\boldsymbol{\tau}} = \hat{\mathbf{r}} + \epsilon f(\hat{\mathbf{r}}). \quad (3)$$

Note that for  $\epsilon$  to be dimensionless, unit vectors have to be introduced. The initial BVP now reads as

$$(\mathcal{L}w)(\mathbf{r}, \epsilon) = \phi(\mathbf{r}, \epsilon), \quad \mathbf{r} \in D, \quad (4)$$

$$(\mathcal{B}w)(\mathbf{r}, \epsilon) = \psi(\mathbf{r}, \epsilon), \quad \mathbf{r} \in \partial D, \quad (5)$$

from which it is clear that any perturbation of the boundary is transferred onto terms involving powers of  $\epsilon$ . In general, perturbations affect the functions  $w(\boldsymbol{\tau})$ ,  $\phi(\boldsymbol{\tau})$ ,  $\psi(\boldsymbol{\tau})$  involved, as well as the operators  $\mathcal{L}$  and  $\mathcal{B}$ . However, if the operator  $\mathcal{L}$  describes the physical properties of the problem it remains unaltered. The sought perturbation terms are evaluated by expanding involved quantities in powers of  $\epsilon$ , given that the perturbation parameter  $\epsilon$  is considered small and the coefficients of the expansion are independent of  $\epsilon$ . These coefficients, on the other hand, are evaluated solving specific expressions, produced after substituting the expansions into (1), (2) and collecting coefficients of  $\epsilon$ .

For example, replacing the Poincaré type expansion  $h(\mathbf{r}, \epsilon) = \sum_{n=0}^{\infty} \epsilon^n h_n(\mathbf{r})$ , where  $h$  represents any of the functions of interest, into (1) and (2), respectively, gives

$$\sum_{n=0}^{\infty} \epsilon^n \mathcal{L}w_n(\mathbf{r}) = \sum_{n=0}^{\infty} \epsilon^n \phi_n(\mathbf{r}), \quad \mathbf{r} \in D \quad (6)$$

and

$$\sum_{n=0}^{\infty} \epsilon^n \sum_{i=0}^n (\mathcal{B}_n w_{n-i})(\mathbf{r}) = \sum_{n=0}^{\infty} \epsilon^n \psi_n(\mathbf{r}), \quad \mathbf{r} \in \partial D, \quad (7)$$

respectively, the latter indicating the Cauchy product of  $(\sum_{n=0}^{\infty} \epsilon^n \mathcal{B}_n(\mathbf{r})) (\sum_{n=0}^{\infty} \epsilon^n w_n(\mathbf{r}))$ .

Collecting now coefficients of  $\epsilon$ , furnishes the BVPs which have to be solved, e.g.

$$(\mathcal{L}w_0)(\mathbf{r}) = \phi_0(\mathbf{r}), \quad \mathbf{r} \in D, \quad (\mathcal{B}_0 w_0)(\mathbf{r}) = \psi_0(\mathbf{r}), \quad \mathbf{r} \in \partial D, \quad (8)$$

$$(\mathcal{L}w_1)(\mathbf{r}) = \phi_1(\mathbf{r}), \quad \mathbf{r} \in D, \quad (\mathcal{B}_0 w_1 + \mathcal{B}_1 w_0)(\mathbf{r}) = \psi_1(\mathbf{r}), \quad \mathbf{r} \in \partial D, \quad (9)$$

$$(\mathcal{L}w_2)(\mathbf{r}) = \phi_2(\mathbf{r}), \quad \mathbf{r} \in D, \quad (\mathcal{B}_0 w_2 + \mathcal{B}_1 w_1 + \mathcal{B}_2 w_0)(\mathbf{r}) = \psi_2(\mathbf{r}), \quad \mathbf{r} \in \partial D, \quad (10)$$

...

...

providing, the unperturbed instance ( $n = 0$ ) (8), the first (9) and second (10) correction and so on. It is evident that for each correction, the solution to the previous one must be obtained in order to proceed. Further details can be found in [23, 24].

## 2 Mathematical Formulation of Medical Monitoring Techniques. Electroencephalography and Magnetoencephalography

Consider a homogeneous conductor modeling the brain denoted by  $\Omega \subset \mathbb{R}^3$  and let  $\partial\Omega$  be its boundary. Let the exterior, not conductive, to the brain region denoted by  $\Omega^c$ . The activity of the brain, represented by a neuronal current  $\mathbf{J}^p$ , can be recorded either as electrical impulses along the surface (EEG), or measuring the magnetic fields generated by  $\mathbf{J}^p$ , the neuronal current a few centimeters above the surface (MEG). The basic assumption is that neuronal current  $\mathbf{J}^p$  is represented either as a discrete or as a continuous distribution of dipoles with specified moments  $\mathbf{Q}$ .

In 1967, Plonsey and Heppner [13] illustrated that the electromagnetic activity of the brain is governed by the quasi-static theory of Maxwell's equations, namely

$$\nabla \times \mathbf{E} = 0, \tag{11}$$

$$\nabla \times \mathbf{B} = \mu_0(\mathbf{J}^p + \sigma \mathbf{E}), \tag{12}$$

$$\nabla \cdot \mathbf{B} = 0, \tag{13}$$

where  $\mathbf{E}$  and  $\mathbf{B}$  denote the electric and magnetic field, respectively, corresponding to the neuronal current  $\mathbf{J}^p$ ,  $\sigma$  is the conductivity of the medium occupying  $\Omega$  and  $\mu_0$  is the magnetic permeability, assumed constant everywhere in  $\mathbb{R}^3$  and equal to the magnetic permeability of the free space. After some straightforward calculations (for details see [25]) one can show that, if the medium is homogeneous, the electric potential  $U$ , introduced as  $\mathbf{E} = -\nabla U$ , is related to the primary current  $\mathbf{J}^p$  by Poisson's equation

$$\Delta U = \frac{1}{\sigma} \nabla \cdot \mathbf{J}^p. \tag{14}$$

Similar, in the current free, non-conductive domain  $\Omega^c$  the magnetic field, due to (13), can be represented as a gradient of a harmonic function, namely

$$\mathbf{B}(\boldsymbol{\tau}, \mathbf{r}_0) = \frac{\mu_0}{4\pi} \nabla W(\boldsymbol{\tau}, \mathbf{r}_0), \tag{15}$$

where  $W(\boldsymbol{\tau}, \mathbf{r}_0)$  denotes the scalar magnetic potential, vanishing at infinity.

Another way is connected to Geselowitz's integral formula [26], which is geometry independent, but relies on knowledge of the electrical potential  $U$

$$\mathbf{B}(\boldsymbol{\tau}; \mathbf{r}_0) = \frac{\mu_0}{4\pi} \mathbf{Q}(\mathbf{r}_0) \times \frac{\boldsymbol{\tau} - \mathbf{r}_0}{|\boldsymbol{\tau} - \mathbf{r}_0|^3} - \frac{\mu_0\sigma}{4\pi} \int_{\Omega} \nabla_{\mathbf{r}'} U(\mathbf{r}'; \mathbf{r}_0) \times \frac{\boldsymbol{\tau} - \mathbf{r}'}{|\boldsymbol{\tau} - \mathbf{r}'|^3} d\nu(\mathbf{r}'). \quad (16)$$

Since

$$\nabla_{\mathbf{r}_0} \frac{1}{|\boldsymbol{\tau} - \mathbf{r}_0|} = -\nabla_{\boldsymbol{\tau}} \frac{1}{|\boldsymbol{\tau} - \mathbf{r}_0|} = \frac{\boldsymbol{\tau} - \mathbf{r}_0}{|\boldsymbol{\tau} - \mathbf{r}_0|^3}, \quad (17)$$

(16) becomes

$$\mathbf{B}(\boldsymbol{\tau}; \mathbf{r}_0) = \frac{\mu_0}{4\pi} \mathbf{Q}(\mathbf{r}_0) \times \nabla_{\mathbf{r}_0} \frac{1}{|\boldsymbol{\tau} - \mathbf{r}_0|} - \frac{\mu_0\sigma}{4\pi} \int_{\Omega} \nabla_{\mathbf{r}'} U(\mathbf{r}'; \mathbf{r}_0) \times \nabla_{\mathbf{r}'} \frac{1}{|\boldsymbol{\tau} - \mathbf{r}'|} d\nu(\mathbf{r}'). \quad (18)$$

In view of the identity

$$\int_{\Omega} (\nabla f) \times (\nabla g) d\nu = \int_{\Omega} \nabla \times (f \nabla g) d\nu = \oint_{\partial\Omega} \hat{\mathbf{v}} \times f \nabla g dS, \quad (19)$$

the integral on the RHS of equation (18) simplifies as

$$\begin{aligned} \int_{\Omega} \nabla_{\mathbf{r}'} U(\mathbf{r}'; \mathbf{r}_0) \times \nabla_{\mathbf{r}'} \frac{1}{|\boldsymbol{\tau} - \mathbf{r}'|} d\nu(\mathbf{r}') &= \int_{\Omega} \nabla_{\mathbf{r}'} \times \left( U(\mathbf{r}'; \mathbf{r}_0) \nabla_{\mathbf{r}'} \frac{1}{|\boldsymbol{\tau} - \mathbf{r}'|} \right) d\nu(\mathbf{r}') \\ &= \oint_{\partial\Omega} U(\mathbf{r}'; \mathbf{r}_0) \hat{\mathbf{v}} \times \frac{\boldsymbol{\tau} - \mathbf{r}'}{|\boldsymbol{\tau} - \mathbf{r}'|^3} dS(\hat{\mathbf{r}}'). \end{aligned} \quad (20)$$

$$(21)$$

## 2.1 The Influence of Geometric Variations on the Forward Problem

### 2.1.1 EEG

Adopt that the neuronal current is represented by a single equivalent dipole at the point  $\mathbf{r}_0$  with moment  $\mathbf{Q}$ , which is sufficient when modeling smaller cortical neuronal sources [27]. Then, the primary current can be represented as  $\mathbf{J}^p = \mathbf{Q}\delta(\boldsymbol{\tau} - \mathbf{r}_0)$ ,  $\delta(\mathbf{r})$  denoting the Dirac measure. Equation (14) then simplifies as

$$\Delta_{\boldsymbol{\tau}} U(\boldsymbol{\tau}, \mathbf{r}_0) = \frac{1}{\sigma} \mathbf{Q}(\mathbf{r}_0) \cdot \nabla_{\boldsymbol{\tau}} \delta(\boldsymbol{\tau} - \mathbf{r}_0), \quad \boldsymbol{\tau} \in \Omega, \quad (22)$$

accompanied by the condition

$$\hat{\mathbf{v}} \cdot \nabla_{\boldsymbol{\tau}} U(\boldsymbol{\tau}, \mathbf{r}_0) = 0, \quad \boldsymbol{\tau} \in \partial\Omega, \quad (23)$$

where  $\hat{\mathbf{v}}$  denotes the unit vector normal to the boundary. The solution to BVP (22) and (23) is

$$U(\boldsymbol{\tau}; \mathbf{r}_0) = \frac{1}{4\pi\sigma} (\mathbf{Q} \cdot \nabla_{\mathbf{r}_0}) \left( \frac{1}{|\boldsymbol{\tau} - \mathbf{r}_0|} + u(\boldsymbol{\tau}; \mathbf{r}_0) \right), \quad (24)$$

where  $-\frac{1}{4\pi} |\boldsymbol{\tau} - \mathbf{r}_0|^{-1}$  is the fundamental solution and  $u(\boldsymbol{\tau}; \mathbf{r}_0)$  is known as we will see shortly.

Due to the arbitrary shape of the boundary, arriving at analytic solutions for the problems (22) and (23) is in general not feasible. Notwithstanding, introducing an auxiliary boundary  $\partial D$  say, in the shape of a sphere of radius  $R$ , it is possible to transform (22) and (23) to approximate problems with reference to the specific symmetrical case. Following the steps outlined in Section 1.1, we connect the boundaries as

$$\boldsymbol{\tau} = R\hat{\mathbf{r}} + \epsilon f(\hat{\mathbf{r}}) \quad (25)$$

with  $f(\hat{\mathbf{r}})$  explicitly carrying units of [L]. The corresponding electric potentials are then straightforwardly computed by replacing

$$U(\boldsymbol{\tau}) = \sum_{n=0}^{\infty} \epsilon^n U_n(\mathbf{r}), \quad (26)$$

into expression (22). After collecting coefficients of  $\epsilon$  we obtain

$$\sigma \Delta_{\mathbf{r}} U_0(\mathbf{r}, \mathbf{r}_0) = \mathbf{Q} \cdot \nabla_{\mathbf{r}} \delta(\mathbf{r} - \mathbf{r}_0), \quad \mathbf{r} \in D, \quad (27)$$

$$\Delta_{\mathbf{r}} U_n(\mathbf{r}, \mathbf{r}_0) = 0, \quad n \geq 1 \quad \mathbf{r} \in D. \quad (28)$$

In order to identify the implication of the deformed boundary on the solution process, expansion (26) is assumed to be valid all the way to the boundary. The Neumann condition (23) now reads

$$\tau^2 \frac{\partial U(\boldsymbol{\tau})}{\partial \tau} - \frac{\partial \tau}{\partial \theta} \frac{\partial U(\boldsymbol{\tau})}{\partial \theta} - \frac{1}{\sin^2 \theta} \frac{\partial \tau}{\partial \phi} \frac{\partial U(\boldsymbol{\tau})}{\partial \phi} = 0, \quad (29)$$

combined with expressions (25), (26) and bearing in mind that  $\partial_r f = 0$ , gives rise to the system

$$\hat{\mathbf{r}} \cdot \nabla U_n(R\hat{\mathbf{r}}, \mathbf{r}_0) = \sum_{k=0}^{n-1} (-1)^k (k+1) \left(\frac{f}{R}\right)^k \left(\nabla f \cdot \nabla U_{n-k-1}(R\hat{\mathbf{r}}, \mathbf{r}_0)\right), \quad (30)$$

the latter holding for every  $n$ , where by convention  $\sum_{n=0}^{-1} = 0$ .

The electric potential of a deformed spherical conductor due to a single dipole excitation can be presented in terms of spherical harmonics as

$$U(\boldsymbol{\tau}, \mathbf{r}_0) = \sum_{k=1}^{\infty} \sum_{\ell=-k}^k \mathbb{B}_k^{\ell}(\boldsymbol{\tau}, \mathbf{r}_0, \mathbf{Q}) Y_k^{\ell}(\hat{\boldsymbol{\tau}}), \quad \boldsymbol{\tau} \in \Omega. \quad (31)$$

In a similar fashion, each  $n$ th-order correction specified by relation (26) is given by

$$U_n(\mathbf{r}, \mathbf{r}_0) = \sum_{k=1}^{\infty} \sum_{\ell=-k}^k \mathbb{A}_{n,k}^{\ell}(r, \mathbf{r}_0, \mathbf{Q}) Y_k^{\ell}(\hat{\mathbf{r}}), \quad \mathbf{r} \in D. \quad (32)$$

Combining relations (26) and (32) yields

$$\mathbb{B}_k^{\ell}(r, \mathbf{r}_0, \mathbf{Q}; \epsilon) = \sum_{n=0}^{\infty} \epsilon^n \mathbb{A}_{n,k}^{\ell}(r, \mathbf{r}_0, \mathbf{Q}). \quad (33)$$

The electric potential in regard to the spherical conductor  $D \subset \mathbb{R}^3$  is then provided via the zeroth-order correction (32), a solution to (27) together with  $\partial_r U_0(R\hat{\mathbf{r}}, \mathbf{r}_0) = 0$ . It is not hard to show that on the surface of the non-deformed conductor we have

$$\mathbb{A}_{0,k}^{\ell}(R, \mathbf{r}_0, \mathbf{Q}) = \frac{\mathbf{Q}}{k\sigma R^{k+1}} \cdot \nabla_{\mathbf{r}_0} r_0^k \bar{Y}_k^{\ell}(\hat{\mathbf{r}}_0), \quad (34)$$

where  $Y$  represent the spherical harmonics and an over-line denotes complex conjugation and we made use of the fact that in spherical coordinates the function  $U$  in (24) enjoys the expansion

$$U(\mathbf{r}; \mathbf{r}_0) = 4\pi \sum_{n=1}^{\infty} \sum_{m=-n}^n \frac{n+1}{n(2n+1)} \frac{r^n r_0^n}{R^{2n+1}} Y_n^m(\hat{\mathbf{r}}) \bar{Y}_n^m(\hat{\mathbf{r}}_0). \quad (35)$$

Conveniently, in conjecture with the generating function

$$\frac{1}{|\mathbf{r} - \mathbf{r}_0|} = \sum_{n=0}^{\infty} \sum_{m=-n}^n \frac{4\pi}{2n+1} \frac{r_0^n}{r^{n+1}} Y_n^m(\hat{\mathbf{r}}) \bar{Y}_n^m(\hat{\mathbf{r}}_0), \quad (36)$$

it is possible to express the surface potential  $U_0$  in closed form (see [28] for details) as

$$U_0(R\hat{\mathbf{r}}, \mathbf{r}_0) = \frac{\mathbf{Q}}{4\pi\sigma} \cdot \left( 2 \frac{R\hat{\mathbf{r}} - \mathbf{r}_0}{|R\hat{\mathbf{r}} - \mathbf{r}_0|^3} + \frac{1}{R|R\hat{\mathbf{r}} - \mathbf{r}_0|} \frac{|R\hat{\mathbf{r}} - \mathbf{r}_0|\hat{\mathbf{r}} + (R\hat{\mathbf{r}} - \mathbf{r}_0)}{|R\hat{\mathbf{r}} - \mathbf{r}_0| + \hat{\mathbf{r}} \cdot (R\hat{\mathbf{r}} - \mathbf{r}_0)} \right). \quad (37)$$

On the other hand, the first-order correction  $U_1$  is derived as the solution to (28) for  $n = 1$  joint by (30), namely

$$R^2 \frac{\partial U_1(R\hat{\mathbf{r}}, \mathbf{r}_0)}{\partial r} = \frac{\partial f(\hat{\mathbf{r}})}{\partial \theta} \frac{\partial U_0(R\hat{\mathbf{r}}, \mathbf{r}_0)}{\partial \theta} + \frac{1}{\sin^2 \theta} \frac{\partial f(\hat{\mathbf{r}})}{\partial \phi} \frac{\partial U_0(R\hat{\mathbf{r}}, \mathbf{r}_0)}{\partial \phi}. \quad (38)$$

Replacing into the latter the expressions for  $U_0$  and  $U_1$  through (32) and integrating the resulting relation over the unit ball, determines the coefficients  $\mathbb{A}_{1,k}^\ell(R, \mathbf{r}_0, \mathbf{Q})$  as

$$\mathbb{A}_{1,k}^\ell = \frac{1}{kR} \sum_{n=1}^{\infty} \sum_{m=-n}^n \mathbb{A}_{0,n}^m \left\{ \oint_{S^2} \left[ \frac{1}{\sin \theta} \frac{\partial f}{\partial \theta} (n j_{n+1}^m Y_{n+1}^m(\hat{\mathbf{r}}) - (n+1) j_n^m Y_{n-1}^m(\hat{\mathbf{r}})) \right. \right. \\ \left. \left. + i \frac{m}{\sin^2 \theta} \frac{\partial f}{\partial \phi} Y_n^m(\hat{\mathbf{r}}) \right] Y_k^\ell(\hat{\mathbf{r}}) dS(\hat{\mathbf{r}}) \right\}, \quad (39)$$

where the coefficients  $\mathbb{A}_{0,n}^m$  are given by (34) and we further employed the recurrence relations

$$\sin \theta \frac{\partial}{\partial \theta} Y_n^m(\hat{\mathbf{r}}) = n j_{n+1}^m Y_{n+1}^m(\hat{\mathbf{r}}) - (n+1) j_n^m Y_{n-1}^m(\hat{\mathbf{r}}), \quad j_n^m = \sqrt{\frac{n^2 - m^2}{4n^2 - 1}}, \quad (40)$$

$$\frac{\partial}{\partial \phi} Y_n^m(\hat{\mathbf{r}}) = i m Y_n^m(\hat{\mathbf{r}}). \quad (41)$$

By the fact that the geometry is fixed and parameters corresponding to the geometry of choice (here the radius  $R$ ) and medium (here the conductivity  $\sigma$ ) are known, the surface potential  $U_0$  is calculated via (37), which leads to knowledge of the coefficients  $\mathbb{A}_{0,k}^\ell$  through (32) for  $n = 0$ , since

$$\mathbb{A}_{0,k}^\ell = \oint_{S^2} U_0(R\hat{\mathbf{r}}, \mathbf{r}_0) \bar{Y}_k^\ell(\hat{\mathbf{r}}) dS(\hat{\mathbf{r}}). \quad (42)$$

After an approximate function  $f(\hat{\mathbf{r}})$  has been derived capable to fit the perturbed surface, the coefficients  $\mathbb{A}_{n,k}^\ell$  corresponding to the  $n$ th-order correction ( $n \geq 1$ ) are computed as demonstrated above. Subsequently, either the corresponding potentials

$U_n, n \geq 1$  via (32) or the coefficients (33) are estimated, leading to the perturbed surface potential  $U(\boldsymbol{\tau})$ .

The electric potential due to a dipolar source can be represented by Geselowitz's integral representation [29]

$$U(\mathbf{r}) = \frac{1}{4\pi\sigma} \mathbf{Q} \cdot \frac{\mathbf{r} - \mathbf{r}_0}{|\mathbf{r} - \mathbf{r}_0|^3} - \frac{1}{4\pi} \oint_{S^2} U(\boldsymbol{\tau}) \hat{\boldsymbol{\nu}} \cdot \frac{\mathbf{r} - \boldsymbol{\tau}}{|\mathbf{r} - \boldsymbol{\tau}|^3} dS(\hat{\boldsymbol{\tau}}), \quad \mathbf{r} \notin S, \quad (43)$$

where  $\mathbf{r}$  is the point of observation (measurement site). In the case of a spherical conductor of radius  $R$ , the integral becomes

$$U(\mathbf{r}) = \frac{1}{4\pi\sigma} \mathbf{Q} \cdot \frac{\mathbf{r} - \mathbf{r}_0}{|\mathbf{r} - \mathbf{r}_0|^3} - \frac{R^2}{4\pi} \oint_{\partial D} U(\mathbf{r}') \frac{\hat{\mathbf{r}}' \cdot \mathbf{r} - R}{|\mathbf{r} - R\hat{\mathbf{r}}'|^3} d\Omega(\hat{\mathbf{r}}'), \quad \mathbf{r} \notin S, \quad (44)$$

where  $\Omega$  denotes the solid angle. The latter is a Fredholm integral equation of the second kind and can be solved accordingly. By expanding all relevant expressions in (43) an integral representation for the perturbed instance is obtained leading again to Fredholm equations of the second kind for both the zeroth  $U_0$  and first order correction  $U_1$ . This scenario will be presented elsewhere.

### 2.1.2 MEG

From (13) and (15) it is evident that the scalar magnetic potential  $W$  is harmonic and can be expanded as

$$W(\mathbf{r}; \mathbf{r}_0) = \sum_{n=1}^{\infty} \sum_{m=-n}^n \mathbb{C}_n^m \frac{Y_n^m(\hat{\mathbf{r}})}{r^{n+1}}, \quad (45)$$

where

$$\mathbb{C}_n^m = \frac{4\pi}{(n+1)(2n+1)} (\mathbf{Q} \times \mathbf{r}_0) \cdot \nabla_{\mathbf{r}_0} \mathbf{r}_0^n \bar{Y}_n^m(\hat{\mathbf{r}}_0), \quad (46)$$

or

$$\frac{\mathbb{C}_n^m}{r^{n+1}} = \oint_{S^2} W(\mathbf{r}; \mathbf{r}_0) \bar{Y}_n^m(\hat{\mathbf{r}}_0) dS(\hat{\mathbf{r}}). \quad (47)$$

Clearly, knowledge of  $W$  provides the values of  $\mathbb{C}_n^m$  via (47) and thus the solution to the forward MEG problem.

The scalar magnetic potential  $W$  introduced in (15) can be computed, using integration along a ray from  $\mathbf{r}$  to infinity as

$$W(\mathbf{r}; \mathbf{r}_0) = -\frac{4\pi}{\mu_0} \int_r^{+\infty} \hat{\boldsymbol{\tau}} \cdot \mathbf{B}(\boldsymbol{\tau}; \mathbf{r}_0) d\tau, \quad (48)$$



which, combined with expression

$$\hat{\mathbf{r}} \cdot \mathbf{B}(\mathbf{r}; \mathbf{r}_0) = -\frac{\mu_0}{4\pi} \mathbf{Q} \times \frac{\mathbf{r}_0 \cdot \hat{\mathbf{r}}}{|\mathbf{r} - \mathbf{r}_0|^3}, \tag{49}$$

becomes

$$W(\mathbf{r}; \mathbf{r}_0) = \mathbf{Q} \times \mathbf{r}_0 \cdot \hat{\mathbf{r}} \int_{\mathbf{r}}^{+\infty} \frac{d\tau}{|\tau \hat{\mathbf{r}} - \mathbf{r}_0|^3}. \tag{50}$$

Performing straightforward calculations (details are provided in [28]) the magnetic potential assumes the form

$$W(\mathbf{r}; \mathbf{r}_0) = \mathbf{Q} \cdot \frac{\hat{\mathbf{r}} \times \hat{\mathbf{P}}}{P(1 + \hat{\mathbf{r}} \cdot \hat{\mathbf{P}})}, \quad \mathbf{P} = \mathbf{r} - \mathbf{r}_0, \tag{51}$$

implying that

$$\mathbf{B}(\mathbf{r}; \mathbf{r}_0) = \frac{\mu_0}{4\pi} \frac{\mathbf{Q} \times (\mathbf{r} - \mathbf{P})}{r P^2 (1 + \hat{\mathbf{r}} \cdot \hat{\mathbf{P}})} - \frac{\mu_0}{4\pi} \frac{U(\mathbf{r}; \mathbf{r}_0)}{P} \left[ \hat{\mathbf{P}} + \frac{r - P}{r(1 + \hat{\mathbf{r}} \cdot \hat{\mathbf{P}})} (\hat{\mathbf{r}} - \hat{\mathbf{P}}) \right]. \tag{52}$$

Based on Geselowitz’s formula (16), the magnetic field can also be computed as

$$\mathbf{B}(\mathbf{r}; \mathbf{r}_0) = \frac{\mu_0}{4\pi} \mathbf{Q}(\mathbf{r}_0) \times \nabla_{\mathbf{r}_0} \frac{1}{|\mathbf{r} - \mathbf{r}_0|} - \frac{\mu_0}{(4\pi)^2} (\mathbf{Q} \cdot \nabla_{\mathbf{r}_0}) \mathbf{K}(\mathbf{r}; \mathbf{r}_0), \tag{53}$$

where

$$\begin{aligned} \mathbf{K}(\mathbf{r}; \mathbf{r}_0) = \oint_{\partial D} \left( \frac{1}{|\mathbf{r}' - \mathbf{r}_0|} + \sum_{n=1}^{\infty} \sum_{m=-n}^n 4\pi \frac{n+1}{n(2n+1)} \frac{r^n r_0^n}{R^{2n+1}} \bar{Y}_n^m(\hat{\mathbf{r}}) Y_n^m(\hat{\mathbf{r}}_0) \right) \mathbf{r}' \\ \times \nabla_{\mathbf{r}'} \frac{1}{|\mathbf{r} - \mathbf{r}'|} ds(\mathbf{r}'). \end{aligned} \tag{54}$$

It can be shown that the magnetic field  $\mathbf{B}$  for the perturbed sphere equals [30]

$$\begin{aligned} \mathbf{B}(\boldsymbol{\tau}; \mathbf{r}_0) &= \sum_{n=0}^{\infty} \epsilon^n \mathbf{B}_n(\mathbf{r}; \mathbf{r}_0) \\ &= \frac{\mu_0}{4\pi} \mathbf{Q} \times \frac{\mathbf{r} - \mathbf{r}_0}{|\mathbf{r} - \mathbf{r}_0|^3} - \frac{\mu_0 \sigma}{4\pi} \oint_{\partial D} \boldsymbol{\alpha}_0(\hat{\mathbf{r}}, \mathbf{r}) U_0(R\hat{\mathbf{r}}, \mathbf{r}) \frac{d\Omega(\hat{\mathbf{r}})}{|\mathbf{r} - R\hat{\mathbf{r}}|^3} \end{aligned}$$

$$\begin{aligned}
& -\epsilon \frac{\mu_0 \sigma}{4\pi} \oint_{\partial D} \left( \alpha_0(\hat{\mathbf{r}}, \mathbf{r}) U_1(R\hat{\mathbf{r}}, \mathbf{r}) + \alpha_1(\hat{\mathbf{r}}, \mathbf{r}) U_0(R\hat{\mathbf{r}}, \mathbf{r}) \right) \frac{d\Omega(\hat{\mathbf{r}})}{|\mathbf{r} - R\hat{\mathbf{r}}|^3} \\
& + \mathcal{O}(\epsilon^2), \tag{55}
\end{aligned}$$

given that

$$\alpha_0(\hat{\mathbf{r}}, \mathbf{r}) = R^2(\hat{\mathbf{r}} \times \mathbf{r}), \tag{56}$$

$$\alpha_1(\hat{\mathbf{r}}, \mathbf{r}) = \left[ 2|\mathbf{r} - R\hat{\mathbf{r}}|^2 - 3R(R - \mathbf{r} \cdot \hat{\mathbf{r}}) \right] \frac{f(\hat{\mathbf{r}})\alpha_0(\hat{\mathbf{r}}, \mathbf{r})}{R|\mathbf{r} - R\hat{\mathbf{r}}|^2} - R\nabla f(\hat{\mathbf{r}}) \times (\mathbf{r} - R\hat{\mathbf{r}}). \tag{57}$$

## 2.2 The Influence of Geometric Variations on the Inversion Algorithm

### 2.2.1 EEG

In order to pinpoint the position  $\mathbf{r}_0 = (x_{01} \ x_{02} \ x_{03})^\top$  and moment  $\mathbf{Q} = (Q_1 \ Q_2 \ Q_3)^\top$  of the dipole, six equations are required in general. These relations are acquired based on formulas translating the spherical harmonics into Cartesian form (for example [21, pp. 155–156]). In the unperturbed case the position and moment of the equivalent dipole are expressed through (34) for  $k$  up to two as

$$\mathbf{r}_0 = \frac{R}{\sqrt{5}} \left( \frac{\mathbb{A}_{0,2}^{-2}}{\mathbb{A}_{0,1}^{-1}} - \frac{\mathbb{A}_{0,2}^2}{\mathbb{A}_{0,1}^1} - i \left( \frac{\mathbb{A}_{0,2}^{-2}}{\mathbb{A}_{0,1}^{-1}} + \frac{\mathbb{A}_{0,2}^2}{\mathbb{A}_{0,1}^1} \right) \frac{1}{\mathbb{A}_{0,1}^{\pm 1}} \left( 2\mathbb{A}_{0,2}^{\pm 1} - \sqrt{2} \frac{\mathbb{A}_{0,1}^0 \mathbb{A}_{0,2}^{\pm 2}}{\mathbb{A}_{0,1}^{\pm 1}} \right) \right)^\top \tag{58}$$

and

$$\mathbf{Q} = \sigma R^2 \sqrt{\frac{\pi}{3}} \left( \sqrt{2} \left( \mathbb{A}_{0,1}^{-1} - \mathbb{A}_{0,1}^1 \right) - i\sqrt{2} \left( \mathbb{A}_{0,1}^{-1} + \mathbb{A}_{0,1}^1 \right) 2\mathbb{A}_{0,1}^0 \right)^\top, \tag{59}$$

respectively.

In order to analyze the effect of perturbations on the inversion algorithm, our starting point is boundary condition (38). Substituting the matching potentials from (32) and expanding the known function  $f(\hat{\mathbf{r}})$  in terms of conjugate spherical harmonics

$$f(\hat{\mathbf{r}}) = \sum_{p=0}^{\infty} \sum_{q=-p}^p \bar{\mathbb{C}}_p^q \bar{\mathbb{Y}}_p^q(\hat{\mathbf{r}}) \tag{60}$$

we obtain

$$\sum_{k=1}^{\infty} \sum_{\ell=-k}^k k R^{k+1} \mathbb{A}_{1,k}^{\ell} \mathbb{Y}_k^{\ell}(\hat{\mathbf{r}}) = \sum_{k=1}^{\infty} \sum_{\ell=-k}^k \sum_{p=0}^{\infty} \sum_{q=-p}^p \mathbb{A}_{0,k}^{\ell} \bar{\mathbb{C}}_p^q \left( \frac{\partial \bar{\mathbb{Y}}_p^q(\hat{\mathbf{r}})}{\partial \theta} \frac{\partial \mathbb{Y}_k^{\ell}(\hat{\mathbf{r}})}{\partial \theta} + \frac{\ell q}{\sin^2 \theta} \bar{\mathbb{Y}}_p^q(\hat{\mathbf{r}}) \mathbb{Y}_k^{\ell}(\hat{\mathbf{r}}) \right). \tag{61}$$

Integrating above first with respect to  $\phi \in (0, 2\pi]$  and then with respect to  $\theta \in [0, \pi]$  and having in mind that

$$\int_{-1}^1 \left( (1-x^2) \frac{dP_n^m(x)}{dx} \frac{dP_k^{\ell}(x)}{dx} + m^2 \frac{P_n^m(x) P_k^{\ell}(x)}{1-x^2} \right) dx = \frac{2n(n+1)}{2n+1} \frac{(n+m)!}{(n-m)!} \delta_{n,k}, \tag{62}$$

where we make use of the Kronecker symbol  $\delta_{n,n} = 1$ , else 0, yields

$$\sum_{k=1}^{\infty} \sum_{\ell=-k}^k \frac{k(k+1)}{R^{k+1}} \left( \mathbf{Q} \cdot \nabla_{\mathbf{r}_0} r_0^k \bar{\mathbb{Y}}_k^{\ell}(\hat{\mathbf{r}}_0) \right) \bar{\mathbb{C}}_k^{\ell} = 0, \tag{63}$$

providing a criterion for the solvability of the BVP at hand. A crucial aspect is the fact that due to relation (63), which is an immediate consequence of the compatibility condition regarding the Neumann problem for the first-order correction, a restriction between the surface deformation  $f(\hat{\mathbf{r}})$ , the dipole’s position  $\mathbf{r}_0$  and moment  $\mathbf{Q}$  is imposed. The fine points will be demonstrated in the Example section.

### 2.2.2 MEG

As aforementioned, the key in calculating the sought position and moments of the dipole is to translate the spherical harmonics in Cartesian coordinates. For example, in the unperturbed case the coefficients are given by (46) and for  $k = 2$  and  $\ell = 0$  we have

$$\mathbb{C}_2^0 = \frac{4\pi}{15} (\mathbf{Q} \times \mathbf{r}_0) \cdot \nabla_{\mathbf{r}_0} r_0^2 \bar{\mathbb{Y}}_2^0(\hat{\mathbf{r}}_0) = \frac{2\pi}{15} \sqrt{\frac{5}{\pi}} (\mathbf{Q} \times \mathbf{r}_0) \cdot (-x_{01} \hat{\mathbf{x}}_1 - x_{02} \hat{\mathbf{x}}_2 + 2x_{03} \hat{\mathbf{x}}_3). \tag{64}$$

Repeating the procedure and after some algebra, we obtain the following  $3 \times 3$  system

$$\begin{pmatrix} \mathbb{C}_1^1 - \mathbb{C}_1^{-1} & i(\mathbb{C}_1^1 + \mathbb{C}_1^{-1}) & 2\sqrt{2}\mathbb{C}_1^0 \\ 0 & i\sqrt{2}\mathbb{C}_1^0 & \mathbb{C}_1^1 + \mathbb{C}_1^{-1} \\ -\sqrt{2}\mathbb{C}_1^0 & 0 & \mathbb{C}_1^{-1} - \mathbb{C}_1^1 \end{pmatrix} \begin{pmatrix} x_{01} \\ x_{02} \\ x_{03} \end{pmatrix} = \begin{pmatrix} \sqrt{\frac{15}{2}}\mathbb{C}_2^0 \\ \frac{\sqrt{5}}{2}\mathbb{C}_2^{-1} + \mathbb{C}_2^1 \\ \frac{\sqrt{5}}{2}\mathbb{C}_2^{-1} - \mathbb{C}_2^1 \end{pmatrix}, \quad (65)$$

from which we can pinpoint the source  $\mathbf{r}_0 = (x_{01} \ x_{02} \ x_{03})^\top$ . On the other hand, using higher-order coefficients, we obtain  $\mathbf{r}_0$  in an easier fashion as

$$\mathbf{r}_0 = \frac{\sqrt{5}}{4} \left( \frac{\mathbb{C}_2^{-2}}{\mathbb{C}_1^{-1}} - \frac{\mathbb{C}_2^2}{\mathbb{C}_1^1} - i \left( \frac{\mathbb{C}_2^{-2}}{\mathbb{C}_1^1} + \frac{\mathbb{C}_2^2}{\mathbb{C}_1^1} \right) \sqrt{3} \frac{\mathbb{C}_2^0}{\mathbb{C}_1^0} - \frac{1}{\sqrt{2}} \left( \frac{\mathbb{C}_2^{-2}\mathbb{C}_1^1}{\mathbb{C}_1^{-1}\mathbb{C}_1^0} + \frac{\mathbb{C}_2^2\mathbb{C}_1^{-1}}{\mathbb{C}_1^1\mathbb{C}_1^0} \right) \right)^\top. \quad (66)$$

Clearly, equating the relations for  $x_{0i}$ ,  $i = 1, 2, 3$ , uniqueness conditions emerge as we will see shortly. Similar, by properly rearranging coefficients and by the fact that

$$\mathbf{Q} \times \mathbf{r}_0 = \hat{\mathbf{i}}(x_{03} Q_2 - x_{02} Q_3) - \hat{\mathbf{j}}(x_{03} Q_1 - x_{01} Q_3) + \hat{\mathbf{k}}(x_{02} Q_1 - x_{01} Q_2), \quad (67)$$

we find for the moments

$$\begin{pmatrix} 0 & x_{03} & -x_{02} \\ -x_{03} & 0 & x_{01} \\ x_{02} & -x_{01} & 0 \end{pmatrix} \begin{pmatrix} Q_1 \\ Q_2 \\ Q_3 \end{pmatrix} = \begin{pmatrix} \sqrt{\frac{3}{2\pi}} (\mathbb{C}_1^{-1} - \mathbb{C}_1^1) \\ -i\sqrt{\frac{3}{2\pi}} (\mathbb{C}_1^{-1} + \mathbb{C}_1^1) \\ \sqrt{\frac{3}{\pi}} \mathbb{C}_1^0 \end{pmatrix}. \quad (68)$$

However, the determinant of the latter vanishes and thus more information is needed.

Before we proceed calculating the moments of the dipole, we note that if  $\mathbb{C}_1^0 = 0$  and one of  $\mathbb{C}_1^\pm = 0$ , results in  $\mathbf{Q} \times \mathbf{r}_0 = \mathbf{0}$ . In the presence of neuronal activity the latter implies  $\mathbf{Q}$  is parallel to  $\mathbf{r}_0$ , i.e. the moment  $\mathbf{Q}$  is always radial in the spherical coordinate system. This leads to that all radial dipoles inside a homogeneous sphere are not visible by MEG and localization of the dipole  $\mathbf{r}_0$ ,  $\mathbf{Q}$  is restrained finding the position  $\mathbf{r}_0$  as well as the tangential components  $Q_{\theta_0}$  and  $Q_{\phi_0}$  as  $\mathbf{Q} = Q_{r_0}\hat{\mathbf{r}}_0 + Q_{\theta_0}\hat{\boldsymbol{\theta}}_0 + Q_{\phi_0}\hat{\boldsymbol{\phi}}_0$ . Following the same procedure, we have

$$Q_{\theta_0} = \sqrt{\frac{3}{2\pi}} \frac{\mathbb{C}_1^{-1} - \mathbb{C}_1^1}{r_0 \sin \phi_0} + \sqrt{\frac{3}{\pi}} \frac{\mathbb{C}_1^0}{r_0 \tan \theta_0 \tan \phi_0}, \quad (69)$$

$$Q_{\phi_0} = -\sqrt{\frac{3}{\pi}} \frac{\mathbb{C}_1^0}{r_0 \sin \theta_0}. \quad (70)$$

Moreover, employing a different combination of coefficients would lead, for example, to

$$Q_{\theta_0} = -\sqrt{\frac{3}{\pi}} \frac{\tan \phi_0}{r_0 \tan \theta_0} \mathbb{C}_1^0 - i\sqrt{\frac{3}{2\pi}} \frac{1}{r_0 \cos \phi_0} (\mathbb{C}_1^{-1} + \mathbb{C}_1^1). \quad (71)$$

Evidently, above relations regarding  $Q_{\theta_0}$  must coincide and by equating them we obtain

$$\frac{i}{\sqrt{2}} \left( e^{-i\phi_0} \mathbb{C}_1^{-1} + e^{i\phi_0} \mathbb{C}_1^1 \right) = -\frac{1}{\tan \theta_0} \mathbb{C}_1^0. \quad (72)$$

namely the condition for  $Q_{\theta_0}$  and  $Q_{\phi_0}$  to be unique and is consistent with the system (68). Additional details can be found in [31].

In the perturbed case the scalar magnetic potential can be express as

$$W(\boldsymbol{\tau}) = \sum_{n=0}^{\infty} \epsilon^n W_n(\mathbf{r}) \quad (73)$$

and expanding both  $W(\boldsymbol{\tau})$  and  $W_n(\mathbf{r})$  as in (45), gives

$$\mathbb{D}_k^\ell(\boldsymbol{\tau}, \mathbf{r}_0, \mathbf{Q}; \epsilon) = \sum_{n=0}^{\infty} \epsilon^n \mathbb{E}_{n,k}^\ell(\mathbf{r}, \mathbf{r}_0, \mathbf{Q}), \quad \mathbb{E}_{0,k}^\ell(\mathbf{r}, \mathbf{r}_0, \mathbf{Q}) = \mathbb{C}_k^\ell(\mathbf{r}, \mathbf{r}_0, \mathbf{Q}). \quad (74)$$

The coefficients  $\mathbb{E}_{n,k}^\ell$  are computed with the aid of (48) by properly expanding the inner product  $\hat{\mathbf{r}}' \cdot \mathbf{B}(\hat{\mathbf{r}}', \mathbf{r}_0)$ , yielding

$$\begin{aligned} W(\mathbf{r}, \mathbf{r}_0; \epsilon) &= \int_r^{+\infty} \frac{\mathbf{Q} \times \mathbf{r}_0 \cdot \hat{\mathbf{r}}'}{|\mathbf{r}' - \mathbf{r}_0|^3} dr' \\ &+ \epsilon \sigma R^2 \int_r^{+\infty} \left[ \oint_{\partial D} U_1(\mathbf{v}) \frac{\hat{\mathbf{r}}' \cdot \nabla f \times \hat{\mathbf{v}}}{|\mathbf{r}' - R\hat{\mathbf{v}}|^3} d\Omega(\hat{\mathbf{v}}) \right] dr'. \end{aligned} \quad (75)$$

By the fact that

$$\nabla_{\mathbf{v}} \frac{1}{|r'\hat{\mathbf{r}}' - \mathbf{v}|} = \frac{r'\hat{\mathbf{r}}' - \mathbf{v}}{|r'\hat{\mathbf{r}}' - \mathbf{v}|^3}, \quad \mathbf{v} = v\hat{\mathbf{v}} \quad (76)$$

and

$$\int_r^\infty \frac{1}{|r'\hat{\mathbf{r}}' - \mathbf{v}|} \frac{dr'}{r'} = \sum_{k=0}^{\infty} \sum_{\ell=-k}^k \frac{4}{(k+1)(2k+1)} \frac{v^k}{r^{k+1}} Y_k^\ell(\hat{\mathbf{r}}') \bar{Y}_k^\ell(\hat{\mathbf{v}}), \quad (77)$$

the first order coefficients, in view of (40) and (41), are

$$\begin{aligned} \mathbb{E}_{1,k}^\ell &= \frac{4\pi\sigma R^{k+1}}{(k+1)(2k+1)} \sum_{p=1}^{\infty} \sum_{q=-p}^p \mathbb{B}_p^q \oint_{\partial D} \left[ \frac{\partial f}{\partial \phi} \frac{1}{\sin \theta} (k j_{k+1}^\ell \bar{Y}_{k+1}^\ell(\hat{\mathbf{r}}) Y_p^q(\hat{\mathbf{r}}) \right. \\ &\quad \left. - (k+1) j_k^\ell \bar{Y}_{k-1}^\ell(\hat{\mathbf{r}}) Y_p^q(\hat{\mathbf{r}}) + ik \frac{\partial f}{\partial \theta} \bar{Y}_k^\ell(\hat{\mathbf{r}}) Y_p^q(\hat{\mathbf{r}}) \right] d\theta d\phi \end{aligned} \quad (78)$$

where  $\mathbb{B}_p^q$  are the coefficients corresponding to the surface electric potential and are given by (33). The coefficients  $\mathbb{E}_{0,k}^\ell$  for the unperturbed case are given by (47).

The analysis so far holds for a single dipole located at  $\mathbf{r}_0$ . Following Fokas [32], let us assume that the neuronal activity is distributed inside the spherical brain as  $\mathbf{J}^p$  at  $\mathbf{r}_0$ . This instance is analyzed replacing  $\mathbf{Q}(\mathbf{r}_0)$  by  $\int_D \mathbf{J}^p(\mathbf{r}_0) d\nu(\mathbf{r}_0)$  in Geselowitz's formula (53), i.e.

$$\frac{4\pi}{\mu_0} \mathbf{B}(\mathbf{r}; \mathbf{r}_0) = \int_D \mathbf{J}^p(\mathbf{r}_0) \times \nabla_{\mathbf{r}_0} \frac{1}{|\mathbf{r} - \mathbf{r}_0|} d\nu(\mathbf{r}_0) - \frac{1}{4\pi} \int_D \mathbf{J}^p(\mathbf{r}_0) \cdot \nabla_{\mathbf{r}_0} \mathbf{K}(\mathbf{r}; \mathbf{r}_0) d\nu(\mathbf{r}_0). \quad (79)$$

Although the analysis is presented in [32], we will reproduce part of it providing all necessary details in order to facilitate the transition to the perturbed case. Utilizing the identities  $\nabla \times (f\mathbf{g}) = \nabla f \times \mathbf{g} + f\nabla \times \mathbf{g}$  and  $\nabla \cdot (\mathbf{f} \otimes \mathbf{g}) = (\nabla \cdot \mathbf{f}) \otimes \mathbf{g} + \mathbf{f} \cdot \nabla \otimes \mathbf{g}$ ,  $\otimes$  denoting the tensor product, gives

$$\mathbf{J}^p(\mathbf{r}_0) \times \nabla_{\mathbf{r}_0} \frac{1}{|\mathbf{r} - \mathbf{r}_0|} = \frac{1}{|\mathbf{r} - \mathbf{r}_0|} \nabla_{\mathbf{r}_0} \times \mathbf{J}^p(\mathbf{r}_0) - \nabla_{\mathbf{r}_0} \times \frac{\mathbf{J}^p(\mathbf{r}_0)}{|\mathbf{r} - \mathbf{r}_0|} \quad (80)$$

and

$$\nabla_{\mathbf{r}_0} \cdot (\mathbf{J}^p(\mathbf{r}_0) \otimes \mathbf{K}(\mathbf{r}; \mathbf{r}_0)) = (\nabla_{\mathbf{r}_0} \cdot \mathbf{J}^p(\mathbf{r}_0)) \otimes \mathbf{K}(\mathbf{r}; \mathbf{r}_0) + \mathbf{J}^p(\mathbf{r}_0) \cdot \nabla_{\mathbf{r}_0} \otimes \mathbf{K}(\mathbf{r}; \mathbf{r}_0). \quad (81)$$

Replacing (80) and (81) into (79) yields

$$\int_D \mathbf{J}^p(\mathbf{r}_0) \times \nabla_{\mathbf{r}_0} \frac{1}{|\mathbf{r} - \mathbf{r}_0|} d\nu(\mathbf{r}_0) = \int_D (\nabla_{\mathbf{r}_0} \times \mathbf{J}^p(\mathbf{r}_0)) \frac{1}{|\mathbf{r} - \mathbf{r}_0|} d\nu(\mathbf{r}_0), \quad (82)$$

since, according to (19)

$$\int_D \nabla_{\mathbf{r}_0} \times \left( \frac{1}{|\mathbf{r} - \mathbf{r}_0|} \mathbf{J}^p(\mathbf{r}_0) \right) d\nu(\mathbf{r}_0) = \int_{\partial D} \hat{\mathbf{n}} \times \left( \frac{1}{|\mathbf{r} - \mathbf{r}_0|} \mathbf{J}^p(\mathbf{r}_0) \right) ds(\mathbf{r}_0) = 0, \quad (83)$$

which vanishes bearing in mind that  $\mathbf{J}^p(\mathbf{r}_0) = \mathbf{0}$  on the boundary  $\partial D$ . Similar,

$$\begin{aligned} \int_D \mathbf{J}^p(\mathbf{r}_0) \cdot \nabla_{\mathbf{r}_0} \mathbf{K}(\mathbf{r}; \mathbf{r}_0) d\nu(\mathbf{r}_0) &= \int_D \nabla_{\mathbf{r}_0} \cdot (\mathbf{J}^p(\mathbf{r}_0) \mathbf{K}(\mathbf{r}; \mathbf{r}_0)) d\nu(\mathbf{r}_0) \\ &\quad - \int_D (\nabla_{\mathbf{r}_0} \cdot \mathbf{J}^p(\mathbf{r}_0)) \mathbf{K}(\mathbf{r}; \mathbf{r}_0) d\nu(\mathbf{r}_0). \end{aligned} \quad (84)$$

Applying the divergence theorem  $\int_D \nabla \otimes \mathbf{f} d\nu = \int_{\partial D} \mathbf{f} \otimes \hat{\mathbf{v}} dS$ , we get

$$\int_D \nabla_{\mathbf{r}_0} \cdot (\mathbf{J}^p(\mathbf{r}_0) \mathbf{K}(\mathbf{r}; \mathbf{r}_0)) d\nu(\mathbf{r}_0) = \int_{\partial D} \hat{\mathbf{v}} \cdot (\mathbf{J}^p(\mathbf{r}_0) \mathbf{K}(\mathbf{r}; \mathbf{r}_0)) ds(\mathbf{r}_0). \quad (85)$$

Therefore, (79) reads as

$$\frac{4\pi}{\mu_0} \mathbf{B}(\mathbf{r}) = \int_D (\nabla_{\mathbf{r}_0} \times \mathbf{J}^p(\mathbf{r}_0)) \frac{d\nu(\mathbf{r}_0)}{|\mathbf{r} - \mathbf{r}_0|} + \frac{1}{4\pi} \int_D (\nabla_{\mathbf{r}_0} \cdot \mathbf{J}^p(\mathbf{r}_0)) \mathbf{K}(\mathbf{r}; \mathbf{r}_0) d\nu(\mathbf{r}_0). \quad (86)$$

The fact that the conductor  $D$  is star shaped and letting the neuronal current  $\mathbf{J}^p$  be  $\mathcal{C}^1$  allows for  $\mathbf{J}^p$  to be represented by Helmholtz's decomposition, namely

$$\mathbf{J}^p(\mathbf{r}_0) = \nabla \Psi(\mathbf{r}_0) + \nabla \times \mathbf{A}(\mathbf{r}_0), \quad \nabla \cdot \mathbf{A}(\mathbf{r}_0) = 0, \quad \mathbf{r}_0 \in D \quad (87)$$

Above in mind, we have  $\nabla_{\mathbf{r}_0} \times \mathbf{J}^p(\mathbf{r}_0) = -\Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0)$  and  $\nabla_{\mathbf{r}_0} \cdot \mathbf{J}^p(\mathbf{r}_0) = \Delta_{\mathbf{r}_0} \Psi(\mathbf{r}_0)$  and equation (86) becomes

$$\frac{4\pi}{\mu_0} \mathbf{B}(\mathbf{r}) = - \int_D (\Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0)) \frac{d\nu(\mathbf{r}_0)}{|\mathbf{r} - \mathbf{r}_0|} + \frac{1}{4\pi} \int_D (\Delta_{\mathbf{r}_0} \Psi(\mathbf{r}_0)) \mathbf{K}(\mathbf{r}; \mathbf{r}_0) d\nu(\mathbf{r}_0). \quad (88)$$

Now, the first integral of (88) writes as

$$\mathbf{r} \cdot \int_D \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) \frac{d\nu(\mathbf{r}_0)}{|\mathbf{r} - \mathbf{r}_0|} = \int_D \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) \cdot \frac{\mathbf{r} - \mathbf{r}_0}{|\mathbf{r} - \mathbf{r}_0|} d\nu(\mathbf{r}_0) + \int_D \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) \cdot \frac{\mathbf{r}_0}{|\mathbf{r} - \mathbf{r}_0|} d\nu(\mathbf{r}_0), \quad (89)$$

where

$$\int_D \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) \cdot \frac{\mathbf{r} - \mathbf{r}_0}{|\mathbf{r} - \mathbf{r}_0|} d\nu(\mathbf{r}_0) = \int_D |\mathbf{r} - \mathbf{r}_0|^2 \left( \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) \cdot \nabla_{\mathbf{r}_0} \frac{1}{|\mathbf{r} - \mathbf{r}_0^3|} \right) d\nu(\mathbf{r}_0) \quad (90)$$

and

$$\begin{aligned} \int_D \nabla_{\mathbf{r}_0} \cdot \left[ \left( |\mathbf{r} - \mathbf{r}_0|^2 \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) \right) \frac{1}{|\mathbf{r} - \mathbf{r}_0|} \right] d\nu(\mathbf{r}_0) = \\ \int_D \left[ \nabla_{\mathbf{r}_0} \cdot \left( |\mathbf{r} - \mathbf{r}_0|^2 \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) \right) \right] \frac{1}{|\mathbf{r} - \mathbf{r}_0|} d\nu(\mathbf{r}_0) + \int_D \left( |\mathbf{r} - \mathbf{r}_0|^2 \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) \right) \\ \cdot \nabla_{\mathbf{r}_0} \frac{1}{|\mathbf{r} - \mathbf{r}_0|} d\nu(\mathbf{r}_0) \end{aligned} \quad (91)$$

so that

$$\begin{aligned} \int_D \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) \cdot \frac{\mathbf{r} - \mathbf{r}_0}{|\mathbf{r} - \mathbf{r}_0|} d\nu(\mathbf{r}_0) = \oint_{\partial D} \hat{\nu}(\mathbf{r}_0) \cdot |\mathbf{r} - \mathbf{r}_0| \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) ds(\mathbf{r}_0) \\ - \int_D \left[ \nabla_{\mathbf{r}_0} \cdot \left( |\mathbf{r} - \mathbf{r}_0|^2 \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) \right) \right] \frac{1}{|\mathbf{r} - \mathbf{r}_0|} d\nu(\mathbf{r}_0) \end{aligned} \quad (92)$$

By the fact that  $\mathbf{J}^P(\mathbf{r}_0)$  vanishes on  $\partial D$  so does  $\mathbf{A}(\mathbf{r}_0)$ . Thus

$$\begin{aligned} \int_D \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) \cdot \frac{\mathbf{r} - \mathbf{r}_0}{|\mathbf{r} - \mathbf{r}_0|} d\nu(\mathbf{r}_0) &= - \int_D \left( \nabla_{\mathbf{r}_0} |\mathbf{r} - \mathbf{r}_0|^2 \right) \cdot \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) \frac{1}{|\mathbf{r} - \mathbf{r}_0|} d\nu(\mathbf{r}_0) \\ &\quad - \int_D |\mathbf{r} - \mathbf{r}_0|^2 \left( \nabla_{\mathbf{r}_0} \cdot \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) \right) \frac{1}{|\mathbf{r} - \mathbf{r}_0|} d\nu(\mathbf{r}_0). \end{aligned} \quad (93)$$

Also  $\nabla_{\mathbf{r}_0} \cdot \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) = \Delta_{\mathbf{r}_0} (\nabla_{\mathbf{r}_0} \cdot \mathbf{A}(\mathbf{r}_0)) = 0$ ,  $\nabla_{\mathbf{r}_0} |\mathbf{r} - \mathbf{r}_0|^2 = -2(\mathbf{r} - \mathbf{r}_0)$  and therefore

$$\int_D \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) \cdot \frac{\mathbf{r} - \mathbf{r}_0}{|\mathbf{r} - \mathbf{r}_0|} d\nu(\mathbf{r}_0) = 2 \int_D (\mathbf{r} - \mathbf{r}_0) \cdot \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) \frac{1}{|\mathbf{r} - \mathbf{r}_0|} d\nu(\mathbf{r}_0) \quad (94)$$

so that

$$\mathbf{r} \cdot \int_D \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) \frac{d\nu(\mathbf{r}_0)}{|\mathbf{r} - \mathbf{r}_0|} = \int_D \mathbf{r}_0 \cdot \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) \frac{d\nu(\mathbf{r}_0)}{|\mathbf{r} - \mathbf{r}_0|}. \quad (95)$$

In the sequel, employing the identity

$$\Delta(\mathbf{f} \cdot \mathbf{g}) = (\Delta \mathbf{f}) \cdot \mathbf{g} + \mathbf{f} \cdot (\Delta \mathbf{g}) + 2(\nabla \otimes \mathbf{f})^\top : (\nabla \otimes \mathbf{g}), \quad (96)$$

$\top$  denoting transposition and  $:$  the double dot product, yields

$$\begin{aligned} \Delta_{\mathbf{r}_0} (\mathbf{r}_0 \cdot \mathbf{A}(\mathbf{r}_0)) &= \Delta_{\mathbf{r}_0} \mathbf{r}_0 \cdot \mathbf{A}(\mathbf{r}_0) + \mathbf{r}_0 \cdot \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) + 2(\nabla \otimes \mathbf{r}_0)^\top : (\nabla_{\mathbf{r}_0} \otimes \mathbf{A}(\mathbf{r}_0)) \\ &= \mathbf{0} \cdot \mathbf{A}(\mathbf{r}_0) + \mathbf{r}_0 \cdot \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) + 2\tilde{\mathbf{I}} : \nabla_{\mathbf{r}_0} \otimes \mathbf{A}(\mathbf{r}_0) \\ &= \mathbf{r}_0 \cdot \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) + 2\nabla_{\mathbf{r}_0} \cdot \mathbf{A}(\mathbf{r}_0) \\ &= \mathbf{r}_0 \cdot \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0). \end{aligned} \quad (97)$$

Combining (97) and (95) gives

$$\mathbf{r} \cdot \int_D \Delta_{\mathbf{r}_0} \mathbf{A}(\mathbf{r}_0) \frac{d\nu(\mathbf{r}_0)}{|\mathbf{r} - \mathbf{r}_0|} = \int_D \Delta_{\mathbf{r}_0} (\mathbf{r}_0 \cdot \mathbf{A}(\mathbf{r}_0)) \frac{d\nu(\mathbf{r}_0)}{|\mathbf{r} - \mathbf{r}_0|} \quad (98)$$

and finally, (88) reads

$$\frac{4\pi}{\mu_0} \mathbf{r} \cdot \mathbf{B}(\mathbf{r}) = - \int_D \Delta_{\mathbf{r}_0} (\mathbf{r}_0 \cdot \mathbf{A}(\mathbf{r}_0)) \frac{d\nu(\mathbf{r}_0)}{|\mathbf{r} - \mathbf{r}_0|} + \frac{1}{4\pi} \int_D (\Delta_{\mathbf{r}_0} \Psi(\mathbf{r}_0)) \mathbf{r} \cdot \mathbf{K}(\mathbf{r}; \mathbf{r}_0) d\nu(\mathbf{r}_0). \quad (99)$$



Next, consider the quantity  $\mathbf{r} \cdot \mathbf{K}(\mathbf{r}; \mathbf{r}_0)$  present in the second integral of the RHS of (99). In view of (54) we have

$$\mathbf{r} \cdot \mathbf{r}' \times \nabla_{\mathbf{r}'} \frac{1}{|\mathbf{r} - \mathbf{r}'|} = \mathbf{r} \cdot \mathbf{r}' \times \frac{\mathbf{r} - \mathbf{r}'}{|\mathbf{r} - \mathbf{r}'|^3} = \mathbf{r} \cdot \frac{\mathbf{r}' \times \mathbf{r}}{|\mathbf{r} - \mathbf{r}'|^3} = 0 \quad (100)$$

and thus  $\mathbf{r} \cdot \mathbf{K}(\mathbf{r}; \mathbf{r}_0) = 0$  so that (99) simplifies as

$$\mathbf{r} \cdot \mathbf{B}(\mathbf{r}) = -\frac{\mu_0}{4\pi} \int_{\partial D} \Delta_{\mathbf{r}_0} (\mathbf{r}_0 \cdot \mathbf{A}(\mathbf{r}_0)) \frac{d\nu(\mathbf{r}_0)}{|\mathbf{r} - \mathbf{r}_0|}. \quad (101)$$

Introducing the expansion

$$\frac{1}{|\mathbf{r} - \mathbf{r}_0|} = \sum_{n=0}^{\infty} \sum_{m=-n}^n \frac{4\pi}{2n+1} \frac{r_0^n}{r^{n+1}} Y_n^m(\hat{\mathbf{r}}) \bar{Y}_n^m(\hat{\mathbf{r}}_0) \quad (102)$$

and (101) reads

$$\mathbf{r} \cdot \mathbf{B}(\mathbf{r}) = -\mu_0 \sum_{n=0}^{\infty} \sum_{m=-n}^n \frac{1}{2n+1} \frac{Y_n^m(\hat{\mathbf{r}})}{r^{n+1}} \int_{r_0 < R} \Delta_{\mathbf{r}_0} (\mathbf{r}_0 \cdot \mathbf{A}(\mathbf{r}_0)) r_0^n \bar{Y}_n^m(\hat{\mathbf{r}}_0) d\nu(\mathbf{r}_0). \quad (103)$$

Further, assume

$$\mathbf{r}_0 \cdot \mathbf{A}(\mathbf{r}_0) = \sum_{n=0}^{\infty} \sum_{m=-n}^n r_0^2 a_n^m(r_0) Y_n^m(\hat{\mathbf{r}}_0), \quad r_0 \in [0, R] \quad (104)$$

so that

$$\Delta_{\mathbf{r}_0} (\mathbf{r}_0 \cdot \mathbf{A}(\mathbf{r}_0)) = \left( \frac{\partial^2}{\partial r_0^2} + \frac{2}{r_0} \frac{\partial}{\partial r_0} + \frac{1}{r_0^2} \mathcal{B}(\hat{\mathbf{r}}_0) \right) (\mathbf{r}_0 \cdot \mathbf{A}(\mathbf{r}_0)), \quad (105)$$

where the Beltrami or surface Laplacian operator  $\mathcal{B}(\hat{\mathbf{r}}_0)$  is

$$\mathcal{B}(\hat{\mathbf{r}}_0) = \frac{\partial^2}{\partial \theta_0^2} + \frac{\cos \theta_0}{\sin \theta_0} \frac{\partial}{\partial \theta_0} + \frac{1}{\sin^2 \theta_0} \frac{\partial^2}{\partial \phi_0^2}, \quad (106)$$

fulfilling

$$\mathcal{B}(\hat{\mathbf{r}}_0) Y_n^m(\hat{\mathbf{r}}_0) = -n(n+1) Y_n^m(\hat{\mathbf{r}}_0). \quad (107)$$

Hence,

$$\Delta_{\mathbf{r}_0}(\mathbf{r}_0 \cdot \mathbf{A}(\mathbf{r}_0)) = \sum_{n=0}^{\infty} \sum_{m=-n}^n \left[ r_0 \frac{d^2 a_n^m}{dr_0^2} + 4 \frac{da_n^m}{dr_0} - \frac{1}{r_0} (n-1)(n+2) a_n^m \right] Y_n^m(\hat{\mathbf{r}}_0). \quad (108)$$

Replacing above into (103) yields

$$\begin{aligned} \mathbf{r} \cdot \mathbf{B}(\mathbf{r}) &= -\mu_0 \sum_{n=1}^{\infty} \sum_{m=-n}^n \frac{1}{2n+1} \frac{Y_n^m(\hat{\mathbf{r}})}{r^{n+1}} \\ &\quad \times \int_0^R \left( r_0^{n+3} \frac{d^2 a_n^m}{dr_0^2} + 4r_0^{n+2} \frac{da_n^m}{dr_0} - r_0^{n+1} (n-1)(n+2) a_n^m \right) dr_0, \end{aligned} \quad (109)$$

since  $d\nu(\mathbf{r}_0) = r_0^2 dr_0 d\Omega(\hat{\mathbf{r}}_0)$  and

$$\int_{S^2(\hat{\mathbf{r}}_0)} Y_p^q(\hat{\mathbf{r}}_0) \bar{Y}_n^m(\hat{\mathbf{r}}_0) d\Omega(\hat{\mathbf{r}}_0) = \delta_{pn} \delta_{qm}. \quad (110)$$

Note, that since  $\mathbf{B}(\mathbf{r}) = \mathcal{O}(r^{-2})$  as  $r \rightarrow \infty$ ,  $n$  must begin with one. Simple calculations lead to

$$\mathbf{r} \cdot \mathbf{B}(\mathbf{r}) = -\mu_0 \sum_{n=1}^{\infty} \sum_{m=-n}^n \frac{1}{2n+1} \frac{R^{n+2}}{r^{n+1}} \left( R \frac{da_n^m}{dr_0} \Big|_{r_0=R} - (n-1) a_n^m(R) \right) Y_n^m(\hat{\mathbf{r}}). \quad (111)$$

In what follows, consider the quantity  $\mathbf{r} \cdot \mathbf{B}(\mathbf{r})$  known from measurements, i.e. if we expand

$$\mathbf{r} \cdot \mathbf{B}(\mathbf{r}) = \sum_{n=1}^{\infty} \sum_{m=-n}^n d_n^m(r) Y_n^m(\hat{\mathbf{r}}), \quad (112)$$

then the coefficients  $d_n^m(r)$  are known as well. Combining (111) and (112) we have

$$R \frac{da_n^m}{dr_0} \Big|_{r_0=R} - (n-1) a_n^m(R) = -\frac{2n+1}{\mu_0} \frac{r^{n+1}}{R^{n+2}} d_n^m(r) \quad (113)$$

from which the coefficients  $a_n^m$  are calculated.

The remaining two components of  $\mathbf{A}(\mathbf{r})$ , namely  $A_\theta(\mathbf{r})$  and  $A_\phi(\mathbf{r})$ , are computed with the aid of  $\nabla \cdot \mathbf{A}(\mathbf{r}) = 0$ , i.e.

$$\frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 A_r(\mathbf{r}) \right) + \frac{1}{r \sin \theta} \frac{\partial}{\partial \theta} (\sin \theta A_\theta(\mathbf{r})) + \frac{1}{r \sin \theta} \frac{\partial}{\partial \phi} A_\phi(\mathbf{r}) = 0, \quad (114)$$

since

$$\mathbf{A}(\mathbf{r}) = A_r(\mathbf{r})\hat{\mathbf{r}} + A_\theta(\mathbf{r})\hat{\boldsymbol{\theta}} + A_\phi(\mathbf{r})\hat{\boldsymbol{\phi}}. \quad (115)$$

Further, considering that  $\mathbf{r} \cdot \mathbf{A}(\mathbf{r}) = r A_r(\mathbf{r})$ , we find, with the help of (104), that

$$r^2 A_r(\mathbf{r}) = \sum_{n=1}^{\infty} \sum_{m=-n}^n r^2 a_n^m(r) Y_n^m(\hat{\mathbf{r}}), \quad (116)$$

which is known. Moreover, let

$$\frac{\partial}{\partial \theta} (\sin \theta A_\theta(\mathbf{r})) = \sin \theta \sum_{n=1}^{\infty} \sum_{m=-n}^n \zeta_n^m(r) Y_n^m(\hat{\mathbf{r}}), \quad (117)$$

$$\frac{\partial}{\partial \phi} A_\phi(\mathbf{r}) = \sin \theta \sum_{n=1}^{\infty} \sum_{m=-n}^n \eta_n^m(r) Y_n^m(\hat{\mathbf{r}}), \quad (118)$$

so that (114) gives

$$\zeta_n^m + \eta_n^m = \frac{1}{r} \frac{d}{dr} (r^2 a_n^m). \quad (119)$$

Switching to the perturbed case, we see from equation (54) that the quantity  $\mathbf{K}(\boldsymbol{\tau}; \mathbf{r}_0)$  is the only one affected by boundary deformations. Therefore, on the deformed boundary,

$$\hat{\mathbf{v}}(\boldsymbol{\tau}) dS(\boldsymbol{\tau}) = (R + \epsilon f) [(R + \epsilon f) \sin \theta \hat{\mathbf{r}} - \epsilon \sin \theta \nabla f] d\theta d\phi. \quad (120)$$

Also, employing Taylor's expansion gives

$$\hat{\mathbf{n}}(\boldsymbol{\tau}) \times \frac{\mathbf{r} - \boldsymbol{\tau}}{|\mathbf{r} - \boldsymbol{\tau}|^3} dS(\boldsymbol{\tau}) = \left( \sum_{n=0}^{\infty} \epsilon^n \boldsymbol{\alpha}_n(\hat{\mathbf{r}}, \mathbf{r}) \right) \frac{dS(\hat{\mathbf{r}})}{|\mathbf{r} - R\hat{\mathbf{r}}|^3}, \quad (121)$$

where the first two coefficients  $\boldsymbol{\alpha}$  are given by (56) and (57), respectively and further

$$\frac{1}{|\boldsymbol{\tau} - \mathbf{r}_0|} = \frac{1}{|R\hat{\mathbf{r}} - \mathbf{r}_0|} - \epsilon f \frac{R - \hat{\mathbf{r}} \cdot \mathbf{r}_0}{|R\hat{\mathbf{r}} - \mathbf{r}_0|^3} + \dots \quad (122)$$

On the other hand,  $\mathbf{r} \cdot \mathbf{K}(\mathbf{r}; \mathbf{r}_0)$  implies  $\mathbf{r} \cdot \boldsymbol{\alpha}_n(\hat{\mathbf{r}}, \mathbf{r})$ ,  $n \in \mathbb{N}$ , namely

$$\mathbf{r} \cdot \boldsymbol{\alpha}_0(\hat{\mathbf{r}}, \mathbf{r}) = 0, \quad (123)$$

$$\mathbf{r} \cdot \boldsymbol{\alpha}_1(\hat{\mathbf{r}}, \mathbf{r}) = R^2 \nabla_s f(\hat{\mathbf{r}}) \times \hat{\mathbf{r}} \cdot \mathbf{r}. \quad (124)$$

Combining all results, finally gives

$$\mathbf{r} \cdot \mathbf{K}(\mathbf{r}; \mathbf{r}_0) = \sum_{n=1}^{\infty} \epsilon^n \mathbf{I}_n(\mathbf{r}; \mathbf{r}_0), \quad (125)$$

provided that

$$\mathbf{I}_1(\mathbf{r}; \mathbf{r}_0) = R^2 \oint_{S^2} \left( \frac{1}{|R\hat{\mathbf{r}} - \mathbf{r}_0|} + U_0(R\hat{\mathbf{r}}; \mathbf{r}_0) \right) \frac{\nabla_s f(\hat{\mathbf{r}}) \times \hat{\mathbf{r}} \cdot \mathbf{r}}{|\mathbf{r} - R\hat{\mathbf{r}}|^3} dS(\hat{\mathbf{r}}). \quad (126)$$

Note that the integrals  $\mathbf{I}_n(\mathbf{r}; \mathbf{r}_0)$  do not depend on  $\mathbf{J}^p$ . Replacing everything back, gives

$$\begin{aligned} \frac{4\pi}{\mu_0} \mathbf{r} \cdot \mathbf{B}(\mathbf{r}) &= - \int_D \Delta_{\mathbf{r}_0} (\mathbf{r}_0 \cdot \mathbf{A}(\mathbf{r}_0)) \frac{d\nu(\mathbf{r}_0)}{|\mathbf{r} - \mathbf{r}_0|} \\ &+ \frac{1}{4\pi} \sum_{n=1}^{\infty} \epsilon^n \int_D (\Delta_{\mathbf{r}_0} \Psi(\mathbf{r}_0)) \mathbf{I}_n(\mathbf{r}; \mathbf{r}_0) d\nu(\mathbf{r}_0). \end{aligned} \quad (127)$$

Following previous analysis we have

$$\begin{aligned} &\int_D \Delta_{\mathbf{r}_0} (\mathbf{r}_0 \cdot \mathbf{A}(\mathbf{r}_0)) \frac{d\nu(\mathbf{r}_0)}{|\mathbf{r} - \mathbf{r}_0|} \\ &= 4\pi \sum_{n=1}^{\infty} \sum_{m=-n}^n \frac{1}{2n+1} \frac{R^{n+2}}{r^{n+1}} \left( R \frac{da_n^m}{dr_0} \Big|_{r_0=R} - (n-1)a_n^m(R) \right) Y_n^m(\hat{\mathbf{r}}). \end{aligned} \quad (128)$$

On the other hand, letting

$$\Psi(\mathbf{r}_0) = \sum_{n=0}^{\infty} \sum_{m=-n}^n \psi_n^m(r_0) Y_n^m(\hat{\mathbf{r}}_0), \quad (129)$$

we find

$$\Delta_{\mathbf{r}_0} \Psi(\mathbf{r}_0) = \sum_{n=0}^{\infty} \sum_{m=-n}^n \left( \ddot{\psi}_n^m(r_0) + \frac{2}{r_0} \dot{\psi}_n^m(r_0) - \frac{n(n+1)}{r_0^2} \psi_n^m(r_0) \right) Y_n^m(\hat{\mathbf{r}}_0). \quad (130)$$

The surface integral  $\mathbf{I}_1$  provided by (126) is evaluated as follows. First we note that the quantity inside the parenthesis enjoys the expansion

$$\frac{1}{|R\hat{\mathbf{r}} - \mathbf{r}_0|} + u(R\hat{\mathbf{r}}; \mathbf{r}_0) = 4\pi \sum_{n=1}^{\infty} \sum_{m=-n}^n \frac{n+1}{n(2n+1)} \frac{r_0^n}{R^{n+1}} Y_n^m(\hat{\mathbf{r}}) \bar{Y}_n^m(\hat{\mathbf{r}}_0). \quad (131)$$

Moreover, by the fact that  $\nabla f(\hat{\mathbf{r}}) \times \hat{\mathbf{r}} \cdot \mathbf{r}' = 0$ , we have

$$\frac{\nabla f(\hat{\mathbf{r}}) \times \hat{\mathbf{r}} \cdot \mathbf{r}'}{|\mathbf{r}' - R\hat{\mathbf{r}}|^3} = (\nabla f(\hat{\mathbf{r}}) \times \hat{\mathbf{r}}) \cdot \nabla_{\mathbf{r}'} \sum_{n=0}^{\infty} \sum_{m=-n}^n \frac{r'^m}{r'^{n+1}} Y_n^m(\hat{\mathbf{r}}) \bar{Y}_n^m(\hat{\mathbf{r}}'). \tag{132}$$

Expand the functions in terms of spherical harmonics as

$$\frac{1}{\sin \theta} \frac{\partial f(\hat{\mathbf{r}})}{\partial \phi} = \sum_{n=0}^{\infty} \sum_{m=-n}^n A_n^m Y_n^m(\hat{\mathbf{r}}), \quad \frac{\partial f(\hat{\mathbf{r}})}{\partial \theta} = \sum_{n=0}^{\infty} \sum_{m=-n}^n B_n^m Y_n^m(\hat{\mathbf{r}}) \tag{133}$$

to find

$$\nabla f(\hat{\mathbf{r}}) \times \hat{\mathbf{r}} = \sum_{n=0}^{\infty} \sum_{m=-n}^n \mathbf{C}_n^m Y_n^m(\hat{\mathbf{r}}), \quad \mathbf{C}_n^m = A_n^m \hat{\boldsymbol{\theta}} - B_n^m \hat{\boldsymbol{\phi}}, \tag{134}$$

given that

$$A_n^m = \oint_{S^2} \frac{1}{\sin \theta} \frac{\partial f(\hat{\mathbf{r}})}{\partial \phi} \bar{Y}_n^m(\hat{\mathbf{r}}) dS(\hat{\mathbf{r}}), \tag{135}$$

$$B_n^m = \oint_{S^2} \frac{\partial f(\hat{\mathbf{r}})}{\partial \theta} \bar{Y}_n^m(\hat{\mathbf{r}}) dS(\hat{\mathbf{r}}). \tag{136}$$

Replacing (134) into (132), gives

$$\begin{aligned} & \sum_{n=0}^{\infty} \sum_{m=-n}^n \mathbf{C}_n^m Y_n^m(\hat{\mathbf{r}}) \cdot \nabla_{\mathbf{r}'} \sum_{n=0}^{\infty} \sum_{m=-n}^n \frac{r'^m}{r'^{n+1}} Y_n^m(\hat{\mathbf{r}}) \bar{Y}_n^m(\hat{\mathbf{r}}) \\ &= \sum_{n=0}^{\infty} \sum_{m=-n}^n \frac{R^{n-1}}{r'^{n+1}} Y_n^m(\hat{\mathbf{r}}) Y_n^m(\hat{\mathbf{r}}) \left( A_n^m \frac{\partial \bar{Y}_n^m(\hat{\mathbf{r}})}{\partial \theta} - \frac{B_n^m}{\sin \theta} \frac{\partial \bar{Y}_n^m(\hat{\mathbf{r}})}{\partial \phi} \right) \end{aligned} \tag{137}$$

Substituting (131), (137) into the surface integral (126) yields

$$\begin{aligned} I_1(\mathbf{r}; \mathbf{r}_0) &= R^2 \oint_{S^2} \left( 4\pi \sum_{n=1}^{\infty} \sum_{m=-n}^n v_n^m(r_0) Y_n^m(\hat{\mathbf{r}}) \bar{Y}_n^m(\hat{\mathbf{r}}_0) \right) \\ &\quad \times \left[ \sum_{p=1}^{\infty} \sum_{q=-p}^p \frac{\tau^{p-1}}{r'^{p+1}} Y_p^q(\hat{\mathbf{r}}) Y_p^q(\hat{\mathbf{r}}) \left( A_p^q \frac{\partial \bar{Y}_p^q(\hat{\mathbf{r}})}{\partial \theta} - \frac{B_p^q}{\sin \theta} \frac{\partial \bar{Y}_p^q(\hat{\mathbf{r}})}{\partial \phi} \right) \right] dS(\hat{\mathbf{r}}), \end{aligned} \tag{138}$$

where the second expansion now begins with  $p = 1$ , since for  $p = 0$  the corresponding spherical harmonic is a constant and the derivatives present vanish. Employing (40) and (41), we find

$$A_p^q \frac{\partial \bar{Y}_p^q(\hat{\mathbf{r}})}{\partial \theta} - \frac{B_p^q}{\sin \theta} \frac{\partial \bar{Y}_p^q(\hat{\mathbf{r}})}{\partial \phi'} = \frac{1}{\sin \theta'} \left[ p j_{p+1}^q A_p^q \bar{Y}_{p+1}^q(\hat{\mathbf{r}}) + i q B_p^q \bar{Y}_p^q(\hat{\mathbf{r}}) - (p+1) j_p^q A_p^q \bar{Y}_{p-1}^q(\hat{\mathbf{r}}) \right] \quad (139)$$

so that the integral now reads

$$I_1(\mathbf{r}; \mathbf{r}_0) = 4\pi R^2 \sum_{n,m} \sum_{p,q} \frac{R^{p-1}}{r^{p+1}} \frac{n+1}{n(2n+1)} \frac{r_0^n}{R^{2n+1}} Y_p^q(\hat{\mathbf{r}}) \bar{Y}_n^m(\hat{\mathbf{r}}_0) \mathbb{J}_{n,p}^{m,q}, \quad (140)$$

where the indices  $n$  and  $p$  start with 1 and

$$\mathbb{J}_{n,p}^{m,q} = p j_{p+1}^q A_p^q \mathcal{S}_{1,n,p}^{m,q} + i q B_p^q \mathcal{S}_{2,n,p}^{m,q} - (p+1) j_p^q A_p^q \mathcal{S}_{3,n,p}^{m,q}, \quad (141)$$

while

$$\mathcal{S}_{1,n,p}^{m,q} = \oint_{S^2} \frac{1}{\sin \theta} Y_n^m(\hat{\mathbf{r}}) Y_p^q(\hat{\mathbf{r}}) \bar{Y}_{p+1}^q(\hat{\mathbf{r}}) dS(\hat{\mathbf{r}}), \quad (142)$$

$$\mathcal{S}_{2,n,p}^{m,q} = \oint_{S^2} \frac{1}{\sin \theta} Y_n^m(\hat{\mathbf{r}}) Y_p^q(\hat{\mathbf{r}}) \bar{Y}_p^q(\hat{\mathbf{r}}) dS(\hat{\mathbf{r}}), \quad (143)$$

$$\mathcal{S}_{3,n,p}^{m,q} = \oint_{S^2} \frac{1}{\sin \theta} Y_n^m(\hat{\mathbf{r}}) Y_p^q(\hat{\mathbf{r}}) \bar{Y}_{p-1}^q(\hat{\mathbf{r}}) dS(\hat{\mathbf{r}}). \quad (144)$$

With above relations, as well as (130), the second integral on the RHS of (127) in the linear regime becomes

$$\begin{aligned} & \int_D (\Delta_{\mathbf{r}_0} \Psi(\mathbf{r}_0)) I_1(\mathbf{r}; \mathbf{r}_0) d\nu(\mathbf{r}_0) = \\ & \int_D \left[ \sum_{n=0}^{\infty} \sum_{m=-n}^n \left( \ddot{\psi}_n^m(r_0) + \frac{2}{r_0} \dot{\psi}_n^m(r_0) - \frac{n(n+1)}{r_0^2} \psi_n^m(r_0) \right) Y_n^m(\hat{\mathbf{r}}_0) \right] \\ & \times 4\pi R^2 \sum_{k,l} \sum_{p,q} \frac{R^{p-1}}{r^{p+1}} \frac{k+1}{k(2k+1)} \frac{r_0^k}{R^{2k+1}} Y_p^q(\hat{\mathbf{r}}) \bar{Y}_k^l(\hat{\mathbf{r}}_0) \mathbb{J}_{k,p}^{l,q} d\nu(\mathbf{r}_0). \end{aligned} \quad (145)$$

Moreover, due to orthogonality the integral corresponding to  $n = 0$  vanishes, namely

$$\oint_{S^2} Y_0^0(\hat{\mathbf{r}}_0) Y_k^l(\hat{\mathbf{r}}_0) dS(\hat{\mathbf{r}}_0) = 0, \quad \forall k \geq 1 \quad (146)$$

and we are left with

$$\int_D (\Delta_{\mathbf{r}_0} \Psi(\mathbf{r}_0)) I_1(\mathbf{r}; \mathbf{r}_0) dV(\mathbf{r}_0) = 4\pi \sum_{n,m} \sum_{p,q} \frac{p+1}{p(2p+1)} \left(\frac{R}{r}\right)^{n+1} \mathbb{J}_{p,n}^{q,m} (R\dot{\psi}_p^q(R) - p\psi_p^q(R)) Y_n^m(\hat{\mathbf{r}}). \tag{147}$$

Replacing everything back into (127) and following the procedure demonstrated, yields

$$\sum_{n,m} d_n^m Y_n^m(\hat{\mathbf{r}}) = \frac{\mu_0}{4\pi} \sum_{n,m} \left\{ -\frac{4\pi}{2n+1} \frac{R^{n+2}}{r^{n+1}} (R\dot{a}_n^m(R) - (n-1)a_n^m(R)) + \epsilon \sum_{p,q} \frac{p+1}{p(2p+1)} \left(\frac{R}{r}\right)^{n+1} \mathbb{J}_{p,n}^{q,m} \times (R\dot{\psi}_p^q(R) - p\psi_p^q(R)) \right\} Y_n^m(\hat{\mathbf{r}}). \tag{148}$$

Since  $a_n^m$  are known via (113), we deduce that the coefficients  $\psi_n^m$  are calculated as

$$d_n^m = \mu_0 \left(\frac{R}{r}\right)^{n+1} \left[ -\frac{R}{2n+1} (R\dot{a}_n^m(R) - (n-1)a_n^m(R)) + \frac{\epsilon}{4\pi} \sum_{p=1}^{\infty} \sum_{q=-p}^p \frac{p+1}{p(2p+1)} \mathbb{J}_{p,n}^{q,m} (R\dot{\psi}_p^q(R) - p\psi_p^q(R)) \right]. \tag{149}$$

The remaining two components are given by equations (119).

### 3 Example

Let us now showcase the particulars of the algorithms presented in an analytical fashion by referring to the simplest example possible, namely that of an univariate function  $f(\theta)$ . In order to facilitate calculations, but without loss of generality, let  $f(\hat{\mathbf{r}}) = \cos \theta$ . The procedure regarding complex functions  $f(\hat{\mathbf{r}})$  can be found in [33].

### 3.1 EEG

When  $f(\hat{\mathbf{r}}) = \cos \theta$ , the corresponding coefficients  $\bar{C}_n^m$  are computed from (60) by integrating over the unit ball, i.e.  $\bar{C}_n^m = \oint_{S^2} \cos \theta Y_n^m(\hat{\mathbf{r}}) dS(\hat{\mathbf{r}})$ . Forasmuch as,  $\cos \theta = \sqrt{4\pi/3} Y_1^0(\hat{\mathbf{r}})$  only coefficients with  $m = 0$  will survive, namely  $\bar{C}_n^0 = \sqrt{4\pi/(2n+1)} \delta_{1,n}$  and thus (63) simplifies to  $\mathbf{Q} \cdot \nabla_{\mathbf{r}_0} r_0 \bar{Y}_1^0(\hat{\mathbf{r}}_0) = 0$ , implying  $Q_3 = 0$ . Hence, for if  $f(\hat{\mathbf{r}}) = \cos \theta$ , we cannot retrieve the  $z$ -coordinate of the moment for the first correction.

Expanding the dipoles position and moment in powers of  $\epsilon$ , we obtain the following approximations (keeping only linear corrections)

$$\mathbf{r}_0 = \mathbf{r}_{0,0} + \epsilon \mathbf{r}_{0,1}, \quad (150)$$

$$\mathbf{Q} = \mathbf{Q}_0 + \epsilon \mathbf{Q}_1, \quad (151)$$

respectively. The zeroth-order corrections  $\mathbf{r}_{0,0}$  and  $\mathbf{Q}_0$  are provided via relations (58) and (59), respectively. On the other hand, the first-order corrections  $\mathbf{r}_{0,1}$  and  $\mathbf{Q}_1$  depend upon the surface deformation, namely the function  $f(\hat{\mathbf{r}})$ .

Nevertheless, relation (39) simplifies accordingly as

$$\mathbb{A}_{1,k}^\ell = -\frac{4\pi}{2k+1} \frac{N_k^\ell}{kR} \frac{(k+\ell)!}{(k-\ell)!} \sum_{n=1}^{\infty} \mathbb{A}_{0,n}^\ell \left( n j_{n+1}^\ell N_{n+1}^\ell \delta_{n+1,k} - (n+1) j_n^\ell N_{n-1}^\ell \delta_{n-1,k} \right), \quad (152)$$

where  $N_k^\ell$  denote the corresponding normalization constants and  $\delta_{n,m}$  the Kronecker symbol. Computing the latter for  $k = 1, 2$ , furnishes the following eight relations

$$\mathbb{A}_{1,1}^\ell = \frac{3j_2^\ell}{R^2} \mathbb{A}_{0,2}^\ell, \quad \ell = -1, 0, 1, \quad (153)$$

$$\mathbb{A}_{1,2}^\ell = \frac{1}{2R^3} \left( 4j_3^\ell \mathbb{A}_{0,3}^\ell - j_2^\ell \mathbb{A}_{0,1}^\ell \right), \quad \ell = -2, -1, 0, 1, 2. \quad (154)$$

However, only six of them are needed in order to identify  $\mathbf{r}_{0,1}$  and  $\mathbf{Q}_1$ , those being

$$\mathbb{A}_{1,1}^{\pm 1} = \mp \frac{3}{4\sigma R^5} \sqrt{\frac{3}{2\pi}} x_{1,3} (Q_{1,1} \mp i Q_{1,2}), \quad (155)$$

$$\mathbb{A}_{1,1}^0 = -\frac{3}{2\sigma R^5} \sqrt{\frac{1}{3\pi}} (Q_{1,1} x_{1,1} + Q_{1,2} x_{1,2}), \quad (156)$$

$$\mathbb{A}_{1,2}^{\pm 2} = \mp \frac{8\sqrt{5}}{9R^2} \mathbb{B}_1^{\pm 1} (x_{1,1} \mp i x_{1,2}), \quad (157)$$

$$\mathbb{A}_{1,2}^0 = \frac{2}{R^2} \mathbb{B}_1^0 x_{1,3}. \quad (158)$$



Solving above system (155)–(158), yields

$$\mathbf{r}_{0,1} = \frac{R^2}{2} \left( \frac{9}{8\sqrt{5}} \left( \frac{\mathbb{A}_{1,2}^{-2}}{\mathbb{A}_{1,1}^{-1}} - \frac{\mathbb{A}_{1,2}^2}{\mathbb{A}_{1,1}^1} \right) - i \frac{9}{8\sqrt{5}} \left( \frac{\mathbb{A}_{1,2}^{-2}}{\mathbb{A}_{1,1}^{-1}} + \frac{\mathbb{A}_{1,2}^2}{\mathbb{A}_{1,1}^1} \right) \frac{\mathbb{A}_{1,2}^0}{\mathbb{A}_{1,1}^0} \right)^\top, \quad (159)$$

$$\mathbf{Q}_1 = \frac{4\sigma R^3}{3} \sqrt{\frac{2\pi}{3}} \frac{\mathbb{A}_{1,1}^0}{\mathbb{A}_{1,2}^0} \left( \mathbb{A}_{1,1}^{-1} - \mathbb{A}_{1,1}^1 - i(\mathbb{A}_{1,1}^{-1} + \mathbb{A}_{1,1}^1) \mathbf{0} \right)^\top. \quad (160)$$

### 3.2 MEG

Substituting  $f(\hat{\mathbf{r}}) = \cos \theta$  into (78) we find

$$\mathbb{E}_{1,k}^\ell = -i4\pi\sigma R^{k+1} \frac{k}{(k+1)(2k+1)} \mathbb{B}_k^\ell. \quad (161)$$

As seen, forming a  $3 \times 3$  system from six equations for  $k = 1, 2$  and utilizing the corresponding Cartesian forms, furnishes

$$x_{1,1} = i \frac{\sqrt{5}\pi}{6} \sigma R^2 \left[ \frac{\mathbb{B}_1^{-1} \mathbb{E}_{1,2}^{-2}}{(\mathbb{E}_{1,1}^{-1})^2} + \frac{\mathbb{B}_1^1 \mathbb{E}_{1,2}^2}{(\mathbb{E}_{1,1}^1)^2} - \frac{4R}{5} \left( \frac{\mathbb{B}_2^2}{\mathbb{E}_{1,1}^1} + \frac{\mathbb{B}_2^{-2}}{\mathbb{E}_{1,1}^{-1}} \right) \right], \quad (162)$$

$$x_{1,2} = \frac{\sqrt{5}\pi}{6} \sigma R^2 \left[ \frac{\mathbb{B}_1^{-1} \mathbb{E}_{1,2}^{-2}}{(\mathbb{E}_{1,1}^{-1})^2} - \frac{\mathbb{B}_1^1 \mathbb{E}_{1,2}^2}{(\mathbb{E}_{1,1}^1)^2} + \frac{4R}{5} \left( \frac{\mathbb{B}_2^2}{\mathbb{E}_{1,1}^1} - \frac{\mathbb{B}_2^{-2}}{\mathbb{E}_{1,1}^{-1}} \right) \right], \quad (163)$$

$$x_{1,3} = i \frac{\sqrt{5}\pi}{3\sqrt{2}} \sigma R^2 \frac{1}{\mathbb{E}_{1,1}^0} \left[ \frac{\mathbb{B}_1^{-1} \mathbb{E}_{1,2}^{-2}}{\mathbb{E}_{1,1}^{-1}} - \frac{\mathbb{B}_1^1 \mathbb{E}_{1,2}^2}{\mathbb{E}_{1,1}^1} + \frac{2R}{5} \left( \mathbb{B}_2^2 - \mathbb{B}_2^{-2} \right) \right]. \quad (164)$$

The moment equals

$$Q_{r_0} = 0, \quad (165)$$

$$Q_{\theta_0} = -i \sqrt{\frac{2\pi}{3}} \sigma R^2 \frac{\mathbb{B}_1^1 - \mathbb{B}_1^{-1}}{r_0 \sin \theta_0}, \quad (166)$$

$$Q_{\phi_0} = -\sqrt{\frac{3}{\pi}} \frac{\mathbb{E}_{1,1}^0}{r_0 \sin \theta_0}. \quad (167)$$

## 4 Conclusions and Discussion

We presented a short review on the application of boundary perturbations for the forward and inverse problems in EEG and MEG. In all four cases, the surface of reference has been the homogeneous sphere, leading to solutions with corrections incorporating the deviation from the spherical geometry, included up to the first correction. Higher order corrections are feasible, but mathematically tedious to obtain and seldom provide more accuracy. A major impediment are the constrains, which are brought in by the Neumann condition (30), demonstrated in the example (Section 3), leading to the fact that components of either the moment of the dipole or its location or both, remain concealed. As a rule of thumb, complex functions introduce complex constrains. However, since the Neumann condition, rising from Green's theorem, cannot be avoided in the framework of the specific BVP, it is an integral part of the presented analysis. Nevertheless, one can still gain valuable insight into the problem at hand, namely if and how strong surface deformations affect measurements.

Advanced analytical or semi-analytical solutions and formulae attaining closed-type forms in EEG and MEG have quite important advantages compared with the pure numerical methods. Indeed, the validity of numerical solutions can be verified by such techniques, nevertheless the basic outcome gained is the facility in incorporating with the inverse problem, knowing the mathematical tools of the forward one itself. On the other hand, bearing in mind that very important physical laws can be derived from analytical methods, we can understand the necessity of tackling with a confident mathematical basis before getting involved with an algorithmic procedure. Therefore, even nowadays, there is always room for such kind of methods that coexist with pure numerical codes and aim to the solution of boundary value problems in physical applications of major importance.

Mathematical and computational work is currently in progress and involves research into several directions, such as the introduction of more complicated geometries for representing the head's shape or the accomplishment of actually difficult inversion algorithms, taking profit from the proposed framework.

## References

1. R.J. Ilmoniemi, R.J. Näätänen, Magnetoencephalography, in *International Encyclopedia of the Social & Behavioral Sciences*, ed. by N.J. Smelser, P.B. Baltes (2001), pp. 9131–9137
2. I.M. Turan, Qualitative and quantitative EEG findings in Schizophrenia. *Schizophrenia Bull.* **3**(1), 61–79 (1977)
3. S. Noachtar, J. Rémi, The role of EEG in epilepsy: a critical review. *Epilepsy Behav.* **15**(1), 22–33 (2009)
4. A.M. Beres, Time is of the essence: a review of electroencephalography (EEG) and event-related brain potentials (ERPs) in language research. *Appl. Psychophysiol. Biofeedback* **42**(4), 247–255 (2017)

5. S. Baillet, Magnetoencephalography for brain electrophysiology and imaging. *Nat. Neurosci.* **20**, 327–339 (2017)
6. T. Kirschstein, R. Köhling, What is the Source of the EEG? *Clin. EEG Neurosci.* **40**(3), 146–149 (2009)
7. M. Hämmäinen, R. Hari, R. Ilmoniemi, J. Knuutila, O. Lounasmaa, Magnetoencephalography: theory, instrumentation and applications to the noninvasive study of human brain function. *Rev. Mod. Phys.* **65**, 413–497 (1993)
8. S. Baillet, Forward and Inverse Problems of MEG/EEG, in *Encyclopedia of Computational Neuroscience*, ed. by D. Jaeger, R. Jung (Springer, New York, 2014)
9. F.N. Wilson, R.H. Bayley, The electric field of an eccentric dipole in a homogeneous spherical conducting medium. *Circulation* **1**(1), 84–92 (1950)
10. H. Hallez, B. Vanrumste, R. Grech, J. Muscat, W. De Clercq, A. Vergult, Y. D’Asseler, K.P. Camilleri, S.G. Fabri, S. Van Huffel, I. Lemahieu, Review on solving the forward problem in EEG source analysis. *J. NeuroEng. Rehabil.* **4**, 46 (2007)
11. R. Grech, T. Cassar, J. Muscat, K.P. Camilleri, S.G. Fabri, M. Zervakis, P. Xanthopoulos, V. Sakkalis, B. Vanrumste, Review on solving the inverse problem in EEG source analysis. *J. NeuroEng. Rehabil.* **5**, 25 (2008)
12. D. Cohen, Magnetoencephalography: evidence of magnetic fields produced by alpha-rhythm currents. *Science* **161**(3843), 784–786 (1968)
13. R. Plonsey, D.B. Heppner, Considerations of quasi-stationarity in electrophysiological systems. *Bull. Math. Biophys.* **29**, 657–664 (1967)
14. J. Sarvas, Basic mathematical and electromagnetic concepts of the biomagnetic inverse problems. *Phys. Med. Biol.* **32**, 11–22 (1987)
15. R. Albanese, P.B. Monk, The inverse source problem for Maxwell’s equations. *Inverse Problems* **22**(3), 1023–1035 (2006)
16. G. Dassios, A.S. Fokas, Electro-magneto-encephalography for a three-shell model: dipoles and beyond for the spherical geometry. *Inverse Problems* **25**(3), 035001 (2009)
17. G. Dassios, G. Fragoyiannis, K. Satrazemi, On the inverse EEG problem for a 1D current distribution. *J. Appl. Math.* **2014**, 715785 (2014)
18. G. Dassios, K. Satrazemi, Inversion of electroencephalography data for a 2-D current distribution. *Math. Methods Appl. Sci.* **38**(6), 1098–1105 (2014)
19. G. Dassios, A.S. Fokas, D. Hadjiloizi, On the complementarity of electroencephalography and magnetoencephalography. *Inverse Problems* **23**, 2541 (2007)
20. G. Dassios, A.S. Fokas, The definite non-uniqueness results for deterministic EEG and MEG data. *Inverse Problems* **29**(6), 065012 (2013)
21. D.A. Varshalovich, A.N. Moskalev, V.K. Khersonskii, *Quantum Theory of Angular Momentum: Irreducible Tensors, Spherical Harmonics, Vector Coupling Coefficients, 3nj Symbols* (1. repr. ed.) (World Scientific, Singapore, 1988)
22. F.W.J. Olver, D.W. Lozier, R.F. Boisvert, C.W. Clark (eds.), *NIST Handbook of Mathematical Functions* (Cambridge University Press, Cambridge, 2010)
23. R. Bellman, *Perturbation Techniques in Mathematics, Engineering and Physics* (Holt, Rinehart and Winston, New York, 1966)
24. N.G. De Bruijn, *Asymptotic Methods in Analysis* (North-Holland Publishing, Amsterdam, 1958)
25. M. Doschoris, F. Kariotou, Mathematical foundation of electroencephalography, in *Electroencephalography*, ed. by P. Sittiprapaporn (IntechOpen, Rijeka, 2017)
26. D.B. Geselowitz, On the magnetic field generated outside an inhomogeneous volume conductor by internal current sources. *IEEE Trans. Mag.* **6**, 346–347 (1970)
27. Z.Z. Acar, S. Makeig, Effects of forward model errors on EEG source localization. *Brain Topogr.* **26**, 378–396 (2013)
28. G. Dassios, Electric and magnetic activity of the brain in spherical and ellipsoidal geometry, in *Mathematical Modeling in Biomedical Imaging I. Lecture Notes in Mathematics*, ed. by H. Ammari, vol. 1983 (Springer, Berlin, 2009)

29. D.B. Geselowitz, On bioelectric potentials in an inhomogeneous volume conductor. *Biophys. J.* **7**(1), 1–11 (1967)
30. M. Doschoris, P. Vafeas, G. Fragoyiannis, The influence of surface deformations on the forward magnetoencephalographic problem. *SIAM J. Appl. Math.* **78**(2), 963–976 (2018)
31. G. Dassios, M. Doschoris and K. Satrazemi, On the resolution of synchronous dipolar excitations via MEG measurements. *Q. Appl. Math.* **76**, 39–45 (2018)
32. A.S. Fokas, Electro-magneto-encephalography for a three-shell model: distributed current in arbitrary, spherical and ellipsoidal geometries. *J. R. Soc. Interface* **6**(34), 479–488 (2009)
33. G. Dassios, M. Doschoris, G. Fragoyannis, Sensitivity analysis of the forward EEG problem depending on head shape variations. *Math. Probl. Eng.* **2015**, Article ID 612528 (2015)

# Poynting–Robertson and Oblateness Effects on the Equilibrium Points of the Perturbed R3BP: Application on Cen X-4 Binary System



Aguda Ekele Vincent and Angela E. Perdiou

**Abstract** We examine the dynamical effects of Poynting–Robertson (P–R) drag and oblateness together with small perturbations in the Coriolis and centrifugal forces on the existence, location and stability of equilibrium points in the photogravitational restricted three-body problem. It is found that under constant P–R drag effect, collinear equilibrium points cease to exist numerically and of course analytically. The problem admits five non-collinear equilibrium points and it is found that the positions of these points depend on all the system parameters except small perturbation in the Coriolis force. Finally, we justify the relevance of the model in astronomy by applying it to Cen X-4 binary system, for which all the equilibrium points have been seen to be unstable.

**MSC** 70F07, 70F15, 70M20, 70K42

## 1 Introduction

The restricted three-body problem (R3BP) consists of two finite bodies, known as primaries which rotate in circular orbits around their common center of mass and a massless body which moves in the plane of motion of the primaries under their gravitational attraction and does not affect their motion. The study of the R3BP is still an active field of research because of its applications in dynamics of the solar and stellar systems, artificial satellites and lunar theory. The circular restricted three-body problem (CR3BP) has been the well known studied problem in celestial

---

A. E. Vincent

Department of Mathematics, School of Basic Sciences, Nigeria Maritime University,  
Okerenkoko, Delta State, Nigeria  
e-mail: [vincentekele@yahoo.com](mailto:vincentekele@yahoo.com)

A. E. Perdiou (✉)

Department of Civil Engineering, School of Engineering, University of Patras, Patras, Greece  
e-mail: [aperdiou@upatras.gr](mailto:aperdiou@upatras.gr)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,  
Springer Optimization and Its Applications 173,  
[https://doi.org/10.1007/978-3-030-72563-1\\_7](https://doi.org/10.1007/978-3-030-72563-1_7)

mechanics. In this problem there are five equilibrium points; three of them lie on the  $x$ -axis and are called collinear while the other two are away from the  $x$ -axis and are called triangular equilibrium points. The three collinear points are generally unstable while the triangular points are generally stable for the mass ratio  $\mu \leq 0.03850\dots$  [30]. These equilibrium points are extensively used in space mission (see [1, 3, 12, 31] and references therein).

In celestial mechanics, many scientists and astronomers over the years have made modifications to the classical CR3BP (e.g. [6, 8, 13, 17–19, 22, 23, 29, 32–34]). Some of the modifications made, include the consideration of one or both primaries as oblate spheroids and/or radiation sources with small change in Coriolis and centrifugal forces and/or under the effect of different kinds of dissipation (Stokes and/or Poynting–Robertson drags). The studying of these issues enable us to get real and accurate data about the dynamical features of the system. For example, Oberti and Vienne [15] showed that the addition of oblateness effects leads to improved approximations of real orbits of certain satellites in the Solar System. Singh [28] examined out-of-plane equilibria by considering effect of a small change in Coriolis and centrifugal forces, when the primaries are both radiating and oblate spheroids. Chernikov [4] studied the existence and stability of equilibrium points under the influence of radiation and Poynting–Robertson drag. He found that six equilibrium points exist at most and pointed out that the collinear points are not positioned on the axis connecting the primaries any more while the triangular points are not symmetrical with respect to this axis. It was found that the triangular points are unstable for P–R effect. Schuerman [21] studied the triangular points of the problem and found that the points are unstable due to P–R effect. Furthermore, Ragos and Zafropoulos [20] extended the problem to the case that both main bodies are radiation sources and studied the existence and stability of the equilibrium points. The P–R effect renders unstable those equilibrium points which are conditionally stable in the classical case. Murray [14] discussed the dynamical effect of different kinds of dissipation (nebular drag, gas drag, and P–R drag) in the circular restricted three body problem and found the collinear points are not positioned on the axis joining the two masses while the displaced triangular points  $L_4$  and  $L_5$  are asymptotically stable for certain classes of drag forces.

Kushvah [11] studied numerically the existence of equilibrium points of the perturbed R3BP, where the bigger and smaller primaries are considered radiation sources and oblate spheroids, respectively, and discussed the P–R effect which is caused due to the radiation pressure. They observed that the collinear points deviate from the axis joining the two primaries, while the triangular points are not symmetrical due to radiation pressure. The P–R effect ruins the stability of equilibrium points known to be conditionally stable in the gravitational case. When the primaries are radiation sources, Singh and Aminu [24] investigated the influences of small perturbations in the Coriolis and centrifugal forces together with P–R drag from both primaries on the triangular points. They found that the positions of these points are affected from the radiation pressure, P–R drag and small perturbation in the centrifugal force. They also discovered that these perturbing forces do not influence the nature of the stability of the points in the presence of P–R

drag as they remain unstable for the binary systems Luyten 726–8 and Kruger 60.1. In the same vein, Singh and Amuda [25] studied the triangular equilibrium points when the effect of radiation pressure from the smaller primary and its Poynting–Robertson (P–R) drag are taken into account and the bigger primary as an oblate spheroid. They found numerically that the equilibrium points of the binary RXJ 0450.1–5856 are unstable. Later, Singh and Amuda [26] investigated the three dimensional case of the problem studied in [25] and they pointed out that the out-of-plane equilibria of the binary Cen X-4 system are unstable. By taking into consideration the P–R effect and Stellar wind drag, Chakraborty and Narayan [5] investigated the photogravitational elliptic restricted three-body problem and found that the equilibrium points are unstable due to the effect of the drag. Recently, Kalantonis et al. [10] studied the stability of the triangular equilibrium points in the elliptic R3BP with radiation and oblateness and showed that the positions of the triangular equilibrium points are given by an analytical formulae in which the parameters of the problem are only involved.

In this work, we aim to make an extension to the work of Singh and Amuda [25] by also taking small perturbations in the Coriolis and centrifugal forces and continue to study numerically the existence and location of the equilibrium points. As an application in this study, we consider the Cen X-4 binary system. The paper is organized as follows: In Section 2, the dynamical equations that involve the parameters of the infinitesimal particle in the binary system under consideration are obtained. In Section 3, we determine the existence and locations of the equilibrium points numerically and verify them graphically for values of the parameters of the problem, while their linear stability is analyzed in Section 4. A numerical application of these results is given in Section 5 while Section 6 summarizes the discussion and conclusion of our study.

## 2 Equations of Motion

The dynamical system consists of two bodies (known as the primaries) which move on circular orbits. We consider a barycentric coordinate system  $Oxyz$  rotating relative to an inertial reference system with angular velocity  $\omega$  about a common  $z$ -axis. The two finite bodies  $P_1$  (bigger primary) and  $P_2$  (smaller primary) have masses  $m_1 = 1 - \mu$  and  $m_2 = \mu$  ( $0 < \mu \leq 1/2$ ), respectively, with  $\mu$  being the mass ratio parameter while the test particle  $P$  is considered to have a mass  $m$ , which is significantly smaller than the masses of the primaries and therefore it does not affect their motion. Also, the bigger primary body is considered to be an oblate spheroid while the smaller one is a source of radiation with its P–R drag. The equations of motion of the test particle in the three-dimensional restricted three-body problem with the origin resting at the center of mass, in a barycentric rotating coordinate system take the form [25]:

$$\begin{aligned}
\ddot{x} - 2n\dot{y} &= n^2x - \frac{(1-\mu)(x+\mu)}{r_1^3} - \frac{\mu(x+\mu-1)q_2}{r_2^3} - \frac{3(1-\mu)(x+\mu)A_1}{2r_1^5} - \\
&\quad \frac{W_2}{r_2^2} \left\{ \frac{x+\mu-1}{r_2^2} [(x+\mu-1)\dot{x} + y\dot{y} + z\dot{z}] + \dot{x} - ny \right\}, \\
\ddot{y} + 2n\dot{x} &= n^2y - \frac{(1-\mu)y}{r_1^3} - \frac{\mu q_2 y}{r_2^3} - \frac{3(1-\mu)A_1 y}{2r_1^5} - \\
&\quad \frac{W_2}{r_2^2} \left\{ \frac{y}{r_2^2} [(x+\mu-1)\dot{x} + y\dot{y} + z\dot{z}] + \dot{y} + n(x+\mu-1) \right\}, \\
\ddot{z} &= -\frac{(1-\mu)z}{r_1^3} - \frac{\mu q_2 z}{r_2^3} - \frac{3A_1(1-\mu)z}{2r_1^5} - \\
&\quad \frac{W_2}{r_2^2} \left\{ \frac{z}{r_2^2} [(x+\mu-1)\dot{x} + y\dot{y} + z\dot{z}] + \dot{z} \right\},
\end{aligned} \tag{1}$$

with

$$r_1^2 = (x+\mu)^2 + y^2 + z^2, \quad r_2^2 = (x+\mu-1)^2 + y^2 + z^2, \quad W_2 = \frac{\mu(1-q_2)}{c_d}, \tag{2}$$

where  $r_i$ ,  $i = 1, 2$  are the distances of the test particle from the bigger and smaller primaries, respectively,  $q_2 \in (0, 1]$ ,  $W_2 \ll 1$  stand for radiation pressure and P-R drag of the smaller body, respectively,  $c_d$  is the dimensional velocity of light which depends on the physical masses of the two bodies and the distance between them, chosen to the value  $c_d = 299792458$  (see [25]) while the dots denote differentiation with respect to time  $t$ . Also,  $A_1$  is the oblateness coefficient of the bigger primary body defined by the formula  $A_1 = (A_E^2 - A_P^2)/5R^2 \ll 1$  where  $A_E$  and  $A_P$  are the equatorial and polar radii of the said primary body, respectively, and  $R$  is the distance between the primaries. On account of the oblateness of the primary body  $m_1$ , the mean perturbed motion  $n$  is defined by  $n^2 = 1 + \frac{3}{2}A_1$ . Additionally, perturbations on the Coriolis and centrifugal forces are included with the help of the parameters  $\alpha$  and  $\beta$ , respectively, such that  $\alpha = 1 + \varepsilon_1$ ,  $\beta = 1 + \varepsilon_2$ ,  $|\varepsilon_i| \ll 1$ ,  $i = 1, 2$ . The unperturbed value of each is taken as unity. Restricting ourselves to the plane  $Oxy$  and following the work of Singh and Aminu [24], the pertinent equations of motion (1) are finally written in the form:



**Table 1** Numerical data for the binary Cen X-4 system

| Binary system | Mass ( $M_{\odot}$ ) |        | Radiation pressure | Binary separation | Dimensionless<br>speed of light | Mass ratio |
|---------------|----------------------|--------|--------------------|-------------------|---------------------------------|------------|
|               | $m_1$                | $m_2$  | $q_2$              | $a$               | $c_d$                           | $\mu$      |
| Cen X-4       | 1.9996               | 0.0801 | 0.993              | 4.31              | 988.323                         | 0.038515   |

$$\begin{aligned}
 \ddot{x} - 2n\alpha\dot{y} &= n^2\beta x - \frac{(1-\mu)(x+\mu)}{r_1^3} - \frac{\mu(x+\mu-1)q_2}{r_2^3} - \frac{3(1-\mu)(x+\mu)A_1}{2r_1^5} - \\
 &\quad \frac{W_2}{r_2^2} \left\{ \frac{x+\mu-1}{r_2} [(x+\mu-1)\dot{x} + y\dot{y}] + \dot{x} - ny \right\}, \\
 \ddot{y} + 2n\alpha\dot{x} &= n^2\beta y - \frac{(1-\mu)y}{r_1^3} - \frac{\mu q_2 y}{r_2^3} - \frac{3(1-\mu)A_1 y}{2r_1^5} - \\
 &\quad \frac{W_2}{r_2^2} \left\{ \frac{y}{r_2} [(x+\mu-1)\dot{x} + y\dot{y}] + \dot{y} + n(x+\mu-1) \right\},
 \end{aligned}
 \tag{3}$$

while now

$$r_1^2 = (x + \mu)^2 + y^2, \quad r_2^2 = (x + \mu - 1)^2 + y^2. \tag{4}$$

The physical parameters of the binary Cen X-4 system are shown in Table 1 (see [2, 25, 26]).

### 3 Existence and Positions of Equilibrium Points

The equilibrium (or Lagrangian) points are obtained when the acceleration ( $\ddot{x}, \ddot{y}$ ) and velocity ( $\dot{x}, \dot{y}$ ) components of the test particle are zero. So, we obtain the coordinates  $(x_0, y_0)$  of equilibrium points as solutions of the equations:

$$\begin{aligned}
 n^2\beta x - \frac{(1-\mu)(x+\mu)}{r_1^3} - \frac{\mu(x+\mu-1)q_2}{r_2^3} - \frac{3(1-\mu)(x+\mu)A_1}{2r_1^5} + \frac{W_2ny}{r_2^2} &= 0, \\
 n^2\beta y - \frac{(1-\mu)y}{r_1^3} - \frac{\mu y q_2}{r_2^3} - \frac{3y(1-\mu)A_1}{2r_1^5} - \frac{W_2n(x+\mu-1)}{r_2^2} &= 0.
 \end{aligned}
 \tag{5}$$

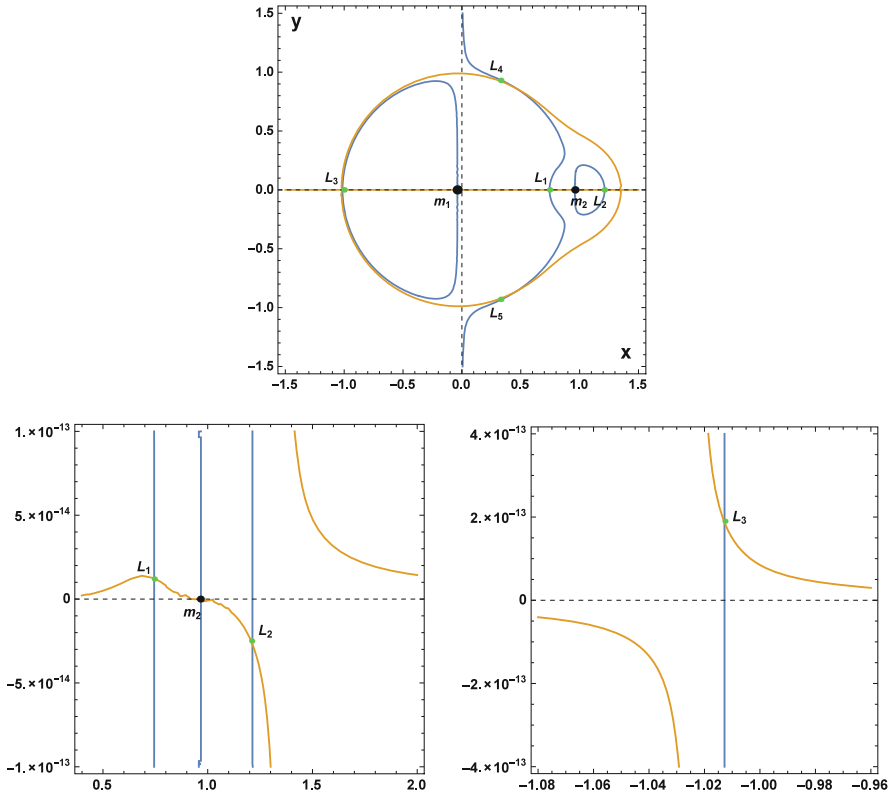
It is interesting to note that for  $A_1 = 0, q_2 = \beta = 1$ , the classical case of the R3BP is recovered while the case  $\beta = 1$  leads to the equations of motion presented in [25]. It is well known that in the classical R3BP there are two types of equilibria or solutions, depending on whether  $y = 0$  or  $y \neq 0$ . Points for which  $y = 0$  are called collinear equilibrium points and they lie on the line connecting the primaries, the  $x$ -axis of the synodic system, while points for which  $y \neq 0$  are called triangular

(non-collinear) equilibrium points and they lie away from the  $x$ -axis of the synodic system.

In the perturbed R3BP where the radiation pressure coupled with P–R drag terms appear, the existence of collinear equilibria as well as the total number of the equilibrium points depend on the particular values of the radiation pressure (see, for example, [4, 14]). Ragos and Zafiroopoulos [20] have shown numerically that in the photogravitational CR3BP including the P–R effect there are at most five equilibrium points (with no collinear points), depending on the values of radiation factors  $q_1$  and  $q_2$ . Following the lead of above paper, we resort to a numerical study in this case of the problem since the system of Equations (5) which provides the  $(x_0, y_0)$  coordinates of the points of equilibrium cannot be solved analytically. In this premise, the equilibrium points are obtained by solving Equations (5) simultaneously using any well-known iterative method for finding roots of non-linear algebraic equations. The aforementioned method has been successfully applied in [16, 27] and [7] (see also references therein) for the determination of equilibrium points in a different model problem of Celestial Mechanics. We observe that our problem admits five non-collinear equilibrium points,  $L_i, i = 1, 2, \dots, 5$ , which positions are independent of the Coriolis force but dependent upon the centrifugal force and the remaining involved parameters.

Generally, to obtain the positions for the collinear equilibrium points we solve Equations (5) for  $y = 0$  but due to the existence of the dissipative term defined by the P–R drag, it is obvious that collinear equilibrium solution does not exist anymore. This is also easy to show geometrically by plotting the contours of the two implicit functions presented in system (5) (see Figure 1). We observe from this figure that the  $y$  components of the equilibrium points  $L_{1,2,3}$  are close to zero but not zero. Moreover, this can be easily seen from bottom-left and right frames in Figure 1 where we enlarge the area close to  $L_{1,2}$  and  $L_3$ , respectively. Therefore, we can conclude that under the effect of P–R drag, induced by the radiation pressure of the smaller primary, there are no equilibrium points that lie exactly on the  $x$ -axis, called collinear equilibrium points. This result agrees with [14, 20] and [11].

So, for the non-collinear equilibrium points, the second Equation (5) holds and the equilibria are obtained by solving both Equations (5) simultaneously. Figure 1 depicts the five non-collinear equilibrium points,  $L_i, i = 1, 2, \dots, 5$  of the problem in the  $xy$ -plane, along with the associated primaries, which have been found by solving numerically the aforementioned system for assumed values of  $\mu = 0.03852$ ,  $\beta = 1.01$ ,  $A_1 = 0.0005$ ,  $q_2 = 0.9999$  and  $c_d = 299792458$ . We denote here that the equilibria in the  $xy$ -plane are given by the mutual intersections of the two coloured curves where blue and brown lines in the figure correspond to the first and second equation of (5), respectively. Here we also note that the intersection points of these curves show the coordinates  $(x_0, y_0)$  of the equilibria on the  $xy$ -plane. It is seen that under the combined effects of the parameters, there exist five non-collinear equilibrium points for which the ordinates of  $L_1$ ,  $L_2$  and  $L_3$  are close to zero but not zero. Therefore, from Figure 1, it is observed that under the combined effects of radiating smaller primary with it P–R drag, and oblateness of



**Fig. 1** The five non-collinear equilibrium points and the position of the primary bodies for  $\mu = 0.03852$ ,  $\beta = 1.01$ ,  $A_1 = 0.0005$ ,  $q_2 = 0.9999$  and  $c_d = 299792458$ . Bottom frames depict zoomed images of  $L_{1,2}$  and  $L_3$ , respectively, with intersections of the curves. Black dots indicate the positions of the bodies  $m_i$ ,  $i = 1, 2$  while the positions of the equilibrium points  $L_i$ ,  $i = 1, 2, \dots, 5$  are denoted by green dots

the bigger primary, the equilibria positions are different from those of the classical R3BP. All these results tally with [20].

### 4 Stability of the Non-collinear Equilibrium Points

To study analytically the solutions in the neighborhood of the non-collinear equilibrium points  $L_i, i = 1, 2, \dots, 5$ , following Ragos and Zafropoulos [20] as well as Singh and Amuda [25], we consider small displacements  $\xi$  and  $\eta$  given to the coordinates of an equilibrium point  $(x_0, y_0)$  such that  $\xi = x - x_0$ ,  $\eta = y - y_0$  and denote the right-hand side of equations of motion (3) by  $\Omega_x = \partial\Omega/\partial x$  and  $\Omega_y = \partial\Omega/\partial y$ , respectively. Then the variational form of the equations of motion is

derived as:

$$\begin{aligned}\ddot{\xi} - 2n\alpha\dot{\eta} &= \Omega_{x\dot{x}}^{(0)}\dot{\xi} + \Omega_{x\dot{y}}^{(0)}\dot{\eta} + \Omega_{xx}^{(0)}\xi + \Omega_{xy}^{(0)}\eta, \\ \ddot{\eta} + 2n\alpha\dot{\xi} &= \Omega_{y\dot{x}}^{(0)}\dot{\xi} + \Omega_{y\dot{y}}^{(0)}\dot{\eta} + \Omega_{yx}^{(0)}\xi + \Omega_{yy}^{(0)}\eta,\end{aligned}\quad (6)$$

where the dots are the derivatives with respect to time  $t$  and only the linear terms in  $\xi$  and  $\eta$  have been taken. Now, we assume solutions of the variational equations of the form:

$$\xi = B_1 e^{\lambda t}, \quad \eta = B_2 e^{\lambda t}, \quad (7)$$

where  $B_i$ ,  $i = 1, 2$ , are arbitrary constants and  $\lambda$  is a parameter. Substituting Equations (7) in Equations (6) and simplifying, we obtain:

$$\begin{aligned}(\lambda^2 - \lambda\Omega_{x\dot{x}}^{(0)} - \Omega_{xx}^{(0)})B_1 + (-2n\alpha\lambda - \lambda\Omega_{x\dot{y}}^{(0)} - \Omega_{xy}^{(0)})B_2 &= 0, \\ (2n\alpha\lambda - \lambda\Omega_{y\dot{x}}^{(0)} - \Omega_{yx}^{(0)})B_1 + (\lambda^2 - \lambda\Omega_{y\dot{y}}^{(0)} - \Omega_{yy}^{(0)})B_2 &= 0.\end{aligned}\quad (8)$$

Now, for the nontrivial solution the determinant of the coefficients matrix of the above system must be zero, namely:

$$\begin{vmatrix} \lambda^2 - \lambda\Omega_{x\dot{x}}^{(0)} - \Omega_{xx}^{(0)} & -2n\alpha\lambda - \lambda\Omega_{x\dot{y}}^{(0)} - \Omega_{xy}^{(0)} \\ 2n\alpha\lambda - \lambda\Omega_{y\dot{x}}^{(0)} - \Omega_{yx}^{(0)} & \lambda^2 - \lambda\Omega_{y\dot{y}}^{(0)} - \Omega_{yy}^{(0)} \end{vmatrix} = 0. \quad (9)$$

Simplifying Equation (9) we obtain the characteristic polynomial corresponding to the system (6) as:

$$\lambda^4 + a\lambda^3 + b\lambda^2 + c\lambda + d = 0, \quad (10)$$

with

$$\begin{aligned}a &= -(\Omega_{y\dot{y}}^{(0)} + \Omega_{x\dot{x}}^{(0)}), \\ b &= 4n^2\alpha^2 + \Omega_{x\dot{x}}^{(0)}\Omega_{y\dot{y}}^{(0)} - \Omega_{xx}^{(0)} - \Omega_{yy}^{(0)} - [\Omega_{x\dot{y}}^{(0)}]^2, \\ c &= \Omega_{x\dot{x}}^{(0)}\Omega_{y\dot{y}}^{(0)} + \Omega_{xx}^{(0)}\Omega_{y\dot{y}}^{(0)} + 2n\alpha\Omega_{xy}^{(0)} - 2n\alpha\Omega_{yx}^{(0)} - \Omega_{y\dot{x}}^{(0)}\Omega_{xy}^{(0)} - \Omega_{yx}^{(0)}\Omega_{x\dot{y}}^{(0)}, \\ d &= \Omega_{xx}^{(0)}\Omega_{yy}^{(0)} - \Omega_{yx}^{(0)}\Omega_{xy}^{(0)},\end{aligned}\quad (11)$$

and the obtained eigenvalues determine the stability or instability of the respective equilibrium point. The second order partial derivative of  $\Omega$  are denoted by subscripts while the superscript "0" means that the corresponding derivatives have been evaluated at the equilibrium points  $(x_0, y_0)$  and are given by the following analytical formulas:

$$\begin{aligned} \Omega_{xx}^{(0)} = & n^2\beta - \frac{(1-\mu)}{r_{10}^3} - \frac{q_2\mu}{r_{20}^3} - \frac{3A_1(1-\mu)}{2r_{10}^5} + \frac{3(1-\mu)(x_0+\mu)^2}{r_{10}^5} \\ & + \frac{3q_2\mu(x_0+\mu-1)^2}{r_{20}^5} + \\ & \frac{15A_1(1-\mu)(x_0+\mu)^2}{2r_{10}^7} - \frac{2nW_2y_0(x_0+\mu-1)}{r_{20}^4}, \end{aligned} \quad (12)$$

$$\begin{aligned} \Omega_{yy}^{(0)} = & n^2\beta - \frac{(1-\mu)}{r_{10}^3} - \frac{q_2\mu}{r_{20}^3} + \frac{3(1-\mu)y_0^2}{r_{10}^5} - \frac{3(1-\mu)A_1}{2r_{10}^5} + \frac{3q_2\mu y_0^2}{r_{20}^5} + \\ & \frac{15A_1(1-\mu)y_0^2}{2r_{10}^7} + \frac{2nW_2y_0(x_0+\mu-1)}{r_{20}^4}, \end{aligned} \quad (13)$$

$$\begin{aligned} \Omega_{xy}^{(0)} = & \frac{nW_2}{r_{20}^2} - \frac{2nW_2y_0^2}{r_{20}^4} + \frac{3(1-\mu)(x_0+\mu)y_0}{r_{10}^5} + \frac{3q_2\mu(x_0+\mu-1)y_0}{r_{20}^5} + \\ & \frac{15A_1(1-\mu)(x_0+\mu)y_0}{2r_{10}^7}, \end{aligned} \quad (14)$$

$$\begin{aligned} \Omega_{yx}^{(0)} = & -\frac{nW_2}{r_{20}^2} + \frac{2nW_2(x_0+\mu-1)^2}{r_{20}^4} + \frac{3(1-\mu)(x_0+\mu)y_0}{r_{10}^5} + \\ & \frac{3q_2\mu(x_0+\mu-1)y_0}{r_{20}^5} + \frac{15A_1(1-\mu)(x_0+\mu)y_0}{2r_{10}^7}, \end{aligned} \quad (15)$$

$$\Omega_{x\dot{x}}^{(0)} = -\frac{W_2}{r_{20}^2}\left(1 + \frac{1}{r_{20}}\right) + \frac{W_2x_0}{r_{20}^4}(2-x_0) + \frac{W_2\mu}{r_{20}^4}(2(1-x_0)-\mu), \quad (16)$$

$$\Omega_{y\dot{y}}^{(0)} = -\frac{W_2}{r_{20}^2}\left(1 + \frac{y_0^2}{r_{20}^2}\right), \quad \Omega_{x\dot{y}}^{(0)} = \frac{W_2y_0}{r_{20}^4}(1-(x_0+\mu)) = \Omega_{y\dot{x}}^{(0)}, \quad (17)$$

with

$$r_{10}^2 = (x_0+\mu)^2 + y_0^2, \quad r_{20}^2 = (x_0+\mu-1)^2 + y_0^2. \quad (18)$$

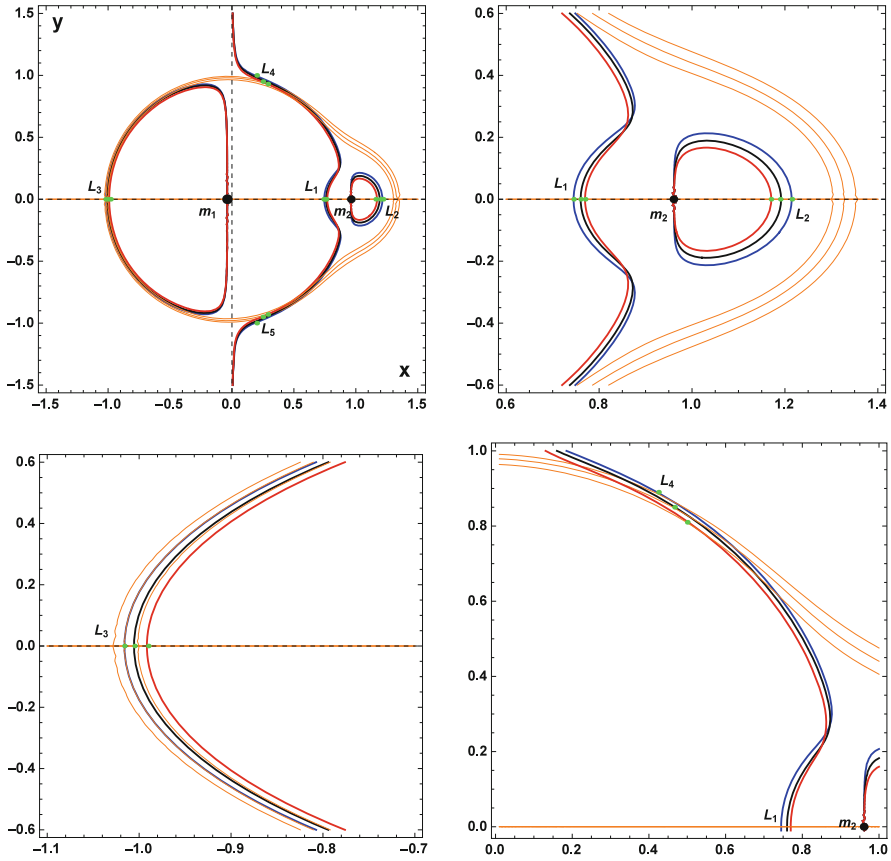
An equilibrium point  $(x_0, y_0)$  is said to be stable in the sense of Lyapunov if and only if all the four roots of the characteristic polynomial, given by Equation (10), are either negative real numbers or distinct imaginary; asymptotically stable if roots are complex with negative real parts and unstable, otherwise.

## 5 Numerical Application

In this section, we compute and examine graphically and numerically the positions of the non-collinear equilibrium points for the binary Cen X-4 system using the astrophysical parameters presented in Table 1 for some assumed oblateness and centrifugal force parameters. As pointed out in Section 2, the adjective non-collinear is due to the fact that  $L_1$ ,  $L_2$  and  $L_3$  do not lie exactly on the  $x$ -axis. In order to visualize the evolution of the equilibria we consider region of the oblateness coefficient  $A_1$  is  $[0, 0.2]$  (see [9]). The investigated region for the values of the Coriolis and centrifugal forces are  $\alpha, \beta \in [1, 1.2]$  (see, e.g. [28]) while the value of the dimensional velocity of light is kept fixed to  $c_d = 988.323$  for all numerical calculations.

Solving Equations (5), using parameters in Table 1, we present in Figures 2 and 3 the positions of the equilibria for the binary system as the two parameters  $A_1$  and  $\beta$  vary in the absence and presence of the P–R drag effect, respectively. For better understanding the evolution of the equilibria, in both figures, we use colour codes to indicate the set of pairs  $(A_1, \beta)$ , while green dots signify the positions of the equilibria. So, the intersections of blue-magenta, black-magenta, and red-magenta curves correspond to three specific pairs of values of  $A_1$  and  $\beta$ ; particularly to  $(0, 1)$ ,  $(0.1, 1.04)$  and  $(0.2, 1.1)$ , respectively. It is necessary to note that, although the curves are identical, their behaviours are different as we observe completely different results regarding the movement of the equilibrium points. From Figure 2 it can be observed that for varying oblateness factor and varying centrifugal force we have five equilibrium points (as in the classical restricted problem), three collinear  $L_{1,2,3}$  and two triangular  $L_{4,5}$ , where equilibria  $L_1$  and  $L_2$  both approach the radiating primary  $m_2$ , while  $L_3$  moves toward the oblate primary  $m_1$  and point  $L_4$  (the situation is same at the symmetric point  $L_5$ ) moves closer to the point  $L_1$ . For clarity purposes, the top-right, bottom-left, and bottom-right frames are enlargements of the top-left frame of Figure 2 (first frame) close to  $L_{1,2}$ ,  $L_3$ , and  $L_{4(5)}$  points, respectively.

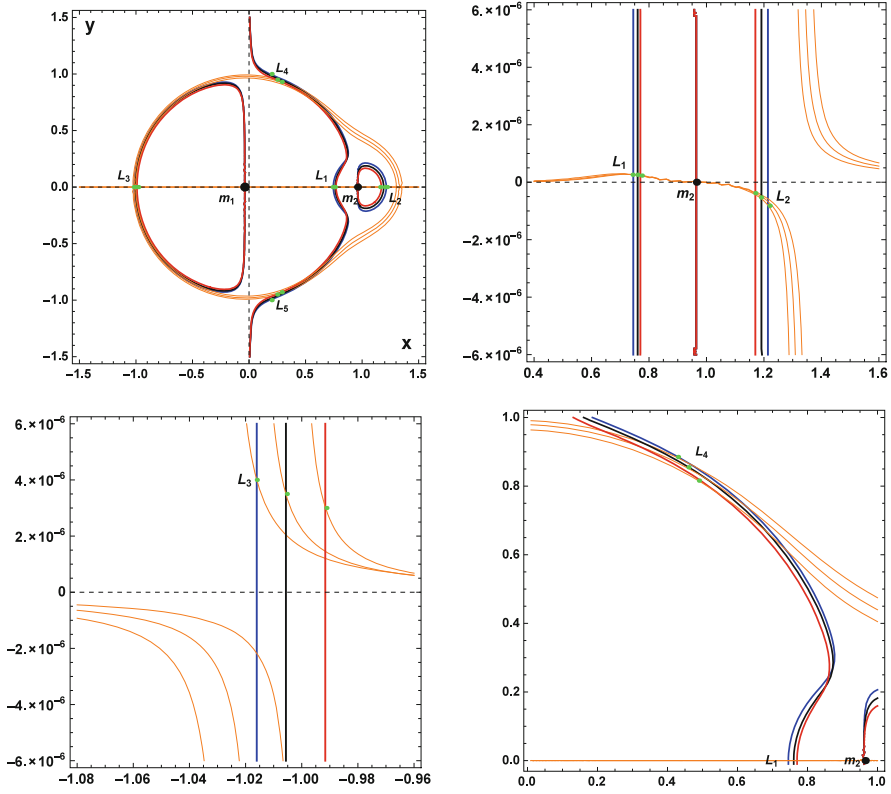
In Table 2, we have evaluated numerically the coordinates of the five equilibrium points for different values of the parameters  $A_1$  and  $\beta$  for the binary system. One can observe in this table that the variational trend of the equilibria is similar to the scenario presented in Figure 2. However, the situation is different in the presence of P–R drag effect as we observe that for increasing values of the oblateness and centrifugal force parameters, there exist five non-collinear equilibrium points positioned off the  $Ox$ -axis. In addition  $L_{1,3,4}$  have  $y > 0$ , while points  $L_{2,5}$  have  $y < 0$ . It can be observed that the equilibria  $L_1$  and  $L_2$  approach the radiating



**Fig. 2** Effect of oblateness and centrifugal force parameters on the collinear  $L_{1,2,3}$  and the triangular  $L_{4,5}$  points of Cen X-4 system without P-R effect (i.e.,  $q_2 = 1$ ,  $W_2 = 0$ ) for  $A_1 = 0$ ,  $\beta = 1$  (blue, magenta);  $A_1 = 0.1$ ,  $\beta = 1.04$  (black, magenta) and  $A_1 = 0.2$ ,  $\beta = 1.1$  (red, magenta). Top-right, bottom-left and right frames: Zoomed areas close to  $L_{1,2}$ ,  $L_3$  and  $L_4$  points, respectively. Black dots represent the primaries while green dots represent the positions of the equilibria

primary body  $m_2$  in opposite directions while  $L_3$  approaches the oblate primary  $m_1$ , and the two non symmetric equilibria  $L_4$  and  $L_5$  approach the displaced  $L_1$  in opposite directions. Tables 3 and 4 provide the locations of the equilibrium points  $L_i$ ,  $i = 1, 2, \dots, 5$  for varying oblateness and centrifugal force parameters in the presence of P-R drag for same fixed values of the parameters. One can see from these tables that the variational trend of the equilibria is similar to the behaviour presented in Figure 3.

Next, since we have already found the coordinates  $(x_0, y_0)$  of the equilibrium points (presented in Tables 2, 3, and 4), we can insert them into the characteristic Equation (10) and thus derive their linear stability numerically. In Tables 5 and 6,



**Fig. 3** Effect of oblateness and centrifugal force parameters on the non-collinear equilibrium points of Cen X-4 system with P-R effect (i.e.,  $q_2 = 0.993$ ,  $W_2 = 2.72790 \times 10^{-7}$ ) for  $A_1 = 0$ ,  $\beta = 1$  (blue, magenta);  $A_1 = 0.1$ ,  $\beta = 1.04$  (black, magenta) and  $A_1 = 0.2$ ,  $\beta = 1.1$  (red, magenta). Top-right, bottom-left and right frames: Zoomed areas close to  $L_{1,2}$ ,  $L_3$  and  $L_4$  points, respectively. Black dots represent the primaries while green dots represent the positions of the equilibria

**Table 2** Positions of the five equilibrium points for varying oblateness and varying centrifugal force in the absence of P-R (i.e.  $q_2 = 1, W_2 = 0$ ) for the binary Cen X-4 system

| $(A_1, \beta)$ | $L_1$         | $L_2$        | $L_3$         | $L_{4,5}$                   |
|----------------|---------------|--------------|---------------|-----------------------------|
| (0, 1)         | (0.744951, 0) | (1.21443, 0) | (-1.01604, 0) | (0.461485, $\pm 0.866025$ ) |
| (0.025, 1.025) | (0.748593, 0) | (1.20566, 0) | (-1.00853, 0) | (0.473601, $\pm 0.849585$ ) |
| (0.05, 1.05)   | (0.751774, 0) | (1.19750, 0) | (-1.00156, 0) | (0.484977, $\pm 0.833890$ ) |
| (0.075, 1.075) | (0.754572, 0) | (1.18988, 0) | (-0.99508, 0) | (0.495655, $\pm 0.818874$ ) |
| (0.1, 1.1)     | (0.757045, 0) | (1.18275, 0) | (-0.98900, 0) | (0.505677, $\pm 0.804479$ ) |

we show the nature of the stability of the equilibrium points for various values of oblateness, Coriolis and centrifugal forces in the absence and presence of the P-R effect, respectively, for the binary Cen X-4 system. Analysis of Tables 5 and 6



**Table 3** Positions of  $L_{1,2,3}$  non-collinear equilibrium points for varying oblateness and varying centrifugal force in the presence of P–R for the binary Cen X-4 system

| $(A_1, \beta)$ | $L_1$                                 | $L_2$                                 | $L_3$                                     |
|----------------|---------------------------------------|---------------------------------------|---|
| (0, 1)         | (0.745409, $2.63737 \times 10^{-7}$ ) | (1.21380, $-5.78026 \times 10^{-7}$ ) | ( $-1.01602$ , $4.01687 \times 10^{-6}$ ) |
| (0.025, 1.025) | (0.749036, $2.60924 \times 10^{-7}$ ) | (1.20505, $-5.44399 \times 10^{-7}$ ) | ( $-1.00851$ , $3.80502 \times 10^{-6}$ ) |
| (0.05, 1.05)   | (0.752204, $2.58813 \times 10^{-7}$ ) | (1.19690, $-5.13966 \times 10^{-7}$ ) | ( $-1.00154$ , $3.61170 \times 10^{-6}$ ) |
| (0.075, 1.075) | (0.754990, $2.57277 \times 10^{-7}$ ) | (1.18930, $-4.86306 \times 10^{-7}$ ) | ( $-0.99501$ , $3.43467 \times 10^{-6}$ ) |
| (0.1, 1.1)     | (0.757452, $2.56225 \times 10^{-7}$ ) | (1.18218, $-4.61069 \times 10^{-7}$ ) | ( $-0.98899$ , $3.27203 \times 10^{-6}$ ) |

**Table 4** Positions of  $L_{4,5}$  non-collinear equilibrium points for varying oblateness and varying centrifugal force in the presence of P–R for the binary Cen X-4 system (continuation of Table 3)

| $(A_1, \beta)$ | $L_4$               | $L_5$                 |
|----------------|---------------------|-----------------------|
| (0, 1)         | (0.46382, 0.864673) | (0.463823, -0.864672) |
| (0.025, 1.025) | (0.47584, 0.848230) | (0.475845, -0.845228) |
| (0.05, 1.05)   | (0.48713, 0.832534) | (0.487133, -0.832533) |
| (0.075, 1.075) | (0.49773, 0.817519) | (0.497729, -0.817517) |
| (0.1, 1.1)     | (0.50767, 0.803125) | (0.507675, -0.803124) |

**Table 5** Stability of Cen X-4 system for small assumed values of oblateness and perturbations in Coriolis and centrifugal forces in the absence of P–R drag effect (see Table 2)

| $L_i, i = 1, 2, \dots, 5$                       | $(x_0, y_0)$                | $\lambda_{1,2}$           | $\lambda_{3,4}$          |
|---|-----------------------------|---------------------------|--------------------------|
| Case: $A_1 = 0, \beta = 1, \alpha = 1$          |                             |                           |                          |
| $L_1$   | (0.744951, 0)               | $\pm 3.145064$            | $\pm 2.469515i$          |
| $L_2$   | (1.21443, 0)                | $\pm 2.002264$            | $\pm 1.772113i$          |
| $L_3$   | (-1.01604, 0)               | $\pm 0.314525$            | $\pm 1.031797i$          |
| $L_{4,5}$                                       | (0.461485, $\pm 0.866025$ ) | $\pm 0.711480i$           | $\pm 0.702705i$          |
| Case: $A_1 = 0.05, \beta = 1.05, \alpha = 1.04$ |                             |                           |                          |
| $L_1$   | (0.751774, 0)               | $\pm 3.338678$            | $\pm 2.591204i$          |
| $L_2$   | (1.19750, 0)                | $\pm 2.224148$            | $\pm 1.956125i$          |
| $L_3$   | (-1.00156, 0)               | $\pm 0.351533$            | $\pm 1.083626i$          |
| $L_{4,5}$                                       | (0.484977, $\pm 0.83389$ )  | $-0.130094 \pm 0.755847i$ | $0.130094 \pm 0.755847i$ |
| Case: $A_1 = 0.1, \beta = 1.1, \alpha = 1.08$   |                             |                           |                          |
| $L_1$   | (0.757045, 0)               | $\pm 3.497120$            | $\pm 2.700120i$          |
| $L_2$   | (1.182750, 0)               | $\pm 2.452064$            | $\pm 2.146047i$          |
| $L_3$   | (-0.989004, 0)              | $\pm 0.386942$            | $\pm 1.141295i$          |
| $L_{4,5}$                                       | (0.505677, $\pm 0.804479$ ) | $-0.180429 \pm 0.806787i$ | $0.180429 \pm 0.806787i$ |

reveals the non existence of pure imaginary roots except in the classical case (i.e.  $q_2 = 1, W_2 = 0, \alpha = \beta = 1$ ). In all cases for all the assumed values of oblateness and perturbations in Coriolis and centrifugal forces with and without P–R effect, there exists at least a positive real root and/or a complex root with positive real part. Consequently the motion of the infinitesimal body is unbounded and thus unstable around all these equilibrium points.

## 6 Discussion and Conclusion

The location and stability of the equilibrium points in the photogravitational restricted three-body problem that accounts for Poynting–Robertson (P–R) drag force with oblateness of the first primary together with small perturbations in the Coriolis and centrifugal forces were studied. It was found both analytically and numerically that in the presence of P–R drag effect the well-known collinear

**Table 6** Stability of Cen X-4 system for small assumed values of oblateness and perturbations in Coriolis and centrifugal forces in the presence of P-R drag effect (see Tables 3 and 4)

| $L_i, i = 1, 2, \dots, 5$                       | $(x_0, y_0)$                         | $\lambda_{1,2}$          | $\lambda_{3,4}$                     |
|---|--------------------------------------|--------------------------|-------------------------------------|
| Case: $A_1 = 0, \beta = 1, \alpha = 1$          |                                      |                          |                                     |
| $L_1$   | $(0.745409, 2.63737 \times 10^{-7})$ | $-3.143139, 3.143130$    | $-3.892 \times 10^{-6} \pm 2.4683i$ |
| $L_2$   | $(1.2138, -5.78026 \times 10^{-7})$  | $-2.003189, 2.003184$    | $-3.721 \times 10^{-6} \pm 1.7726i$ |
| $L_3$   | $(-1.01602, 4.01687 \times 10^{-6})$ | $-0.314652, 0.314652$    | $-1.857 \times 10^{-7} \pm 1.0318i$ |
| $L_4$   | $(0.46382, 0.864673)$                | $-0.01326 \pm 0.707256i$ | $0.013264 \pm 0.707207i$            |
| $L_5$   | $(0.463823, -0.864672)$              | $-0.01336 \pm 0.707257i$ | $0.013358 \pm 0.707209i$            |
| Case: $A_1 = 0.05, \beta = 1.05, \alpha = 1.04$ |                                      |                          |                                     |
| $L_1$   | $(0.752204, 2.58813 \times 10^{-7})$ | $-3.336270, 3.336260$    | $-4.186 \times 10^{-6} \pm 2.5898i$ |
| $L_2$   | $(1.1969, -5.13966 \times 10^{-7})$  | $-2.225351, 2.225345$    | $-4.183 \times 10^{-6} \pm 1.957i$  |
| $L_3$   | $(-1.00154, 3.6117 \times 10^{-6})$  | $-0.351708, 0.351708$    | $-1.932 \times 10^{-7} \pm 1.0837i$ |
| $L_4$   | $(0.48713, 0.832534)$                | $-0.13090 \pm 0.755987i$ | $0.130890 \pm 0.755983i$            |
| $L_5$   | $(0.487133, -0.832533)$              | $-0.13090 \pm 0.75599i$  | $0.130902 \pm 0.755983i$            |
| Case: $A_1 = 0.1, \beta = 1.1, \alpha = 1.08$   |                                      |                          |                                     |
| $L_1$   | $(0.757452, 2.56225 \times 10^{-7})$ | $-3.494287, 3.494280$    | $-4.462 \times 10^{-6} \pm 2.6985i$ |
| $L_2$   | $(1.18218, -4.61069 \times 10^{-7})$ | $-2.453508, 2.453500$    | $-4.671 \times 10^{-6} \pm 2.1468i$ |
| $L_3$   | $(-0.98899, 3.27203 \times 10^{-6})$ | $-0.387111, 0.387111$    | $-1.993 \times 10^{-7} \pm 1.1413i$ |
| $L_4$   | $(0.507673, 0.803125)$               | $-0.18105 \pm 0.806928i$ | $0.181047 \pm 0.806923i$            |
| $L_5$   | $(0.507675, -0.803124)$              | $-0.18105 \pm 0.806929i$ | $0.181051 \pm 0.806924i$            |

equilibrium points of the circular restricted three-body problem cease to exist while the respective triangular equilibrium points do not form equilateral triangles.

Using the astrophysical parameters of the Cen X-4 binary system we performed a numerical study for its equilibrium points and showed that in the case where P–R drag was considered five non-collinear equilibrium points exist whereas in the absence of P–R drag force there are also five equilibrium points but three of them are located on the axis joining the primaries and the rest two form in the plane of motion equilateral triangles with the primaries, as in the circular restricted three-body problem. It was also found that the equilibrium points are independent of the effect of small perturbation in the Coriolis force but are affected by the small perturbation in centrifugal force. For the stability of the five equilibria, the four roots of the characteristic polynomial were determined numerically and found that are unstable due to the existence of at least one positive real root or a complex root with positive real part. The instability of the equilibrium points agrees with the results existing in the literature when the primaries are not oblate spheroids and small perturbations in the Coriolis and centrifugal forces are not considered (for details we refer to [4, 21])

## References

1. E.I. Abouelmagd, F. Alzahrani, A. Hobiny, J.L.G. Guirao, M. Alhothuali, Periodic orbits around the collinear libration points. *J. Nonlinear Sci. Appl.* **9**, 1716–1727 (2016)
2. N. Bello, A. Umar, On the stability of triangular points in the relativistic R3BP when the bigger primary is oblate and the smaller one radiating with application on Cen X-4 binary system. *Results Phys.* **9**, 1067–1076 (2018)
3. L.R. Capdevila, K.C. Howell, A transfer network linking Earth, Moon, and the triangular libration point regions in the Earth-Moon system. *Adv. Space. Res.* **62**, 1826–1852 (2018)
4. J.A. Chernikov, The photogravitational restricted three-body problem. *Sov. Astron.* **14**, 176–181 (1970)
5. A. Chakraborty, A. Narayan, Effect of stellar wind and Poynting–Robertson drag on photogravitational elliptic restricted three-body problem. *Sol. Syst. Res.* **52**, 168–179 (2018)
6. S.M. Elshaboury, E.I. Abouelmagd, V.S. Kalantonis, E.A. Perdios, The planar restricted three-body problem when both primaries are triaxial rigid bodies: Equilibrium points and periodic orbits. *Astrophys. Space Sci.* **361**, 315 (2017)
7. D. Fakis, T. Kalvouridis, The Copenhagen problem with a quasi-homogeneous potential. *Astrophys. Space Sci.* **362**, 102 (2017)
8. F. Gao, R. Wang, Bifurcation analysis and periodic solutions of the HD 191408 system with triaxial and radiative perturbations. *Universe* **6**, 35 (2020)
9. V.S. Kalantonis, C.N. Douskos, E.A. Perdios, Numerical determination of homoclinic and heteroclinic orbits at collinear equilibria in the restricted three-body problem with oblateness. *Celest. Mech. Dyn. Aston.* **94**, 135–153 (2006)
10. V.S. Kalantonis, A.E. Perdiou, E.A. Perdios, On the stability of the triangular equilibrium points in the elliptic restricted three-body problem with radiation and oblateness, in *Mathematical Analysis and Applications*, ed. by T. Rassias, P. Pardalos. Springer Optim. Its Appl., vol. 154 (Springer, Cham, 2019), pp. 273–286
11. B.S. Kushvah, The effect of radiation pressure on the equilibrium points in the generalized photogravitational restricted three body problem. *Astrophys. Space Sci.* **315**, 231–241 (2008)

12. G. Mingotti, F. Toppato, F. Bernelli-Zazzera, Earth-Mars transfers with ballistic escape and low-thrust capture. *Celest. Mech. Dyn. Astr.* **110**, 169–188 (2011)
13. Z.E. Musielak, B. Quarles, The three-body problem. *Rep. Progr. Phys.* **77**, 065901 (2014)
14. C.D. Murray, Dynamical effects of drag in the circular restricted three-body problem. I: location and stability of the Lagrangian equilibrium points. *Icarus* **112**, 465–484 (1994)
15. P. Oberti, A. Vienne, An upgraded theory for Helene, Telesto, and Calypso. *Astron. Astrophys.* **397**, 353–359 (2003)
16. J.P. Papadouris, K.E. Papadakis, Equilibrium points in the photogravitational restricted four-body problem. *Astrophys. Space Sci.* **344**, 21–38 (2013)
17. E.A. Perdios, V.S. Kalantonis, Critical periodic orbits in the restricted three-body problem with oblateness. *Astrophys. Space Sci.* **305**, 331–336 (2006)
18. E.A. Perdios, V.S. Kalantonis, Self-resonant bifurcations of the Sitnikov family and the appearance of 3D isolas in the restricted three-body problem. *Celest. Mech. Dyn. Astr.* **113**, 377–386 (2011)
19. E.A. Perdios, V.S. Kalantonis, C.N. Douskos, Critical periodic orbits in the restricted three-body problem with oblateness. *Astrophys. Space Sci.* **314**, 199–208 (2008)
20. O. Ragos, F.A. Zafiroopoulos, A numerical study of the influence of the Poynting–Robertson effect on the equilibrium points of the photogravitational restricted three-body problem. I. Coplanar case. *Astron. Astrophys.* **300**, 568–578 (1995)
21. D.W. Schuerman, Influence of the Poynting–Robertson effect on triangular points of the photogravitational restricted three-body problem. *Astrophys. J.*, **238**, 337–342 (1980)
22. J. Singh, A.E. Perdiou, J.M. Gyegwe, V.S. Kalantonis, Periodic orbits around the collinear equilibrium points for binary Sirius, Procyon, Luhman 16,  $\alpha$ -Centuari and Luyten 726–8 systems: the spatial case. *J. Phys. Commun.* **1**, 025008 (2017)
23. J. Singh, A.E. Perdiou, J.M. Gyegwe, E.A. Perdios, Periodic solutions around the collinear equilibrium points in the perturbed restricted three-body problem with triaxial and radiating primaries for binary HD 191408, Kruger 60 and HD 155876 systems. *Appl. Math. Comput.* **325**, 358–374 (2018)
24. J. Singh, A. Aminu, Instability of triangular libration points in the perturbed photogravitational R3BP with Poynting–Robertson (P–R) drag. *Astrophys. Space Sci.* **351**, 473–482 (2014)
25. J. Singh, T.O. Amuda, Poynting–Robertson (P–R) drag and oblateness effects on motion around the triangular equilibrium points in the photogravitational R3BP. *Astrophys. Space Sci.* **350**, 119–126 (2014)
26. J. Singh, T.O. Amuda, Out-of-plane equilibrium points in the photogravitational CR3BP with oblateness with P–R drag. *J. Astrophys. Astr.* **36**, 291–305 (2015)
27. J. Singh, A.E. Vincent, Equilibrium points in the restricted fourbody problem with radiation pressure. *Few-Body Syst.* **57**, 83–91 (2016)
28. J. Singh, Motion around the out of plane equilibrium points of the perturbed R3BP. *Astrophys. Space Sci.* **342**, 13–19 (2013)
29. Md.S. Suraj, R. Aggarwal, A. Mittal, Md C.A., The perturbed restricted three-body problem with angular velocity: Analysis of basins of convergence linked to the libration points. *Int. J. Nonlinear Mech.* **123**, 103494 (2020)
30. V. Szebehely, *Theory of Orbits. The Restricted Problem of Three-Bodies* (Academic Press, New York, 1967)
31. P. Verrier, T. Waters, J. Sieber, Evolution of the  $L_1$  halo family in the radial solar sail circular restricted three-body problem. *Celest. Mech. Dyn. Astron.* **120**, 373–400 (2014)
32. X.Y. Zeng, H.X. Baoyin, J.F. Li, Updated rotating mass dipole with oblateness of one primary (I): equilibria in the equator and their stability. *Astrophys. Space Sci.* **361**, 14 (2015)
33. E.E. Zotos, F.L. Dubeibe, Orbital dynamics in the post Newtonian planar circular restricted Sun-Jupiter system. *Int. J. Mod. Phys. D* **27**, 1850036 (2018)
34. E.E. Zotos, D. Veras, The grain size survival threshold in one-planet post-main-sequence exoplanetary systems. *Astron. Astrophys.* **637**, A14 (2020)

# Localization and Perturbation of Complex Zeros of Solutions to Second Order Differential Equations with Polynomial Coefficients. A Survey



Michael Gil<sup>\*</sup>

**Abstract** This paper is a survey of the recent results of the author on the complex zeros of solutions to linear homogeneous second order ordinary differential equations with polynomial coefficients. In particular, estimates for the sums and products of the zeros are derived. These estimates give us bounds for the function counting the zeros of solutions and information about the zero-free domain. Some other applications of the obtained estimates for the sums and products of the zeros are also discussed. In addition, we investigate the variation of the zeros of solutions under perturbations of the coefficients. Illustrative examples are also presented. A part of the results presented in the paper is new.

**AMS (MOS) Subject Classification** 34C10, 34A30

## 1 Introduction

This paper is a survey of the recent results of the author on the zeros of solutions to a linear ordinary differential equation (ODE) with polynomial coefficients in the complex domain.

The literature devoted to the zeros of the solutions of such equations is very rich. Besides, the main tool is the Nevanlinna theory. The excellent exposition of the Nevanlinna theory and its applications to differential equations is given in the book [32]. In that book, in particular, the well-known results of Banks [4–6], Brüggemann [7, 8], Hellerstein and Rossi [26–28, 36], and other mathematicians are reflected. The classical comparison principle for zeros of ODEs in the complex plane is presented in [29].

The real zeros of solutions to equations with polynomial coefficients were investigated in the papers by Gundersen [25], Eremenko and Merenkov [11], and

---

M. Gil<sup>\*</sup> (✉)

Department of Mathematics, Ben Gurion University of the Negev, Beer-Sheva, Israel  
e-mail: [gilmi@bezeqint.net](mailto:gilmi@bezeqint.net)

by C.Z. Huang [30]. The paper [37] studies the convergence of the zeros of a non-trivial (entire) solution to the linear differential equation

$$f'' + \{Q_1(z)e^{P_1(z)} + Q_2(z)e^{P_2(z)} + Q_3(z)e^{P_3(z)}\}f = 0,$$

where  $P_j$  are polynomials of degree  $n \geq 1$  and  $Q_j (\neq 0)$  are entire functions of order less than  $n$  ( $j = 1, 2, 3$ ). In the paper [13], by certain separation and comparison results, estimates for the counting functions of the zeros of solutions to  $n$ th-order linear differential equations are deduced. These estimates generalize known results for the zeros of solutions to third- and fourth-order linear differential equations. The remarkable results on the zeros of a wide class of ordinary differential equations with polynomial coefficients, whose solutions are classical orthogonal polynomials, have been established by N. Anghel [2]. In addition, in the paper [3] N. Anghel investigated the following question: when is an entire function of finite order, the solution to a complex second order homogeneous linear differential equation with polynomial coefficients? He gives two (equivalent) answers to this question, one of which involves certain Stieltjes-like relations for the zeros of solutions, the second one requires the vanishing of all but finitely many suitable expressions constructed via the relations of the sums of the zeros of the function derived in [17].

In connection with the recent results on the complex zeros of solutions to ODEs see also the papers [12, 13, 31, 34], and references given therein. Certainly, we could not survey the whole subject here and refer the reader to the listed publications.

It should be noted that in the above cited works mainly the asymptotic distributions of zeros are investigated. At the same time, bounds for the zeros of solutions are very important in various applications. But to the best of our knowledge, they have been investigated considerably less than the asymptotic distributions. In the paper [18] the author has established bounds for the sums of the zeros of solutions for the second order homogeneous equations with polynomial coefficients. In the interesting paper [9], some results from [18] have been extended to the equation  $u^{(m)} = P(z)u$ , where  $P$  is a polynomial and  $m > 2$ . In the papers [19, 22] and [21] the main result from [18] have been extended to the second order ODEs with non-polynomial coefficients, to ODEs having singular points and to non-homogeneous second order ODEs, respectively. In addition, in the paper [20] the author has derived a bound for the products of the zeros of solutions to ODEs with polynomial coefficients.

It should be also noted that to the best of our knowledge, perturbations of the zeros of solutions were almost not investigated in the available literature. Here we can mention only the paper [19] on perturbations of the zeros of solutions to second order differential equations with polynomial coefficients.

The present paper reflects some results from the just pointed papers of the author. Besides, the proofs are considerably simplified. Our main tool is the recent results for the zeros of entire functions established in [14, 15] (see also [17]).

A few words about the contents. The paper consists of 13 sections.

Section 2 contains solution estimates for non-homogeneous ODEs.

In Section 3 we recall the basic properties of singular values of compact operators, which we need for the proofs of our main results.

In Sections 4 and 5 we obtain bounds for the sums and products of zeros of finite order entire functions via Taylor coefficients and via orders of the functions. As it was above mentioned, bounds for the zeros of entire functions are our main tool.

In Sections 6 we derive bounds for the sums and products of zeros of solutions to the equation  $u'' = P(z)u$  ( $z \in \mathbf{C}$ ), where  $P(z)$  is a polynomial. In Section 7 we discuss some applications of these bounds.

In Sections 8, 9, and 10 we obtain perturbation results for the zeros of entire functions, which are used in Section 12 to obtain perturbation bounds for the zeros of ODEs.

In Section 13 we present an example which illustrates the results obtained in Section 12.

## 2 Solution Estimates for ODEs

Consider the equation

$$\frac{d^2u}{dz^2} = Q(z)u + F(z) \quad (z \in \mathbf{C}, u(0) = u_0 \in \mathbf{C}, u'(0) = u_1 \in \mathbf{C}), \tag{1}$$

where  $Q(z)$  and  $F(z)$  are entire functions. A solution of (1) is a twice continuously differentiable function  $u(z)$  defined for all  $z \in \mathbf{C}$  and satisfying the given initial conditions. Since the equation is linear, the existence and uniqueness of solutions is well known, cf. [32]. About the recent results on solution estimates for ordinary differential equations see for instance the books [1, 35] and references given therein. For an entire function  $f$  and a positive number  $r$  put  $M_f(r) = \sup_{|z| \leq r} |f(z)|$ .

**Lemma 1** *A solution  $u(z)$  of Equation (1) satisfies the inequality*

$$M_u(r) \leq \left( |u_0| + r|u_1| + \int_0^r (r-s)M_F(s)ds \right) \cosh \left( r\sqrt{M_Q(r)} \right) \\ (r \geq 0, \cosh x = (e^x + e^{-x})/2, x \geq 0).$$

**Proof** For a fixed  $t \in [0, 2\pi)$  and  $z = re^{it}$  from (1) we have

$$\frac{1}{e^{2it}} \frac{d^2u(re^{it})}{dr^2} = Q(re^{it})u(re^{it}) + F(re^{it}).$$

Integrating twice this equation in  $r$ , we obtain

$$u(re^{it}) = e^{2it}(u_0 + ru_1) + e^{2it} \int_0^r (r-s)[Q(se^{it})u(se^{it}) + F(se^{it})]ds.$$



Hence,

$$\begin{aligned} M_u(r) &\leq |u_0| + r|u_1| + \int_0^r (r-s)(M_Q(s)M_u(s) + M_F(s))ds \\ &= \int_0^r (r-s)M_Q(s)M_u(s)ds + H(r), \end{aligned}$$

where

$$H(r) = |u_0| + r|u_1| + \int_0^r (r-s)M_F(s)ds.$$

Due to the comparison lemma [10, Lemma III.2.1], we have  $M_u(r) \leq v(r)$ , where  $v(r)$  is a solution of the equation

$$v(r) = H(r) + \int_0^r (r-s)M_Q(s)v(s)ds = H(r) + (Vv)(r).$$

Here  $V$  is the Volterra operator defined by

$$(Vv)(r) = \int_0^r (r-s)M_Q(s)v(s)ds = \int_0^r \int_0^{r_1} M_Q(r_2)v(r_2)dr_2dr_1,$$

and therefore,

$$v = \sum_{k=0}^{\infty} V^k H. \quad (2)$$

But for any positive nondecreasing function  $h(r)$  we have

$$(Vh)(r) = \int_0^r \int_0^{r_1} M_Q(r_2)h(r_2)dr_2dr_1 \leq h(r)M_Q(r) \int_0^r \int_0^{r_1} dr_2dr_1 = h(r)M_Q(r)r^2/2.$$

Similarly,

$$\begin{aligned} (V^2h)(r) &= \int_0^r \int_0^{r_1} M_Q(r_2) \int_0^{r_2} \int_0^{r_3} M_Q(r_4)h(r_4)dr_4 dr_3 dr_2 dr_1 \\ &\leq h(r)M_Q^2(r) \int_0^r \int_0^{r_1} \int_0^{r_2} \int_0^{r_3} dr_4 dr_3 dr_2 dr_1 = h(r)M_Q^2(r)r^4/4! \end{aligned}$$

Continuously this process we obtain

$$(V^m h)(r) \leq h(r)M_Q^m(r) \frac{r^{2m}}{(2m)!} \quad (m = 1, 2, \dots).$$

Thus from (2) it follows

$$M_u(r) \leq v(r) \leq H(r) \sum_{k=0}^{\infty} \frac{M_Q^k(r)r^{2k}}{(2k)!}.$$

But

$$\sum_{k=0}^{\infty} \frac{M_Q^k(r)r^{2k}}{(2k)!} = \cosh(r\sqrt{M_Q(r)}).$$

This implies the required result.

Q.E.D.

Since  $\cosh(x) \leq e^x, x \geq 0$ , we obtain the following corollary.

**Corollary 1** *A solution  $u(z)$  of Equation (1) satisfies the inequality*

$$M_u(r) \leq \left( |u_0| + r|u_1| + \int_0^r (r-s)M_F(s)ds \right) \exp\left(r\sqrt{M_Q(r)}\right) \quad (r \geq 0).$$

Now let  $Q(z) = P(z)$ , where

$$P(z) = \sum_{k=0}^n c_k z^k \quad (c_n \neq 0; c_k \in \mathbf{C}, k = 0, \dots, n < \infty).$$

Consider the equation

$$u'' = P(z)u + F(z), \quad u(0) = u_0, u'(0) = u_1 \quad (u_0, u_1 \in \mathbf{C}). \tag{3}$$

In the considered case  $M_Q(r) \leq p(r)$ , where

$$p(r) = \sum_{k=0}^n |c_k|r^k.$$

According to Corollary 1, a solution  $u(z)$  of (3) satisfies the inequality

$$M_u(r) \leq (|u_0| + r|u_1|) \exp\left(r\sqrt{p(r)}\right) \quad (r \geq 0). \tag{4}$$

Recall the Young inequality

$$ab \leq \frac{a^m}{m} + \frac{a^l}{l} \quad \left(\frac{1}{m} + \frac{1}{l} = 1; m > 1; a, b > 0\right).$$

Then with  $m = (n + 2)/(k + 2)$  we have

$$r^k \leq \frac{kr^{n+2}}{n+2} + 1 - \frac{k}{n+2} = \frac{(k+2)r^{n+2}}{n+2} + \frac{n+2-(k+2)}{n+2}.$$

Hence

$$r^2 p(r) = \sum_{k=0}^n |c_k| r^{k+2} \leq \sum_{k=0}^n |c_k| \left( \frac{(k+2)r^{n+2}}{n+2} + \frac{n-k}{n+2} \right) = b_0 r^{n+2} + b_1,$$

where

$$b_0 = \frac{1}{n+2} \sum_{k=0}^n |c_k| (k+2)$$

and

$$b_1 = \frac{1}{n+2} \sum_{k=0}^n |c_k| (n-k).$$

Since

$$\sqrt{a+b} \leq \sqrt{a} + \sqrt{b} \quad (a, b \geq 0),$$

we get

$$r\sqrt{p(r)} = \sqrt{r^2 p(r)} \leq \sqrt{b_0} r^{n/2+1} + \sqrt{b_1}.$$

Thus

$$\exp[r\sqrt{M_P(r)}] \leq \exp[r\sqrt{p(r)}] \leq \eta_P \exp[\gamma_P r^{n/2+1}],$$

where

$$\gamma_P = \sqrt{\frac{1}{n+2} \sum_{k=0}^n |c_k| (k+2)} \quad \text{and} \quad \eta_P = \exp \left[ \sqrt{\frac{1}{n+2} \sum_{k=0}^n |c_k| (n-k)} \right].$$

Now (4) implies

**Corollary 2** *A solution  $u(z)$  of Equation (3) satisfies the inequality*

$$M_u(r) \leq \eta_P \left( |u_0| + r|u_1| + \int_0^r (r-s) M_F(s) ds \right) \exp[\gamma_P r^{n/2+1}] \quad (r \geq 0).$$

### 3 Singular Values of Compact Operators

Let  $\mathcal{H}$  be a complex separable Hilbert space, cf. [24].

For a compact operator  $A$  acting in  $\mathcal{H}$ , by  $\lambda_k(A)$  ( $k = 1, 2, \dots$ ) we denote the eigenvalues of  $A$  taken with their multiplicities and ordered in the non-increasing way of their absolute values;  $s_k(A)$  ( $k = 1, 2, \dots$ ) are the singular numbers (i.e. the eigenvalues of  $(A^*A)^{1/2}$ ), taken with their multiplicities and ordered in the non-increasing way. Here  $A^*$  is the operator adjoint to  $A$ . In the sequel we need the following well-known results, cf. [24, Section II.4.2], [23, Section IV.4].

**Lemma 2** *Let  $A$  and  $B$  be compact operators in  $\mathcal{H}$ . Then*

$$\sum_{k=1}^j s_k(A+B) \leq \sum_{k=1}^j (s_k(A) + s_k(B)),$$

$$\sum_{k=1}^j s_k^p(AB) \leq \sum_{k=1}^j s_k^p(A)s_k^p(B) \quad (p \geq 1)$$

and

$$\prod_{k=1}^j s_k(AB) \leq \prod_{k=1}^j s_k(A)s_k(B) \quad (j = 1, 2, \dots).$$

In addition, if  $C$  and  $D$  are bounded linear operators in  $\mathcal{H}$ , then

$$s_k(CAD) \leq \|C\| \|D\| s_k(A) \quad (k \geq 1),$$

where  $\|C\|$  means the operator norm of  $C$ .

According to Corollary 2.2 from the book [23] the operators  $A$  and  $A^*$  have the same singular values.

Recall that  $A$  is said to be normal if  $AA^* = A^*A$ .

**Lemma 3 (Weyl's Inequalities)** *The inequalities*

$$\prod_{j=1}^k |\lambda_j(A)| \leq \prod_{j=1}^k s_j(A)$$

and

$$\sum_{j=1}^k |\lambda_j(A)| \leq \sum_{j=1}^k s_j(A) \quad (k = 1, 2, \dots)$$

are true. They become equalities if and only if  $A$  is normal.

For the proof see Theorem IV.3.1 and Corollary IV.3.4 from [23], or Section II.3.1 from [24, Section II.4.2].

### 4 Bounds for Zeros of Entire Functions via Taylor Coefficients

Let us consider the entire function

$$h(z) = \sum_{k=0}^{\infty} \frac{a_k z^k}{(k!)^\alpha} \quad (0 < \alpha \leq 1, z \in \mathbf{C}, a_0 = 1, a_k \in \mathbf{C}, k \geq 1). \tag{5}$$

Enumerate the zeros  $z_k(h)$  of  $h$  with the multiplicities in non-decreasing order of their absolute values:  $|z_k(h)| \leq |z_{k+1}(h)|$  ( $k = 1, 2, \dots$ ) and assume that

$$\theta(h) := \left[ \sum_{k=1}^{\infty} |a_k|^2 \right]^{1/2} < \infty. \tag{6}$$

The aim of this section is to prove the following theorem.

**Theorem 1** *Let  $h$  be defined by (5) and condition (6) hold. Then*

$$\sum_{k=1}^j \frac{1}{|z_k(h)|} \leq \theta(h) + \sum_{k=1}^j \frac{1}{(k+1)^\alpha} \quad (j = 1, 2, \dots)$$

and

$$\prod_{k=1}^j \frac{1}{|z_k(h)|} \leq (\theta(h) + \frac{1}{2^\alpha}) \prod_{k=2}^j \frac{1}{(k+1)^\alpha} \quad (j = 2, 3, \dots).$$

To prove this theorem introduce the polynomial

$$f_n(z) = \sum_{k=0}^n \frac{a_k z^{n-k}}{(k!)^\alpha} \quad (0 < \alpha \leq 1, z \in \mathbf{C}, a_0 = 1, a_k \in \mathbf{C}, k > 1)$$

and the  $n \times n$ -matrix

$$F_n = \begin{pmatrix} -a_1 & -a_2 & \dots & -a_{n-1} & -a_n \\ 1/(2^\alpha) & 0 & \dots & 0 & 0 \\ 0 & 1/(3^\alpha) & \dots & 0 & 0 \\ \cdot & \cdot & \dots & \cdot & \cdot \\ 0 & 0 & \dots & 1/(n^\alpha) & 0 \end{pmatrix}.$$

Let  $z_k(f_n)$  ( $k = 1, \dots, n$ ) be the zeros of  $f_n$  with their multiplicities enumerated in non-increasing order of their absolute values. and  $\lambda_k(F_n)$  be the eigenvalues of  $F_n$  taken with the multiplicities enumerated in the non-increasing order of their absolute values.

**Lemma 4** *One has  $\lambda_k(F_n) = z_k(f_n)$  ( $k = 1, \dots, n$ ).*

*Proof* Clearly,  $f_n$  is the characteristic polynomial of the matrix

$$B = \begin{pmatrix} -a_1 - \frac{a_2}{2^\alpha} \cdots - \frac{a_{n-1}}{((n-1)!)^\alpha} - \frac{a_n}{(n!)^\alpha} & & & & \\ 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \cdot & \cdot & \dots & \cdot & \cdot \\ 0 & 0 & \dots & 1 & 0 \end{pmatrix}.$$

Following [17, Lemma 5.2.1, p. 117], put

$$m_k = \frac{1}{k^\alpha} \text{ and } \psi_k = \frac{1}{(k!)^\alpha} = m_1 m_2 \cdots m_k \quad (k = 1, \dots, n).$$

Then

$$F_n = \begin{pmatrix} -a_1 & -a_2 & \dots & -a_{n-1} & -a_n \\ m_2 & 0 & \dots & 0 & 0 \\ 0 & m_3 & \dots & 0 & 0 \\ \cdot & \cdot & \dots & \cdot & \cdot \\ 0 & 0 & \dots & m_n & 0 \end{pmatrix}$$

and

$$B = \begin{pmatrix} -a_1 & -a_2 \psi_2 & \dots & -a_{n-1} \psi_{n-1} & -a_n \psi_n \\ 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \cdot & \cdot & \dots & \cdot & \cdot \\ 0 & 0 & \dots & 1 & 0 \end{pmatrix}.$$

Let  $\mu$  be an eigenvalue of  $B$ , i.e. for the eigenvector  $(x_k)_{k=1}^n \in \mathbf{C}^n$ , we have

$$-a_1 x_1 - a_2 \psi_2 x_2 - \dots - a_{n-1} \psi_{n-1} x_{n-1} - a_n \psi_n x_n = \mu x_1,$$

$$x_k = \mu x_{k+1} \quad (k = 1, \dots, n - 1).$$

Put  $x_k = y_k / \psi_k$ . Since  $\psi_1 = 1$ , we obtain

$$-a_1 y_1 - a_2 y_2 \dots - a_{n-1} y_{n-1} - a_n y_n = \mu y_1$$

and

$$\frac{y_k}{\psi_k} = \mu \frac{y_{k+1}}{\psi_{k+1}} \quad (k = 1, \dots, n-1).$$

Or

$$m_{k+1}y_k = \frac{y_k\psi_{k+1}}{\psi_k} = \mu y_{k+1} \quad (k = 1, \dots, n-1).$$

These equalities are equivalent to the equality  $F_n y = \mu y$  with  $y = (y_k)$ . In other words  $TBT^{-1} = F_n$ , where  $T = \text{diag}(1, \psi_2, \dots, \psi_n)$  and therefore

$$T^{-1} = \text{diag}\left(1, \frac{1}{\psi_2}, \dots, \frac{1}{\psi_n}\right).$$

This proves the lemma.

Q.E.D.

Put

$$h_n(z) = z^n f_n(1/z) = \sum_{k=0}^n \frac{a_k z^k}{(k!)^\alpha}.$$

Then

$$z_k(h_n) = \frac{1}{z_k(f_n)} = \frac{1}{\lambda_k(F_n)} \quad (k = 1, \dots, n). \quad (7)$$

Here  $z_k(h_n)$  are the zeros of  $h_n$  with their multiplicities enumerated in non-decreasing order of their absolute values.

Furthermore, note that  $F_n = M + C$ , where

$$M = \begin{pmatrix} -a_1 & -a_2 & \dots & -a_{n-1} & -a_n \\ 0 & 0 & \dots & 0 & 0 \\ \cdot & \cdot & \dots & \cdot & \cdot \\ 0 & 0 & \dots & 0 & 0 \end{pmatrix} \text{ and } C = \begin{pmatrix} 0 & 0 & \dots & 0 & 0 \\ 1/2^\alpha & 0 & \dots & 0 & 0 \\ 0 & 1/3^\alpha & \dots & 0 & 0 \\ \cdot & \cdot & \dots & \cdot & \cdot \\ 0 & 0 & \dots & 1/n^\alpha & 0 \end{pmatrix}.$$

Therefore, with

$$\theta(h_n) = \left[ \sum_{k=1}^n |a_k|^2 \right]^{1/2},$$

we have

$$MM^* = \begin{pmatrix} \theta^2(h_n) & 0 & \dots & 0 & 0 \\ 0 & 0 & \dots & 0 & 0 \\ \cdot & \cdot & \dots & \cdot & \cdot \\ 0 & 0 & \dots & 0 & 0 \end{pmatrix} \text{ and } CC^* = \begin{pmatrix} 0 & 0 & \dots & 0 & 0 \\ 0 & 1/2^{2\alpha} & \dots & 0 & 0 \\ 0 & 0 & \dots & 0 & 0 \\ \cdot & \cdot & \dots & \cdot & \cdot \\ 0 & 0 & \dots & 0 & 1/n^{2\alpha} \end{pmatrix}.$$

Hence, the singular values of  $M$  are  $s_1(M) = \theta(h_n)$  and  $s_k(M) = 0$  for  $k > 1$ . In addition,

$$s_k(C) = \frac{1}{(k + 1)^\alpha} \quad (k = 1, \dots, n - 1), \quad s_n(C) = 0.$$

Thus, applying Corollary 2.2 from [24], we arrive at the following result.

**Lemma 5** *The singular values  $s_k(F_n)$  of  $F_n$  satisfy the inequalities*

$$s_1(F_n) \leq \theta(h_n) + \frac{1}{2^\alpha} \text{ and } s_k(F_n) \leq \frac{1}{(k + 1)^\alpha} \text{ for } k = 2, \dots, n.$$

Hence due to Weyl’s inequalities (see the previous section), we obtain

$$\sum_{k=1}^j |\lambda_k(F_n)| \leq \theta(h_n) + \sum_{k=1}^j \frac{1}{(k + 1)^\alpha} \quad (j = 1, \dots, n)$$

and

$$\prod_{k=1}^j |\lambda_k(F_n)| \leq (\theta(h_n) + 1/2^\alpha) \prod_{k=2}^j \frac{1}{(k + 1)^\alpha} \quad (j = 2, \dots, n).$$

Now (7) implies

$$\sum_{k=1}^j \frac{1}{|z_k(h_n)|} \leq \theta(h_n) + \sum_{k=1}^j \frac{1}{(k + 1)^\alpha} \quad (j = 1, \dots, n) \tag{8}$$

and

$$\prod_{k=1}^j \frac{1}{|z_k(h_n)|} \leq (\theta(h_n) + 1/2^\alpha) \prod_{k=2}^j \frac{1}{(k + 1)^\alpha} \quad (j = 2, \dots, n). \tag{9}$$



**Proof of Theorem 1** In each compact domain  $\Omega \in \mathbf{C}$ , we have  $h_n(z) \rightarrow h(z)$  ( $n \rightarrow \infty$ ) uniformly in  $\Omega$ . Due to the Hurwitz theorem [33, p. 5]  $z_k(h_n) \rightarrow z_k(h)$  for  $z_k(h) \in \Omega$ . Now (8) and (9) prove the theorem. Q.E.D.

## 5 Bounds for Sums and Products of Zeros of an Entire Function via Its Order

**Lemma 6** For an entire function  $f(z)$  let the inequality

$$M_f(r) \leq \exp[Br^\rho] \quad (B = \text{const} > 0; \rho \geq 1, r > 0) \quad (10)$$

be fulfilled. Then its Taylor coefficients  $f_j$  ( $j = 1, 2, \dots$ ) satisfy the inequality

$$|f_j| \leq \frac{(eB\rho)^{j/\rho}}{(j!)^{1/\rho}} \quad (j \geq 1).$$

**Proof** By the well-known inequality for the coefficients of a power series

$$|f_j| \leq \frac{M_f(r)}{r^j} \leq \frac{e^{Br^\rho}}{r^j}.$$

Employing the usual method for finding extrema it is easy to see that the function in the right-hand side of this inequality takes its smallest value in the range  $r > 0$  for  $r = (\frac{j}{B\rho})^{1/\rho}$  and therefore

$$|f_j| \leq \left(\frac{eB\rho}{j}\right)^{j/\rho}.$$

Since  $j^j \geq j!$ , we obtain

$$|f_j| \leq (eB\rho)^{j/\rho} \frac{1}{(j^j)^{1/\rho}} \leq \frac{(eB\rho)^{j/\rho}}{(j!)^{1/\rho}},$$

as claimed. Q.E.D.

If an entire function  $v(z)$  satisfies the inequality

$$M_v(r) \leq r^m \exp[Br^\rho] \quad (m = 1, 2, \dots),$$

then we can write  $v(z) = f(z)z^m$  with  $f(z)$  satisfying (10). Then due Lemma 6 the Taylor coefficients at zero  $v_j$  of  $v$  satisfy the relations

$$|v_j|=|f_{j-m}| \leq \frac{(eB\rho)^{(j-m)/\rho}}{((j-m)!)^{1/\rho}} \quad (j=m, m+1, \dots) \text{ and } v_j = 0 \quad (j=0, 1, \dots, m-1).$$

Hence we get

**Corollary 3** *Let an entire function  $v(z)$  be subject to the inequality*

$$M_v(r) \leq \exp[B_v r^\rho](D_0 + D_1 r + D_2 r^2 + \dots + D_m r^m) \quad (r > 0; D_j = \text{const} \geq 0, j = 1, \dots, m < \infty). \tag{11}$$

*Then its Taylor coefficients at zero  $v_j$  satisfy the inequalities*

$$|v_j| \leq \sum_{k=0}^m D_k \frac{(eB_v \rho)^{(j-k)/\rho}}{((j-k)!)^{1/\rho}} \quad (j = 0, 1, \dots).$$

Note that

$$\frac{1}{(j-k)!} = 0 \text{ if } k > j \text{ and } 0! = 1.$$

Put

$$d_j = v_j(j)^{1/\rho}, \xi = \xi_v := (eB_v \rho)^{1/\rho}$$

and

$$l_j := (j!)^{1/\rho} \sum_{k=0}^m D_k \frac{\xi_v^{j-k}}{((j-k)!)^{1/\rho}} \quad (j = 1, \dots).$$

Then due to Corollary 3  $|d_j| \leq l_j$ . Assume that

$$eB_v \rho < 1. \tag{12}$$

Then  $\xi_v < 1$  and for  $j < m$  we have

$$l_j \leq (j!)^{1/\rho} \sum_{k=0}^j D_k \xi_v^{j-k} \leq C_0(j!)^{1/\rho},$$

where

$$C_0 := \sum_{k=0}^m D_k.$$

If  $j \geq m$ , then

$$l_j \leq (j!)^{1/\rho} \sum_{k=0}^m D_k \frac{\xi^{j-k}}{((j-k)!)^{1/\rho}} \leq C_0 j^m \xi^{j-m}.$$

Thus

$$\begin{aligned} \psi_1(v) &:= \sum_{j=1}^{\infty} l_j = \sum_{j=1}^{\infty} (j!)^{1/\rho} \sum_{k=0}^m D_k \frac{\xi^{j-k}}{((j-k)!)^{1/\rho}} \\ &\leq C_0 \left( \sum_{j=0}^{m-1} (j!)^{1/\rho} + \sum_{j=m}^{\infty} \xi^{j-m} j^m \right) < \infty. \end{aligned}$$

The following quantity plays an essential role hereafter:

$$\psi(v) := \left( \sum_{j=1}^{\infty} l_j^2 \right)^{1/2} = \left[ \sum_{j=1}^{\infty} \left( (j!)^{1/\rho} \sum_{k=0}^m D_k \frac{\xi^{j-k}}{((j-k)!)^{1/\rho}} \right)^2 \right]^{1/2}.$$

Since  $\psi(v) \leq \psi_1(v)$ , we have  $\psi(v) < \infty$ . We thus have proved the following lemma.

**Lemma 7** *Let an entire function  $v(z)$  be subject to inequalities (11) and (12). Then it is representable as*

$$v(z) = \sum_{k=0}^{\infty} \frac{d_k z^k}{(k!)^{1/\rho}},$$

and

$$\theta(v) := \left[ \sum_{j=1}^{\infty} |d_j|^2 \right]^{1/2} \leq \psi(v) < \infty.$$

Making use of Theorem 1 and the latter lemma we get our next result.

**Theorem 2** *Let  $v$  be an entire function satisfying inequalities (11) and (12), and let  $v(0) = 1$ . Then*

$$\sum_{k=1}^j \frac{1}{|z_k(v)|} \leq \psi(v) + \sum_{k=1}^j \frac{1}{(k+1)^{1/\rho}} \quad (j = 1, 2, \dots)$$

and

$$\prod_{k=1}^j \frac{1}{|z_k(v)|} \leq (\psi(v) + \frac{1}{2^{1/\rho}}) \prod_{k=2}^j \frac{1}{(k+1)^{1/\rho}} \quad (j = 2, 3, \dots).$$

*Remark 1* Condition (12) is not very restrictive. If  $eB\rho \geq 1$ , then we can apply the substitution  $z = wc$  with

$$c^\rho < \frac{1}{eB\rho}.$$

Indeed, put  $v_c(w) = v(cw)$ . Then condition (11) implies.

$$M_{v_c}(r) \leq (D_0 + D_1cr + \dots + D_mcr^m) \exp[Bc^\rho r^\rho] \quad (r > 0).$$

Condition (12) takes the form

$$eBc^\rho \rho < 1.$$

So if

$$c^\rho < \frac{1}{eB\rho},$$

then Theorem 2 can be applied to the function  $v_c(z)$ . Besides,  $z_k(v) = cz_k(v_c)$ .

## 6 Sums and Products of Zeros of Solutions to Homogeneous ODEs with Polynomial Coefficients

Again consider the equation

$$u''(z) = P(z)u(z) \quad (u(0) = 1, u'(0) = u_1 \in \mathbf{C}, z \in \mathbf{C}), \tag{13}$$

where

$$P(z) = \sum_{k=0}^n c_k z^k \quad (c_n \neq 0; c_k \in \mathbf{C}, k = 0, \dots, n).$$

As is well-known [32], a solution  $u(z)$  of (13) is an entire function whose order is no more than  $n/2 + 1$ , and the zeros  $z_k(u)$  of  $u(z)$  are simple. Enumerate the zeros of  $u$  in the nondecreasing order of their absolute values:  $|z_k(u)| \leq |z_{k+1}(u)|$  ( $k = 1, 2, \dots$ ).

According to Corollary 2  $u(t)$  satisfies inequality (11) with

$$B = \gamma_P, \rho = n/2 + 1, D_0 = \eta_P, D_1 = \eta_P|u_1|, D_k = 0 \quad (k \geq 2).$$

Recall that

$$\gamma_P = \sqrt{\frac{1}{n+2} \sum_{k=0}^n |c_k|(k+2)} \text{ and } \eta_P = \exp \left[ \sqrt{\frac{1}{n+2} \sum_{k=0}^n |c_k|(n-k)} \right].$$

Condition (12) takes the form

$$e\gamma_P(n+1/2) < 1. \tag{14}$$

Hence,

$$\xi_P := (e\gamma_P(n+1/2))^{2/(2n+1)} < 1.$$

Simple calculations show that

$$\psi(u) = \eta_P \left[ \sum_{j=1}^{\infty} \xi_P^{2j} \left( 1 + |u_1| \frac{j^{2/(n+2)}}{\xi_P} \right)^2 \right]^{1/2} < \infty.$$

Now Theorem 2 implies

**Theorem 3** *If condition (14) holds, then the zeros of a solution  $u(z)$  to Equation (13) satisfy the inequalities*

$$\sum_{k=1}^j \frac{1}{|z_k(u)|} \leq \psi(u) + \sum_{k=1}^j \frac{1}{(k+1)^{2/(n+2)}} \quad (j = 1, 2, \dots) \tag{15}$$

and

$$\prod_{k=1}^j \frac{1}{|z_k(u)|} \leq (\psi(u) + \frac{1}{2^{2/(n+2)}}) \prod_{k=2}^j \frac{1}{(k+1)^{2/(n+2)}} \quad (j = 2, 3, \dots).$$

*Remark 2* Condition (14) is not very restrictive. If it does not hold:  $e\gamma_P(n+1/2) \geq 1$ , then one can apply the substitution  $z = wb$  with

$$0 < b = \text{const} < \frac{1}{e\gamma_P(n/2+1)}. \tag{16}$$

Indeed, substituting  $z = wb$  into (13), with  $u_b(w) = u(bw)$  we have

$$\frac{1}{b^2} \frac{d^2 u_b(w)}{dw^2} = P(bw)u_b(w).$$

Or

$$\frac{d^2u_b(w)}{dw^2} = P_b(w)u_b(w), \tag{17}$$

where

$$P_b(w) = b^2 P(bw) = \sum_{k=0}^n c_k b^{k+2} w^k.$$

If  $e\gamma_P(n + 1/2) \geq 1$ , then due to (16) we have  $b < 1$  and therefore

$$\gamma_{P_b} = \sqrt{\frac{1}{n+2} \sum_{k=0}^n b^{k+2} |c_k| (k+2)} \leq b \sqrt{\frac{1}{n+2} \sum_{k=0}^n |c_k| (k+2)} = \gamma_P b.$$

According to (16) condition (14) for Equation (17) is fulfilled:

$$e\gamma_{P_b}(n/2 + 1) < e b \gamma_P(n + 1/2) < 1.$$

Therefore we can apply Theorem 3. Besides  $z_k(u) = b z_k(u_b)$  ( $k = 1, 2, \dots$ ).

*Example 1* To estimate the sharpness of Theorem 3, consider the equation

$$u'' + a^2 u = 0, u(0) = 1, u'(0) = 0 \quad (a = \text{const} > 0).$$

Then  $u(z) = \cos(az)$  and its zeros are

$$\frac{\pi}{a}(m + 1/2) \quad (m = 0, \pm 1, \pm 2, \dots).$$

So

$$z_{2k}(\cos(az)) = \frac{\pi}{a}(k - 1/2), z_{2k-1} = \frac{\pi}{a}(-k + 1/2) \quad (k = 1, 2, \dots)$$

and

$$\sum_{k=1}^{2j} \frac{1}{|z_k(\cos(az))|} = \frac{2a}{\pi} \sum_{k=1}^j \frac{1}{k - 1/2} \quad (j = 1, 2, \dots). \tag{18}$$

In the considered case  $n = 0, \gamma_P = a$  and  $\eta_P = 1$ . Thus,

$$e\gamma_P(n + 1/2) = \frac{ea}{2}.$$

So condition (14) holds, provided

$$a < \frac{2}{e}.$$

This condition yields  $\xi_P := ea/2 < 1$  and  $\psi(u) = \psi(\cos(az))$ , where

$$\psi(\cos(az)) = \left[ \sum_{k=1}^{\infty} \xi_P^{2k} \right]^{1/2} = \left[ \sum_{k=1}^{\infty} (ea/2)^{2k} \right]^{1/2} = \frac{1}{[(2/ea)^2 - 1]^{1/2}}.$$

Now Theorem 3 implies

$$\sum_{k=1}^{2j} \frac{1}{|z_k(\cos(az))|} \leq \psi(\cos(az)) + \sum_{k=1}^{2j} \frac{1}{k+1} \quad (j = 1, 2, \dots) \quad (19)$$

and

$$\prod_{k=1}^{2j} \frac{1}{|z_k(\cos(az))|} \leq (\psi(\cos(az)) + \frac{1}{2}) \prod_{k=2}^{2j} \frac{1}{(k+1)} \quad (j = 2, 3, \dots).$$

We can see that (19) and (18) are asymptotically equivalent.

## 7 Applications of Theorem 3

Again  $u(z)$  is a solution of Equation (13). Recall that  $u(0) = 1$ . Assume that  $n \geq 1$  and for the brevity put

$$\alpha = \frac{2}{n+2}.$$

Since  $|z_k(u)| \leq |z_{k+1}(u)|$ , Theorem 3 implies that

$$\frac{j}{|z_j(u)|} \leq \psi(u) + \sum_{k=1}^j \frac{1}{(k+1)^\alpha} \quad (j = 1, 2, \dots),$$

provided condition (14) holds. But

$$\sum_{k=1}^j (k+1)^{-\alpha} \leq \int_1^{j+1} \frac{dx}{x^\alpha} = \frac{(1+j)^{1-\alpha} - 1}{1-\alpha} \quad (0 < \alpha < 1).$$

Thus,

$$\frac{j}{|z_j(u)|} \leq \psi(u) + \frac{(1+j)^{1-\alpha} - 1}{1-\alpha}$$

and therefore,

$$|z_j(u)| \geq \chi_j(u), \tag{20}$$

where

$$\chi_j(u) = \frac{j}{\psi(u) + \frac{(1+j)^{1-\alpha} - 1}{1-\alpha}}.$$

If  $|z_j(u)| \geq a$  ( $a > 0$ ), then  $u(z)$  has in  $\Omega(a) := \{z \in \mathbf{C} : |z| < a\}$  no more than  $j - 1$  zeros. Denote by  $\nu(f, a)$  the number of the zeros of an entire function  $f$  inside  $\Omega(a)$ , i.e.  $\nu(f, a)$  is the counting function of the zeros of  $f$ . In particular, due to (20)  $\nu(u, a) = 0$  for any positive

$$a < \chi_1(u) = \frac{1}{\psi(u) + \frac{2^{1-\alpha} - 1}{1-\alpha}}.$$

We thus arrive at

**Corollary 4** *Let condition (14) hold and  $n \geq 1$ . Then the counting function of the zeros of a solution  $u(z)$  of (13) satisfies the inequality  $\nu(u, a) \leq j - 1$  for any positive  $a \leq \chi_j(u)$  ( $j = 1, 2, \dots$ ).*

*Moreover the circle  $\{z \in \mathbf{C} : |z| < \chi_1(u)\}$  is a zero-free domain of  $u(z)$ .*

To consider additional applications of Theorem 3 recall some inequalities for convex functions. The following result is classical, cf. [24, Lemma II.3.4], [23, p. 53].

**Lemma 8** *Let  $\phi(x)$  ( $-\infty \leq x \leq \infty$ ) be a convex continuous function, such that*

$$\phi(-\infty) = \lim_{x \rightarrow -\infty} \phi(x) = 0,$$

*and  $a_j, b_j$  ( $j = 1, 2, \dots, l \leq \infty$ ) be two non-increasing sequences of real numbers, such that*

$$\sum_{k=1}^j a_k \leq \sum_{k=1}^j b_k \quad (j = 1, 2, \dots, l).$$

*Then*

$$\sum_{k=1}^j \phi(a_k) \leq \sum_{k=1}^j \phi(b_k) \quad (j = 1, 2, \dots, l).$$



The next result is also well known, cf. [24, Chapter II], [23, p. 53].

**Lemma 9** *Let a scalar-valued function  $\Phi(t_1, t_2, \dots, t_j)$  with an integer  $j$  be defined on the domain*

$$-\infty < t_j \leq t_{j-1} \dots \leq t_2 \leq t_1 < \infty$$

*and have continuous partial derivatives, satisfying the condition*

$$\frac{\partial \Phi}{\partial t_1} > \frac{\partial \Phi}{\partial t_2} > \dots > \frac{\partial \Phi}{\partial t_j} > 0 \text{ for } t_1 > t_2 > \dots > t_j,$$

*and  $a_k, b_k$  ( $k = 1, 2, \dots, j$ ) be two non-increasing sequences of real numbers satisfying the condition*

$$\sum_{k=1}^m a_k \leq \sum_{k=1}^m b_k \quad (m = 1, 2, \dots, j).$$

*Then  $\Phi(a_1, \dots, a_j) \leq \Phi(b_1, \dots, b_j)$ .*

Furthermore, put

$$\vartheta_1 = \psi(u) + \frac{1}{2^\alpha} \text{ and } \vartheta_k = \frac{1}{(k+1)^\alpha} \quad (k = 2, 3, \dots).$$

Inequality (15) and Lemma 8 yield

**Corollary 5** *Let  $\phi(t)$  ( $0 \leq t < \infty$ ) be a continuous convex function, such that  $\phi(0) = 0$ . Then*

$$\sum_{k=1}^j \phi\left(\frac{1}{|z_k(u)|}\right) \leq \sum_{k=1}^j \phi(\vartheta_k) \quad (j = 1, 2, \dots).$$

*In particular, for any  $p \geq 1$  and  $j = 2, 3, \dots$ , we have*

$$\sum_{k=1}^j \frac{1}{|z_k(u)|^p} \leq \sum_{k=1}^j \vartheta_k^p$$

*and therefore, if  $p > 1/\alpha = n/2 + 1$ , then*

$$\sum_{k=1}^{\infty} \frac{1}{|z_k(u)|^p} \leq \vartheta_1^p + \zeta(p\alpha) - 1 - \frac{1}{2^{p\alpha}} < \infty,$$

where  $\zeta(z)$  is the zeta Riemann function:

$$\zeta(z) = \sum_{k=1}^{\infty} \frac{1}{k^z} \quad (\text{Re } z > 1).$$

In addition, making use of (15) and Lemma 9, we arrive at

**Corollary 6** Let  $\Phi(t_1, t_2, \dots, t_j)$  satisfy the hypothesis of Lemma 9. Then

$$\Phi\left(\frac{1}{|z_1(u)|}, \dots, \frac{1}{|z_j(u)|}\right) \leq \Phi(\vartheta_1, \dots, \vartheta_j).$$

In particular, let  $\{m_k\}_{k=1}^{\infty}$  be a decreasing sequence of positive numbers with  $m_1 = 1$ . Then the previous corollary yields the inequality

$$\sum_{k=1}^j \frac{m_k}{|z_k(u)|} \leq \psi(u) + \sum_{k=1}^j \frac{m_k}{(k+1)^\alpha} \quad (j = 1, 2, \dots).$$

## 8 A Perturbation Bound for the Zeros of Entire Functions in Terms of Taylor Coefficients

**Definition 1** Let  $z_j(h)$  and  $z_j(\tilde{h})$  ( $j = 1, 2, \dots$ ) be the zeros of entire functions  $h$  and  $\tilde{h}$ , respectively, enumerated with the multiplicities in the non-decreasing order of their absolute values. Then the quantity

$$\text{rv}_h(\tilde{h}) = \sup_j \inf_k \left| \frac{1}{z_k(h)} - \frac{1}{z_j(\tilde{h})} \right|$$

will be called the relative variation of the zeros of  $\tilde{h}$  with respect to the zeros of  $h$ .

In this section we consider entire functions of the form

$$h(z) = \sum_{k=0}^{\infty} \frac{a_k z^k}{(k!)^\alpha} \text{ and } \tilde{h}(z) = \sum_{k=0}^{\infty} \frac{\tilde{a}_k z^k}{(k!)^\alpha} \quad (a_0 = \tilde{a}_0 = 1, 0 < \alpha \leq 1) \tag{21}$$

with complex coefficients  $a_k, \tilde{a}_k$  ( $k = 1, 2, \dots$ ), assuming that

$$\theta(h) = \left(\sum_{k=1}^{\infty} |a_k|^2\right)^{1/2} < \infty \text{ and } \theta(\tilde{h}) = \left(\sum_{k=1}^{\infty} |\tilde{a}_k|^2\right)^{1/2} < \infty. \tag{22}$$

To investigate the relative variation of the zeros take an integer  $p$  satisfying the inequality

$$p > \frac{1}{2\alpha} \quad (23)$$

and put

$$\varpi_p(h) := 2 \left[ \left( \theta(h) + \frac{1}{2\alpha} \right)^{2p} + \zeta(2\alpha p) - 1 - \frac{1}{2^{2\alpha p}} \right]^{1/2p}.$$

Recall that  $\zeta(z)$  is the Riemann zeta function. Finally, denote

$$q := \left[ \sum_{k=1}^{\infty} |a_k - \tilde{a}_k|^2 \right]^{1/2}$$

and

$$\chi_p(h, y) := \sum_{k=0}^{p-1} \frac{\varpi_p^k(h)}{y^{k+1}} \exp \left[ \frac{1}{2} + \frac{\varpi_p^{2p}(h)}{2y^{2p}} \right] \quad (y > 0).$$

**Theorem 4** *Let  $h$  and  $\tilde{h}$  be defined by (21) and conditions (22) be fulfilled. Then for any integer  $p$  satisfying inequality (23) we have*

$$\text{rv}_h(\tilde{h}) \leq y_p(q, h),$$

where  $y_p(q, h)$  is the unique positive root of the equation

$$q \chi_p(h, y) = 1.$$

The proof of this theorem is presented in the next section. Since  $\chi_p(h, \cdot)$  is a monotonically decreasing function, the latter theorem implies

**Corollary 7** *Let  $h$  and  $\tilde{h}$  be defined by (21), and conditions (22) and (23) be fulfilled. Then all the zeros of  $\tilde{h}$  lie in the union of the sets  $W_k(p, h)$  ( $k = 1, 2, 3, \dots$ ), where*

$$W_k(p, h) := \left\{ \lambda \in \mathbf{C} : q \chi_p \left( h, \left| \frac{1}{z_k(h)} - \frac{1}{\lambda} \right| \right) \geq 1 \right\}.$$

*In particular, if  $h$  has  $l < \infty$  zeros (we do not take into account the zero limits at infinity), then all the zeros of  $\tilde{h}$  lie in the set*

$$\bigcup_{k=0}^l W_k(p, h),$$

where

$$W_0(p, h) = \{\lambda \in \mathbf{C} : q\chi_p(h, \frac{1}{|\lambda|}) \geq 1\}.$$

Note that

$$\chi_p\left(h, \frac{1}{|\lambda|}\right) = \sum_{k=0}^{p-1} \varpi_p^k(h) |\lambda|^{k+1} \exp\left[\frac{1}{2}(1 + \varpi_p^{2p}(h) |\lambda|^{2p})\right].$$

Furthermore, we need the following lemma.

**Lemma 10** *The unique positive root  $z_a$  of the equation*

$$\sum_{j=0}^{p-1} \frac{1}{y^{j+1}} \exp\left[\frac{1}{2}\left(1 + \frac{1}{y^{2p}}\right)\right] = a \quad (a = \text{const} > 0)$$

satisfies the inequality  $z_a \leq \delta_p(a)$ , where

$$\delta_p(a) := \begin{cases} pe/a & \text{if } a \leq pe, \\ [\ln(a/p)]^{-1/2p} & \text{if } a > pe \end{cases}.$$

For the proof see [16, Lemma 8.3.3] or Lemma 1.6.4 from [17]. Substitute the equality  $y = x\varpi_p(h)$  into the equation  $q\chi_p(h, y) = 1$ , and apply the latter lemma. Then, we have the inequality

$$y_p(q, h) \leq \delta_p(q, h), \tag{24}$$

where

$$\delta_p(q, h) := \begin{cases} epq & \text{if } \varpi_p(h) \leq epq, \\ \varpi_p(h) [\ln(\varpi_p(h)/qp)]^{-1/2p} & \text{if } \varpi_p(h) > epq. \end{cases}$$

Now Theorem 4 yields the inequality

$$\text{rv}_f(\tilde{h}) \leq \delta_p(q, h).$$

If  $h$  has an infinite set of zeros, then according to Theorem 4 for any zero  $z(\tilde{h})$  of  $\tilde{h}$ , there is a zero  $z(h)$  of  $h$ , such that

$$|z(\tilde{h}) - z(h)| \leq y_p(q, h) |z(\tilde{h})z(h)|$$

Hence,

$$|z(h)| \leq |z(\tilde{h})|(1 + y_p(q, h)|z(h)|)$$

and therefore, for any zero  $z(\tilde{h})$  of  $\tilde{h}$ , there is a zero  $z(h)$  of  $h$ , such that

$$|z(\tilde{h})| \geq \frac{|z(h)|}{y_p(q, h)|z(h)| + 1}.$$

Now (24) implies

$$|z(\tilde{h})| \geq \frac{|z(h)|}{\delta_p(q, h)|z(h)| + 1}.$$

## 9 Proof of Theorem 4

For an integer  $n \geq 2$ , consider the polynomials

$$h_n(z) = \sum_{k=0}^n \frac{a_k z^k}{(k!)^\alpha} \text{ and } \tilde{h}_n(z) = \sum_{k=0}^n \frac{\tilde{a}_k z^k}{(k!)^\alpha}.$$

Put

$$\theta(h_n) := \left[ \sum_{k=1}^n |a_k|^2 \right]^{1/2}, \quad \varpi_p(h_n) := 2 \left[ \left( \theta(h_n) + \frac{1}{2^{p\alpha}} \right)^{2p} + \sum_{k=2}^n \frac{1}{(k+1)^{2p\alpha}} \right]^{1/2p},$$

$$\chi_p(h_n, y) := \sum_{k=0}^{p-1} \frac{\varpi_p^k(h_n)}{y^{k+1}} \exp \left[ \frac{1}{2} + \frac{\varpi_p^{2p}(h_n)}{2y^{2p}} \right] \quad (y > 0),$$

and

$$q_n := \left[ \sum_{k=1}^n |a_k - \tilde{a}_k|^2 \right]^{1/2}.$$

Let

$$f_n(\lambda) = \lambda^n h_n(1/\lambda) = \sum_{k=0}^n \frac{a_k}{k^\alpha} \lambda^{n-k} \text{ and } \tilde{f}_n(\lambda) = \lambda^n \tilde{h}_n(1/\lambda) = \sum_{k=0}^n \frac{\tilde{a}_k}{k^\alpha} \lambda^{n-k}.$$

**Lemma 11** For any zero  $z(\tilde{f}_n)$  of  $\tilde{f}_n$ , there is a zero  $z(f_n)$  of  $f_n$ , such that

$$|z(f_n) - z(\tilde{f}_n)| \leq y_p(q_n, f_n),$$

where  $y_p(q_n, f_n)$  is the unique (positive) root of the equation

$$q_n \chi_p(h_n, y) = 1.$$

**Proof** Introduce the matrices

$$F_n = \begin{pmatrix} -a_1 & -a_2 & \dots & -a_{n-1} & -a_n \\ \frac{1}{2^\alpha} & 0 & \dots & 0 & 0 \\ 0 & \frac{1}{3^\alpha} & \dots & 0 & 0 \\ \cdot & \cdot & \dots & \cdot & \cdot \\ 0 & 0 & \dots & \frac{1}{n^\alpha} & 0 \end{pmatrix}$$

and

$$\tilde{F}_n = \begin{pmatrix} -\tilde{a}_1 & -\tilde{a}_2 & \dots & -\tilde{a}_{n-1} & -\tilde{a}_n \\ \frac{1}{2^\alpha} & 0 & \dots & 0 & 0 \\ 0 & \frac{1}{3^\alpha} & \dots & 0 & 0 \\ \cdot & \cdot & \dots & \cdot & \cdot \\ 0 & 0 & \dots & \frac{1}{n^\alpha} & 0 \end{pmatrix}.$$

Thanks to Lemma 4, we have

$$\lambda_k(F_n) = z_k(f_n), \quad \lambda_k(\tilde{F}_n) = z_k(\tilde{f}_n) \quad (k = 1, 2, \dots, n),$$

where  $\lambda_k(A)$ ,  $k = 1, \dots, n$  are the eigenvalues of an  $n \times n$  matrix  $A$  with their multiplicities. Besides,

$$\tilde{F}_n - F_n = \begin{pmatrix} \tilde{a}_1 - a_1 & \tilde{a}_2 - a_2 & \dots & \tilde{a}_{n-1} - a_{n-1} & \tilde{a}_n - a_n \\ 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & \dots & 0 & 0 \\ \cdot & \cdot & \dots & \cdot & \cdot \\ 0 & 0 & \dots & 0 & 0 \end{pmatrix}$$

and

$$q_n = \|\tilde{F}_n - F_n\|,$$

where  $\|\cdot\|$  is the spectral norm, i.e. the operator norm with respect to the Euclidean vector norm. Making use of [17, Theorem 2.12.4], we can assert that for any  $\lambda_j(\tilde{F}_n)$  there is a  $\lambda_k(F_n)$ , such that

$$|\lambda_j(\tilde{F}_n) - \lambda_k(F_n)| \leq y_p(q_n), \quad (25)$$

where  $y_p(q_n)$  is the unique (positive) root of the equation

$$q_n \sum_{k=0}^{p-1} \frac{(2N_{2p}(F_n))^k}{y^{k+1}} \exp \left[ \frac{1}{2} + \frac{(2N_{2p}(F_n))^{2p}}{2y^{2p}} \right] = 1.$$

Here

$$N_{2p}(F_n) = \left[ \sum_{k=1}^n s_k^{2p}(F_n) \right]^{1/2p}$$

is the Schatten-von Neumann norm of  $F_n$ . Due to Lemma 5,

$$s_1(F_n) \leq \theta(h_n) + \frac{1}{2^\alpha}$$

and

$$s_k(F_n) \leq \frac{1}{(k+1)^\alpha}$$

for  $k > 1$ . Thus,

$$N_{2p}^{2p}(F_n) \leq \left( \theta(h_n) + \frac{1}{2^\alpha} \right)^{2p} + \sum_{k=2}^n \frac{1}{(k+1)^{2p\alpha}} = \varpi_p^{2p}(h_n).$$

Now (25) implies the required result.

Q.E.D.

**Proof of Theorem 4** Since

$$z_k(h_n) = \frac{1}{z_k(f_n)} \text{ and } z_k(\tilde{h}_n) = \frac{1}{z_k(f_n)},$$

taking into account that  $\varpi_p(h_n) \leq \varpi_p(h)$ , and the roots continuously depend on coefficients, we have the required result, letting in the previous lemma  $n \rightarrow \infty$ . Q.E.D.

## 10 A Perturbation Bound for the Zeros of an Entire Function via Its Order

Let  $h(z)$  and  $\tilde{h}(z)$  be entire functions with  $h(0) = \tilde{h}(0) = 1$ . Put  $v(z) = h(z) - \tilde{h}(z)$  and assume that

$$M_h(r) \leq \exp[B_h r^\rho](C_0 + C_1 r + C_2 r^2 + \dots + C_s) \tag{26}$$

$(r > 0; B_h, C_j = \text{const} \geq 0, j = 0, 1, \dots, s < \infty),$

and

$$M_v(r) \leq \exp[B_v r^\rho](D_1 r + D_2 r^2 + \dots + D_m r^m)$$

$$(r > 0; B_v, D_j = \text{const} \geq 0, j = 1, \dots, m < \infty). \tag{27}$$

By Corollary 3 the Taylor coefficients at zero  $v_j$  of  $v$  satisfy the inequalities

$$|v_j| \leq \sum_{k=1}^m D_k \frac{\xi_v^{j-k}}{((j-k)!)^{1/\rho}} \quad (j = 1, \dots),$$

where  $\xi_v = (eB_v \rho)^{1/\rho}$ . Similarly, the Taylor coefficients at zero  $h_j$  of  $h$  satisfy the inequalities

$$|h_j| \leq \sum_{k=0}^s C_k \frac{\xi_h^{j-k}}{((j-k)!)^{1/\rho}} \quad (j = 1, 2, \dots),$$

where  $\xi_h := (eB_h \rho)^{1/\rho}$ . Recall that  $\frac{1}{(j-k)!} = 0$  if  $k > j$  and  $0! = 1$ . If

$$B_h \leq B_v \text{ and } eB_v \rho < 1, \tag{28}$$

then  $\xi_v < 1$  and  $\xi_h < 1$ . Besides,

$$\psi(v) = \left[ \sum_{j=1}^{\infty} (j!)^{2/\rho} \left( \sum_{k=1}^m D_k \frac{\xi_v^{j-k}}{((j-k)!)^{1/\rho}} \right)^2 \right]^{1/2} < \infty$$

and

$$\psi(h) = \left[ \sum_{j=1}^{\infty} (j!)^{2/\rho} \left( \sum_{k=0}^s C_k \frac{\xi_h^{j-k}}{((j-k)!)^{1/\rho}} \right)^2 \right]^{1/2} < \infty$$

(see Section 5). Due to Lemma 7  $v(z)$  is representable as

$$v(z) = \sum_{k=0}^{\infty} \frac{d_k z^k}{(k!)^{1/\rho}}, \tag{29}$$

with

$$\theta^2(v) = \sum_{j=1}^{\infty} |d_j|^2 \leq \psi^2(v)$$



and  $h(z)$  is representable as

$$h(z) = \sum_{k=0}^{\infty} \frac{a_k z^k}{(k!)^{1/\rho}},$$

with

$$\theta^2(h) = \sum_{j=1}^{\infty} |a_j|^2 \leq \psi^2(h).$$

Since  $v(z) = h(z) - \tilde{h}(z)$ , the Taylor coefficients  $\tilde{h}_k$  of  $\tilde{h}$  are equal to  $h_k - v_k$  and due to (29) we can write

$$\tilde{h}(z) = \sum_{k=0}^{\infty} \frac{\tilde{a}_k z^k}{(k!)^{1/\rho}},$$

with

$$\theta^2(\tilde{h}) = \sum_{j=1}^{\infty} |\tilde{a}_j|^2 < \infty.$$

Besides, in (29) we have  $d_j = a_j - \tilde{a}_j$ . So

$$q = \left[ \sum_{k=1}^{\infty} |a_k - \tilde{a}_k|^2 \right]^{1/2} \leq \psi(v).$$

Furthermore, according to (23) for an integer  $p$  satisfying the inequality

$$p > \rho/2 \tag{30}$$

we have

$$\varpi_p(h) = 2 \left[ (\theta(h) + \frac{1}{2^{1/\rho}})^{2p} + \zeta(2p/\rho) - 1 - \frac{1}{2^{2p/\rho}} \right]^{1/2p} \leq \pi_p(h),$$

where

$$\pi_p(h) := 2 \left[ (\psi(h) + \frac{1}{2^{1/\rho}})^{2p} + \zeta(2p/\rho) - 1 - \frac{1}{2^{2p/\rho}} \right]^{1/2p}.$$

Hence  $\chi_p(h, y) \leq \tau_p(h, y)$ , where

$$\tau_p(h, y) = \sum_{k=0}^{p-1} \frac{\pi_p^k(h)}{y^{k+1}} \exp \left[ \frac{1}{2} + \frac{\pi_p^{2p}(h)}{2y^{2p}} \right] \quad (y > 0).$$

Now Lemma 10 implies

**Theorem 5** *Let entire functions  $h$  and  $\tilde{h}$  satisfy the conditions  $h(0) = \tilde{h}(0) = 1$ , (26), (27) and (28). Then for any integer  $p$  satisfying inequality (30) we have  $rv_h(\tilde{h}) \leq x_p(v, h)$ , where  $x_p(v, h)$  is the unique positive root of the equation*

$$\psi(v)\tau_p(h, y) = 1.$$

Since  $\tau_p(h, \cdot)$  is a monotonically decreasing function, the latter theorem implies that all the zeros of  $\tilde{h}$  lie in the union of the sets  $\hat{W}_k(h, \psi(v))$  ( $k = 1, 2, 3, \dots$ ), where

$$\hat{W}_k := \left\{ \lambda \in \mathbf{C} : \psi(v)\tau_p \left( h, \left| \frac{1}{z_k(h)} - \frac{1}{\lambda} \right| \right) \geq 1 \right\}.$$

In particular if  $h$  has a finite number  $l$  of zeros, then all the zeros of  $\tilde{h}$  lie in the set

$$\cup_{k=0}^l \hat{W}_k(h, \psi(v)),$$

where

$$\hat{W}_0(h, \psi(v)) = \left\{ \lambda \in \mathbf{C} : \psi(v)\tau_p \left( h, \frac{1}{|\lambda|} \right) \geq 1 \right\}.$$

Note that

$$\tau_p \left( h, \frac{1}{|\lambda|} \right) = \sum_{k=0}^{p-1} \pi_p^k(h) |\lambda|^{k+1} \exp \left[ \frac{1}{2} (1 + \pi_p^{2p}(h) |\lambda|^{2p}) \right].$$

Furthermore, according to (24) we have the inequality

$$x_p(v, h) \leq \omega_p(v, h), \tag{31}$$

where

$$\omega_p(v, h) := \begin{cases} ep\psi(v) & \text{if } \pi_p(h) \leq ep\psi(v), \\ \pi_p(h) [\ln (\frac{\pi_p(h)}{\psi(v)^p})]^{-1/2p} & \text{if } \pi_p(h) > ep\psi(v) \end{cases}.$$

Now Theorem 5 yields the inequality  $rv_f(\tilde{h}) \leq \omega_p(v, h)$ .

If  $h$  has an infinite set of zeros, then according to Theorem 5 and (31), for any zero  $z(\tilde{h})$  of  $\tilde{h}$ , there is a zero  $z(h)$  of  $h$ , such that

$$|z(\tilde{h}) - z(h)| \leq x_p(v, h)|z(\tilde{h})z(h)| \leq \omega_p(v, h)|z(\tilde{h})z(h)|.$$

These relations imply the inequalities

$$|z(h)| - |z(\tilde{h})| \leq x_p(v, h)|z(\tilde{h})z(h)| \leq \omega_p(v, h)|z(\tilde{h})z(h)|.$$

Hence,

$$|z(\tilde{h})| \geq \frac{|z(h)|}{x_p(v, h)|z(h)| + 1} \geq \frac{|z(h)|}{\omega_p(v, h)|z(h)| + 1}.$$

## 11 An Estimate for the Difference of Two Solutions

Consider the equations

$$u'' = P(z)u, \tag{32}$$

and

$$\tilde{u}'' = \tilde{P}(z)\tilde{u} \quad (z \in \mathbf{C}), \tag{33}$$

where

$$P(z) = \sum_{k=0}^n c_k z^k \text{ and } \tilde{P}(z) = \sum_{k=0}^n \tilde{c}_k z^k \quad (c_n \neq 0; \tilde{c}_n \neq 0)$$

are polynomials with complex coefficients.

In the sequel it is assumed that

$$\tilde{u}(0) = u(0) = 1, \quad \tilde{u}'(0) = u'(0) = u_1 \quad (u_1 \in \mathbf{C}). \tag{34}$$

Recall that numbers  $\eta_P$  and  $\gamma_P$  are defined in Section 2.

**Lemma 12** *Let  $u(z)$  and  $\tilde{u}(z)$  be the solutions of (32) and (33), respectively. Let the initial conditions (34) hold. Then*

$$|u(z) - \tilde{u}(z)| \leq \sum_{k=0}^n \frac{|c_k - \tilde{c}_k| r^{k+2}}{(k+1)(k+2)} (1 + r|u_1|) \eta_P \eta_{\tilde{P}} \exp[(\gamma_P + \gamma_{\tilde{P}}) r^{n/2+1}]$$

$$(|z| \leq r, r > 0).$$

**Proof** Put  $w(z) = u(z) - \tilde{u}(z)$ . Then

$$w'' = P(z)w + (P(z) - \tilde{P}(z))\tilde{u}.$$

We have

$$|P(z) - \tilde{P}(z)| \leq \hat{v}(r) \quad (|z| = r, r > 0),$$

where

$$\hat{v}(r) := \sum_{k=0}^n |c_k - \tilde{c}_k| r^k.$$

Due to Corollary 2

$$M_{\tilde{u}}(r) \leq \eta_{\tilde{p}}(1 + r|u_1|) \exp[\gamma_{\tilde{p}} r^{n/2+1}].$$

Thus,

$$|(P(z) - \tilde{P}(z))\tilde{u}(z)| \leq \beta(r) \quad (|z| \leq r),$$

where

$$\beta(r) := \hat{v}(r)\eta_{\tilde{p}}(1 + r|u_1|) \exp[\gamma_{\tilde{p}} r^{n/2+1}].$$

Since  $w(0) = w'(0) = 0$ , due to Corollary 1,

$$\begin{aligned} M_w(r) &\leq \eta_P \exp[\gamma_P r^{n/2+1}] \int_0^r (r-s)\beta(s)ds \\ &\leq (1 + r|u_1|)\eta_P \eta_{\tilde{p}} \exp[(\gamma_P + \gamma_{\tilde{p}})r^{n/2+1}] \int_0^r (r-s)\hat{v}(s)ds. \end{aligned}$$

But

$$\int_0^r (r-s)\hat{v}(s)ds \leq \sum_{k=0}^n |c_k - \tilde{c}_k| \int_0^r (r-s)s^k ds = \sum_{k=0}^n |c_k - \tilde{c}_k| \frac{r^{k+2}}{(k+1)(k+2)}.$$

Thus  $w(z)$  satisfies the inequality

$$M_w(r) \leq \sum_{k=0}^n |c_k - \tilde{c}_k| \frac{r^{k+2}}{(k+1)(k+2)} (1 + r|u_1|)\eta_P \eta_{\tilde{p}} \exp[(\gamma_P + \gamma_{\tilde{p}})r^{n/2+1}],$$

as claimed.

Q.E.D.

Put  $b_0 = \eta_P \eta_{\tilde{P}}$ . Then due to the latter lemma we can write

$$|w(z)| = |u(z) - \tilde{u}(z)| \leq \sum_{k=2}^{n+3} g_k r^k \exp[(\gamma_P + \gamma_{\tilde{P}})r^{n/2+1}] \quad (|z| \leq r, r > 0), \quad (35)$$

where

$$g_2 = \frac{b_0}{2} |c_0 - \tilde{c}_0|, g_3 = b_0 \left( \frac{|c_1 - \tilde{c}_1|}{6} + |u_1| \frac{|c_0 - \tilde{c}_0|}{2} \right),$$

$$g_4 = b_0 \left( \frac{|c_2 - \tilde{c}_2|}{12} + |u_1| \frac{|c_1 - \tilde{c}_1|}{6} \right), \dots, g_{n+1} = b_0 \left( \frac{|c_n - \tilde{c}_n|}{(n+1)(n+2)} + |u_1| \frac{|c_{n-1} - \tilde{c}_{n-1}|}{(n-1)n} \right),$$

$$g_{n+3} = \frac{b_0}{(n+1)(n+2)} |u_1| |c_n - \tilde{c}_n|.$$

We thus arrive at our next result.

**Corollary 8** *Let  $u(z)$  and  $\tilde{u}(z)$  be the solutions of (32) and (33), respectively. Let the initial condition (34) hold. Then inequality (35) is valid.*

## 12 Perturbations of the Zeros of Solutions to ODEs

As above  $z_k(u)$  and  $z_k(\tilde{u})$  ( $k = 1, 2, \dots$ ) are the zeros of the solutions  $u$  and  $\tilde{u}$  to (32) and (33), respectively, enumerated in the non-decreasing order of their absolute values. To apply Theorem 5 to perturbations of the zeros of solutions to the considered equations, note that due to Corollary 2,

$$M_u(r) \leq \eta_P (1 + r|u_1|) \exp[\gamma_P r^{n/2+1}].$$

So taking in Theorem 5  $h(z) = u(z)$  and  $v(z) = w(z)$ , where  $w(z) = u(z) - \tilde{u}(z)$ , we have

$$B_u = \gamma_P, \rho = n/2 + 1 \text{ and } \xi_u = \xi_P,$$

where

$$\xi_P := (e\gamma_P(n/2 + 1))^{2/(n+2)}.$$

Besides,  $C_0 = \eta_P$ ,  $C_1 = \eta_P |u_1|$  and  $C_k = 0$  otherwise. Similarly, due to (35) we have

$$B_w = \gamma_P + \gamma_{\tilde{P}} \text{ and } \xi_w = \xi(P, \tilde{P}),$$

where

$$\xi(P, \tilde{P}) := (e(\gamma_P + \gamma_{\tilde{P}})(n/2 + 1))^{2/(n+2)}.$$

In addition,  $D_k = g_k$  ( $k = 2, \dots, n + 3$ ) and  $D_k = 0$  otherwise. If

$$e(\gamma_P + \gamma_{\tilde{P}})(n/2 + 1) < 1, \tag{36}$$

then  $\xi(P, \tilde{P}) < 1$  and  $\xi_P < 1$ , and therefore,  $\psi(w) = \Delta$ , where

$$\Delta := \left[ \sum_{j=1}^{\infty} (j!)^{4/(n+2)} \left( \sum_{k=2}^{n+1} g_k \frac{\xi^{j-k}(P, \tilde{P})}{((j-k)!)^{4/(n+2)}} \right)^2 \right]^{1/2} < \infty$$

and

$$\psi(u) = \eta_P \left[ \sum_{j=1}^{\infty} (j!)^{4/(n+2)} \left( \frac{\xi_P^j}{(j!)^{4/(n+2)}} + |u_1| \frac{\xi_P^{j-1}}{(j-1!)^{4/(n+2)}} \right)^2 \right]^{1/2} < \infty.$$

Furthermore, for an integer  $p$  satisfying the inequality

$$p > \frac{n+2}{4}, \tag{37}$$

we can write

$$\pi_p(u) = 2 \left[ (\psi(u) + \frac{1}{2^{2/(n+2)}})^{2p} + \zeta(4p/(n+2)) - 1 - \frac{1}{2^{4p/(n+2)}} \right]^{1/2p}$$

and

$$\tau_p(u, y) = \sum_{k=0}^{p-1} \frac{\pi_p^k(u)}{y^{k+1}} \exp \left[ \frac{1}{2} + \frac{\pi_p^{2p}(u)}{2y^{2p}} \right] \quad (y > 0).$$

Now Theorem 5 implies

**Theorem 6** *Let  $u(z)$  and  $\tilde{u}(z)$  be the solutions to Equations (32) and (33), respectively. Let conditions (34) and (36) hold. Then for an integer  $p$  satisfying the inequality (37) one has*

$$\text{rv}_u(\tilde{u}) \leq x_p(\Delta, u),$$

where  $x_p(\Delta, u)$  is the unique positive root of the equation

$$\Delta \cdot \tau_p(u, y) = 1.$$

According to Remark 2 condition (36) is not very restrictive. If it does not hold, i.e. if  $e(\gamma_p + \gamma_{\bar{p}})(n + 1/2) \geq 1$ , then one can apply the substitution  $z = w\hat{b}$  with

$$0 < \hat{b} = \text{const} < \frac{1}{e(\gamma_p + \gamma_{\bar{p}})(n/2 + 1)}.$$

Furthermore, according to (24) we have the inequality

$$x_p(\Delta, u) \leq \omega_p(\Delta, w), \tag{38}$$

where

$$\omega_p(\Delta, w) := \begin{cases} ep\Delta & \text{if } \pi_p(u) \leq ep\Delta, \\ \pi_p(u) [\ln(\pi_p(u)/(p\Delta))]^{-1/2p} & \text{if } \pi_p(u) > ep\Delta \end{cases}.$$

Now Theorem 6 yields the inequality

$$\text{rv}_u(\tilde{u}) \leq \omega_p(u, w).$$

Since  $\tau_p(u, \cdot)$  is a monotonically decreasing function, the latter theorem implies that all the zeros of  $\tilde{u}$  lie in the union of the sets  $\hat{W}_k(u, \Delta)$  ( $k = 1, 2, 3, \dots$ ), where

$$\hat{W}_k(u, \Delta) := \left\{ \lambda \in \mathbf{C} : \Delta \cdot \tau_p\left(u, \left| \frac{1}{z_k(u)} - \frac{1}{\lambda} \right| \right) \geq 1 \right\}.$$

In particular, if  $u$  has  $l < \infty$  zeros, then all the zeros of  $\tilde{u}$  lie in the set

$$\cup_{k=0}^l \hat{W}_k(u, \Delta),$$

where

$$\hat{W}_0(u, \Delta) = \{ \lambda \in \mathbf{C} : \Delta \cdot \tau_p(u, \frac{1}{|\lambda|}) \geq 1 \}.$$

Note that

$$\tau_p\left(u, \frac{1}{|\lambda|}\right) = \sum_{k=0}^{p-1} \pi_p^k(u) |\lambda|^{k+1} \exp\left[\frac{1}{2}(1 + \pi_p^{2p}(u) |\lambda|^{2p})\right].$$

### 13 Example to Theorem 6

To illustrate Theorem 6, consider the equations

$$u'' + a^2u = 0, u(0) = 1, u'(0) = 0 \quad (a = \text{const} > 0). \tag{39}$$

and

$$\tilde{u}'' + \tilde{a}^2\tilde{u} = 0, \tilde{u}(0) = 1, \tilde{u}'(0) = 0 \quad (\tilde{a} = \text{const} > 0). \tag{40}$$

Then  $u(z) = \cos(az)$ ,  $\tilde{u}(z) = \cos(\tilde{a}z)$  and their zeros are

$$\frac{\pi}{a}(m + 1/2) \text{ and } \frac{\pi}{\tilde{a}}(m + 1/2) \quad (m = 0, \pm 1, \pm 2, \dots),$$

respectively. Consequently,

$$\begin{aligned} \text{rv}_u(\tilde{u}) &= \sup_j \inf_k \left| \frac{1}{z_k(h)} - \frac{1}{z_j(\tilde{h})} \right| \\ &= \sup_{j=0, \pm 1, \pm 2, \dots} \inf_{k=0, \pm 1, \pm 2, \dots} \left| \frac{a}{\pi(k + 1/2)} - \frac{\tilde{a}}{\pi(j + 1/2)} \right| \\ &= \frac{|\tilde{a} - a|}{\pi} \sup_{j=0, \pm 1, \pm 2, \dots} \frac{1}{j + 1/2}. \end{aligned}$$

Hence,

$$\text{rv}_u(\tilde{u}) = 2 \frac{|\tilde{a} - a|}{\pi}. \tag{41}$$

In the considered case  $n = 0$ ,  $\gamma_P = a$ ,  $\gamma_{\tilde{P}} = \tilde{a}$ , and  $\eta_{\tilde{P}} = \eta_P = 1$ . In addition,  $\xi(P, \tilde{P}) = e(a + \tilde{a})$ , and  $\xi_P = ea$ . So condition (36) holds, provided

$$e(a + \tilde{a}) < 1. \tag{42}$$

This condition yields  $\psi(u) = \psi(\cos(az)) = \psi_1(a)$ , where

$$\psi_1(a) = \left[ \sum_{k=1}^{\infty} \xi_P^{2k} \right]^{1/2} = \left[ \sum_{k=1}^{\infty} (ae)^{2k} \right]^{1/2} = \frac{1}{[1/(ae)^2 - 1]^{1/2}}.$$

Moreover,  $g_2 = |a - \tilde{a}|/2$ ,  $g_k = 0$  for  $k \geq 3$ , and



$$\Delta = g_2 \left[ \sum_{k=2}^{\infty} (k(k-1)\xi^{k-2}(P, \tilde{P}))^2 \right]^{1/2} \leq g_2 \sum_{k=2}^{\infty} k(k-1)\xi^{k-2}(P, \tilde{P}).$$

Since

$$\sum_{k=2}^{\infty} x^{k-2} k(k-1) = \frac{d^2}{dx^2} \sum_{k=0}^{\infty} x^k = \frac{d^2}{dx^2} (1-x)^{-1} = \frac{6}{(1-x)^3} \quad (0 < x < 1),$$

we have  $\Delta \leq \Delta_1$ , where

$$\Delta_1 = \frac{6g_2}{(1-\xi(P, \tilde{P}))^3} = \frac{3|a-\tilde{a}|}{(1-(a+\tilde{a})e)^3}.$$

With  $p = 1$  the inequality (37) holds and  $\pi_1(u) = \hat{\pi}(a)$ , where

$$\hat{\pi}(a) := 2 \left[ \left( \psi_1(a) + \frac{1}{2} \right)^2 + \zeta(2) - 1 - \frac{1}{2^2} \right]^{1/2},$$

and  $\tau_1(u, y) = \hat{\tau}(a, y)$ , where

$$\hat{\tau}(a, y) := \frac{1}{y} \exp \left[ \frac{1}{2} + \frac{\hat{\pi}^2(a)}{2y^2} \right] \quad (y > 0).$$

Since  $\cos 0 = 1$  and  $\sin 0 = 0$ , making use of Theorem 6, we can assert the following result: *if condition (42) holds, then for solutions of (39) and (40) we have  $rv_u(\tilde{u}) \leq x(\Delta, a)$ , where  $x(\Delta, a)$  is the unique positive root of the equation*

$$\Delta_1 \cdot \hat{\tau}(a, y) = 1.$$

By (38)

$$x_1(\Delta, q) \leq \omega_1(\Delta_1),$$

where

$$\omega_1(\Delta_1) := \begin{cases} e\Delta_1 & \text{if } \hat{\pi}(a) \leq e\Delta_1, \\ \hat{\pi}(a) [\ln(\hat{\pi}(a)/\Delta_1)]^{-1/2} & \text{if } \hat{\pi}(a) > e\Delta_1. \end{cases}$$

So if  $\hat{\pi}(a) \leq e\Delta_1$ , then

$$rv_u(\tilde{u}) \leq e\Delta_1 = \frac{3e|a-\tilde{a}|}{(1-e(a+\tilde{a}))^3}.$$

For sufficiently small  $\tilde{a}$  and  $a$  this inequality is “close” to (41).

## References

1. D. Andrica, Th.M. Rassias (eds.), *Differential and Integral Inequalities*. Springer Optimization and Its Applications, vol. 151 (Springer, Switzerland, 2019)
2. N. Anghel, Stieltjes-Calogero-Gil' relations associated to entire functions of finite order. *J. Math. Phys.* **51**(5), 251–262 (2010)
3. N. Anghel, Entire functions of finite order as solutions to certain complex differential equations. *Proc. Am. Math. Soc.* **140**, 2319–2332 (2012)
4. S. Bank, A note on the zeros of solutions  $w'' - P(z)w = 0$  where  $P$  is a polynomial. *Appl. Anal.* **25**, 29–41 (1987)
5. S. Bank, A note on the location of complex zeros of solutions of linear differential equations. *Complex Variables Theory Appl.* **12**, 159–167 (1989)
6. S. Bank, On the complex zeros of solutions of linear differential equations. *Ann. Mat. Pura Appl.* **161**, 83–112 (1992)
7. F. Brügghemann, On the zeros of fundamental systems of linear differential equations with polynomial coefficients. *Complex Variables Theory Appl.* **15**, 159–166 (1990)
8. F. Brügghemann, On the solutions of linear differential equations with real zeros; proof of a conjecture of Hellerstein and Rossi. *Proc. Am. Math. Soc.* **113**, 371–379 (1991)
9. T.B. Cao, K. Liu, H.-Y. Xu, Bounds for the sums of zeros of solutions of  $u(m) = P(z)u$  where  $P$  is a polynomial. *Electron. J. Qual. Theory Differ. Equ.* **60**, 10 (2011)
10. J.L. Daleckii, M.G. Krein, *Stability of Solutions of Differential Equations in Banach Space* (American Mathematical Society, Providence, 1971)
11. A. Eremenko, S. Merenkov, Nevanlinna functions with real zeros. *Ill. J. Math.* **49**(4), 1093–1110 (2005)
12. L. Gao, On the growth of solutions of higher-order algebraic differential equations. *Acta Math. Scientia (B)* **22**(4), 459–465 (2002)
13. M. Gaudenzi, On the number of zeros of solutions of a linear differential equation. *J. Math. Anal. Appl.* **221**(1), 306–325 (1998)
14. M.I. Gil', Approximations of zeros of entire functions by zeros of polynomials. *J. Approx. Theory* **106**, 66–76 (2000)
15. M.I. Gil', Perturbations of zeros of a class of entire functions. *Complex Var.* **42**, 97–106 (2000)
16. M.I. Gil', *Operator Functions and Localization of Spectra*. Lectures Notes in Mathematics vol. 1830 (Springer, Berlin, 2003)
17. M.I. Gil', *Localization and Perturbation of Zeros of Entire Functions* (CRC Press, Taylor and Francis Group, New York, 2010)
18. M.I. Gil', Bounds for zeros of solutions of second order differential equations with polynomial coefficients. *Results Math.* **59**, 115–124 (2011)
19. M.I. Gil', Sums of zeros of solutions to second order ODE with non-polynomial coefficients. *Electron. J. Diff. Equ.* **2012**(107), 1–8 (2012)
20. M.I. Gil', Bounds for products of zeros of solutions to nonhomogeneous ODE with polynomial coefficients. *Int. J. Differ. Equ.* **2015**, 690519 (2015)
21. M.I. Gil', Sums of zeros of solutions to non-homogeneous ODE with polynomial coefficients. *J. Math. Anal. Appl.* **421**(2), 1917–1924 (2015)
22. M.I. Gil', Inequalities for zeros of solutions to second order ODE with one singular point. *Differ. Equ. Appl.* **8**(1), 69–76 (2016)
23. I.C. Gohberg, S. Goldberg, N. Krupnik, *Traces and Determinants of Linear Operators* (Birkhäuser, Basel, 2000)
24. I.C. Gohberg, M.G. Krein, *Introduction to the Theory of Linear Nonselfadjoint Operators*. Trans. Mathem. Monographs, vol. 18 (American Mathematical Society, Providence, 1969)
25. G.G. Gundersen, On the real zeros of solutions of  $f'' + A(z)f = 0$  where  $A(z)$  is entire. *Ann. Acad. Sci. Fenn. Ser. A.I. Math* **11**, 275–294 (1986)
26. S. Hellerstein, J. Rossi, Zeros of meromorphic solutions of second order linear differential equations. *Math. Z.* **192**, 603–612 (1986)

27. S. Hellerstein, J. Rossi, On the distribution of zeros of solutions of second-order differential equations. *Complex Var. Theory Appl.* **13**, 99–109 (1989)
28. S. Hellerstein, J. Rossi, Schwarzian derivatives and zeros of solutions to second order linear differential equations. *Proc. Am. Math. Soc.* **113**, 741–746 (1991)
29. E. Hille, *Lectures on Ordinary Differential Equations* (Addison-Wesley, Ontario, 1969)
30. C.Z. Huang, Real zeros of solutions of second order linear differential equations. *Kodai Math. J.* **14**, 113–122 (1991)
31. C.-H. Lin, Y. Sibuya, T. Tabara, Zeros of solutions of a second order linear differential equation with polynomial coefficients. *Funkc. Ekvacioj* **36**(2), 375–384 (1993)
32. I. Laine, *Nevanlinna Theory and Complex Differential Equations* (Walter de Gruyter Berlin, 1993)
33. M. Marden, *Geometry of Polynomials* (American Mathematical Society, Providence, 1985)
34. G.M. Muminov, On the zeros of solutions of the differential equation  $\omega^{(2m)} + p(z)\omega = 0$ . *Demonstr. Math.* **35**(1), 41–48 (2002)
35. P.M. Pardalos, T.M. Rassias (eds.), *Mathematics Without Boundaries. Surveys in Interdisciplinary Research*, vol. VIII (Springer, New York, 2014)
36. J. Rossi, The Tsuji characteristic and real zeros of second order ordinary differential equations. *J. Lond. Math. Soc.* **36**(2), 490–500 (1987)
37. J. Tu, Z.X. Chen, Zeros of solutions of certain second order linear differential equation. *J. Math. Anal. Appl.* **332**, 1279–291 (2007)

# Dynamics of a Higher-Order Ginzburg–Landau-Type Equation



Theodoros P. Horikis, Nikos I. Karachalios, and Dimitrios J. Frantzeskakis

**Abstract** We study possible dynamical scenarios associated with a higher-order Ginzburg–Landau-type equation. In particular, first we discuss and prove the existence of a limit set (attractor), capturing the long-time dynamics of the system. Then, we examine conditions for finite-time collapse of the solutions of the model at hand, and find that the collapse dynamics is chiefly controlled by the linear/nonlinear gain/loss strengths. Finally, considering the model as a perturbed nonlinear Schrödinger equation, we employ perturbation theory for solitons to analyze the influence of gain/loss and other higher-order effects on the dynamics of bright and dark solitons.

## 1 Introduction

In this work, our aim is to study the dynamics of a higher-order Ginzburg–Landau type equation. In particular, the model under consideration has the form of a higher-order nonlinear Schrödinger (NLS) equation incorporating gain and loss. The origin of our motivation is the following dimensionless higher-order NLS equation:

$$\partial_t u + \frac{i s}{2} \partial_x^2 u - i |u|^2 u = \beta \partial_x^3 u + \mu \partial_x (|u|^2 u) + (v - i \sigma_R) u \partial_x (|u|^2), \quad (1)$$

---

T. P. Horikis (✉)

Department of Mathematics, University of Ioannina, Ioannina, Greece  
e-mail: [horikis@uoi.gr](mailto:horikis@uoi.gr)

N. I. Karachalios

Department of Mathematics, University of Thessaly, Lamia, Greece  
e-mail: [karan@uth.gr](mailto:karan@uth.gr)

D. J. Frantzeskakis

Department of Physics, University of Athens, Athens, Greece  
e-mail: [dfrantz@phys.uoa.gr](mailto:dfrantz@phys.uoa.gr)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,  
Springer Optimization and Its Applications 173,  
[https://doi.org/10.1007/978-3-030-72563-1\\_9](https://doi.org/10.1007/978-3-030-72563-1_9)

187

where  $u(x, t)$  is a complex field,  $\beta$ ,  $\mu$ ,  $\nu$  and  $\sigma_R$  are real constants, while  $s = \pm 1$  stands for normal ( $s = +1$ ) or anomalous ( $s = -1$ ) group-velocity dispersion. Note that Equation (1) can be viewed as a perturbed NLS equation, with the perturbation—in case of small values of relevant coefficients—appearing in the right-hand side (see, e.g., Refs. [1–3]).

Variants of Equation (1) appear in a variety of physical contexts, where they are derived at higher-order approximations of perturbation theory [the lowest-order nonlinear model is simply the NLS equation in the left-hand side of Equation (1)]. The most prominent example is probably that of nonlinear optics [1–3]. In this case,  $t$  and  $x$  denote propagation distance and retarded time (in a reference frame moving with the group velocity), respectively, while  $u(x, t)$  is the complex electric field envelope.

For  $\beta = \mu = \nu = \sigma_R = 0$ , Equation (1) reduces to the unperturbed equation, i.e., the completely integrable NLS [4], which supports bright soliton solutions (for  $s = -1$ ) [5], or dark soliton solutions (for  $s = +1$ ) [6]. As concerns the origin of the higher-order terms, we mention the following. While the unperturbed NLS equation is sufficient to describe optical pulse propagation, for ultra-short pulses, third-order dispersion and self-steepening (characterized by coefficients  $\beta$ ,  $\mu$  and  $\nu$ , respectively) become important and have to be incorporated in the model. Similar situations also occur in other contexts and, thus, corresponding versions of Equation (2) have been derived and used, e.g., in the context of nonlinear metamaterials [7–9], but also in the problem of water waves in finite depth [10–12]. Moreover, in the context of optics, and for relatively long propagation distances, higher-order nonlinear dissipative effects, such as the stimulated Raman scattering (SRS) effect, of strength  $\sigma_R > 0$ , are also important [1–3].

In addition to the above mentioned effects, our aim is to investigate the dynamics in the framework of Equation (1), but also incorporating linear or nonlinear gain and loss. This way, we are going to analyze the following model:

$$\partial_t u + \frac{is}{2} \partial_x^2 u - i|u|^2 u = \gamma u + \delta |u|^2 u + \mu \partial_x (|u|^2 u) + \beta \partial_x^3 u + (\nu - i\sigma_R) u \partial_x (|u|^2), \quad (2)$$

which also incorporates dissipative effects, such as linear loss (for  $\gamma < 0$ ) [or gain (for  $\gamma > 0$ )]. These effects are physically relevant in the context of nonlinear optics [1–3, 13]: indeed, nonlinear gain ( $\delta > 0$ ) [or loss ( $\delta < 0$ )] may be used to counterbalance the effects from the linear loss/gain mechanisms, which may potentially lead to the stabilization of optical solitons—see, e.g., Refs. [14, 15]. Notice that it is the presence of gain/loss that renders Equation (2) a higher-order cubically nonlinear Ginzburg–Landau-type equation (see recent studies [16–18] on such models), featuring zero diffusion.

In this work, we will discuss various possible dynamical scenarios associated with Equation (2). In particular, the organization of the presentation and main results of this work can be described as follows.

In Section 2, first we show that the incorporation of gain and loss terms in the model gives rise to the existence of an attractor, capturing the long-time dynamics of the system. A rigorous proof is provided, based on the interpretation of the energy balance equation and properties of the functional (phase) space in which the problem defines an infinite-dimensional flow. It will also be discussed that although the gain/loss effects are pivotal for the dissipative nature of the infinite-dimensional flow that will be defined below, the structure of the attractor is basically determined by the other higher-order effects. In the same Section (Section 2), we also examine conditions for finite-time collapse of the solutions of the model. In particular, upon using energy arguments, we find that the collapse dynamics is chiefly controlled by the linear/nonlinear gain/loss strengths. We also identify a critical value of the linear gain, separating the possible decay of solutions to the trivial zero-state, from collapse.

In addition, considering the higher-order Ginzburg–Landau-type equation as a perturbed NLS equation, in Section 3 we study the dynamics of bright and dark solitons under the influence of the higher-order effects. The analysis is based on various perturbative techniques, relying on general aspects of the perturbation theory for bright and dark solitons. Specifically, we adopt the so-called adiabatic approximation, according to which the soliton form does not change due to the (small) perturbation, but its characteristics (center, amplitude, velocity, etc.) become unknown functions of time. We derive relevant evolution equations for the soliton characteristics and describe the pertinent soliton dynamics. We also briefly discuss still another method to analyze soliton dynamics, namely a multiscale expansion technique that asymptotically reduces the model to a Korteweg-de Vries–Burgers (KdV-B) equation. This way, we discuss various other nonlinear wave structures that can be supported by the higher-order effects, namely anti-dark solitons, as well as shock waves and rarefaction waves.

## 2 Limit Set and Collapse

### 2.1 Existence of the Limit Set

Let us consider the case  $s = -1$ , and supplement Equation (2) with periodic boundary conditions for  $u$  and its spatial derivatives up to the-second order, namely:

$$\begin{aligned} u(x + 2L, t) &= u(x, t), \quad \text{and} \\ \partial_x^j(x + 2L, t) &= \partial_x^j(x, t), \quad j = 1, 2, \end{aligned} \tag{3}$$

$\forall (x, t) \in \mathbb{R} \times [0, T]$ , for some  $T > 0$ , where  $L > 0$  is given. The initial condition

$$u(x, 0) = u_0(x), \quad \forall x \in \mathbb{R}, \tag{4}$$

also satisfies the periodicity conditions (3).

As shown in Ref. [19], all possible regimes except  $\gamma > 0, \delta < 0$ , are associated with finite-time collapse or decay. Furthermore, a critical value  $\gamma^*$  can be identified in the regime  $\gamma < 0, \delta > 0$ , which separates finite-time collapse from the decay of solutions. On the other hand, for  $\gamma > 0, \delta < 0$ , below we prove the existence of a limit set (attractor)  $\omega(\mathcal{B})$ , attracting all bounded orbits initiating from arbitrary, appropriately smooth initial data  $u_0$  (considered elements of a bounded set  $\mathcal{B}$  in a suitable Sobolev space). Notice that, as shown numerically in Ref. [20], the attractor  $\omega(\mathcal{B})$  captures the full route, ranging from Poincaré–Bendixson limit-cycle dynamics to quasiperiodic or chaotic dynamics.

The starting point of our proof is the power balance equation:

$$\frac{d}{dt} \int_{-L}^L |u|^2 dx = 2\gamma \int_{-L}^L |u|^2 dx + 2\delta \int_{-L}^L |u|^4 dx, \tag{5}$$

satisfied by any local solution  $u \in C([0, T], H_{per}^k(\Omega))$ , which initiates from sufficiently smooth initial data  $u_0 \in H_{per}^k(\Omega)$ , for fixed  $k \geq 3$ . Here,  $H_{per}^k(\Omega)$  denotes the Sobolev spaces of periodic functions  $H_{per}^k$  [21], in the fundamental interval  $\Omega = [-L, L]$ . Analysis of (5), results in the asymptotic estimate:

$$\limsup_{t \rightarrow \infty} \frac{1}{2L} \int_{-L}^L |u(x, t)|^2 dx \leq -\frac{\gamma}{\delta}, \tag{6}$$

implying that local in time solutions  $u \in C([0, T], H_{per}^k(\Omega))$  are uniformly bounded in  $L^2(\Omega)$ . This allows for the definition of the extended dynamical system:

$$\varphi(t, u_0) : H_{per}^k(\Omega) \rightarrow L^2(\Omega), \quad \varphi(t, u_0) = u,$$

whose orbits are bounded  $\forall t \geq 0$ . Moreover, from the above asymptotic estimate, we derive the following: if  $L^2(\Omega)$  is endowed with the equivalent averaged norm

$$\|u\|_\alpha^2 = \frac{1}{2L} \int_{-L}^L |u|^2 dx,$$

then its ball:

$$\mathcal{B}_\alpha(0, \rho) = \left\{ u \in L^2(\Omega) : \|u\|_\alpha^2 \leq \rho^2, \quad \rho^2 > -\frac{\gamma}{\delta} \right\},$$

attracts all bounded sets  $\mathcal{B} \in H_{per}^k(\Omega)$ . That is, there exists  $T^* > 0$ , such that  $\varphi(t, \mathcal{B}) \subset \mathcal{B}_\alpha$ , for all  $t \geq T^*$ . Thus, we may define for any bounded set  $\mathcal{B} \in H_{per}^k(\Omega)$ ,  $k \geq 3$ , its  $\omega$ -limit set in  $L^2(\Omega)$ ,

$$\omega(\mathcal{B}) = \bigcap_{s \geq 0} \overline{\bigcup_{t \geq s} \varphi(t, \mathcal{B})}.$$

The closures are taken with respect to the weak topology of  $L^2(\Omega)$ . Then, the standard (embedding) properties of Sobolev spaces imply that the attractor  $\omega(\mathcal{B})$  is at least weakly compact in  $L^2(\Omega)$ , or relatively compact in the dual space  $H_{per}^{-1}(\Omega)$ .

In the direct numerical simulations of Ref. [20], it was found that apart from the gain/loss parameters  $\gamma$  and  $\delta$ , the other higher-order effects play also important role on the dynamics. In particular, the competition between the third-order dispersion (characterized by the coefficient  $\beta$ ) and SRS effect (characterized by the coefficient  $\sigma_R$  gives rise to rich dynamics (briefly mentioned above): indeed, the dynamics ranges from Poincaré–Bendixson-type scenarios, in the sense that bounded solutions may converge either to distinct equilibria via orbital connections or to space-time periodic solutions, to the emergence of almost periodic and chaotic behavior. A main result is that third-order dispersion has a dominant role in the development of such complex dynamics, since it can be chiefly responsible (even in the absence of other higher-order effects) for the existence of periodic, quasiperiodic, and chaotic spatiotemporal structures.

We conclude by illustrating some representative results illustrating the richness of these dynamics.

Figure 1 depicts the birth of a space time periodic solution emerging from the modulation instability of the continuous wave (cw) steady-state solution of amplitude  $|\phi_b|^2 = -\frac{\gamma}{\delta}$  for the choice of parameters  $\beta = 0.55$ ,  $\sigma_R = 0.01$ ,  $\gamma = 1.5$ ,  $\delta = -1$ ,  $\sigma_R = 0.3$ ,  $\mu = \nu = 0.01$ . The initial condition is a single-mode cw of the form

$$u_0(x) = \epsilon e^{-i \frac{K\pi x}{L}}, \quad K > 0. \tag{7}$$

with  $K = 5$  and  $\epsilon = 0.01$ . This is one of the examples showing the Poincaré–Bendixson type dynamics when the instability of a steady state gives rise to the birth of a limit-cycle. The results visualise the asymptotic behavior in the  $2D$ -finite dimensional subspace

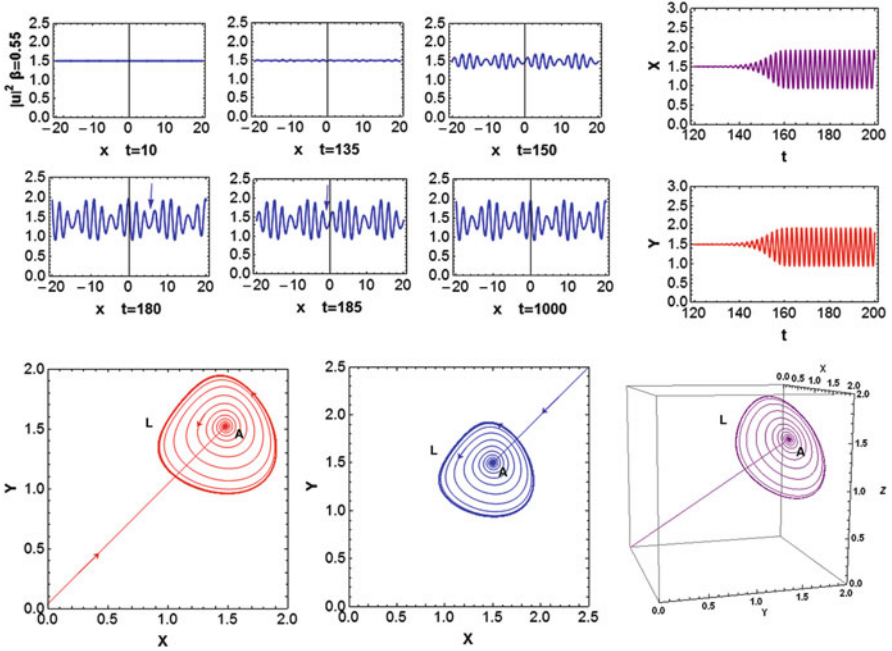
$$\mathcal{P}_2 = \{(X, Y) \in \mathbb{R}^2 : (X(t), Y(t)) = (|u(x_1, t)|^2, |u(x_2, t)|^2), \quad x_1, x_2 \in \Omega, t \geq 0\},$$

for some arbitrarily chosen fixed spatial coordinates  $x_1, x_2$ . In this subspace, the steady-state  $\phi_b$  is marked by the point  $\mathbf{A} = (|\phi_b|^2, |\phi_b|^2) = (-\frac{\gamma}{\delta}, -\frac{\gamma}{\delta})$ . The  $3D$ -counterpart is defined as

$$\mathcal{P}_3 = \{(X, Y, Z) \in \mathbb{R}^3 : (X(t), Y(t), Z(t)) = (|u(x_1, t)|^2, |u(x_2, t)|^2, |u(x_3, t)|^2), \quad x_1, x_2, x_3 \in \Omega, t \geq 0\}. \tag{8}$$

The emergence of limit-cycles characterizing the global attractor persists up to certain thresholds for the parameter  $\beta$ .



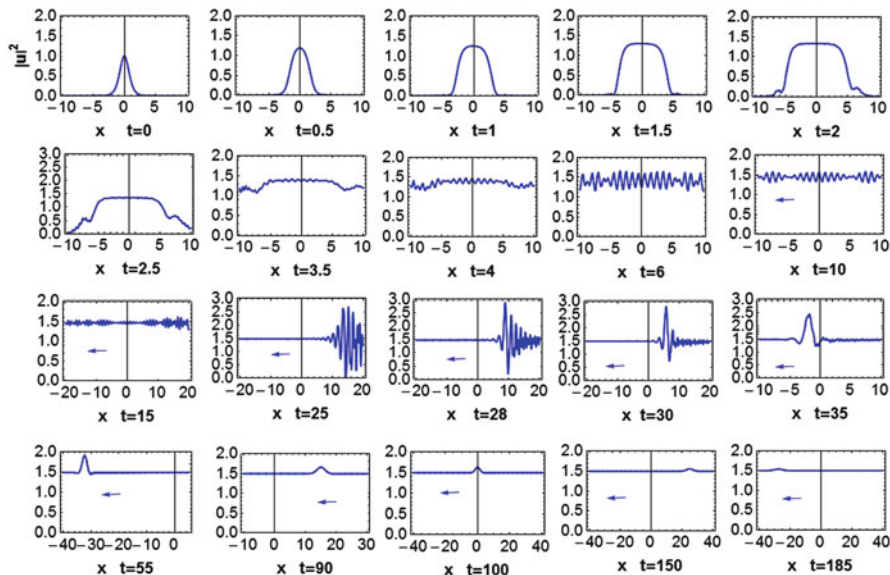


**Fig. 1** (Color Online) Upper left panel: The birth of a space-time periodic solution from the instability of the cw-steady state of amplitude  $|\phi_b|^2 = -\gamma/\delta$ , for cw-initial data (7) of  $K = 5$  and  $\epsilon = 0.01$ . Parameters  $\beta = 0.55$ ,  $\sigma_R = 0.01$ ,  $\gamma = 1.5$ ,  $\delta = -1$ ,  $\sigma_R = 0.3$ ,  $\mu = \nu = 0.01$ . Upper right panel: Integral curves  $X(t) = |u(x_1, t)|^2$  (upper fig.-purple curve) and  $Y(t) = |u(x_2, t)|^2$  (bottom fig.-red curve), for the spatial coordinates  $x_1 = 0$  and  $x_2 = 5$  respectively: Convergence to a periodic solution, for the set of parameters of the upper left panel. Bottom left panel: The space-time periodic solution of the left upper panel, as a limit cycle in the 2D-phase space  $\mathcal{P}_2$  for the spatial coordinates  $x_1 = 0$  and  $x_2 = 5$ . Bottom middle panel: Convergence to the limit cycle of the bottom left panel for the cw-initial condition of  $K = 5$  and  $\epsilon = \sqrt{3}$ . Bottom right panel: The space-time periodic solution of the upper right panel as a limit cycle on the 3D-phase space  $\mathcal{P}_3$  (defined by (8)), for the choice of spatial coordinates  $x_1 = 0$ ,  $x_2 = 5$  and  $x_3 = 10$

On the other hand, even when the steady state  $\phi_b$  is asymptotically stable, the convergence may include highly non-trivial transient dynamics. Figure 2 depicts an example of the evolution of the initial condition

$$u_0(x) = \epsilon \operatorname{sech} x. \quad (9)$$

Such initial data correspond to the profile of a “bright soliton” as an initial state. The example is for  $\epsilon = 1$  and parameters  $\sigma_R = \mu = \nu = \beta = 0.01$ . The gain/loss strengths are  $\gamma = 1.5$ ,  $\delta = -1$ . We observe formation of a “shock-wave” transitioning to an unstable periodic solution and then, the formation of a decaying travelling pulse, prior to the ultimate convergence to the asymptotically stable state  $\phi_b$ .



**Fig. 2** (Color Online) Snapshots of the evolution of the density  $|u|^2$  for initial data (9) with  $\epsilon = 1$ , when  $\sigma_R = \mu = \nu = \beta = 0.01$  and  $\gamma = 1.5, \delta = -1$ . Formation of a decaying “bright” traveling solitary pulse, prior to the ultimate convergence to the steady state  $\phi_b$ . The array indicates the direction of the travelling pulses

## 2.2 Conditions for Collapse

The question of collapse concerns sufficiently smooth (weak) solutions of Equation (2). The existence of such solutions, is guaranteed by a local existence result associated with the initial-boundary value problem (2)–(4). In particular, the methods which are used in order to prove such a local existence result in the Sobolev spaces of periodic functions  $H^k_{per}$  [22, 23], are based on the lines of approach of [24–27]. The application of these methods to establish local existence for the problem (2)–(4), although involving lengthy computations, is now considered as standard. Thus, we refrain from showing the details here, and we just state it in:

**Theorem 1** *Let  $u_0 \in H^k_{per}(\Omega)$  for any integer  $k \geq 2$ , and  $\beta, \gamma, \delta, \mu, \nu, \sigma_R \in \mathbb{R}$ . Then there exists  $T > 0$ , such that the problem (2)–(4), has a unique solution*

$$u \in C([0, T], H^k_{per}(\Omega)) \quad \text{and} \quad u_t \in C([0, T], H^{k-3}_{per}(\Omega)).$$

Moreover, the solution  $u \in H^k_{per}(\Omega)$  depends continuously on the initial data  $u_0 \in H^k_{per}(\Omega)$ , i.e., the solution operator

$$\mathcal{S}(t) : H^k_{per}(\Omega) \mapsto H^k_{per}(\Omega), \quad t \in [0, T], \tag{10}$$

$$u_0 \mapsto \mathcal{S}(t)u_0 = u,$$

is continuous.

Here, for the sake of completeness, let us recall the definition of  $H_{per}^k(\Omega)$ :

$$H_{per}^k(\Omega) = \{u : \Omega \rightarrow \mathbb{C}, \quad u \text{ and } \frac{\partial^j u}{\partial x^j} \in L^2(\Omega), \quad j = 1, 2, \dots, k; \\ u(x), \text{ and } \frac{\partial^j u}{\partial x^j}(x) \text{ for } j = 1, 2, \dots, k - 1, \text{ are } 2L\text{-periodic}\}. \quad (11)$$

Since our analytical energy arguments for examining collapse require sufficiently smoothness of local-in-time solutions, we shall implement Theorem 1 by assuming that  $k = 3$ , at least. As it follows from the definition of the Sobolev spaces (2.2), this assumption means that the initial condition  $u_0(x)$ ,  $x \in \Omega$ , and its spatial (weak) derivatives, at least up to the 2nd-order, are  $2L$ -periodic. Then, it turns out from Theorem 1, that the unique, local-in-time solution  $u(x, t)$  of (2) satisfies the periodic boundary conditions (3) for  $t \geq 0$ , and is sufficiently (weakly) smooth. Such periodicity and smoothness properties of the local-in-time solutions are sufficient for our purposes (see Theorem 2 below).

Next, we adopt the method of deriving a differential inequality for the functional

$$M(t) = \frac{e^{-2\gamma t}}{2L} \int_{-L}^L |u(x, t)|^2 dx, \quad (12)$$

and then, showing that its solution diverges in finite-time under appropriate assumptions on its initial value at time  $t = 0$ ; see [22, 23, 28–30] and references therein. Note that the choice of this functional is not arbitrary; in fact, it is a direct consequence of the conservation laws of the NLS model. For a discrete counterpart of this argument applied in discrete Ginzburg–Landau-type equations, we refer to [31]. For applications of these types of arguments in the study of escape dynamics for Klein–Gordon chains, we refer to [32].

**Theorem 2** For  $u_0 \in H_{per}^k(\Omega)$ ,  $k \geq 3$ , let  $\mathcal{S}(t)u_0 = u \in C([0, T_{\max}), H_{per}^k(\Omega))$  be the local- in- time solution of the problem (2)–(4), with  $[0, T_{\max})$  be its maximal interval of existence. Assume that the parameter  $\delta > 0$  and that the initial condition  $u_0(x)$  is such that  $M(0) > 0$ . Then,  $T_{\max}$  is finite, in the following cases:

$$(i) \quad T_{\max} \leq \frac{1}{2\gamma} \log \left[ 1 + \frac{\gamma}{\delta M(0)} \right], \quad (13)$$

$$\text{for } \gamma \neq 0, \text{ and } \gamma > -\delta M(0). \quad (14)$$

$$(ii) \quad T_{\max} \leq \frac{1}{2\delta M(0)}, \text{ for } \gamma = 0. \quad (15)$$

**Proof** For any  $T < T_{\max}$ , since  $k \geq 3$ , due to the continuous embedding [22]:

$$C([0, T], H_{per}^k(\Omega)) \subset C([0, T], L^2(\Omega)),$$

the solution  $\mathcal{S}(t)u_0 = u \in C([0, T], L^2(\Omega))$ . Furthermore, it follows from Theorem 1, that  $u_t \in C([0, T], L^2(\Omega))$ . Then, by differentiating  $M(t)$  with respect to time, we find that:

$$\frac{dM(t)}{dt} = -\gamma \frac{e^{-2\gamma t}}{L} \int_{-L}^L |u|^2 dx + \frac{e^{-2\gamma t}}{L} \operatorname{Re} \int_{-L}^L u_t \bar{u} dx. \tag{16}$$

In the second term on the right-hand side of (16), we substitute  $u_t$  by the right-hand side of Equation (2). Then, after some computations, Equation (16) results in:

$$\frac{dM(t)}{dt} = \delta \frac{e^{-2\gamma t}}{L} \int_{-L}^L |u|^4 dx. \tag{17}$$

Next, by the Cauchy-Schwarz inequality, we have

$$\int_{-L}^L |u|^2 dx \leq \sqrt{2L} \left( \int_{-L}^L |u|^4 dx \right)^{1/2}. \tag{18}$$

Therefore, for the functional  $M(t)$  defined in (12), we get the inequality

$$M(t) \leq \frac{e^{-2\gamma t}}{\sqrt{2L}} \left( \int_{-L}^L |u|^4 dx \right)^{1/2}, \tag{19}$$

which in turns, implies the estimate

$$M(t)^2 \leq \frac{e^{-4\gamma t}}{2L} \int_{-L}^L |u|^4 dx, \tag{20}$$

for all  $t \in [0, T_{\max})$ . On the other hand, from (17) we have that

$$\int_{-L}^L |u|^4 dx = e^{2\gamma t} \frac{L}{\delta} \frac{dM(t)}{dt},$$

and hence, we may rewrite (20) as

$$[M(t)]^2 \leq \frac{e^{-2\gamma t}}{2\delta} \frac{dM(t)}{dt}, \quad \text{or} \quad \frac{\frac{dM(t)}{dt}}{[M(t)]^2} \geq 2\delta e^{2\gamma t}. \tag{21}$$

Since

$$\frac{d}{dt} \left[ \frac{1}{M(t)} \right] = - \frac{\frac{dM(t)}{dt}}{[M(t)]^2},$$

we get from (21) the differential inequality

$$\frac{d}{dt} \left[ \frac{1}{M(t)} \right] \leq -2\delta e^{2\gamma t}. \quad (22)$$

Integration of (22) with respect to time, implies that

$$\frac{1}{M(t)} \leq \frac{1}{M(0)} - 2\delta \int_0^t e^{2\gamma s} ds,$$

and since  $M(t) > 0$ , we see that  $M(0) > 0$  satisfies the inequality

$$2\delta \int_0^t e^{2\gamma s} ds \leq \frac{1}{M(0)}. \quad (23)$$

From (23), we shall distinguish between the following cases for the damping parameter  $\gamma$ :

- We assume that the damping parameter  $\gamma \neq 0$ . In this case, (23) implies that

$$\frac{2\delta}{2\gamma} \left( e^{2\gamma t} - 1 \right) \leq \frac{1}{M(0)}, \quad \text{or} \quad e^{2\gamma t} \leq 1 + \frac{\gamma}{\delta M(0)}.$$

Thus, assuming that

$$\frac{\gamma}{\delta M(0)} > -1,$$

we derive that the maximal time of existence is finite, since

$$t \leq \frac{1}{2\gamma} \log \left[ 1 + \frac{\gamma}{\delta M(0)} \right].$$

The inequality above, proves the estimate of the collapse time (13) under assumption (14), that is, case (i) of the Theorem.

- Assume that the damping parameter is  $\gamma = 0$ . Then, Equation (23) implies that

$$2\delta t \leq \frac{1}{M(0)},$$

or

$$t \leq \frac{1}{2\delta M(0)}.$$

This inequality proves the estimate of the collapse time (15) in the absence of damping, that is, case (ii) of the Theorem.  $\square$

From condition (14) on the definition of the analytical upper bound of the blow-up time

$$T_{\max}[\gamma, \delta, M(0)] = \frac{1}{2\gamma} \log \left[ 1 + \frac{\gamma}{\delta M(0)} \right], \tag{24}$$

given in (13), we define a critical value of the linear gain/loss parameter as

$$\gamma^* = -\delta M(0). \tag{25}$$

We observe that

$$\lim_{\gamma \rightarrow \gamma^*} T_{\max}[\gamma, \delta, M(0)] = +\infty, \tag{26}$$

while  $T_{\max}[\gamma, \delta, M(0)]$  is finite if

$$\gamma > \gamma^*, \tag{27}$$

according to the condition (14). Then, (26) suggests that when  $\delta > 0$ , the critical value  $\gamma^*$  may act as a critical point separating the two dynamical behaviors: blow-up in finite-time for  $\gamma > \gamma^*$  and global existence for  $\gamma < \gamma^*$ .

We also remark that the analytical upper bound for the blow-up time (15) in the case  $\gamma = 0$ ,

$$T_{\max}[\delta, M(0)] = \frac{1}{2\delta M(0)}, \tag{28}$$

is actually the limit of the analytical upper bound (24) for  $\gamma > 0$  as  $\gamma \rightarrow 0$ , e.g.,

$$\lim_{\gamma \rightarrow 0} T_{\max}[\gamma, \delta, M(0)] = T_{\max}[\delta, M(0)]. \tag{29}$$

The analytical estimates for the blow-up time have been proved sharp with respect to their dependence on the several parameters as it was illustrated by the relevant numerical simulations [19].

### 3 Soliton Dynamics: Perturbative Approach

Below, our aim is to consider the higher-order Ginzburg–Landau equation (2) as a perturbed NLS equation. This can be done, upon rewriting Equation (2) in the following form,

$$iu_t - \frac{s}{2}u_{xx} + |u|^2u = \varepsilon F[u], \quad (30)$$

where subscripts denote partial derivatives, and the functional perturbation  $F[u]$  is given by:

$$F[u] = i\gamma u + i\delta|u|^2u + i\beta u_{xxx} + i\mu(|u|^2u)_x + (i\nu + \sigma_R)u(|u|^2)_x. \quad (31)$$

In other, words, we consider the situation where the coefficients of the gain/loss and higher-order terms are small, i.e., of the order of a formal small parameter  $\varepsilon$  (with  $0 < \varepsilon \ll 1$ ). This problem finds applications in long-haul optical fiber communications, where the terms involved in  $F[u]$  can indeed be considered as small perturbations [1].

Based on the fact that, for  $\varepsilon = 0$ , Equation (30) becomes the traditional NLS model that possesses bright or dark solitons for  $s = -1$  and  $s = +1$  respectively, we will study separately these two cases, and investigate how the perturbation (31) alters the soliton dynamics. Our analysis relies on various perturbation techniques that have been developed in the past, both for bright [33–35] and dark [36–38], including the perturbed inverse scattering method, the variational approach (or Lagrangian method), the Lie transform method, and others (see also [1–3] and references therein). Among these techniques, a particularly convenient method to study the soliton dynamics is the so-called adiabatic approach. According to this, an adiabatic relation is the balance between nonlinearity and dispersion, so that (amplitude)×(width)=const. In other words, it is assumed that—in the presence of the perturbations—the functional form of the soliton remains unchanged, but the soliton parameters change (slowly) as the soliton evolves. Thus, the soliton parameters are treated as unknown functions of  $t$ , and their evolution is determined by the evolution of the conserved quantities (integrals of motion) of the unperturbed NLS. Particularly relevant such conserved quantities include the energy:

$$E = \int_{-\infty}^{\infty} |u|^2 dx, \quad (32)$$

the momentum,

$$P = \frac{i}{2} \int_{-\infty}^{\infty} (u\bar{u}_x - \bar{u}u_x) dx, \quad (33)$$

where overbar denotes complex conjugate, and the Hamiltonian:

$$H = \frac{1}{2} \int_{-\infty}^{-\infty} (s|u_x|^2 + |u|^4) dx. \tag{34}$$

In addition, for our considerations below, it is also useful to introduce still another conserved quantity, namely the central position of the soliton(s)—alias “soliton center”—given by:

$$R_{bs} = \int_{-\infty}^{+\infty} x|u|^2 dx, \quad R_{ds} = \int_{-\infty}^{+\infty} x(u_\infty^2 - |u|^2) dx, \tag{35}$$

for the bright and dark solitons respectively.

### 3.1 Perturbation Theory for Bright Solitons ( $s = -1$ )

We start with the case of  $s = -1$ , i.e., the case of bright solitons. First, using Equation (30) and its complex conjugate, it is straightforward to derive the following equations for the evolution of the NLS conserved quantities under the action of the perturbation:

$$\frac{dE}{dt} = \varepsilon \int_{-\infty}^{+\infty} (\bar{u}F + u\bar{F}) dx, \tag{36}$$

$$\frac{dP}{dt} = \varepsilon i \int_{-\infty}^{+\infty} (\bar{u}_x F - u_x \bar{F}) dx, \tag{37}$$

$$\frac{dH}{dt} = 2\varepsilon \int_{-\infty}^{+\infty} \left[ \left( \frac{1}{2} \bar{u}_{xx} + |u|^2 \bar{u} \right) F + \left( \frac{1}{2} u_{xx} + |u|^2 u \right) \bar{F} \right] dx. \tag{38}$$

For sufficiently small perturbation, the form of the bright soliton solution  $u_{bs}(x, t)$  may be assumed to have the following rather general form, where all its parameters are allowed to vary in  $t$  as

$$u_{bs}(x, t) = \eta(t) \operatorname{sech}[\eta(t)(x - x_0(t))] \exp[-i\kappa(t)x + i\phi(t)], \tag{39}$$

where the soliton’s amplitude (and inverse width)  $\eta$ , its central position  $x_0$ , the wavenumber  $\kappa$ , and phase  $\phi$  are unknown functions of  $t$  that have to be determined. Notice that, in the absence of the perturbation,  $x_0$  and  $\phi$  are constants, given by:

$$\frac{dx_0}{dt} = -\kappa, \quad \frac{d\phi}{dt} = \frac{1}{2} (\eta^2 - \kappa^2). \tag{40}$$

Substituting the soliton (39) into Equations (37)–(38) [and (35)], we obtain a set of four ordinary differential equations (ODEs) for the four unknown soliton parameters:



$$\frac{d\eta}{dt} = -\text{Im} \left\{ \int_{-\infty}^{\infty} F[u] \bar{u} dx \right\}, \quad (41)$$

$$\frac{d\kappa}{dt} = \text{Re} \left\{ \int_{-\infty}^{\infty} F[u] \tanh[\eta(x - x_0)] \bar{u} dx \right\}, \quad (42)$$

$$\frac{dt_0}{dt} = -\kappa - \frac{1}{\eta^2} \text{Im} \left\{ \int_{-\infty}^{\infty} F[u] (x - x_0) \bar{u} dx \right\}, \quad (43)$$

$$\begin{aligned} \frac{d\phi}{dt} = & \frac{1}{2}(\eta^2 - \kappa^2) + x_0 \frac{d\kappa}{dt} \\ & - \text{Re} \left\{ \int_{-\infty}^{\infty} F[u] \left[ \frac{1}{\eta} - (x - x_0) \tanh[\eta(x - x_0)] \right] \bar{u} dx \right\}, \end{aligned} \quad (44)$$

where Re and Im stand for the real and imaginary parts, respectively.

Before analyzing the full problem, where the perturbation  $F[u]$  is given by Equation (31), it is relevant to consider at first a simple example. In particular, let the linear loss/gain term be a small perturbation, i.e.,  $F[u] = i\gamma u$ , with  $\gamma \ll 1$ , and assume that  $\delta = \beta = \nu = \sigma_R = 0$ . Then, substituting this form of  $F[u]$  into Equations (41)–(44), and performing the integrations, it is found that the soliton wavenumber  $\kappa$  and the central position  $x_0$  remain unaffected of the perturbation, while the soliton amplitude  $\eta$  and phase  $\phi$  evolve, due to the presence of the loss/gain, as follows:

$$\eta(t) = \exp(2\gamma t), \quad \phi(t) = \phi(0) - \frac{1}{8\gamma} [1 - \exp(4\gamma t)]. \quad (45)$$

To obtain the above result, it was assumed that  $\eta(0) = 1$  and  $\kappa(0) = x_0(0) = 0$  (hence  $\kappa(t) = x_0(t) = 0 \forall t$ ). Thus, in the presence of loss,  $\gamma < 0$  (or gain,  $\gamma > 0$ ) the soliton amplitude decreases (or increases), while its width increases (or decreases), i.e., the soliton broadens (or is compressed) during its evolution.

We now return to the full problem, and study the effect of the perturbation (31) on the dynamics of bright solitons. Following the same procedure, i.e., substituting Equation (31) into Equations (41)–(44), and performing the integrations, we find that the soliton parameters evolve according to the following system:

$$\frac{d\eta}{dt} = \frac{2}{3}\eta(3\gamma + 2\delta\eta^2), \quad (46)$$

$$\frac{d\kappa}{dt} = -\frac{8}{15}\sigma_R\eta^4, \quad (47)$$

$$\frac{dt_0}{dt} = -\kappa + \frac{1}{3}(3\beta - 3\mu - 2\nu)\eta + 3\beta\kappa^2\eta^2, \quad (48)$$

$$\frac{d\phi}{dt} = -\kappa \left[ (\mu - 3\beta)\eta^2 + \beta(\eta^2 - 2)q^2 \right] + \frac{1}{2}(\eta^2 - \kappa^2) - \frac{8}{15}\sigma_R t_0 \eta^4. \quad (49)$$

Although the result in this case is more complicated, it is still possible to arrive at a simple analytical result. Indeed, first observe that Equation (46) can be solved analytically to provide the functional form of  $\eta(z)$ , which is found to be:

$$\eta^2(t) = \frac{3C\gamma e^{4\gamma t}}{1 - 2C\delta e^{4\gamma t}}, \quad C = \frac{\eta^2(0)}{3\gamma + 2\delta\eta^2(0)}. \tag{50}$$

Then, the wavenumber  $\kappa(t)$  can be obtained from Equation (47) by simply integrating the above expression for  $\eta$ . Finally, having found  $\eta(t)$  and  $\kappa(t)$ , integration of Equations (48) and (49) yield, respectively, the functional forms of  $x_0(z)$  and  $\phi(t)$ .

### 3.2 Perturbation Theory for Dark Solitons ( $s = +1$ )

In this section, we consider the case  $s = +1$ , and provide analytical results based on the perturbation theory for dark solitons devised in Ref. [38]. We start by noting that, for  $\varepsilon = 0$ , the unperturbed defocusing NLS Equation (30) possesses the following single dark soliton solution:

$$u_{ds}(x, t) = [A + iB \tanh(BX)]e^{i\sigma_0}, \tag{51}$$

where  $X = x - X_0$ , with  $X_0 = x - \int_0^t A(s)ds - x_0$  being the dark soliton center,  $A^2 + B^2 = u_\infty^2$ , and  $\Delta\phi_0 = 2 \tan^{-1}(B/A)$  is the phase jump across the soliton. The latter, is equal to  $\pi$  for stationary, so-called “black” solitons with  $A = 0$  (moving solitons with  $A \neq 0$  are termed “grey”) [39, 40]; finally,  $A$  and  $B$  depict the velocity and depth of the dark soliton, respectively, while  $x_0$  and  $\sigma_0$  are real parameters. Notice that  $u_\infty$  represents the boundary condition at infinity, i.e.,  $u_\infty = u(x \rightarrow \infty)$ , and sets the amplitude of the soliton background. The dark soliton (51) is, therefore, comprised of a background of constant density, and a density dip that propagates on top, accompanied by a phase jump across the minimum density.

The effect of the perturbation of Equation (31) on the dark soliton dynamics will now be studied upon assuming that the soliton parameters are slowly varying functions of  $t$ . As shown in Ref. [38], dissipative terms—similar to the ones considered here—give rise to a *shelf*, which develops around the soliton; the shelf has a non-trivial contribution to the integrals employed in order to find expressions for the soliton parameters. Thus, this perturbative approach is better suited here, compared to ones merely relying on the adiabatic approximation [36, 37], as they do not take into account this important contribution.

Our analysis starts with the dynamics of the soliton background. Assuming that  $u(x \rightarrow \infty) = u_0(t)$ , we derive from Equation (30) the equation:

$$iu_{0t} - |u_0|^2 u_0 = i\gamma u_0 + i\delta |u_0|^2 u_0. \tag{52}$$

Then, employing the polar decomposition  $u_0 = u_\infty(t) \exp(i\theta(t))$ , we obtain:

$$u'_\infty = (\gamma + \delta u_\infty^2)u_\infty, \quad \theta' = u_\infty^2, \quad (53)$$

where primes denote differentiation with respect to  $t$ . The role of the term of strength  $\delta$  is now more obvious: a nontrivial equilibrium (constant solution), exists iff  $\gamma\delta < 0$  which is  $u_\infty^2 = -\gamma/\delta$ . Note the relevance of the solution  $u_\infty^2$  with the upper bound in the estimate (6). It is also the density of the cw steady-state solution  $\phi_b$  (see Fig. 1). We focus here on these solutions, i.e., solutions that tend to stabilize the soliton, by keeping its parameters constant. The evolution of the rest of the soliton parameters [see (51)] can be found by means of a multiscale boundary layer theory [38]; the resulting evolution equations are:

$$2BA_t = \text{Re} \left\{ \int_{-\infty}^{\infty} F[u_{\text{ds}}](\bar{u}_{\text{ds}})_t dx \right\}, \quad (54)$$

$$Bx_{0t} = \text{Im} \left\{ \int_{-\infty}^{\infty} x(F[u_\infty]u_\infty - F[u_{\text{ds}}]\bar{u}_{\text{ds}}) dx \right\}, \quad (55)$$

$$u_\infty\sigma_{0t} = \text{Im} \left\{ \int_{-\infty}^{\infty} (F[u_\infty]u_\infty - F[u_{\text{ds}}]\bar{u}_{\text{ds}}) dx \right\} + \text{Re} \{F[u_\infty]\}, \quad (56)$$

$$BB_t = u_\infty u_{\infty t} - AA_t, \quad (57)$$

$$u_\infty^2 \Delta\phi_{0t} = 2AB_t - 2BA_t, \quad (58)$$

$$q_1^\pm = \frac{1}{2} \frac{\sigma_{0t} \pm \Delta\phi_{0t}}{u_\infty \mp A}, \quad \phi_{1t}^\pm = \mp 2q_1^\pm. \quad (59)$$

Here, we should mention that  $q_1^\pm$  and  $\phi_1^\pm$  in Equations (59) represent the asymptotics of the shelf, induced by the perturbation  $F[u]$ , as  $x \rightarrow \pm\infty$  respectively; in fact, they are higher-order corrections to the soliton, so that the shelf amplitude is  $u_\infty + q_1^\pm$  and its speed  $u_\infty$ . Integrating the above equations, and using Equation (53), finally yields:

$$u'_\infty = (\gamma + \delta u_\infty^2)u_\infty, \quad (60)$$

$$A' = \frac{4}{15}\sigma_R A^4 + \frac{2}{3}\delta A^3 - \frac{8}{15}\sigma_R u_\infty^2 A^2 + \left(\gamma + \frac{\delta}{3}u_\infty^2\right)A + \frac{4}{15}\sigma_R u_\infty^4, \quad (61)$$

$$x'_0 = \left(2\beta - \mu - \frac{2\nu}{3}\right)A^2 - \left(2\beta + 2\mu + \frac{4\nu}{3}\right)u_\infty^2, \quad (62)$$

$$\sigma'_0 = \frac{B_z}{u_\infty} - \frac{2B}{3u_\infty} \left(3\gamma + 4u_\infty^2\delta + 2\delta A^2\right). \quad (63)$$

These equations show that the evolution of the soliton center, described by the equation  $X'_0 = A + x'_0$ , is affected by all parameters of Equation (30) [directly or indirectly from  $A(t)$ ]. On the other hand, the rest of the soliton characteristics, i.e., the background, the dip and the shelf, only depend on  $\gamma$ ,  $\delta$  and  $\sigma_R$ . This implies

that soliton stabilization can be targeted accordingly. Indeed, stable fixed points of this system correspond to stable solitons traveling on top of a constant background with a constant speed. It is possible to identify two such solitons, namely a grey and a black one, supported in the presence ( $\sigma_R \neq 0$ ) and in the absence ( $\sigma_R = 0$ ) of the SRS effect, respectively. In both cases, the background assumes the same form: this can be obtained by means of Equation (60), which depicts the nontrivial stationary solution  $u_\infty^2 = -\gamma/\delta$  for  $\gamma\delta < 0$ , i.e., for linear loss and nonlinear gain, or vice versa.

We start with the case  $\sigma_R \neq 0$ . Substituting the above mentioned constant background in Equation (61), and seeking stationary solutions for the soliton velocity, we arrive at a 4th-order algebraic equation for  $A$ . Solving this equation, we find that there exists only one root, namely  $A = (4\delta\sigma_R)^{-1}(-5\delta^2 + \sqrt{25\delta^4 - 16\gamma\delta\sigma_R^2})$ , which does not violate the relationship  $A^2 + B^2 = u_\infty^2$ . Thus, a stable soliton exists for:

$$u_\infty^2 = -\frac{\gamma}{\delta}, \quad A = \frac{-5\delta^2 + \sqrt{25\delta^4 - 16\gamma\delta\sigma_R^2}}{4\delta\sigma_R}. \tag{64}$$

Note that since  $\gamma\delta < 0$  the quantity under the square root is always positive.

In general, the solution of Equation (60) with  $u_\infty(0) = u_0$  is:

$$u_\infty^2(t) = \frac{\gamma u_0^2 e^{2\gamma t}}{\gamma + \delta u_0^2 - \delta u_0^2 e^{2\gamma t}}, \tag{65}$$

which suggests that there is a (finite) time for which the background exhibits blow-up, as it was discussed in Theorem 2. Indeed, the denominator becomes zero when

$$t = t_* \equiv \frac{1}{2\gamma} \ln \left( 1 + \frac{\gamma}{\delta u_0^2} \right). \tag{66}$$

The unexpected feature here is that the addition of the term  $i\delta|u|^2u$  which compensates the effect of the linear loss term  $i\gamma u$  may result in blow-up of the background in finite time, even when the other soliton parameters remain finite. In addition, Equation (65) indicates that an equilibrium can also be reached in finite time when the denominator is a multiple of the numerator. Nevertheless, while under this requirement the background will be stabilized, this does not necessarily guarantees the stabilization of the other soliton parameters.

Next, we consider the case of  $\sigma_R = 0$ . In this case, Equations (60) and (61) lead to the following equations for the background and soliton velocity:

$$u_\infty^2 = -\frac{\gamma}{\delta}, \quad A' = \frac{2}{3}(\delta A^2 - \gamma)A, \tag{67}$$

Obviously, the above equation for the velocity depicts a stationary solution  $A = 0$  (recall that  $\gamma\delta < 0$ ), that corresponds to a black soliton. Hence, when SRS is absent (which would give a frequency downshift causing the soliton to move), a stable black soliton can exist.

An important comment is in order here. While for these specific choice of  $u_\infty$  and  $A$  the soliton gets stabilized, this does not mean that the shelf is no longer present. In fact, the shelf is always present in the perturbed NLS, even though its amplitude is small, since it appears as a higher-order correction in the perturbation theory [38]. Thus, the shelf does not affect the soliton propagation but it does, however, affect soliton interactions (see Ref. [41] for a relevant study, but in the framework of another dissipative NLS model). Notice, also, that the shelf can be suppressed with counter effect the destabilization of the soliton.

Finally, we briefly consider the case where gain/loss terms are absent, i.e.,  $\gamma = \delta = 0$ . In this particular case, the dark soliton dynamics is merely driven by the SRS effect. Indeed, now the evolution of the background and soliton velocity is described by the following equations:

$$u'_\infty = 0, \quad A' = \frac{4}{15}\sigma_R(A^2 - u_\infty^2)^2, \quad (68)$$

which recover the results obtained in Refs. [36, 37]. The soliton dynamics in this case can be understood as follows. Since  $A^2 \neq u_\infty^2$ , the right-hand-side of the second equation is always positive and, thus, the soliton becomes shallower and faster, i.e.,  $B \rightarrow 0$  and  $A \rightarrow u_\infty$ , so that the condition  $A^2 + B^2 = u_\infty^2$  is satisfied. Thus, the dark soliton eventually decays to the stationary background. It is therefore clear that no stable dark soliton (in the sense of the existence of stationary soliton parameters) exists in this case.

### 3.3 Solitons and Shock Waves in an Effective KdV-Burgers Picture

Finally, for completeness, it is relevant to briefly mention the following. Apart from the direct perturbation theory for solitons, there exists still another method to analyze the dynamics of dark solitons in the framework of Equation (30). Indeed, as shown in Ref. [42] for the special case of  $\gamma = \delta = 0$ , it is possible to employ a multiscale expansion method and asymptotically reduce the higher-order NLS equation to a Korteweg-de Vries–Burgers (KdV-B) equation. This can be done upon seeking solutions of the form:

$$u(x, t) = [u_\infty + U(x, t)] \exp[iu_\infty^2 t + i\phi(x, t)], \quad (69)$$

where  $U(x, t)$  and  $\phi(x, t)$  are unknown real functions (to be determined) representing an amplitude and a phase modulation of the background wave  $u_b = u_\infty \exp(iu_\infty^2 t)$ . Then, it is assumed that these functions are presented in the form of formal asymptotic series,

$$U = \epsilon^2 U_1 + \epsilon^4 U_2 + \dots, \quad \phi = \epsilon \phi_1 + \epsilon^3 \phi_2 + \dots, \quad (70)$$

where the unknown functions  $U_j$  and  $\phi_j$  ( $j = 1, 2, \dots$ ) depend on the slow variables  $X = \epsilon(x - vt)$  and  $T = \epsilon^3 t$ , with  $v$  being an unknown velocity, and  $\epsilon$  being a formal small parameter. Substituting Equations (69)–(70) into Equation (30), we obtain the following results. First, to the lowest-order of approximation in  $\epsilon$  of the perturbation technique, we derive the unknown velocity  $v$  and an equation connecting the functions  $\phi_{1T}$  and  $U_1$ . Second, to the next order of approximation, we derive the following KdV-B equation:

$$U_{1T} + c_1 U_1 U_{1X} + c_2 U_{1XXX} = c_3 U_{1XX}, \quad (71)$$

where,  $U_1$  obviously represents the soliton amplitude. The coefficients of the underlying KdV equation,  $c_1$  and  $c_2$ , depend on the coefficients of the pNLS,  $\beta$  and  $\nu$ , as well as on the background amplitude  $u_\infty$ , while the diffusion coefficient  $c_3$  depends on the SRS parameter,  $\sigma_R$ .

Importantly, the relevant asymptotic reduction to the KdV-B equation can be performed for both normal and anomalous dispersion cases, i.e., for both  $s = \pm 1$ . Of course, in the case  $s = -1$  it is known [1–3] that the soliton’s background plane wave is prone to the modulational instability (MI), but this long-wavelength instability may be suppressed: indeed, in applications, one expects periodic or other boundary conditions in the  $x$ -direction, meaning that the admitted wavenumbers are quantized, hence they are limited from below by a minimum wavenumber,  $k_{\min}$ , which corresponds to the transverse size of the system. In such a case, if  $k_{\min} > K_{\max}$  (where  $K$  is the perturbation wavenumber characterizing the MI, and  $K_{\max}$  defines the width of the instability band,  $0 \leq K \leq K_{\max}$ ), no quantized wavenumber can get into the instability band and, hence, the MI is eliminated.

The effective KdV-B description of the soliton dynamics offers a number of interesting results. First, in the absence of the SRS effect ( $\sigma_R = 0$ ), dark solitons small-amplitude dark solitary wave solutions can exist for both the normal and anomalous dispersion regimes. This result is in sharp contrast with the conventional form and certain perturbed versions of the NLS equation, where dark solitons solely exist for the normal dispersion regime ( $s = +1$ ). In addition, in this latter regime, there exists another type of solution, namely an *anti-dark soliton*, having the form of a hump, rather than a dip, on top of the background plane wave. Notice that the transformation from the dark to the anti-dark soliton is possible (see details in Ref. [42]).

When the SRS effect is present ( $\sigma_R \neq 0$ ), the soliton dynamics is governed by a KdV-B equation. In this case, the evolution of solitons can be studied by means of the perturbation theory for solitons [33, 35]. The results that can be obtained in this case show that the solitons experience a decrease in their amplitudes and/or their velocities, depending on the direction of propagation and the dispersion region ( $s = -1$  or  $s = +1$ ). In particular, right-going solitons experience a decrease in both their amplitudes and their velocities, while the evolution of left-going solitons depends on

$s$ : for  $s = -1$ , they increase their amplitudes and decrease their velocities, while for  $s = +1$ , they decrease their amplitudes and increase their velocities—in accordance with the results presented in the previous section.

Still another nonlinear wave structure that can be predicted to occur in this setting, is the one of a traveling shock wave [43]. In the effective KdV-B picture, the existence of such a structure is not surprising, because the KdV-B equation possesses stable traveling shock wave solutions. The latter, are obviously supported by the SRS effect (recall that if  $\sigma_R = 0$  then  $c_3 = 0$  and the diffusion term in Equation (71) vanishes, as was also found by means of other methods in other studies [44–46]. Notice that, as before, shock wave type structures are possible for both normal and anomalous dispersion cases. In particular, in the case of the normal dispersion ( $s = +1$ ), the structure has the usual shock wave profile, while in the case of the anomalous dispersion ( $s = -1$ ) it has the form of a rarefaction wave.

Finally, based on the analysis of the shock wave structure of the KdV-B equation, one may deduce the relevant profiles in the context of the perturbed NLS equation. Thus, the structure of the front of the shock solutions may be monotonic, in the nonlinearity-dominated regime, or oscillatory in the dispersion-dominated regime. In fact, since the former regime is only accessible for  $s = -1$  [43] the front of the rarefaction wave is monotonic. On the other hand, the profile of the front of the shock wave supported for  $s = +1$ , may be either monotonic in the nonlinearity-dominated regime (resembling the regular stationary solutions of the Burgers equation), or oscillating. It is interesting to mention that the oscillations in the kink front can be studied in the framework of the perturbation theory for solitons of the KdV equation, treating the diffusion term as a small perturbation [33, 35]. This way, it can be deduced that the oscillations of the shock front can be considered as a succession of KdV solitons, a fact that completes the connection between the soliton and shock wave solutions of the perturbed NLS model.

## References

1. A. Hasegawa, Y. Kodama, *Solitons in Optical Communications* (Oxford University Press, Oxford, 1996)
2. G.P. Agrawal, *Nonlinear Fiber Optics* (Academic Press, London, 2012)
3. Yu.S. Kivshar, G.P. Agrawal, *Optical Solitons: From Fibers to Photonic Crystals* (Academic Press, London, 2003)
4. M.J. Ablowitz, H. Segur, *Solitons and Inverse Scattering Transform* (SIAM, 1981)
5. V.E. Zakharov, A.B. Shabat, *Sov. Phys. JETP*. **34**, 62 (1972)
6. V.E. Zakharov, A. B. Shabat, *Sov. Phys. JETP*. **37**, 823 (1973)
7. M. Scalora, M.S. Syrchin, N. Akozbek, E.Y. Poliakov, G. D’Aguanno, N. Mattiucci, M.J. Bloemer, A.M. Zheltikov, *Phys. Rev. Lett.* **95**, 013902 (2005)
8. S. Wen, Y. Xiang, X. Dai, Z. Tang, W. Su, D. Fan, *Phys. Rev. A* **75**, 033815 (2007)
9. N.L. Tsitsas, N. Rompotis, I. Kourakis, P.G. Kevrekidis, D.J. Frantzeskakis, *Phys. Rev. E* **79**, 037601 (2009)
10. R.S. Johnson, *Proc. R. Soc. Lond. A* **357**, 131 (1977)
11. Yu.V. Sedletsky, *J. Exp. Theor. Phys.* **97**, 180 (2003)

12. A.V. Slunyaev, J. Exp. Theor. Phys. **101**, 926 (2005)
13. N.N. Akhmediev, A. Ankiewicz, *Solitons. Nonlinear Pulses and Beams* (Chapman and Hall, London, 1997)
14. H. Ikeda, M. Matsumoto, A. Hasegawa, Opt. Lett. **20**, 1113 (1995)
15. T.P. Horikis, D.J. Frantzeskakis, Opt. Lett. **38**, 5098 (2013)
16. S.C.V. Latas, M.F.S. Ferreira, M.V. Facão, Appl. Phys. B **104**, 131 (2011)
17. S.C.V. Latas, M.F.S. Ferreira, Opt. Lett. **37**, 3897 (2012)
18. I.M. Uzunov, T.N. Arabadzhiev, Z.D. Georgiev, Opt. Fiber Tech. **24**, 15 (2015)
19. V. Achilleos, S. Diamantidis, D.J. Frantzeskakis, T.P. Horikis, N.I. Karachalios, P.G. Kevrekidis, Physica D **316**, 57 (2016)
20. V. Achilleos, A.R. Bishop, S. Diamantidis, D.J. Frantzeskakis, T.P. Horikis, N.I. Karachalios, P.G. Kevrekidis, Phys. Rev. E **94**, 012210 (2016)
21. R. Temam, *Infinite-Dimensional Dynamical Systems in Mechanics and Physics* (Springer, Berlin, 1997)
22. T. Cazenave, *Semilinear Schrödinger Equations*. Courant Lecture Notes, vol. 10 (American Mathematical Society, Providence, 2003)
23. C. Sulem, P.L. Sulem, *The Nonlinear Schrödinger Equation. Self-Focusing and Wave Collapse* (Springer, Berlin, 1999)
24. T. Kato, *Quasilinear Equations of Evolution with Applications to Partial Differential Equations*. Lecture Notes in Mathematics, vol. 448 (Springer, New York, 1975), pp. pp. 25–70
25. T. Kato, *Abstract Differential Equations and Nonlinear Mixed Problems* (Lezione Fermiane Pisa, 1985)
26. T. Kato, C.Y. Lai, J. Funct. Anal. **56**, 15 (1984)
27. T. Kato, *Nonlinear Schrödinger Equations. Schrödinger Operators*. Sønderborg, 1988, 218–263, Lecture Notes in Phys., vol. 345 (Springer, Berlin, 1989)
28. J.M. Ball, Q. J. Math. Oxf. **28**, 473 (1977)
29. T. Ozawa, Y. Yamazaki, Nonlinearity **16**, 2029 (2003)
30. M. Taylor, *Partial Differential Equations III*. Applied Mathematical Sciences, vol. 117 (Springer, New York, 1996)
31. N.I. Karachalios, H. Nistazakis, A. Yannacopoulos, Discrete Cont. Dyn. Syst. A **19**, 711 (2007)
32. V. Achilleos, A. Álvarez, J. Cuevas, D.J. Frantzeskakis, N.I. Karachalios, P.G. Kevrekidis, B. Sánchez-Rey, Phys. D **244**, 1 (2013)
33. V.I. Karpman, E.M. Maslov, Sov. Phys. JETP **46**, 281 (1977)
34. D.J. Kaup, A.C. Newell, Proc. R. Soc. Lond. Ser. A **361**, 413 (1978)
35. Yu.S. Kivshar, B.A. Malomed, Rev. Mod. Phys. **61**, 761 (1989)
36. I.M. Uzunov, V.S. Gerdjikov, Phys. Rev. A **47**, 1582 (1993)
37. Yu.S. Kivshar, X. Yang, Phys. Rev. E **49**, 1657 (1994)
38. M.J. Ablowitz, S.D. Nixon, T.P. Horikis, D.J. Frantzeskakis, Proc. R. Soc. A **467**, 2597(2011)
39. Yu.S. Kivshar, B. Luther-Davies, Phys. Rep. **298**, 81 (1998)
40. D.J. Frantzeskakis, J. Phys. A: Math. Theor. **43**, 213001 (2010)
41. M.J. Ablowitz, S.D. Nixon, T.P. Horikis, D.J. Frantzeskakis, J. Phys. A: Math. Theor. **46**, 095201 (2013)
42. D.J. Frantzeskakis, J. Phys. A: Math. Gen. **29**, 3631 (1996)
43. D.J. Frantzeskakis, J. Opt. Soc. Am. B **9**, 2359 (1997)
44. Yu.S. Kivshar, Phys. Rev. A **42**, 1757 (1990)
45. G.P. Agrawal, C. Headley III, Phys. Rev. A **46**, 1573 (1992)
46. Yu. S. Kivshar and B. A. Malomed, Opt. Lett. **18**, 485 (1993).



# The Role of Differential Equations in Applied Statistics



Christos P. Kitsos and C. S. A. Nisiotis

**Abstract** The target of this paper is to discuss, investigate and present how the differential equations are applied in Statistics. The stochastic orientation of Statistics creates problems to adopt the differential equations as an individual tool, but Applied Statistics is using the differential equations either through applications from other fields, like Chemistry or as a tool to explain “variation” in stochastic processes.

## 1 Introduction

Differential equations (de), [30] among others, support a number of sciences. Quincy Wright in his early work introduced (de) in political science, on a study of War. In such cases, there is a covered uncertainty. In others, like electric circuits, see Appendix 1, there is a solid background for (de), and not any uncertainty. In Astronomy, were (de) are extensively used, this uncertainty is measured through a Probability model, the Gaussian, [5].

Consider two points in the Universe and their distance  $l$ . Due to the expansion of the Universe the distance  $l$  it is not a constant function, but it depends on time  $t$ . Therefore considering a scale factor  $s(t)$  acting as normalizing factor, the distance  $D = D(t)$  is eventually of the form

$$D(t) = s(t)l(t)$$

---

C. P. Kitsos (✉)

Department of Informatics, University of West Attica, Aigaleo, Greece

e-mail: [xkitsos@uniwa.gr](mailto:xkitsos@uniwa.gr)

C. S. A. Nisiotis

Department of Public and Community Health, University of West Attica, Aigaleo, Greece

e-mail: [csnisiotis@uniwa.gr](mailto:csnisiotis@uniwa.gr)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,

Springer Optimization and Its Applications 173,

[https://doi.org/10.1007/978-3-030-72563-1\\_10](https://doi.org/10.1007/978-3-030-72563-1_10)

From  $D(t)$  we can obtain the velocity  $V(t)$  and the acceleration  $\alpha(t)$  as

$$\begin{aligned} V(t) &= \dot{D}(t) = \dot{s}(t)l(t) + s(t)\dot{l}(t) \\ \alpha(t) &= \dot{V}(t) = \ddot{s}l + 2\dot{s}\dot{l} + s\ddot{l} \end{aligned}$$

where:

$\dot{s}(t)l(t)$  = is known as velocity due to the expansion

$s(t)\dot{l}(t)$  = is known as peculiar velocity, usually denoted by  $u$

$2\dot{s}l(t)\dot{l}(t) + s(t)\ddot{l}(t)$  = is known as peculiar acceleration,

Considering the gravitational potential  $G$  the equation of motion is [10],

$$\frac{d^2 D}{dt^2} = \nabla_D G$$

If we set:

$$\Psi(l, t) := G(l, t) + \frac{1}{2}s\ddot{l}l^2$$

It can be proven, [10, 46], that:

$$\frac{\partial u}{\partial t} + \frac{\dot{s}}{s}u = -\frac{1}{s}\nabla\Psi$$

Let  $p = p(x)$  be the density of the Universe at the position  $x$ , and  $\bar{p}$  the mean density of the Universe. Thus, if we define:

$$\delta = \delta(x) = \frac{p(x) - \bar{p}}{\bar{p}}$$

the probability distribution of  $\delta$  is Normal with mean zero and variance  $\sigma_\delta^2$  known i.e.  $\delta \sim N(0, \sigma_\delta^2)$  or:

$$P(\delta) = \frac{1}{\sqrt{2\pi}\sigma_\delta} \exp\left(-\frac{1}{2}\frac{\delta^2}{\sigma_\delta^2}\right)$$

The above relation provides evidence that Astronomy is adopting Statistics on the way that Statistics is adopting (de). As Statistics is working with descriptive data, not continuous, the difference equations are adopted widely, but we shall not refer to this characteristic form of equations. In Section 4 we discuss how probability is adopted the (de) to measure the “rate of change” in a random walk.

Now, although the optimal experimental design theory was early originated [41], it was applied in Chemistry only by the pioneering work [15], who worked with the dilution assessment problem, while [25] investigates the problem for the sequential point of view. In principle, the essential problem with the non-linear experimental design problems is, that there is not a global statistical framework, for all the adopted models. There is a solid statistical theoretical background for any nonlinear function, but not for the particular nonlinear equation, which is arisen though a differential equation approach, in Chemical Kinetics, see the early work in [11–13], either from Chemistry or Biology. So each model should be particularly investigated, in order to calculate the optimal design points, which are functions of the unknown parameters! The proposed solution to overpass the unknown parameter repentance, [25, 26, 28] is to adopt the sequential principle of design.

The differential equations play a vital role to Statistical application in Chemical Kinetics, typical example being the Michaelis–Menten model, [29]. They are applied to create the model, which as it is based on observations the continuation is destroyed. Therefore the observations obtained are coming from the model plus a stochastic error. The usual assumption is the normality of the errors, with mean zero and variance constant  $\sigma^2 > 0$ . Here comes another problem in statistical nonlinear case: the variance depends on the input variables and the parameters, we want to estimate. Although the D-optimality criterion appeared to have an aesthetic appeal, as it requests the minimization of the variance, in Chemical Kinetics, see [11] another problem appears: as nonlinear models does not provide, in Statistical terms, the appropriate information, the Taylor expansion has been applied, see the early work [6, 25], among others. Therefore what is consider in this paper is:

- The differential equation (de) that approaches the phenomenon
- Form the model in general Statistical model
- Consider the non Linear Statistical model
- Provide a general framework to obtain the optimal design.

The use and abuse of regression is also discussed as well as the use of differential equations in Statistics. Therefore the (de) are applied, in principle, in three characteristic lines of thoughts:

1. To provide a theoretical framework and a solid background to physical problems, see Appendix 1.
2. Although a theoretical background is developed through (de) principles there is an underlying uncertainty covered by Statistics.

## 2 Theoretical Framework for the Nonlinear Design

Consider the non-empty set  $U \subseteq \mathcal{R}^k$  in which the  $k$  predictor variables or covariates or explanatory variables or independent variables are  $u = (u_1, u_2, \dots, u_k)$ . We

assume that their values are within the known as experimental region or design space  $U$ . A typical example of an input variable from Chemical Kinetics is time.

The parameter space  $\Theta \subseteq \mathcal{R}^P$  is the set where the  $p$ -term parameter vector  $\theta = (\theta_1, \theta_2, \dots, \theta_p)$  takes values. Where the sequential procedure of design is adopted [26, 31]  $\Theta$  is assumed compact. When the response vector  $y$  is supposed to take any value in the response space,  $\Psi$  we also suppose that a *regression model*, in principle a nonlinear function  $f$ , that links  $u$  and  $\theta$ , and consists of the deterministic portion  $f(u, \theta)$  and the stochastic portion,  $e$ , known as error, linked, eventually, through the relation

$$y = f(u, \theta) + e, \text{ with } E(y) = \eta = f(u, \theta) \quad (1)$$

where  $E(y)$  means the expected value of  $y$ .

For the independent identically distributed errors we suppose that are normally distributed with mean 0 and variance  $\sigma^2 > 0$ . The function  $f(u, \theta)$  (which describes the chemical kinetic model in our scenario) is assumed to be continuous with the second order derivatives of  $f(\cdot)$  with respect to  $\theta$  existing at and near the true value of the parameter, see [16]. For model (1) we introduce the quantity [3, 4]:

$$S_n(\theta) = \sum (y_i - f(u_i, \theta))^2 = \|y - f(u, \theta)\|_2^2 \quad (2)$$

where  $\|\cdot\|_2$  is the  $l_2$ -norm. An estimate  $\hat{\theta}$  will be called the *least squares estimate* (LSE) if  $S_n(\hat{\theta}) = \min\{S_n(\theta); \theta \in \Theta\}$ .

In the non-linear regression problems the variance  $\sigma^2 = \sigma^2(u, \theta)$ . That is  $\sigma^2$  depends on the design point and the parameter vector. In the linear case it is independent of the parameter  $\theta$ .

The concept of the *average-per-observation information matrix* (apoin) is important for the nonlinear design problem, as described [29]. It is defined for  $\xi_n$ , the  $n$ -point design measure [16] (practically: the portion of observations at the optimal design space), to be, for the discrete case, equal to

$$M(\theta, \xi_n) = n^{-1} \sum I(\theta, u_i) \quad (3)$$

while for the continuous case is

$$M(\theta, \xi) = \int_U I(\theta, u) \xi(du). \quad (4)$$

The following partition of matrix  $M$  is usually considered

$$M = \begin{pmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{pmatrix} \quad (5)$$

with  $M_{11} \in \text{Mat}(s, s)$ ,  $M_{12} \in \text{Mat}(s, p - s)$ ,  $M_{22} \in \text{Mat}(p - s, p - s)$ ,  $1 \leq s < p$  and  $\text{Mat}(n, p)$  the set of  $n \times p$  matrices, see also [20, 25].

We now briefly discussed how the Optimal Experimental Design Theory, for the Non-Linear model proceeds, see [16]. The average per observation information matrix is considered, were  $\theta$  takes its true value. Then we can define the following operator  $J_Q$  applied to (apoi)  $M$ , Kitsos (1986), through a considered known matrix  $Q$ :

$$J_Q(M) = QM^-(\theta, \xi)Q^T,$$

$M^-$  a generalized inverse of  $M$  and  $Q^T \in \text{Mat}(p, s)$ .

Given the above notation a real valued convex, decreasing function,  $\omega$  say, on the set of nonnegative definite matrices, say  $N\text{Mat}(s, s)$ , is applied to  $J_Q$  is what we consider as an optimality criterion function.

The design measure  $\xi^*$  is called  $\omega$ -optimal if and only if:

$$\omega \{ J_Q [M(\theta, \xi^*)] \} = \min \left\{ \omega \left\{ QM^-(\theta, \xi)Q^T \right\}, \xi \in \mathcal{E} \right\}. \tag{6}$$

For the special case for  $\omega$  such as  $\omega(\cdot) = \log [\det (QM^-Q^T)]$ , and for  $Q = I \in \text{Mat}(p, p)$  the identity matrix we are referring to  $D$ -optimal design, which is adopted in this paper. Different values of the convex function and the matrix  $Q$  define other optimality criteria (When the trace of the covariance matrix is minimized we are referred to  $A$ -optimality), which are beyond the target of this paper. Now, two theoretical results (DR1 and DR2) are useful to reduce the number of parameters required to calculate  $D$ -optimal designs:

DR1: Let  $f(x, \theta)$  be any nonlinear model of the form:

$$f(x, \theta) = L(x_1, \beta_1) + NL(x_2, \beta_2)$$

where  $\theta = (\theta_1, \theta_2, \dots, \theta_s; \theta_{s+1}, \dots, \theta_p) = (\beta_1; \beta_2)$ ,  $NL(x_2, \beta_2)$  is the non-linear part,  $L(x, \beta_1) = \beta_1^T x_1$ ,  $\beta_1^T$  the transpose of  $\beta_1$  and  $\dim \beta_1 = \dim x_1 = s$ .

Then the  $D$ -optimal design depends only on  $\beta_2$  (the Non-Linearly involved parameters) as well as the  $D_s$ -optimal design for the vector  $\beta_1$ .

DR2: Hill [20] defined the partially nonlinear model for the  $k$  parameters,  $k < p$ , to be the one for which  $\nabla f(u, \theta) = B(\theta)H(u, \beta)$ , where  $B(\theta)$  is a matrix depending on  $\theta = (\theta_1, \theta_2, \dots, \theta_p)$ ,  $\beta$  is the vector of the  $k$  parameters which appear in a nonlinear way and  $H$  is an appropriate matrix. In such a case the design depends only on the parameters of  $\beta$ .

We try to have results reducing the number of the involved parameters, as former information is needed for the involved parameters. So the less the involved parameters, the less the required information, especially when DR1 holds.

### 3 Growth Curves

In Statistics the nonlinear models are mainly produced by a linear system of differential equations, [7–9, 12, 19, 21, 22, 27, 31, 32, 40, 47], . See also the pioneer work [37], and [35, 45] . In principle compartment models are based on a division of the system into compartments. Then it is assumed that the “rates of flow” of whatever is under investigation, with typical example the drugs between compartments follow the first order kinetics, with the relevant (de) and eventually the rate of transfer is proportional to the concentration.

The target is to find a growth curve, as a mathematical model, to describe the plant dry weight data over a pre-defined season. A number of growth curves (gc) have been developed in various fields. The main approach remains the same, either they serve a chemical oriented purpose, [32, 33] (cgc) or a plant gc, (pgc).

We pay attention to pgc, in the presence of its dry weight  $W$ , which varies with time  $t$ ,  $W = W(t)$ . Moreover growth exists at the expense of a substrate  $S$ . The rate of the growth reaction is linearly proportional to  $S$  and  $W$ , so with  $\lambda$  being a constant it holds:

$$\frac{dIV}{dt} = \lambda SW \quad (7)$$

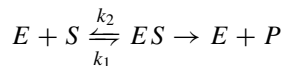
When  $W$  be maximum i.e.  $W = W_m$  say, then  $S = 0$  and in then the equality can be reduced to:

$$\frac{dW}{dt} = \lambda(W_m - W)W \quad (8)$$

The solution of (8) is, with  $W_0 = W(t = 0)$  be the starting observation:

$$W \equiv W_0 \exp(W_m \lambda t) \quad (9)$$

In Statistical terms, as there were introduced in the above section, model (9) is considered as  $\eta = \theta_0 \exp(\theta_1 t)$ , with the parameters definition to be clear, see [33]. The Michaelis–Menten (M-M) model is a typical example, see [29]



with  $E$  = enzyme,  $S$  = substrate and  $P$  = product. Eventually the model is:

$$u = \frac{u_{max}[s]}{k + [s]} \quad (10)$$

with  $u$  = speed of the steady-state reaction [ $\cdot$ ] declares concentrations:

$$k = \frac{k_2 + k_3}{k_1}$$

Equation (10) declare that:

$$\left. \frac{du}{d[s]} \right|_{[s]=0} = \frac{u_{max}}{k} \tag{11}$$

Adopting the general notation in Statistical terms introduced in Section 2 model (12) can be considered as:

$$\eta = \frac{\theta_1 u}{\theta_2 + u} \tag{12}$$

where  $\eta$  is the rate of the enzyme reaction,  $\theta_1$  corresponds to the maximum rate reaction and  $\theta_2$  is the half saturation constant, with  $u$  being the concentration of substrate, see also [34].

Moreover under DR1 proposition considering  $(\theta = \theta_1, \theta_2; \theta_3, \theta_4) = (\beta_1, \beta_2)$  with

$$L(u, \beta_1) = \theta_1 + \theta_2 u$$

$$NL(u, \beta_2) = \frac{\theta_3 u}{\theta_4 + u}$$

then an extending (M-M) model is:

$$f(u, \theta) = \theta_1 + \theta_2 u + \frac{\theta_3 u}{\theta_4 + u} \quad (extM - M)$$

has exactly the same D-optimal design as

$$n_N = \frac{\theta_3 u}{\theta_4 + u}$$

Again under DR1 the design does not depend on  $\theta_3$  so prior information for (ext M-M) only for  $\theta_4$  is needed, while for the design approach see [33].

The Gompertz growth equation for given dry matter is  $D$ , by two plant research first-order (de) (1st de), with  $S$  being the senescence parameter:

$$\frac{dD}{dt} = \lambda D, \quad \frac{d\lambda}{dt} = -\lambda S \tag{13}$$

or by the one:  $\frac{dD}{dt} = \lambda_0 D e^{-St}$  with  $\lambda_0 =$  specific growth at  $t = 0$ .

The Semi-Empirical model, [45], is the model which describes the lightflux densities within the canopy of an isolated plant.

Given a point  $M = M(r)$ , with  $r$  being the radius describing the position of the point  $M$  when the ray of light traverse. This de provides, by integration:

$$I = I_0 \exp \left[ - \int_0^s kF(r)ds \right] \tag{14}$$

with  $I_0 = I(s = 0)$ . If  $k$  and  $F$  does not vary along the chosen path it is known that we can have:

$$I = I_0 e^{-kFS} \tag{15}$$

The estimating photosynthesis rate  $P$  for the presence of light variability is another example of adopting a mathematical approach. Let  $I$  be the flux density falling a leaf then it can be easily proved that when I changes:

$$\Delta P = \frac{\partial P}{\partial I} \Delta I + \frac{1}{2} \frac{\partial^2 P}{\partial I^2} (\Delta I)^2 \tag{16}$$

If we let the  $CO_2$  density to be  $C$  and  $\alpha, \beta$  constants then it is:

$$P = \frac{(\alpha I)(\beta C)}{\alpha I + \beta C} \tag{17}$$

Considering (9) for  $N$  fluctuations with  $N$  being large and taking into account the first order terms are canceled about the mean value variation:

$$\Delta P = \frac{1}{2} I^2 \frac{\partial^2 P}{\partial I^2} (CV)^2 \tag{18}$$

with  $(CV)$  being the coefficient of variation. From (17) and (18) one can obtain:

$$\frac{\Delta P}{P} = - \frac{(\alpha I)(\beta C)}{(\alpha I + \beta C)^2} (CV)^2 \tag{19}$$

Then the maximum values can be evaluated for the absolute correction  $P$ , especially for and the relative correction  $\Delta P/P$  as a function of the coefficient of variation:

$$\left( \frac{\Delta P}{P} \right)_{max} = - \frac{1}{4} (CV)^2 \quad \text{at } \alpha S = \beta C \tag{20}$$

Let us consider the light flux density  $I$  from a steady value  $I_1$  say to a new one  $I_2$ , say. The effect on photosynthesis rate of charging is under consideration.

Let us assume that the centre is  $t = 0$ , where this charge occurs fast. Then for  $t < 0$  let  $I = I_1$  and for  $t > 0$  let  $I = I_2$ , while the density of  $CO_2$  is  $C$ . In such a case the system in photosynthesizing at steady rate  $P_1$  and  $P_2$  is respectively:



$$P_i = \frac{aI_i\beta C}{aI_i + \beta C} \quad i = 1, 2$$

The time constant is in this case:

$$\tau = \frac{k}{aI\beta C}$$

Then as it holds:

$$\frac{dP}{dt} = \frac{1}{\tau}(P_2 - P)$$

one can obtain:

$$\frac{dP}{P_2 - P} = \frac{dt}{\tau}$$

Therefore the final model can be:

$$P = P_2 - (P_2 - P_1)e^{-t/\tau}$$

Thus from the general framework of Section 2 it can be written following the notation on a nonlinear Statistical model as:

$$\eta = \theta_0 - \theta_1 \exp(-\theta_2 u) \tag{21}$$

Notice that under DR1, DR2 the design depends only on  $\theta_2$ .

Notice that Equation (21) hides what is behind, as Statistics are involving of solving the stochastic problem

$$\eta = f(u, \theta) + error$$

This is an important point for our scenario: The use of differential equation is only to form, eventually equations like the (21) one. This is exactly what we present briefly in this paper in this section.

For an irreversible, unimolecular homogeneous reaction of the type  $A \xrightarrow{k} B$  under isothermal conditions, the rate equation is described by:  $\frac{d[A]}{dt} = -k[A]$ , where  $[A]$ , as usually, is the concentration or partial pressure of reagent A,  $t$  is the reaction time and  $k$  the reaction rate constant<sup>[21]</sup>. Integrating, with a given initial concentration of A at time  $t = 0$ ,  $[A]_0$ ,

$$\frac{[A]}{[A]_0} = X_A = \exp(-kt) = \exp \left[ -At \exp \left( -\frac{\Delta E}{RT} \right) \right] \tag{22}$$

where the fraction  $\frac{[A]}{[A]_0}$  resembles the molar fraction of A,  $X_A$ ,  $k$  is the rate constant, which is a function of the absolute temperature  $T$ . Then (22) arises.

As  $A$  and  $\frac{\Delta E}{R}$  are highly correlated, it is desirable to introduce  $k_1$ , the rate of the reaction at some specific reference temperature  $T_0$ .

Thus  $k_1 = A \exp\left(-\frac{\Delta E}{RT}\right)$ , so that Equation (21) is reduced to:

$$X_A = \exp\left\{-k_1 t_1 \exp\left[-\frac{\Delta E}{R}\left(\frac{1}{T} - \frac{1}{T_0}\right)\right]\right\}. \quad (23)$$

Rewriting Equation (21) in terms of the rate  $\theta_1 (= k_1)$ , and letting (under the general framework of a nonlinear model):

$$\theta_2 = \frac{\Delta E}{R}, \quad t_2 = \frac{1}{T} - \frac{1}{T_0},$$

we obtain the suitable model for the developed experimental design framework as in Equation (1), with  $\theta \subseteq \mathcal{R}^2$ , see also [33]:

$$\eta = \exp\left[-\theta_1 t_1 \exp(-\theta_2 t_2)\right], \quad u = (t_1, t_2) \in U \subseteq \mathcal{R}^+ \times \Delta, \quad \Delta = [380, 450]. \quad (24)$$

This is the first order decay law for  $\eta$ , the molar fraction of A ( $X_A$ ).

If it is assumed that  $f(\infty, \theta) - f(0, \theta) = 1$  then the model is reduced to

$$\eta = \theta_1 - \exp(-\theta_2 u). \quad (25)$$

Under DR1 model (24) depends only on  $\theta_2$ , therefore prior information only for  $\theta_2$  is needed.

Recall to verify that the design depends only on  $\theta_2$ .

The reaction network of two irreversible first-order reactions in series proceeds according to the scheme  $A \xrightarrow{k_1} B \xrightarrow{k_2} C$ : a raw material  $A$  reacts to form a partial product  $B$  which, in turn, decomposes to give a substance  $C$ .

We shall refer below to the law of mass action is the proposition that the rate of the chemical reaction is directly proportional to the product of the activities or concentrations of the reactants. This law forms (at least) a differential equation and not only explains but also predicts behaviors of solutions in dynamic equilibrium.

Using the Guldberg-Waage form, of the reaction rates, to describe the network, provide for constant volume the simultaneous equations:

$$\left. \begin{aligned} \frac{d[A]}{dt} &= -k_1[A], & \frac{d[B]}{dt} &= k_1[A] - k_2[B], & \frac{d[C]}{dt} &= k_2[B], \\ [A]_0 + [B]_0 + [C]_0 &= [A] + [B] + [C] \end{aligned} \right\} \quad (26)$$

Integration of the differential equation for the concentration of A,  $[A]$ , with the initial conditions  $[A] = [A]_0$  at  $t = 0$  yields  $[A] = [A]_0 \exp(-k_1 t)$ . Substituting

into the differential equation for the concentration of  $B$ ,  $[B]$ , gives:

$$\frac{d[B]}{dt} + k_2[B] = k_1[A]_0 \exp(-k_1t). \quad (27)$$

With the initial conditions  $[B]_0 = [C]_0 = 0$  at  $t = 0$  the solution of (27) for the concentration of  $B$ , adopting the integrating factor method, is eventually

$$[B] = \frac{k_1[A]_0}{k_2 - k_1} [\exp(-k_1t) - \exp(-k_2t)]. \quad (28)$$

The experimental design model produced between the response  $\eta$  (concentration of  $B$ ,  $[B]$ ) and time  $t$ , will be in terms of Equation (1):

$$\eta = \frac{\theta_1}{\theta_1 - \theta_2} [\exp(-\theta_2 u) - \exp(-\theta_1 u)], \quad u \in U = \mathcal{R}^+, \quad \Theta \subseteq \mathcal{R}^2 \quad (29)$$

where  $\theta_1$  and  $\theta_2$  correspond to  $k_1$  and  $k_2$ , respectively, both functions of the temperature according to the Arrhenius law, [41, 43]. The optimal design for the model that appears in (29) is considered with  $(t, T) \in U \subseteq \mathcal{R}^+ \times \Delta$ ,  $\Delta = [380, 450]$ . For the experimental design of this model, also known as compartmental, see [2]. Kitsos and Kolovos [33] worked on a number of Chemical Models and calculated their D-optimal statistical design points, see following Section 2.

## 4 Differential Equations in Probability

Differential equations can be traced as a useful tool in Probability Theory, [5, 14, 23]. In the sequence we present in a compact form differential equations are essential in Applied Statistical Theory.

Although the difference equations play an important role in Probability theory differential equations are also very useful. The well known Chapman–Kolmogorov differential equations are the most applied either the forward one or the backwards type. In this section we are focused on the following type differential equations applied in Probability Theory.

### 4.1 Brownian Motion

The English botanist Brown, in 1827, introduced the physical phenomenon, known hear after the Brownian motion. In physics, was introduced in 1905 by Einstein. For us in this paper Brownian motion is an example of a continuous time, continues state space, Markov process  $X(t) = X_t$ ,  $t \in T$  with characteristics:

- B1: A process with independent increments. It is reflected then that the changes of  $X_t$  over non-over-lapping time periods are independent random variables (irv.).
- B2: The probability distribution of  $X(t_2) - X(t_1)$ ,  $t_2 > t_1$  depends only on  $t_2 - t_1$
- B3:  $P[X(r) - X(s) \leq x] = \frac{1}{\sqrt{2\pi(r-s)}} \int_{-\infty}^x \exp\left[-\frac{u^2}{2} A(r-s) du\right]$ ,  $r > s$  with  $A$  being a positive constant.

Let assume that  $X_0 = 0$ . Note that  $E(X(t)) = 0$ ,  $Var(X(t)) = \sigma^2(X_t) = At$ . At time  $t_0$  the particle process the position  $X(t_0) = X_0$ . The conditional probability density of  $X(t + t_0)$  given that  $X(t_0) = x_0$  is denoted by  $\phi(x, b|x_0)$ , see also Appendix 2.

It is important that Einstein proved that  $\phi(x, b|x_0)$  satisfies the partial differential equation, known as the diffusion equation:

$$\frac{\partial \phi}{\partial t} = D \frac{\partial^2 \phi}{\partial x^2}$$

with  $D$  being the diffusion coefficient. For the Diffusion equation the requests for the density function  $p(x, t|x_0)$  act as boundary conditions and provide a unique solution. Moreover:

$$D = \frac{2RT}{Nf}$$

where  $R$  is the gas constant,  $T$  the temperature,  $N$  is Avogadro's number and  $f$  is a coefficient of friction. With the appropriate scales  $D = \frac{1}{2}$  and then:

$$\phi(x, t|x_0) = \frac{1}{\sqrt{2\pi t}} \exp\left[-\frac{1}{2t}(x - x_0)^2\right]$$

the well known Gaussian distribution.

Therefore it comes easily that a Brownian motion is a stochastic process  $\{X(t), t \geq 0\}$  with:

- Br1.  $\forall t, s \quad X(t + s) - X(s) \sim N(0, \sigma^2 t)$
- Br2. Consider time intervals  $[t_1, t_2] \cap [t_3, t_4] = \emptyset$ ,  $t_L \leq t_{L+1}$ ,  $L = 1, 2, 3$  then  $X(t_{L+1}) - X(t_L) \sim N(0, \sigma^2 t)$   $L = 1, 2, 3$  and are independent random variables. This is true for  $r$  disjoint time intervals.
- Br3.  $X(0) = 0$ ,  $X(t)$  is continuous as  $t = 0$

The  $X^*(t) = X(t)/\sigma \sim N(0, 1)$  is the standard Brownian motion. Given the interval  $I = (t_0, t_1)$ ,  $X(0) = 0$ , the probability that  $X(t)$  has at least one zero within  $I$  is

$$\gamma = \frac{2}{\pi} \arccos \sqrt{\frac{t_0}{t_1}}$$

Another very important stochastic process in Probability Theory is the Wiener process. The Wiener process,  $W_t$ , is a continuous-time process based on:

- W1 :  $W_0 = 0$
- W2 :  $W_t$  is almost surely continuous
- W3 :  $W_t$  has independent increments
- W4 :  $W_t - W_s \sim N(0, t - s), \quad 0 \leq s < t$

Through the Wiener process the stochastic differential equation (s-de):

$$dX_t = \mu_t dt + \sigma_t dW_t, \quad \mu_t = \mu(X_t, t), \quad \sigma_t = \sigma(X_t, t)$$

and  $\mu$  is the drift and the diffusion coefficient:

$$D = D(X_t, t) = \frac{1}{2} \sigma^2(X_t, t)$$

The Fokker–Plank equation for the probability density  $p(x, t)$  of a given variable (rv)  $X_t$  is:

$$\frac{\partial p}{\partial t} = -\frac{\partial}{\partial x} [\mu(x, t)p(x, t)] + \frac{\partial^2}{\partial x^2} [D(x, t)p(x, t)]$$

When the drift is zero and the diffusion constant we are referred to Brownian motion discussed above.

### 4.2 Pure Birth Process

Suppose we have a system  $E = \{E_1, E_2, \dots, E_j, E_{j+1}, \dots, E_n\}$  and that from state  $E_j$  you can move only to  $E_{j+1}$ . Moreover in state  $E_n$  at time  $t$ , the probability of a jump at a sort time interval  $I = (t, t + h)$  equals  $\lambda_n h + O(h)$ .

The probability of more than one jump with  $I$  is  $O(h)$ .

Let  $P_n(t)$  be the probability that at time  $t$  the system is at stage  $E_n$ . Then,

$$\frac{P_n(t + h) - P_n(t)}{h} = -\lambda P_n(t) + P_{n-1}(t) + \frac{O(h)}{h} \tag{30}$$

with  $h \rightarrow 0$  the above is reduced to:

$$P'_n(t) = -\lambda P_n(t) + \lambda P_{n-1}(t) \tag{31}$$

with:

$$P'_0(t) = -\lambda P_0(t) \tag{32}$$

The Poisson process satisfies the above (de) and it holds

$$P_k(0) = 1, \quad P_n(0) = 0, \quad n \neq k$$

The Poisson process acts as the solution of (31) with form

$$P_n(t) = \frac{(\lambda t)^n}{n!} e^{-\lambda t}, \quad n \geq 1.$$

### 4.3 The Birth-and-Death Process

The assumption about the system  $E = \{E_1, E_2, \dots, E_{j-1}, E_j, E_{j+1}\}$  is now changing, that the system changes only through transitions from states to their nearest ones (i.e., can move backwards). In principle from  $E_n$  moves either to  $E_{n+1}$  or  $E_{n-1}$ , while  $E_0$  give rise only to  $E_1$  movement. Consider the time interval  $I = (t, t + h)$  the probability that the transition is from  $E_n$  to  $E_{n+1}$  equals to  $\lambda_n h + O(h)$ , while the probability that from  $E_n$  there is a movement to  $E_{n-1}$  equals to  $\mu_n h + O(h)$ . The probability that within  $I$  move that one change takes place is  $O(h)$ .

If we let  $v_n = \lambda_n + \mu_n$ , due to the independence of the events, it can be eventually proved:

$$P_n(t + h) = P_n(t)[1 - v_n h] + \lambda_{n-1} h P_{n-1}(t) + \mu_{n+1} h P_{n+1}(t) + O(h) \quad (33)$$

Thus, from the above equation we obtain (divining by  $h, h \rightarrow 0$ )

This theory is extended to more variables under the light of physical applications and will not discussed here. We believe that the birth / death process is very important in applications.

$$\begin{aligned} P'_n(t) &= -v_n P_n(t) + \lambda_{n-1} P_{n-1}(t) + \mu_{n+1} P_{n+1}(t), \quad n \geq 1 \\ P'_0(t) &= -\lambda_0 P_0(t) + \mu_1 P_1(t) \end{aligned} \quad (34)$$

with initial conditions:

$$P_1(0) = 1, \quad P_n(0) = 0$$

If the coefficients are bounded there is a unique solution, under the regularity condition  $\sum P_n(t) = 1$ , while for  $\sum P_n(t) < 1$  there exist infinitely many solutions.

The solution of (34) can be obtained by induction. Let:

$$\pi_0 = 1, \quad \pi_j = \frac{\lambda_0 \lambda_1 \cdots \lambda_{j-1}}{\mu_1 \mu_2 \cdots \mu_j} \quad j \geq 1$$

Then from (34) we obtain:  $P_1 = \lambda_0/\mu_1$        $P_1 = \pi_1 P_0$

Assuming that:  $P_k = \pi_k P_0$        $k = 1, \dots, j$  we have eventually, [23]  $P_{j+1} = \pi_{j+1} P_0$

Trivially the consequence  $P_j, j = 1, 2, \dots$  defines a distribution, if  $\sum P_j = 1$ . Provided that  $\sum \pi_k < \infty$  it holds:

$$P_j = \frac{\pi_j}{\sum \pi_k} \quad j = 0, 1, 2, \dots$$

It can be proven that the limit

$$\lim_{k \rightarrow \infty} P_n(t) = P_n \tag{35}$$

exist and is independent of the initial conditions. Due to (35) for an ordinary Markov chain it holds:

$$P_{n,n+1} = \frac{\lambda_n}{\nu_n}, \quad P_{n,n-1} = \frac{\mu_n}{\nu_n}$$

Let  $a_i, i = 1, 2, \dots, n$  be the probability of absorption into state 0 from the initial state 1. Then we obtain that:

$$a_i = \frac{\lambda_i}{\nu_i} a_{i+1} + \frac{\mu_i}{\nu_i} a_{i-1} \quad i \geq 1$$

where  $a_0 = 0$ . The above relation can be considered through the “embedded random walk” associated to a given  $b-d$  process.

As far as applications concerns see the early work [14] and [23].

The values of birth and death process give rise to a number of applications. The most well known is coming from the queuing theory, where the Kendall’s notation is described by three factors: if “Arrivals” denotes the time between arrivals in queue, “Service” is the time service distribution and “channels” is the number of service channels open at the node, the queue is denoted by Assistance/Service/Channel. If the arrival process and the service time distribution is Markovian (or Memoryless) the M/M/1 system is a typical one. In principle the Arrival process is a Poisson process and the Service time distribution is an exponential distribution, see the pioneering work [24].

The birth and death rates are constant, say  $\lambda_i = \lambda, \mu_i = \mu$  with the average rate of arrivals to be  $\lambda$ , and the average service time to be  $1/\mu$ .

The corresponding (de) for the evaluation of the probability that M/M/1 system is at state  $k$  at the given time  $t$  with  $\nu = \lambda + \mu$  are:

$$P'_0(t) = \mu P_1(t) - \lambda P_0(t)$$

$$P'_k(t) = \lambda P_{k-1}(t) + \mu P_{k+1}(t) - \nu P_k(t)$$

when more than one channel exists a M/M/C queue exists with  $(C) = S$  servers and an infinite buffer, the birth and death process is characterized by  $\lambda_i = \lambda$  as above and:

$$\mu_i = \mu, \text{ for } i \leq S - 1, \mu_i = S\mu, \text{ for } i \geq S$$

We introduce the notation:

$$v_k = \lambda + k\mu, \quad v_s = \lambda + S\mu$$

Then the corresponding (de) are:

$$P'_0(t) = \mu P_1(t) - \lambda P_0(t)$$

$$P'_k(t) = \lambda P_{k-1}(t) + (k+1)\mu P_{k+1}(t) - v_k P_k(t), \text{ for } k = 1, 2, \dots, S-1$$

$$P'_k(t) = \lambda P_{k-1}(t) + S\mu P_{k+1}(t) - v_s P_k(t), k \geq S$$

For different queues different (de) are needed, so (de) are essential in Probability Theory.

## 5 Discussion

The Design theory is widely applied, [36, 39, 48] to a number of different oriented experiments. At the same time different fields of Mathematics are considering Statistics, [1], among others. The target of this paper is to discuss how differential equations (de) are applied to Statistics. The back-bone of Statistics is the optimal Design theory. Although the problem of (de) seems to have with Statistics as intersection the null set this is not the true situation. As the Ca problem as faced [31] seems to have no relation with [38] this is not true: to solve the low dose problems, as Robins-Monro iterations, it is similar to solve a non-linear equation with Newton-Raphson method.

Same story with the solution of the partial (de), [42] and the optimal Design Problem, [16, 44]. Still the hidden theory is based on (de), so does a number of Growth Models, facing by Statistics. As it was proven in Section 3 the Growth Curves in most of the Chemical, Biological cases are based on non-linear problems coming from (de). For the Linear case of Growth Curves see [18] were no (de) are needed.

In Probability theory the use of (de) is clear to birth and death stochastic processes are based on random walks, see Appendix 2. In Physics or Electrical Engineering (de) are applied extensively, see Appendix 1, and their use defines a solid and compact theoretical model, while in Statistics there exists a stochastic orientation



of the use: either the stochastic error in models or the Probabilistic development in random works. The Browian motion, eventually joints three sciences: Math, Stat and Physics.

There are three lines of though for facing a de: the theoretical inside, see [42] the applied—how we can solve a de—see [30] and the Numerical Analysis point of view, see [17].

It is true that no such a deep development of de exists in Statistics. But we would say, that as Statistics serves all Sciences, so do and de. We tried to cover how this so widely used lines of thought are, eventually, communicate.

## Appendix 1

Let us consider the main types of fundamental electrical circuits, so essential to build a differential equations approach to electrical circuits. To avoid any confusion the imaginary unit is denoted by  $j = (0, 1)$  and the current is denoted by  $i = i(t)$ .

### 1. An $RL$ —circuit.

The voltage source, where it is assumed that it is AC, and thus  $E = E_0 \sin \omega t$ . Then based on the 2nd law of Kirchoff:

$$L \frac{di}{dt} + Ri = E_0 \sin \omega t \quad t \geq 0$$

The (general) solution is:

$$i(t) = i_0 e^{-\frac{R}{L}t} + E_0 (R^2 + L^2 \omega^2)^{-1/2} \sin(\omega t - \text{Arctan} \frac{\omega L}{R}) + E_0 \omega (R^2 + L^2 \omega^2)^{-1} e^{-\frac{R}{L}t}$$

### 2. An $RC$ —Circuit.

Then it holds:

$$Ri + \frac{1}{C} \int_0^t i dt_1 = E(t)$$

And if  $E(t) = \text{constant} = E_0$  then:

$$Ri + \frac{1}{C} \int_0^t i dt_1 = E_0$$

And the solution is:

$$i = i(t) = \frac{E_0}{R} e^{-t/RC},$$

3. *RC—circuit* with  $E = E_0 \cos(\omega t)$ 

Then it holds:

$$Ri + \frac{1}{C} \int_0^t i dt_1 = E_0 \cos(\omega t)$$

And if we let:

$$u = \int_0^t i dt_1$$

We obtain that:

$$R \frac{du}{dt} + \frac{u}{C} = E_0 \cos(\omega t), \quad u(0) = 0$$

Thus:

$$\begin{aligned} u &= t e^{-t/RC} \int_0^t \frac{E_0}{R} e^{-t_1/RC} \cos(t) du \\ &= \frac{E_0}{R} \left[ \frac{(1/RC) (\cos(\omega t) - e^{-t/RC}) + \omega \sin(\omega t)}{(1/RC)^2 + \omega^2} \right] = u_1 + u_2 \end{aligned}$$

With:

$$u_1 = \frac{-E_0/R (e^{-t/RC}/RC)}{(1/RC)^2 + \omega^2}, \quad u_2 = \frac{E_0 (1/RC) \cos(\omega t) + \omega \sin(\omega t)}{R ((1/RC)^2 + \omega^2)}$$

Eventually we can evaluate that:

$$\begin{aligned} i_1 &= \frac{du_1}{dt} = \frac{-E_0/R (e^{-t/RC}/RC)^2}{(1/RC)^2 + \omega^2} = \frac{E_0}{R} e^{-t/RC} / (1 + \omega^2 (RC)^2) \\ i_2 &= \frac{du_2}{dt} = \frac{E_0}{R} \left[ \frac{-\omega RC \sin(\omega t) + \omega^2 (RC)^2 \cos(\omega t)}{1 + (RC)^2 \omega^2} \right] \\ &\cong \frac{E_0}{R} \left[ \frac{-\omega RC \sin(\omega t)}{1} \right] = -\omega C E_0 \sin(\omega t) \end{aligned}$$

Notice that with  $t \ll i_1$ ,  $t = RC$  say then:  $i_1 \cong \frac{E_0}{R} \frac{1}{e}$ .

4. *RL—circuit*

Assuming that  $E = E_0$  the (de) concerning this circuit is:

$$L \frac{di}{dt} + Ri = E_0, \quad L(0) = 0$$

With solution:  $i = \frac{E_0}{R} - \frac{E_0}{R} e^{-\frac{R}{L}t}$

5. *LC—circuit* with  $E = E_0$ .

The corresponding (de) is:

$$L \frac{di}{dt} + \frac{1}{C} \int_0^t i dt_1 = E(t) \text{ and the solution is :}$$

$$i = E_0 \sqrt{C/L} \sin \left( t / \sqrt{LC} \right) \text{ with } i(0) = 0 \text{ and } L i'(0) = E_0$$

6. *LC—circuit* with  $E = E_0 \cos(\omega t)$ .

The corresponding equation is:  $L \frac{di}{dt} + \frac{1}{C} \int_0^t i dt = E_0 \cos(\omega t)$ .

With solution eventually:  $i = \frac{e^{j\omega t}}{(-L\omega^2 + j\omega R + 1/C)}$

As it is known that it holds:

$$\sin(\omega t) = \frac{e^{j\omega t} - e^{-j\omega t}}{2j}$$

with the appropriate calculations the general solution eventually is:

$$i = - \frac{E_0 \omega (1/C - L\omega^2) \sin(\omega t + E_0 R \omega^2 \cos(\omega t))}{(-L\omega^2 + 1/C)^2 + \omega^2 R^2}$$

There is an extensive approach which is beyond the target of this appendix.

## Appendix 2: Introduction to Heat Equation

Consider a random walk, where we assume that the probability to move to the closest up, down, backwards, forwards points are equal to 1/4. Let  $P(x, y; t)$  be the probability that a particle at time  $t$  is at  $(x, y) \in R^2$  point. Then we can see that at time  $t + 1$  holds:

$$P(x, y; t + 1) = \frac{1}{4} \{ P(x-1, y; t) + P(x, y-1; t) + P(x+1, y, t) + P(x, y+1; t) \}$$

Thus for the difference:

$$P(x, y; t + 1) - P(x, y; t) = \frac{1}{4} \left\{ P(x+1, y; t) - 2P(x, y; t) + P(x-1, y, t) + P(x, y+1; t) - 2P(x, y; t) + P(x, y-1; t) \right\}$$

This difference equation approximated by the two-dimensional heat equation:

$$\frac{\partial P}{\partial t} = C \left( \frac{\partial^2 P}{\partial x^2} + \frac{\partial^2 P}{\partial y^2} \right)$$

If we assume that the random—walk takes place in a limited domain  $D$ , usually of the form  $D = (l, u) \times [L, U] \subseteq R^2$ . More over it is assumed that the particle is absorbed when it reaches the boundary. Let  $(x_0, y_0)$  be the boundary points and  $W = W(x_b, y_b)$  the associated “profit that is paid out”.

If we denote by  $P(x, y; x_0, y_0)$  the probability that the particle start from the (interior) point  $(x, y)$  to be absorbed at the boundary  $(x_b, y_b)$  the expected “profit” is:

$$u(x, y) = \sum_b P(x, y; x_b, y_b) W(x_b, y_b)$$

and satisfies the difference equation

$$u(x, y) = \frac{1}{4} \{u(x+1, y) + u(x-1, y) + u(x, y+1) + u(x, y-1)\}$$

with  $u(x_b, y_b) = W(x_b, y_b)$ . This equation is a well-known approximation of the Laplace equation:

$$\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = 0$$

Recall that

$$\nabla^2 u = \begin{cases} 0 & \text{Laplace equation} \\ C(x, y) & \text{Poisson equation} \end{cases}$$

Through the Laplace equation three different well-known problems, associated with the boundary are defined: Dirichlet’s problem, Neuman’s problem, Robbin–Churchill’s problem for elliptic equations.

## References

1. S.I. Amari, *Differential Geometrical Methods in Statistics* (Springer, Berlin, 2012)
2. A.C. Atkinson, K. Chaloner, A.M. Herzberg, J. Juritz, Optimum experimental designs for properties of a compartmental model. *Biometrics* **49**(2), 325–337 (1993)
3. Y. Bard, *Nonlinear Parameter Estimation* (Academic Press, New York, 1974)
4. M.D. Bates, G.D. Watts, *Nonlinear Regression Analysis and Its Applications* (Wiley, New York, 1988)

5. B.R. Bhat, *Modern Probability Theory* (Wiley, New Delhi, 1985)
6. G.E.P. Box, H.L. Lucas, Design of experiments in Nonlinear situation. *Biometrika* **49**, 77–90 (1959)
7. G.E.P. Box, W.G. Hunter, The experimental study of physical mechanisms, *Technometrics* **7**, 23–42 (1965)
8. S. Brunauer, P.H. Emmett, E. Teller, Adsorption of gases in multimolecular layers. *J. Am. Chem. Soc.* **60**(2), 309–319 (1938)
9. N.L. Carr, Kinetics of catalytic isomerisation of n-pentane. *Indus. Eng. Chem.* **52**(5), 391–396 (1960)
10. M. Davis, P.J.E. Peebles, On the integration of the BBGKY equations for the development of strongly nonlinear clustering in an expanding universe. *Astrophys. J. Suppl. Ser.* **34**, 425–50 (1977)
11. M.E. Davis, R.J. Davis, *Fundamentals of Chemical Reaction Engineering*, 1st edn. (McGraw Hill, New York, 2003), 368 pp.
12. C.T. De Wit, Resource use efficiency in agriculture. *Agric. Syst.* **40**(1–3), 125–151 (1992)
13. J. Downie, K.A. Shelstad, W.F. Graydon, Kinetics of the vapour-phase oxidation of toluene over a vanadium catalyst. *Can. J. Chem. Eng.* **39**(5), 201–204 (1961)
14. W. Feller, *An Introduction to Probability Theory and Its Applications* (Wiley, New York, 1968)
15. R.A. Fisher, On the Mathematical Foundations of Theoretical Statistics. *Philos. Trans. R. Soc. Lond. Ser. A*, **222**, 309–368 (1922)
16. I. Ford, C.P. Kitsos, D.M. Titterington, Recent advances in nonlinear experimental design. *Technometrics* **31**(1), 49–60 (1989)
17. C.E. Fröberg, *Numerical Mathematics: Theory and Computer Applications* (Benjamin-Cummings, San Francisco, 1985)
18. A.F. Graybill, *Theory and Application of the Linear Model* (Duxbury Press, 1976)
19. K. Harmsen, A modified Mitscherlich equation for rainfed crop production in semi-arid areas: 1. *Theory. Neth. J. Agric. Sci.* **48**(3), 237–250 (2000)
20. P.D.H. Hill, D-optimal Designs for partially Nonlinear Regression Models. *Technometrics* **22**, 275–276 (1980)
21. L.H. Hosten, G.F. Froment, Isomerization of n-pentane. *Ind. Eng. Chem. Process. Des. Dev.* **10**(2), 280–287 (1971)
22. I.S. Jaswal, R.F. Mann, J.A. Juusola, J. Downie, The vapour-phase oxidation of benzene over a vanadium pentoxide catalyst. *Can. J. Chem. Eng.* **47**(3), 284–287 (1969)
23. S. Karlin, M.H. Taylor, *A First Course in Stochastic Processes* (Academic Press, New York, 1975)
24. D.G. Kendall, Stochastic processes occurring in the theory of queues and their analysis by the method of the imbedded Markov chain. *Ann. Math. Stat.* **1**, 338–54 (1953)
25. C.P. Kitsos, Design and inference for Non-linear Problems. Unpublished Ph.D. thesis, Univ. of Glasgow
26. C.P. Kitsos, Fully-sequential procedures in nonlinear design problems. *Comput. Stat. Data Anal.* **8**(1), 13–19 (1989)
27. C.P. Kitsos, On the support points of D-optimal non-linear experimental design for kinetics, in *MODA4: Advances in Model-Oriented Data Analysis*, ed. by C.P. Kitsos, W.G. Muller (Physica Verlag, Heidelberg, 1995), pp. 71–76
28. C.P. Kitsos, Nonlinear: optimal—sequential experiment designs and applications, in *Proceedings of the 51st International Conference on Industrial Statistics: Aims and Computational Aspects*, ed. by C.P. Kitsos, L. Edler, Athens, Greece, August 16–17, 1997 (Physica-Verlag, Heidelberg, 1997), pp. 151–163
29. C.P. Kitsos, Design aspects for the Michaelis Menten model. *Biometrical Lett.* **38**(1), 53–66 (2001)
30. C.P. Kitsos, *Technological Mathematics and Statistics*, vols. I, II (New Technologies Pub., Greece, 2009)
31. C.P. Kitsos, Sequential Approaches for Ca tolerance models. *Biometrie und Medizinische Informatik, Greifswalder Seminarberichte*, vol. 18 (Shaker Verlag, 2011), pp. 87–98

32. C.P. Kitsos, K.G. Kolovos, An Optimal Calibration Design for pH Meters. *Instrumentation Sci. Tech.* **38**(6), 436–447 (2010)
33. C.P. Kitsos, K.G. Kolovos, A compilation of the D-optimal designs in chemical kinetics. *Chem. Eng. Commun.* **200**(2), 185–204 (2013)
34. S. Lopez, J. France, W.J. Gerrits, M.S. Dhanoa, D.J. Humphries, J. Dijkstra, A generalized Michaelis-Menten equation for the analysis of growth. *J. Anim. Sci.* **78**(7), 1816–1828 (2000)
35. P. Mars, D.W. van Krevelen, Oxidations carried out by means of vanadium oxidized catalysts. *Chem. Eng. Sci.* **3**, 41–59 (1954)
36. I. Martinez, I. Ortiz, C. Rodriguez, Optimal design for weighted rational models. *Appl. Math. Lett.* **2009**, **22**(12), 1892–1895
37. E.A. Mitscherlich, The law of the minimum and the law of diminishing soil productivity (In German). *Landwirtschaftliche Jahrbuecher* **38**, 537–552 (1909)
38. J. Ortega, W. Rheinbolt, *Iterative Solution of Equations in Several Variables* (Academic Press, New York 1970)
39. L.J. Rodriguez-Aragon, J. Lopez-Fidalgo, Optimal designs for the Arrhenius equation. *Chemom. Intell. Lab. Syst.* **77**(1–2), 131–138 (2005)
40. G.A.F. Seber, C.J. Wild, *Nonlinear Regression* (Wiley, New York, 1989), 768 pp.
41. K. Smith, On the standard deviation of adjusted and interpolated values of an observed polynomial function and its constants and the guidance they give towards a proper choice of the distribution of observations. *Biometrika* **12**(1–2), 1–85 (1918)
42. G.D. Smith, *Solution of Partial Differential Equations* (Oxford University Press, London, 1965)
43. J.M. Smith, *Chemical Engineering Kinetics* (McGraw Hill, Singapore, 1981)
44. D. Steinberg, W. Hunter, Experimental design: review and comment. *Technometrics* **26**(2), 71–97 (1984)
45. H.M.J. Thorney, *Mathematical Models in Plant Physiology* (Academic Press, London, 1976)
46. G.B. Van Albada, Formation and evolution of clusters of galaxies (Errata: 15 330). *Bull. Bull. Astron. Inst. Neth.* **15**, 165 (1960)
47. Q. Wright, *A Study of War. 1942*, 2 vols. (Univ. of Chicago Press, Chicago, 1942)
48. V. Zarikas, V. Gikas, C.P. Kitsos, Evaluation of the optimal design Cosinor model for enhancing the potential of robotic theodolite kinematic observation. *Measurement* **43**(10), 1416–1424 (2010)

# Geometric Derivation and Analysis of Multi-Symplectic Numerical Schemes for Differential Equations



Odysseas Kosmas, Dimitrios Papadopoulos, and Dimitrios Vlachos

**Abstract** In the current work we present a class of numerical techniques for the solution of multi-symplectic PDEs arising at various physical problems. We first consider the advantages of discrete variational principles and how to use them in order to create multi-symplectic integrators. We then consider the nonstandard finite difference framework from which these integrators derive. The latter is now expressed at the appropriate discrete jet bundle, using triangle and square discretization. The preservation of the discrete multi-symplectic structure by the numerical schemes is shown for several one- and two- dimensional test cases, like the linear wave equation and the nonlinear Klein–Gordon equation.

## 1 Introduction and Motivation

In general, symplectic integrators are robust, efficient and accurate in preserving the long time behavior of the solutions of Hamiltonian ordinary differential equations (ODEs) [1]. The basic feature of a symplectic integrator is that the numerical performance is designed to preserve a physical observable property, i.e., the symplectic form at each time step. Recently, it was shown that many conservative partial differential equations (PDEs) allow for description similar to the symplectic structure of Hamiltonian ODEs, called the multi-symplectic formulation (see, e.g., Refs. [2–5]). For example, in Ref. [2] authors develop the multi-symplectic

---

O. Kosmas

Modelling and Simulation Centre, MACE, University of Manchester, Manchester, UK  
e-mail: [odysseas.kosmas@manchester.ac.uk](mailto:odysseas.kosmas@manchester.ac.uk)

D. Papadopoulos

Delta Pi Systems Ltd., Thessaloniki, Greece  
e-mail: [dimitris@delta-pi-systems.eu](mailto:dimitris@delta-pi-systems.eu)

D. Vlachos (✉)

Department of Informatics & Telecommunications, University of Peloponnese, Tripolis, Greece  
e-mail: [dvlachos@uop.gr](mailto:dvlachos@uop.gr)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,  
Springer Optimization and Its Applications 173,  
[https://doi.org/10.1007/978-3-030-72563-1\\_11](https://doi.org/10.1007/978-3-030-72563-1_11)

231

structure of Hamiltonian PDEs from a Lagrangian formulation, using the variational principle. The wave equation and its multisymplectic structure have been studied by [6–8] from the Hamiltonian viewpoint.

On the other hand, in the past decades, nonstandard finite difference schemes have been well established by Mickens [9–11] to compensate the weaknesses that may be caused by standard finite difference methods as, such as the numerical instabilities. Regarding the positivity, the boundedness, and the monotonicity of solutions, nonstandard finite difference schemes have a better performance than standard ones, due to their flexibility to construct a nonstandard finite difference method. The latter can preserve certain properties and structures, which are obeyed by the original equations.

In the present paper, following our previous work [12] we pay special attention to the geometric structure of multisymplectic integrators through the use of nonstandard finite difference schemes for variational partial differential equations (PDEs). The considered approach comes as a first step towards developing a Veselov type discretization for PDEs in variational form, e.g. [2, 4, 5] and combines it with nonstandard nonstandard finite difference schemes of Mickens [9–11]. The resulting multisymplectic-momentum integrators have very good energy performance in the level of the conservation of a nearby Hamiltonian, under appropriate circumstances, up to exponentially small error [2].

In the following Section 2 we present a short overview of the standard numerical techniques relying on variational integrator schemes and their special case of exponential variational integrators in Section 3. Afterwards, nonstandard finite difference properties are employed for the derivation of nonstandard variational integrators by using a triangle discretization of the spacetime (Section 4.1). Then, in Sections 5 and 6, we demonstrate concrete applications of the proposed integrators, for the numerical solution of the linear wave equation, the Laplace equation and the Poisson equation. In Section 7, we perform dispersion analysis and convergence experiments to further illustrate the numerical properties of the method. Finally, in Section 8, we summarize the main conclusions coming out of our study.

## 2 Review of Variational Integrators

The discrete Euler–Lagrange equations can be derived in correspondence to the steps of derivation of the Euler–Lagrange equations in the continuous formulation of Lagrangian dynamics [3]. Denoting the tangent bundle of the configuration manifold  $Q$  by  $TQ$ , the continuous Lagrangian  $L : TQ \rightarrow \mathbb{R}$  can be defined. In the discrete setting, considering approximate configurations  $q_k \approx q(t_k)$  and  $q_{k+1} \approx q(t_{k+1})$  at the time nodes  $t_k, t_{k+1}$ , with  $h = t_{k+1} - t_k$  being the fixed time step, a discrete Lagrangian  $L_d : Q \times Q \rightarrow \mathbb{R}$  is defined to approximate the action integral along the curve segment between  $q_k$  and  $q_{k+1}$ , i.e.,



$$L_d(q_k, q_{k+1}) \approx \int_{t_k}^{t_{k+1}} L(q(t), \dot{q}(t)) dt. \tag{1}$$

Defining the discrete trajectory  $\gamma_d = (q_0, \dots, q_N)$ ,  $N \in \mathbb{N}$ , one can obtain the action sum

$$S_d(\gamma_d) = \sum_{k=1}^{N-1} L_d(q_k, q_{k+1}). \tag{2}$$

The discrete Hamilton’s principle states that a motion  $\gamma_d$  of the discrete mechanical system extremizes the action sum, i.e.,  $\delta S_d = 0$ . Through differentiation and rearrangement of the terms, holding the end points  $q_0$  and  $q_N$  fixed, the discrete Euler–Lagrange equations are obtained [3]

$$D_2 L_d(q_{k-1}, q_k) + D_1 L_d(q_k, q_{k+1}) = 0, \quad k = 1, \dots, N - 1, \tag{3}$$

where the notation  $D_i L_d$  indicates derivative with respect to the  $i$ -th argument of  $L_d$ , see also [3, 12–16].

The definition of the discrete conjugate momentum at time steps  $k$  and  $k + 1$  reads

$$p_k = -D_1 L_d(q_k, q_{k+1}), \quad p_{k+1} = D_2 L_d(q_k, q_{k+1}), \quad k = 0, \dots, N - 1. \tag{4}$$

The above equations, also known as position-momentum form of a variational integrator, can be used when an initial condition  $(q_0, p_0)$  is known, to obtain  $(q_1, p_1)$ .

To construct high order methods, we approximate the action integral along the curve segment between  $q_k$  and  $q_{k+1}$  using a discrete Lagrangian that depends only on the end points. We obtain expressions for configurations  $q_k^j$  and velocities  $\dot{q}_k^j$  for  $j = 0, \dots, S - 1$ ,  $S \in \mathbb{N}$  at time  $t_k^j \in [t_k, t_{k+1}]$  by expressing  $t_k^j = t_k + C_k^j h$  for  $C_k^j \in [0, 1]$  such that  $C_k^0 = 0$ ,  $C_k^{S-1} = 1$  using

$$q_k^j = g_1(t_k^j)q_k + g_2(t_k^j)q_{k+1}, \quad \dot{q}_k^j = \dot{g}_1(t_k^j)q_k + \dot{g}_2(t_k^j)q_{k+1}, \tag{5}$$

where  $h \in \mathbb{R}$  is the time step. We choose functions

$$g_1(t_k^j) = \sin\left(u - \frac{t_k^j - t_k}{h}u\right) (\sin u)^{-1}, \quad g_2(t_k^j) = \sin\left(\frac{t_k^j - t_k}{h}u\right) (\sin u)^{-1}, \tag{6}$$

to represent the oscillatory behavior of the solution, see [17, 18]. For continuity,  $g_1(t_{k+1}) = g_2(t_k) = 0$  and  $g_1(t_k) = g_2(t_{k+1}) = 1$  is required.

For any different choice of interpolation used, we define the discrete Lagrangian by the weighted sum

$$L_d(q_k, q_{k+1}) = h \sum_{j=0}^{S-1} w^j L(q(t_k^j), \dot{q}(t_k^j)), \quad (7)$$

where it can be easily proved that for maximal algebraic order

$$\sum_{j=0}^{S-1} w^j (C_k^j)^m = \frac{1}{m+1}, \quad (8)$$

where  $m = 0, 1, \dots, S-1$  and  $k = 0, 1, \dots, N-1$  see [17, 18].

Applying the above interpolation technique with the trigonometric expressions of (6), following the phase lag analysis of [13, 14, 17, 18], the parameter  $u$  can be chosen as  $u = \omega h$ . For problems that include a constant and known domain frequency  $\omega$  (such as the harmonic oscillator) the parameter  $u$  can be easily computed. For the solution of orbital problems of the general  $N$ -body problem, where no unique frequency is given, a new parameter  $u$  must be defined by estimating the frequency of the motion of any moving point mass [16, 19–21].

### 3 Exponential Integrators

We now consider the Hamiltonian systems

$$\ddot{q} + \Omega q = g(q), \quad g(q) = -\nabla U(q), \quad (9)$$

where  $\Omega$  is a diagonal matrix (will contain diagonal entries  $\omega$  with large modulus) and  $U(q)$  is a smooth potential function. We are interested in the long time behavior of numerical solutions when  $\omega h$  is not small.

Since  $q_{n+1} - 2 \cos(h\omega)q_n + q_{n-1} = 0$  is an exact discretisation of (9) we can consider the numerical scheme

$$q_{n+1} - 2 \cos(h\omega)q_n + q_{n-1} = h^2 \psi(\omega h) g(\phi(\omega h)q_n), \quad (10)$$

where the functions  $\psi(\omega h)$  and  $\phi(\omega h)$  are even, real-valued functions satisfying  $\psi(0) = \phi(0) = 1$ , see [1]. The resulting methods using the latter numerical scheme are known as exponential integrators (for some examples of those integrators see the appendix).

### 3.1 Exponential High Order Variational Integrators

If we now use the phase fitted variational integrator for the system (9) the result of the discrete Euler–Lagrange equations (3) will be

$$q_{n+1} + \Lambda(u, \omega, h, S)q_n + q_{n-1} = h^2\Psi(\omega h)g(\Phi(\omega h)q_n), \tag{11}$$

where

$$\Lambda(u, \omega, h, S) = \frac{\sum_{j=0}^{S-1} w^j \left[ \dot{g}_1(t_k^j)^2 + \dot{g}_2(t_k^j)^2 - \omega^2(g_1(t_k^j)^2 + g_2(t_k^j)^2) \right]}{\sum_{j=0}^{S-1} w^j \left[ \dot{g}_1(t_k^j)\dot{g}_2(t_k^j) - \omega^2 g_1(t_k^j)g_2(t_k^j) \right]}. \tag{12}$$

Using the above expressions, to obtain exponential variational integrators that use expressions for configurations  $q_k^j$  and velocities  $\dot{q}_k^j$  taken from (5), we get

$$\Lambda(u, \omega, h, S) = -2 \cos(\omega h). \tag{13}$$

In [16] we have proved (using the phase lag analysis of [22]) that exponentially fitted methods using phase fitted variational integrators can be derived when (13) holds. So phase fitted variational integrators using trigonometric interpolation can be considered as exponential integrators, i.e. when using phase fitted variational integrators, keeping the phase lag zero the resulting methods are exponentially fitted methods (exponential integrators). Those methods has been tested on several numerical results in [16].

### 3.2 Frequency Estimation for Mass Points Motion in Three Dimensions

In our previous work [16], we constructed adaptive time step variational integrators using phase fitting techniques and estimated the required frequency through the use of a harmonic oscillator with given frequency  $\omega$ . Here, in solving the general  $N$ -body problem by using a constant time step, a new frequency estimation is necessary in order to find for each body i) the frequency at an initial time  $t_0$  and ii) the frequency at time  $t_k$  for  $k = 1, \dots, N - 1$ .

It is now clear that, by applying the trigonometric interpolation (6), the parameter  $u$  can be chosen as  $u = \omega h$ . For problems for which the domain of frequency  $\omega$  is fixed and known (such as the harmonic oscillator) the parameter  $u$  can be easily computed. For the solution of orbital problems involved in the general  $N$ -body

problem, where no unique frequency is determined, the parameter  $u$  must be defined by estimating the frequency of the motion of any moving material point.

Towards this purpose, we consider the general case of  $N$  masses moving in three dimensions. If  $q_i(t)$  ( $i = 1, \dots, N$ ) denotes the trajectory of the  $i$ -th particle, its curvature can be computed from the known expression

$$k_i(t) = \frac{|\dot{q}_i(t) \times \ddot{q}_i(t)|}{|\dot{q}_i(t)|^3}, \quad (14)$$

where  $\dot{q}_i(t)$  the velocity of the  $i$ -th mass with magnitude  $|\dot{q}_i(t)|$  at a point  $q_i(t)$ . After a short time  $h$ , the angular displacement of that mass is  $h|\dot{q}_i(t) \times \ddot{q}_i(t)|/|\dot{q}_i(t)|^2$ , which for each mass's actual frequency gives the expression

$$\omega_i(t) = \frac{|\dot{q}_i(t) \times \ddot{q}_i(t)|}{|\dot{q}_i(t)|^2}. \quad (15)$$

From (14) and (15) the well known relation  $\omega_i(t) = k_i(t)|\dot{q}_i(t)|$  holds (see also [16]).

For the specific case of many-particle physical problems, that can be described via a Lagrangian of the form  $L(q, \dot{q}) = \frac{1}{2}\dot{q}^T M(q)\dot{q} - V(q)$ , where  $M(q)$  represents a symmetric positive definite mass matrix and  $V$  is a potential function, the continuous Euler–Lagrange equations are  $M(q)\ddot{q} = -\nabla V(q)$ . In this case, the expression for frequency estimation (15), referred to the  $i$ -th body at time  $t_k$ ,  $k = 1, \dots, N - 1$ , takes the form

$$\omega_i(t_k) = h^{-1} \frac{|M^{-1}(q_k)p_k \times (M^{-1}(q_k)p_k - M^{-1}(q_{k-1})p_{k-1})|}{|M^{-1}(q_k)p_k|^2}, \quad (16)$$

where the quantities on the right hand side are the mass matrix, the configuration and the momentum of the  $i$ -th body. Since the frequency  $\omega_i(t_k)$  must be also known at an initial time instant  $t_0$  (in which the initial positions are  $\bar{q}_0$  and initial momenta are  $\bar{p}_0$ ), using the continuous Euler–Lagrange equation at  $t_0$  we obtain

$$\omega_i(t_0) = \frac{|M^{-1}(\bar{q}_0)\bar{p}_0 \times (-M^{-1}(\bar{q}_0)\nabla V(\bar{q}_0))|}{|M^{-1}(\bar{q}_0)\bar{p}_0|^2}. \quad (17)$$

Equations (16) and (17) provide an “estimated frequency” for each mass in the general motion of the  $N$ -body problem. This allows us to derive high order variational integrator methods using trigonometric interpolation where the frequency is estimated at every time step of the integration procedure. These methods show better energy behavior, i.e. smaller total energy oscillation than other methods which employ constant frequency, see [14, 16].

Before closing this section, it should be mentioned that the linear stability of our method is comprehensively analyzed in our previous works [14, 16, 19].

### 4 Triangle and Square Discretization

In order to express the discrete Lagrangian and discrete Hamilton function, we will use the definition of the tangent bundle  $TQ$  and cotangent bundle  $T^*Q$  as in [2] to fields over the higher-dimensional manifold  $X$ . In this way, we also view fields over  $X$  as sections of some fiber bundle  $B \rightarrow X$ , with fiber  $Y$ , and then consider the first jet bundle  $J^1B$  and its dual  $(J^1B^*)$  as the appropriate analogs of the tangent and cotangent bundles.

It is then possible to use the generalization of the Veselov discretization [4, 5] to multisymplectic field theory, by discretizing the spacetime  $X$ . For simplicity reasons we will restrict ourselves to the discrete analogue of  $\dim X = 2$ . Thus, we take  $X = \mathbb{Z} \times \mathbb{Z} = (i, j)$  and the fiber bundle  $Y$  to be  $X \times F$  for some smooth manifold  $F$  [2, 12].

#### 4.1 Triangle Discretization

Assume that we have a uniform quadrangular mesh in the base space, with mesh lengths  $\Delta x$  and  $\Delta t$ . The nodes in this mesh are denoted by  $(i, j) \in \mathbb{Z} \times \mathbb{Z}$ , corresponding to the points  $(x_i, t_j) := (i\Delta x, j\Delta t) \in \mathbb{R}^2$ . We denote the value of the field  $u$  at the node  $(i, j)$  by  $u_i^j$ . We label the triangle at  $(i, j)$  with three ordered triple  $((i, j), (i + 1, j), (i, j + 1))$  as  $\Delta_{ij}$ , and we define  $X_\Delta$  to be the set of all such triangles, see Figure 1.

Then, the discrete jet bundle is defined as follows [2]

$$J^1_\Delta Y := \{(u_i^j, u_{i+1}^j, u_i^{j+1}) \in \mathbb{R}^3 : ((i, j), (i + 1, j), (i, j + 1)) \in X_\Delta\}, \tag{18}$$

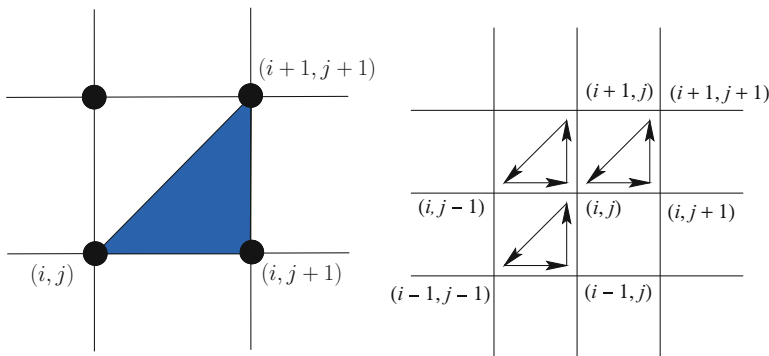


Fig. 1 The triangles which touch  $(i, j)$

which is equal to  $X_{\Delta} \times \mathbb{R}^3$ . The field  $u$  can be now defined by averaging the fields over all vertices of the triangle (see Figure 1a)

$$u \rightarrow \frac{u_i^j + u_i^{j+1} + u_{i+1}^{j+1}}{3}, \tag{19}$$

while the derivatives can be expressed using nonstandard finite differences [9–11]

$$\frac{du}{dt} \rightarrow \frac{u_i^{j+1} - u_i^j}{\phi(\Delta t)}, \quad \frac{du}{dx} \rightarrow \frac{u_{i+1}^{j+1} - u_i^{j+1}}{\psi(\Delta x)}, \tag{20}$$

with [9, 10]

$$\phi(\Delta t) = 2 \sin\left(\frac{\Delta t}{2}\right), \quad \psi(\Delta x) = 2 \sin\left(\frac{\Delta x}{2}\right). \tag{21}$$

Using the latter expressions, we can obtain the discrete Lagrangian at any triangle, which depends on the edges of the triangle, i.e.,  $L_d(u_i^j, u_i^{j+1}, u_{i+1}^{j+1})$ , while the discrete Euler–Lagrange field equations are

$$D_1 L_d(u_i^j, u_i^{j+1}, u_{i+1}^{j+1}) + D_2 L_d(u_i^{j-1}, u_i^j, u_{i+1}^j) + D_3 L_d(u_{i-1}^{j-1}, u_{i-1}^j, u_i^j) = 0, \tag{22}$$

see Figure 1 (right).

### 4.2 Square Discretization

For the cases where square discretization is used, and if we also denote a square at  $(i, j)$  with four ordered quaternion  $((i, j), (i + 1, j), (i + 1, j + 1), (i, j + 1))$  by  $\square_j^i$ , we can consider  $X_{\square}$  to be the set of all such squares, see Figure 1. Then, the discrete jet bundle is defined as (for more details see [2] and references therein)

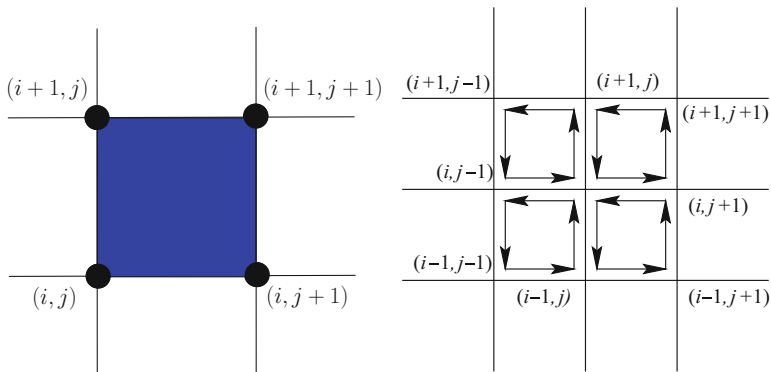
$$J_{\square}^1 Y := \{(u_i^j, u_{i+1}^j, u_{i+1}^{j+1}, u_i^{j+1}) \in \mathbb{R}^4 : ((i, j), (i + 1, j), (i + 1, j + 1), (i, j + 1)) \in X_{\square}\}, \tag{23}$$

which is equal to  $X_{\square} \times \mathbb{R}^4$ .

By averaging the fields over all vertices of the square, the field  $u$  can be now obtained as (see Figure 2 (left))

$$u \rightarrow \frac{u_i^j + u_{i+1}^j + u_i^{j+1} + u_{i+1}^{j+1}}{4}. \tag{24}$$

As above, the expressions for the derivatives can be taken from [9–11] for the discrete Lagrangian, which now depends on the edges of the square, i.e.,



**Fig. 2** The triangles which touch  $(i, j)$

$L_d(u_i^j, u_{i+1}^j, u_{i+1}^{j+1}, u_i^{j+1})$ . As a result, the discrete Euler–Lagrange field equations are

$$\begin{aligned}
 &D_1 L_d(u_i^j, u_{i+1}^{j+1}, u_{i+1}^j, u_i^{j-1}) + D_2 L_d(u_i^{j-1}, u_i^j, u_{i+1}^j, u_{i+1}^{j-1}) + \\
 &D_3 L_d(u_{i-1}^{j-1}, u_{i-1}^j, u_i^j, u_i^{j-1}) + D_4 L_d(u_{i-1}^j, u_{i-1}^{j+1}, u_i^{j+1}, u_i^j) = 0, \tag{25}
 \end{aligned}$$

see Figure 2 (right).

## 5 Numerical Examples Using Triangle Discretization

To illustrate the proposed method, we consider the basic PDEs of three physical problems, i.e., the linear wave equation, the Laplace equation, and the Poisson equation (see [2] and [23, 24]). In the following subsections, for representation requirements, quadrilaterals have been by interpolating the solution on triangles.

### 5.1 Linear Wave Equation

The linear wave equation contains second order partial derivatives of the wavefunction  $u(x, t)$  with respect to time and space, respectively, as (see e.g. [23, 24])

$$\frac{\partial^2 u}{\partial t^2} + c \frac{\partial^2 u}{\partial x^2} = 0. \tag{26}$$

This equation may be considered for the description of the wave function, i.e., the amplitude of oscillation, that is created from a one-dimensional medium (e.g. a

string extended in the  $x$ -direction). For the special case that the velocity of the wave, representing by the parameter  $c$ , is chosen as  $c = -1$ , the corresponding Lagrangian is [12]

$$L(u, u_t, u_x) = \frac{1}{2}u_t^2 - \frac{1}{2}u_x^2, \quad (27)$$

where the derivatives are  $\partial u/\partial t = u_t$  and  $\partial u/\partial x = u_x$ .

If we use triangle discretization, described in Section 4.1, we end up with discrete Lagrangian

$$L_d(u_i^j, u_i^{j+1}, u_{i+1}^{j+1}) = \frac{1}{2}\Delta t \Delta x \left[ \frac{1}{2} \left( \frac{u_i^{j+1} - u_i^j}{\phi(\Delta t)} \right)^2 - \frac{1}{2} \left( \frac{u_{i+1}^{j+1} - u_i^{j+1}}{\psi(\Delta x)} \right)^2 \right], \quad (28)$$

where  $\Delta t$  and  $\Delta x$  are the mesh lengths for time and space, respectively. Applying the above discrete Lagrangian to the discrete Euler–Lagrange field equations (22), we get

$$\frac{u_i^{j+1} - 2u_i^j + u_i^{j-1}}{(\phi(\Delta t))^2} - \frac{u_{i+1}^j - 2u_i^j + u_{i-1}^j}{(\psi(\Delta x))^2} = 0. \quad (29)$$

The latter expression represents the variational integrator for the linear wave Equation (26), resulting through the use of the proposed nonstandard finite difference schemes.

In Figure 3 the solution  $u(x, t)$  of (29) is shown in a 3-D diagram. We have chosen as initial conditions  $0 < x < 1$ ,  $u(x, 0) = 0.5[1 - \cos(2\pi x)]$ ,  $u_t(x, 0) = 0.1$  and as boundary conditions  $u(0, t) = u(1, t)$ ,  $u_x(0, t) = u_x(1, t)$ , the latter being periodic. The grid discretization has been taken to be  $\Delta t = 0.01$  and  $\Delta x = 0.01$ . As seen, the time evolution of the solution  $u(x = \text{const.}, t)$  is a continuous function, while the periodicity is preserved.

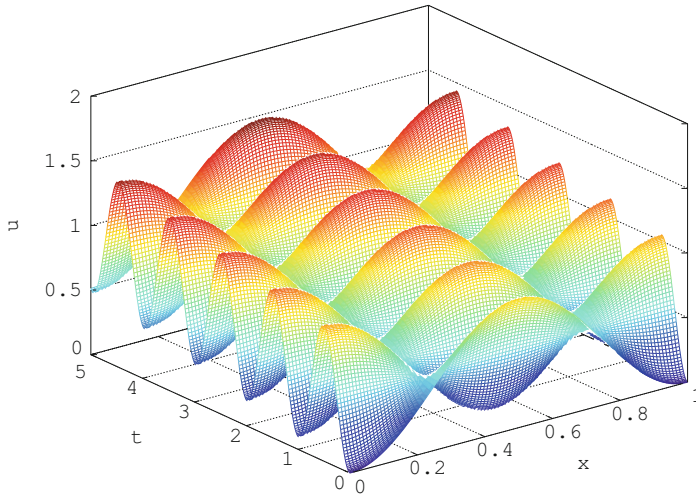
## 5.2 Laplace Equation

As another physical example, we have chosen the Laplace equation over a 2-D scalar field  $u(x, y)$ . It is written as

$$u_{xx} + u_{yy} = 0. \quad (30)$$

The function  $u(x, y)$  may describe a potential in a 2-D medium or a potential inside a 3-D medium, which does not depend on the third coordinate  $z$ . Thus, the 2-dimensional second order PDE (30) governs a variety of equilibrium physical





**Fig. 3** The waveforms of linear wave equation (26)

phenomena such as temperature distribution in solids, electric field in electrostatics, inviscid and irrotational two-dimensional flow (potential flow), groundwater flow, etc.

The corresponding continuous Lagrangian of (30) takes the form

$$L(u, u_x, u_y) = \frac{1}{2}u_x^2 + \frac{1}{2}u_y^2. \tag{31}$$

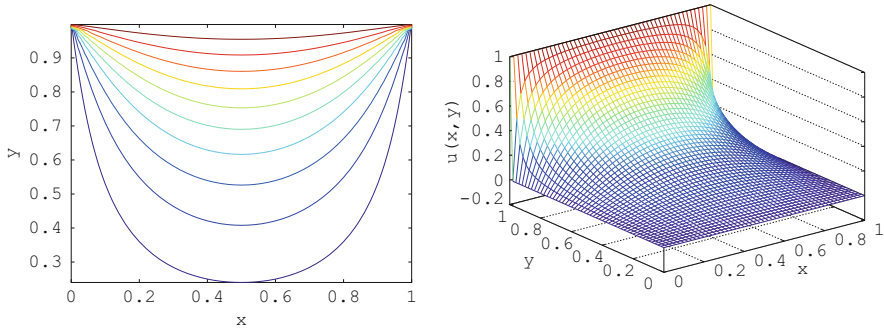
By applying the triangle discretization of Section 4.1, the discrete Lagrangian can be written as

$$L_d(u_i^j, u_i^{j+1}, u_{i+1}^{j+1}) = \frac{1}{2}\Delta x\Delta y \left[ \frac{1}{2} \left( \frac{u_i^{j+1} - u_i^j}{\phi(\Delta x)} \right)^2 + \frac{1}{2} \left( \frac{u_{i+1}^{j+1} - u_i^{j+1}}{\psi(\Delta y)} \right)^2 \right]. \tag{32}$$

From the latter Lagrangian, working in a similar manner to that followed in Section 4.2, results the integrator from the proposed nonstandard finite difference schemes

$$\frac{u_i^{j+1} - 2u_i^j + u_i^{j-1}}{(\phi(\Delta x))^2} + \frac{u_{i+1}^j - 2u_i^j + u_{i-1}^j}{(\psi(\Delta y))^2} = 0. \tag{33}$$

The solution of the above equation, when considering the boundary conditions  $u(x, 0) = 0, u(x, 1) = 1$  and  $u(0, y) = u(1, y) = 0$ , is plotted in Figure 4. The grid discretization has been chosen to be  $\Delta x = 0.02$  and  $\Delta y = 0.02$ .



**Fig. 4** Contour plot (left) and three-dimensional surface plot (right) of the solution of Laplace equation with boundary conditions  $u(x, 0) = 0$ ,  $u(x, 1) = 1$ ,  $u(0, y) = u(1, y) = 0$ , and discretization:  $\Delta x = 0.02$ ,  $\Delta y = 0.02$

### 5.3 Poisson Equation

As a final application to illustrate the advantages of the proposed variational integrator relying on nonstandard finite difference schemes, we examine the Poisson equation, which is an elliptic PDE of the form

$$-u_{xx} - u_{yy} = f(x, y). \quad (34)$$

Obviously, this equation in physical applications presents an additional complexity compared to the Laplace equation (30). Now the right hand side is a non-zero function  $f(x, y)$ , which may be considered as a source (or a load) function defined on some two-dimensional domain denoted by  $\Omega \subset \mathbb{R}^2$  (it could also be a general non-linear function  $f(u, x, y)$ ). A solution  $u$  satisfying (34) will also satisfy specific conditions on the boundaries of the domain  $\Omega$ . For example, for the element  $\partial\Omega$  the general condition holds

$$\alpha u + \beta \frac{\partial u}{\partial n} = g \quad \text{on } \partial\Omega, \quad (35)$$

where  $\partial u/\partial n$  denotes the directional derivative in the direction normal to the boundary  $\partial\Omega$  and  $\alpha$  and  $\beta$  are constants [23, 24].

As it is well known, the system of (34) and (35) is referred to as a boundary value problem for the Poisson equation. If the constant  $\beta$  in Equation (35) is zero, then the boundary condition is of Dirichlet type, and the boundary value problem is referred to as the Dirichlet problem for the Poisson equation. Alternatively, if the constant  $\alpha$  is zero, then we correspondingly have a Neumann boundary condition, and the problem is referred to as a Neumann problem. A third possibility exists when the Dirichlet conditions hold on a part of the boundary  $\partial\Omega_D$ , and Neumann conditions

hold on the remainder  $\partial\Omega \setminus \partial\Omega_D$  (or indeed mixed conditions where  $\alpha$  and  $\beta$  are both nonzero), see [23, 24] and references therein.

Equation (34) can also be obtained by starting from the Lagrangian

$$L(u, u_x, u_y) = \frac{1}{2}u_x^2 + \frac{1}{2}u_y^2 - fu. \tag{36}$$

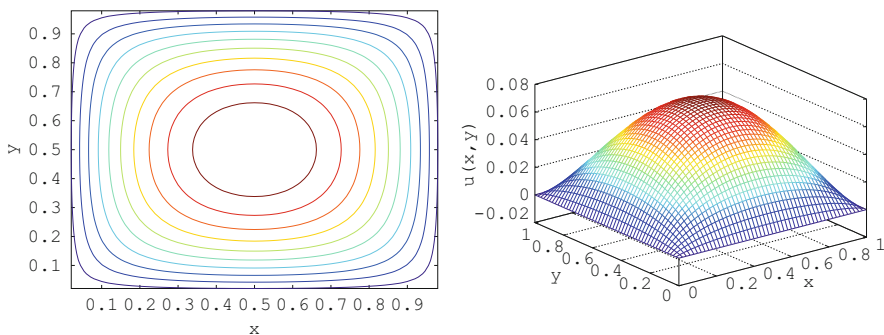
The triangle discretization of Section 4.1 in the Poisson problem defines the discrete Lagrangian

$$L_d(u_i^j, u_i^{j+1}, u_{i+1}^{j+1}) = \frac{1}{2}\Delta x\Delta y \left[ \frac{1}{2} \left( \frac{u_i^{j+1} - u_i^j}{\phi(\Delta x)} \right)^2 + \frac{1}{2} \left( \frac{u_{i+1}^{j+1} - u_i^{j+1}}{\psi(\Delta y)} \right)^2 \right] - \frac{f_i^j u_i^j + f_i^{j+1} u_i^{j+1} + f_{i+1}^{j+1} u_{i+1}^{j+1}}{3}. \tag{37}$$

By inserting the latter discrete Lagrangian into the discrete Euler–Lagrange field equations (22) and elaborating as done in [2], the resulting integrator from the proposed nonstandard finite difference schemes is

$$-\frac{u_i^{j+1} - 2u_i^j + u_i^{j-1}}{(\phi(\Delta x))^2} - \frac{u_{i+1}^j - 2u_i^j + u_{i-1}^j}{(\psi(\Delta y))^2} = f_i^j + \partial f_i^j / \partial u_i^j. \tag{38}$$

As a special case we chose the source term  $f(x, y) \equiv 1$ , so  $\partial f_i^j / \partial u_i^j = 0$  in (38), and the boundary conditions  $u(0, y) = u(1, y) = 0$  and  $u(x, 0) = u(x, 1) = 0$ . Figure 5 shows the numerical results obtained with the discretization  $\Delta x = 0.02$  and  $\Delta y = 0.02$ .



**Fig. 5** Contour plot (left) and three-dimensional surface plot (right) of the solution of Poisson equation, using the variational integrator with nonstandard finite difference schemes. The source term was chosen  $f(x, y) \equiv 1$ , while the boundary conditions  $u(0, y) = u(1, y) = 0$ ,  $u(x, 0) = u(x, 1) = 0$  for discretization  $\Delta x = 0.02$ ,  $\Delta y = 0.02$

## 6 Numerical Examples Using Square Discretization

To illustrate the behavior of the proposed method, we will consider the Klein–Gordon equation, which plays a significant role in many scientific applications such as solid state physics, nonlinear optics and quantum field theory, see for example [25].

### 6.1 Klein–Gordon

For the general case, the initial-value problem of the one-dimensional nonlinear Klein–Gordon equation is given by

$$u_{tt} + \alpha u_{xx} + g(u) = f(x, t), \quad (39)$$

where  $u = u(x, t)$  represents the wave displacement at position  $x$  and time  $t$ ,  $\alpha$  is a known constant and  $g(u)$  is the nonlinear force which, in the physical applications has also other forms [25].

Here we will consider the special case that  $\alpha = -1$ ,  $g(u) = u^3 - u$  and  $f(x, t) = 0$  resulting

$$u_{tt} = u_{xx} - u^3 + u. \quad (40)$$

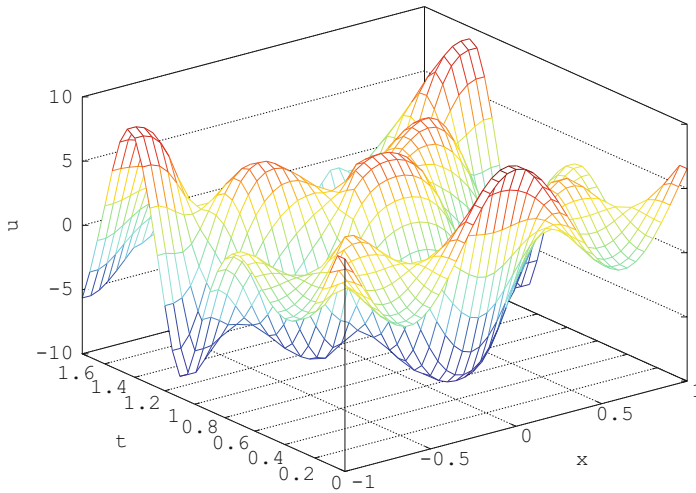
The above equation can be described using the Lagrangian

$$L(u, u_t, u_x) = \frac{1}{2}u_t^2 - \frac{1}{2}u_x^2 - \frac{1}{4}u^4 - \frac{1}{2}u^2.$$

Following Section 4.2 we can obtain the discrete Lagrangian that now uses square discretization as

$$\begin{aligned} L_d(u_i^j, u_{i+1}^j, u_{i+1}^{j+1}, u_i^{j+1}) &= \frac{\Delta t \Delta x}{2} \left( \frac{u_i^{j+1} - u_i^j}{2\phi(\Delta t)} + \frac{u_{i+1}^{j+1} - u_{i+1}^j}{2\phi(\Delta t)} \right)^2 - \\ &\quad \frac{\Delta t \Delta x}{2} \left( \frac{u_{i+1}^{j+1} - u_i^{j+1}}{2\psi(\Delta x)} + \frac{u_{i+1}^j - u_i^j}{2\psi(\Delta x)} \right)^2 - \\ &\quad - \frac{\Delta t \Delta x}{4} u_i^j u_{i+1}^j u_{i+1}^{j+1} u_i^{j+1} + \\ &\quad \frac{\Delta t \Delta x}{2} \left( \frac{u_i^j u_{i+1}^j + u_i^j u_{i+1}^{j+1} + u_i^j u_i^{j+1} + u_{i+1}^j u_{i+1}^{j+1} + u_{i+1}^{j+1} u_i^{j+1} + u_{i+1}^{j+1} u_i^j}{6} \right), \end{aligned}$$

The waveforms of Klein-Gordon equation



**Fig. 6** Numerical solution of the Klein–Gordon equation (40) using square discretization of Section 4.2

which we will consider for the discrete Euler–Lagrange equations (25) in order to derive the resulting integrator from the proposed nonstandard finite difference schemes.

Figure 6 shows the numerical results obtained with the discretization  $\Delta t = 0.05$  and  $\Delta x = 0.05$ . To that we have used initial conditions  $u(x, 0) = A(1 + \cos(\frac{2\pi x}{L}))$  where  $A = 5$  and  $u_t(x, 0) = 0$ , while the boundary conditions were  $u(-1, t) = u(1, t)$  and  $u_x(-1, t) = u_x(1, t)$ .

## 7 Analysis of the Proposed Schemes

A dispersion analysis and mesh convergence experiments are performed in this section in order to show the numerical properties of the proposed method.

### 7.1 Dispersion Analysis

We will now turn our study to the dispersion-dissipation properties of the derived numerical schemes and compare them with the ones of [2]. To that end, similar to [26], we consider the discrete analog of the Fourier mode

$$u_i^j = \hat{u} e^{\mathbf{i}(ik\Delta x + j\omega\Delta t)}, \tag{41}$$

where  $\mathbf{i}^2 = -1$ . Using  $\bar{k} = k\Delta x$  and  $\bar{\omega} = \omega\Delta t$ , the latter equation results in

$$u_i^j = \hat{u} e^{\mathbf{i}(i\bar{k} + j\bar{\omega})}. \tag{42}$$

Following the above, the multi-symplectic scheme of [2], also known as leap frog algorithm, for the case of the linear wave (26) gives

$$\frac{u_i^{j+1} - 2u_i^j + u_i^{j-1}}{(\Delta t)^2} - \frac{u_{i+1}^j - 2u_i^j + u_{i-1}^j}{(\Delta x)^2} = 0. \tag{43}$$

When substituting (42) in the latter equation, we get the discrete dispersion relationship

$$\frac{e^{\mathbf{i}\bar{k}}}{(\Delta t)^2} \left[ e^{2\mathbf{i}\bar{\omega}} - 2e^{\mathbf{i}\bar{\omega}} + 1 \right] - \frac{e^{\mathbf{i}\bar{\omega}}}{(\Delta x)^2} \left[ e^{2\mathbf{i}\bar{k}} - 2e^{\mathbf{i}\bar{k}} + 1 \right] = 0. \tag{44}$$

As a second example we consider the second order implicit Runge–Kutta scheme described in [27, 28] and [29]. This scheme, also known as implicit Crank–Nicolson, is a symplectic time discretization of order two, which for the case of (26) gives

$$4 \left( u_i^{j+2} - 2u_i^{j+1} + u_i^j \right) - \lambda^2 \left( u_{i-1}^{j+2} - 2u_i^{j+2} + u_{i+1}^{j+2} \right) - 2\lambda^2 \left( u_{i-1}^{j+1} - 2u_i^{j+1} + u_{i+1}^{j+1} \right) - \lambda^2 \left( u_{i-1}^j - 2u_i^j + u_{i+1}^j \right) = 0, \tag{45}$$

where

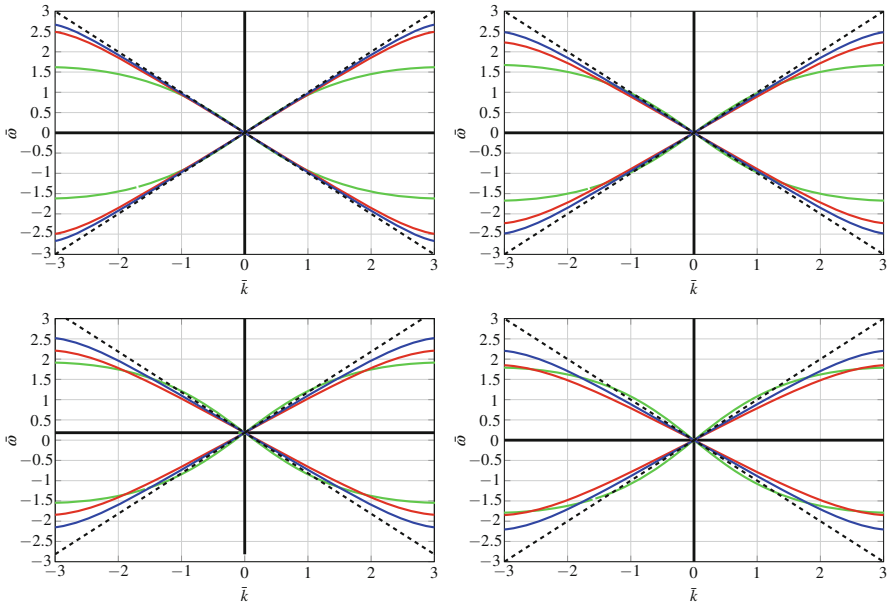
$$\lambda^2 = \left( \frac{\Delta t}{\Delta x} \right)^2. \tag{46}$$

Substituting to the above integrator the form (42) we obtain the discrete dispersion relationship

$$\frac{4e^{\mathbf{i}\bar{k}}}{(\Delta t)^2} \left[ e^{2\mathbf{i}\bar{\omega}} - 2e^{\mathbf{i}\bar{\omega}} + 1 \right] - \frac{\left[ e^{2\mathbf{i}\bar{k}} - 2e^{\mathbf{i}\bar{k}} + 1 \right]}{(\Delta x)^2} \left[ e^{2\mathbf{i}\bar{\omega}} - 2e^{\mathbf{i}\bar{\omega}} + 1 \right] = 0. \tag{47}$$

For the case of the linear wave Equation (26) the integrator with the proposed technique, i.e., (29) for  $u_i^j$  of (42) gives

$$(\cos \bar{\omega} - 1) (1 - \cos \Delta x) - (\cos \bar{k} - 1) (1 - \cos \Delta t) = 0. \tag{48}$$



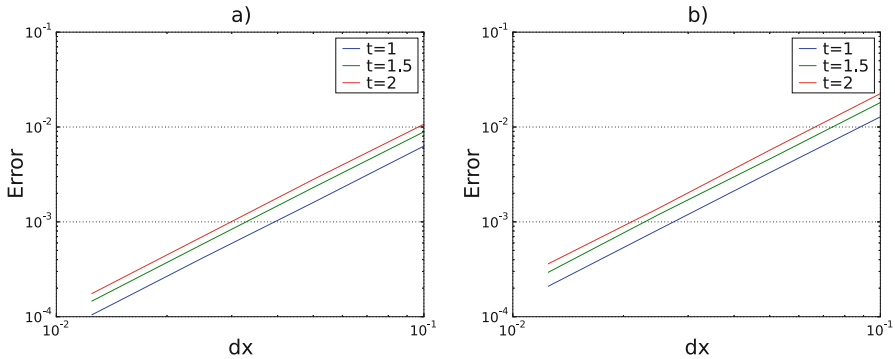
**Fig. 7** Dispersion curves for the linear wave equation with the proposed method (red), the leap frog scheme of [2] (blue), the implicit Runge–Kutta (green) and the analytic one (dashed black) for  $\lambda = \{0.95, 0.9, 0.85, 0.8\}$

For now we will restrict ourselves only to  $\lambda \leq 1$ , but due to symmetry, all other cases can be easily obtained. Figure 7 shows the discrete dispersion relationships for  $\lambda = \{0.95, 0.9, 0.85, 0.8\}$ . Specifically, to each sub plot we can see the dispersion curve of the leap frog scheme, i.e., Equation (44), with blue line, the red line corresponds to the proposed method, described by (48), while the green line is the one for the implicit Runge–Kutta scheme, Equation (47). For all the choices of  $\lambda$  tested the behaviour of the method using nonstandard finite difference schemes is close to the excellent behaviour of the leap frog scheme, and much better than the implicit Runge–Kutta scheme.

### 7.2 Convergence Experiments

In order to show the grid independence of the solution, following the finite element convention, the  $l^\infty$ -norm error is calculated between the solutions on two successive grids according to

$$e_h = \max_i \{|u_i^f - u_i^c|, \dots, |u_{nel}^f - u_{nel}^c|\}, \tag{49}$$



**Fig. 8** Error of numerical solution as a function of grid size  $\Delta x$  for different time steps: (a) triangle discretization (b) square discretization

where  $u_i^f$  is the solution on the fine grid,  $u_i^c$  the solution on a coarse grid interpolated on the fine one. Here,  $n_{el}$  are the total number of elements, where the elements of the mesh are either triangles or squares. A sample convergence of the calculations for the Klein–Gordon case is shown in Figure 8 in a logarithmic plot for triangle and square discretizations and for different time steps. It can be easily seen that by decreasing the space discretization the error is also decreased linearly in the log scale.

## 8 Summary and Conclusions

The derivation of advantageous multisymplectic numerical methods, relying on nonstandard finite difference schemes, is investigated. The numerical solution of the linear wave equation, the 2-D Laplace equation, and the 2-D Poisson equation, which are addressed in this study, show a good energy behavior and the preservation of the discrete multisymplectic structure of the proposed numerical schemes. Moreover, we showed with the help of dispersion analysis and mesh convergence experiments the numerical properties of the proposed method.

Future applications may include the field equation of incompressible fluid dynamics, like that of Cotter et al [30] and Pavlov et al [31], which could be of interest in investigating the properties of 3-D media. For partial differential equations arising in the field of fluid dynamics, dissipative terms should be taken into consideration. These dissipative perturbations necessitate application of techniques similar to [32, 33] but in the case of PDEs. Furthermore, a possible application in complex geometries, as they appear in real world problems, would necessitate the extension of this methodology to non-uniform grids.

The variational method presented in this work can be applied in a variety of physical problems, ranging from magnetic field simulations in NMR [34] to



inverse problems that arise in geophysics [35] and others. Future work may include comparison with other numerical of PDEs, such as the finite element method or the finite volume method.

**Acknowledgments** Dr. Odysseas Kosmas wishes to acknowledge the support of EPSRC via grand EP/N026136/1 “Geometric Mechanics of Solids”.

## Appendix

By denoting  $\text{sinc}(\xi) = \sin(\xi)/\xi$ , special cases of the exponential integrators described using (10) can be obtained, i.e.

- Gautschi type exponential integrators [36] for

$$\psi(\Omega h) = \text{sinc}^2\left(\frac{\Omega h}{2}\right), \quad \phi(\Omega h) = 1$$

- Deuffhard type exponential integrators [37] for

$$\psi(\Omega h) = \text{sinc}(\Omega h), \quad \phi(\Omega h) = 1$$

- García-Archilla et all. type exponential integrators [38] for

$$\psi(\Omega h) = \text{sinc}^2(\Omega h), \quad \phi(\Omega h) = \text{sinc}(\Omega h)$$

Finally, in [1] a way to write the Störmer-Verlet algorithm as an exponential integrators is presenting.

## References

1. E. Hairer C. Lubich, G. Wanner, Geometric numerical integration illustrated by the Störmer-Verlet method. *Acta Numer.* **12**, 399 (2003)
2. J.E. Marsden, G.W. Patrick, S. Shkoller, Multisymplectic geometry, variational integrators, and nonlinear PDEs. *Commun. Math. Phys.* **199**, 351 (1998)
3. J.E. Marsden, M. West, Discrete mechanics and variational integrators. *Acta Numer.* **10**, 357 (2001)
4. A.P. Veselov, Integrable discrete-time systems and difference operators. *Funkts. Anal. Prilozhen.* **22**, 1 (1988)
5. A.P. Veselov, Integrable Lagrangian correspondences and the factorization of matrix polynomials. *Funkts. Anal. Prilozhen.* **25**, 38 (1991)
6. T.J. Bridges, Multi-symplectic structures and wave propagation. *Math. Proc. Camb. Philos. Soc.* **121**, 1 (1997)
7. T.J. Bridges, S. Reich Multi-symplectic integrators: numerical schemes for Hamiltonian PDEs that conserve symplecticity. *Phys. Lett. A.* **284**, 4–5 (2001)

8. T.J. Bridges, S. Reich, Numerical methods for Hamiltonian PDEs. *J. Phys.* **39**, 19 (2006)
9. R.E. Mickens, *Applications of Nonstandard Finite Difference Schemes* (World Scientific Publishing, Singapore, 2000)
10. R.E. Mickens, Nonstandard finite difference schemes for differential equations. *J. Differ. Equ. Appl.* **8**, 823 (2002)
11. R.E. Mickens, Dynamic consistency: a fundamental principle for constructing nonstandard finite difference schemes for differential equations. *J. Differ. Equ. Appl.* **11**, 645 (2005)
12. O.T. Kosmas, D. Papadopoulos, Multisymplectic structure of numerical methods derived using nonstandard finite difference schemes. *J. Phys. Conf. Ser.* **490** (2014)
13. O.T. Kosmas, Charged particle in an electromagnetic field using variational integrators. *ICNAAM Numer. Anal. Appl. Math.* **1389**, 1927 (2011)
14. O.T. Kosmas, S. Leyendecker, Analysis of higher order phase fitted variational integrators. *Adv. Comput. Math.* **42**, 605 (2016)
15. O.T. Kosmas D.S. Vlachos, Local path fitting: a new approach to variational integrators. *J. Comput. Appl. Math.* **236**, 2632 (2012)
16. O.T. Kosmas, S. Leyendecker, Variational integrators for orbital problems using frequency estimation. *Adv. Comput. Math.* **45**, 1–21 (2019)
17. O.T. Kosmas, D.S. Vlachos, Phase-fitted discrete Lagrangian integrators. *Comput. Phys. Commun.* **181**, 562–568 (2010)
18. O.T. Kosmas, S. Leyendecker, Phase lag analysis of variational integrators using interpolation techniques. *PAMM Proc. Appl. Math. Mech.* **12**, 677–678 (2012)
19. O.T. Kosmas, S. Leyendecker, Stability analysis of high order phase fitted variational integrators, in *Proceedings of WCCM XI—ECCM V—ECFD VI*, vol. 1389 (2014), pp. 865–866
20. O.T. Kosmas, S. Leyendecker, Family of high order exponential variational integrators for split potential systems. *J. Phys. Conf. Ser.* **574**, 012002 (2015)
21. O.T. Kosmas, D.S. Vlachos, A space-time geodesic approach for phase fitted variational integrators. *J. Phys. Conf. Ser.* **738**, 012133 (2016)
22. L. Brusca, L. Nigro, A one-step method for direct integration of structural dynamic equations. *Int. J. Numer. Methods Eng.* **15**, 685–699 (1980)
23. L.C. Evans, *Partial Differential Equations* (American Mathematical Society, Providence, 1998)
24. V.I. Arnold, *Lectures on Partial Differential Equations* (Springer, Berlin, 2000)
25. H. Han, Z. Zhang, Split local absorbing conditions for one-dimensional nonlinear Klein–Gordon equation on unbounded domain. *J. Comput. Phys.* **227**, 8992 (2008)
26. J.W. Thomas, *Numerical Partial Differential Equations*, Finite Difference Methods, vol. 1 (Springer, New York, 1995)
27. J.M. Sanz-Serna, M.P. Calvo, *Numerical Hamiltonian Problems* (Chapman & Hall, London, 1994)
28. J.M. Sanz-Serna, *Solving Numerically Hamiltonian Systems*. Proceedings of the International Congress of Mathematicians (Birkhäuser, Basel, 1995)
29. S. Reich, Multi-symplectic Runge–Kutta collocation methods for Hamiltonian wave equations. *J. Comput. Phys.* **157**, 473 (2000)
30. C.J. Cotter, D.D. Holm, P.E. Hydon, Multisymplectic formulation of fluid dynamics using the inverse map. *Proc. R. Soc. A* **463**, 2671 (2007)
31. D. Pavlov, P. Mullen, Y. Tong, E. Kanso, J.E. Marsden, M. Desbrun, Structure-preserving discretization of incompressible fluids. *Phys. D: Nonlinear Phenom.* **240**, 443 (2011)
32. E. Hairer, C. Lubich, Invariant tori of dissipatively perturbed Hamiltonian systems under symplectic discretization. *Appl. Numer. Math.* **29**, 57–71 (1999)
33. D. Stoffer, On the qualitative behaviour of symplectic integrators. III: perturbed integrable systems. *J. Math. Anal. Appl.* **217**, 521–545 (1998)
34. D. Papadopoulos, M.A. Voda, S. Stapf, F. Casanova, M. Behr, B. Blümich, Magnetic field simulations in support of interdiffusion quantification with NMR. *Chem. Eng. Sci.* **63**, 4694 (2008)
35. D. Papadopoulos, M. Herty, V. Rath, M. Behr, Identification of uncertainties in the shape of geophysical objects with level sets and the adjoint method. *Comput. Geosci.* **15**, 737 (2011)

36. W. Gautschi, Numerical integration of ordinary differential equations based on trigonometric polynomials. *Numer. Math.* **3**, 1 (1961)
37. P. Deuffhard, A study of extrapolation methods based on multistep schemes without parasitic solutions. *Z. Angew. Math. Phys.* **30**, 2 (1979)
38. B. García-Archilla, M.J. Sanz-Serna, R.D. Skeel, Long-time-step methods for oscillatory differential equations. *SIAM J. Sci. Comput.* **20**, 3 (1999)

# Non-radial Solutions of a Supercritical Equation in Expanding Domains: The Limit Case



Nikos Labropoulos

**Abstract** In this article, we introduce a new method to prove the existence of an infinite sequence of distinct non-radial but symmetric nodal (i.e. sign changing) solutions for supercritical nonlinear elliptic problems defined in the whole Euclidean space. By ‘symmetric’ we mean that both the domain and the solution remain invariant under the action of a compact subgroup  $G$  of the isometry group  $O(n)$ , without finite subgroup. The key ingredient of the method is a process through which an open symmetric domain of the  $n$ -dimensional space can be extended in an appropriate manner to ‘fill’ eventually the entire space ‘almost everywhere’, remaining symmetric, and giving a sequence of domains where in each of them subsequently we solve an appropriate auxiliary problem. Passing to the limit we obtain the solution of the problem as a limit of the sequence formed by the solutions of the corresponding to the domains sequence of equations.

The base model problem of interest is stated below:

$$(P) \quad \begin{cases} \Delta_p u = |u|^a u, & u \in C^2(\mathbb{R}^n), \quad n \geq 3 \\ 1 < p < n - k, \quad 0 \leq a \leq p^*(k) = \frac{(n-k)p}{n-k-p}, \end{cases}$$

where  $p_G^*$  is the critical exponent of the embedding

$$H_{0,G}^{1,p}(\Omega) \hookrightarrow L^{p_G^*}(\Omega)$$

(being the critical of the supercritical one) and  $k$  is the minimum orbit dimension in  $G$ . However, we will focus on the critical of the supercritical case  $a = p^*(k)$ , since on the one hand it is the most important and on the other hand it covers all the other cases.

---

N. Labropoulos (✉)

Department of Mathematics, National Technical University of Athens, Athens, Greece  
e-mail: [nal@math.ntua.gr](mailto:nal@math.ntua.gr)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,

Springer Optimization and Its Applications 173,

[https://doi.org/10.1007/978-3-030-72563-1\\_12](https://doi.org/10.1007/978-3-030-72563-1_12)

By  $H_{0,G}^{1,p}(\Omega)$  is denoted the closure of the subspace  $C_{0,G}^\infty(\Omega)$  consisting of all  $G$ -invariant functions in  $C_0^\infty(\Omega)$ .

### 1 Introduction

In this article, the main objective is to prove the existence of non-radial nodal (sign-changing) solutions of the above problem (P), in the case where the exponent  $a$  is the critical of supercritical exponent, since the rest of the cases have been studied. Thus, the problem (P) is set out in detail as follows:

$$(P) \quad \begin{cases} \Delta_p u = |u|^{p^*(k)-2}u & \text{in } \mathbb{R}^n, \quad n \geq 3 \\ 1 < p < n - k, \quad p^*(k) = \frac{(n-k)p}{n-k-p}, \end{cases}$$

Here,  $G$  is a group of symmetries that acts on the domains and the functions defined on them together,  $k$  is the minimum dimension orbit of all orbits of  $G$ ,

$$\Delta_p u = -\operatorname{div} \left( |\nabla u|^{p-2} \nabla u \right), \quad 1 < p \neq 2$$

is the  $p$ -Laplacian operator (note that if  $p = 2$ , is the Laplace–Beltrami operator) and  $p^*(k)$  is the critical exponent of the Sobolev embedding

$$H_G^{1,p}(\Omega) \hookrightarrow L^p(\Omega).$$

By  $H_G^{1,p}(\Omega)$  is denoted the subspace of all  $G$ -invariant functions in  $H^{1,p}(\Omega)$ .

In problem (P) the solutions obtained are such that

$$\int_{\mathbb{R}^n} |\nabla u|^p dx \rightarrow \infty.$$

We study both the cases,  $p = 2$  and  $p \in (1, 2) \cup (2, n - k)$ , however, to avoid any confusion we note that throughout the article we denote by

$$\Delta_p u = -\operatorname{div} \left( |\nabla u|^{p-2} \nabla u \right), \quad 1 < p < n - k$$

the  $p$ -Laplacian as well as the Laplace–Beltrami operator but when we refer to other articles the Laplace–Beltrami operator is denoted as in the referred articles, i.e. without the minus conversion.

For the problem (P), we prove the existence and find both the type and the number of the solutions to the problem (P). For this aim we use *the method of expanding domains* which was successfully introduced for the first time in [42]. In

that article this method was used firstly to prove the existence of a solution and secondly to determine the type and the number of the solutions to critical nonlinear elliptic problem:

$$(P_0) \quad -\Delta u = |u|^{\frac{4}{n-2}} u, \quad u \in C^2(\mathbb{R}^n), \quad n \geq 3.$$

Concerning the method itself it seems to have particular value because it can be used and in other types of partial differential equations.

Both cases, i.e.  $p = 2$  and  $p \neq 2$ , are extremely interesting and that is why for several decades now many researchers have been paying attention to them.

Problem  $(P_0)$  consists a special case of  $(P)$  for  $p = 2$  and it owns its origin in many astrophysical and physical contexts and more precisely in the Lane-Emden-Fowler problem,

$$(P'_0) \quad \begin{cases} -\Delta u = u^q \\ u > 0 \text{ in } \Omega, \quad u = 0 \text{ on } \partial\Omega, \end{cases}$$

where  $\Omega$  is a domain with smooth boundary in  $\mathbb{R}^N$  and  $p > 1$ . But its greatest interest lies in its relation to the Yamabe problem (see in [5, 57, 64, 68]) and for a complete and detailed study we refer to [6], nevertheless it has an autonomous presence holding an important place among the most famous nonlinear partial differential equations). We refer, also, to the classical papers [20, 30, 47], which are some of the large number of very good papers that are devoted to the study of this problem.

Gidas, Ni, and Nirenberg, in their celebrated paper [30], proved symmetry and some related properties of positive solutions of a larger class of second order elliptic equations. Concerning the equation

$$-\Delta u = |u|^{\frac{4}{n-2}} u, \quad u \in C^2(\mathbb{R}^n), \quad n \geq 3,$$

they proved that any positive solution of this, which has finite energy, namely

$$\int_{\mathbb{R}^n} |\nabla u|^2 dx < +\infty,$$

is necessarily of the form

$$u(x) = \left( \frac{\lambda \sqrt{n(n-2)}}{\lambda^2 + |x - x_0|^2} \right)^{\frac{n-2}{2}}, \quad \lambda > 0, x_0 \in \mathbb{R}^n.$$

These solutions yield the well-known one-instanton solutions in a regular gauge of the Yang–Mills equation. In addition, since this equation is invariant under the conformal transformations of  $\mathbb{R}^n$ , if  $u(x)$  is a solution, then

$$\lambda^{\frac{n-2}{2}} u\left(\frac{x-x_0}{\lambda}\right), \forall \lambda > 0 \text{ and } x_0 \in \mathbb{R}^n$$

is also a solution. Moreover, all solutions obtained in this way have the same energy and we will say that these solutions are equivalent. In particular, all these solutions are equivalent.

Ding in his also celebrated article [20] using Ambrosetti and Rabinowitz analysis (see in [2]) proved that this problem has infinite distinct solutions  $u_k \in C^2(\mathbb{R}^n)$ ,  $k = 1, 2, \dots$ , which changes sign and such that

$$\lim_{k \rightarrow \infty} \int_{\mathbb{R}^n} |\nabla u_k|^2 dx \rightarrow \infty.$$

Ding showed that it is possible to solve the equation in the whole Euclidean space, reduced the problem to an equivalent problem on  $\mathbb{S}^n$ , the Euclidean  $n$ -sphere throughout a conformal deformation. However, this method cannot be used in the case of the  $p$ -Laplacian operator, because this operator is not a conformal invariant operator.

Mazzeo and Smale in their also celebrated article [47] proved that if  $\Omega$  is an open set in  $\mathbb{R}^n$  and  $u$  is a positive  $C^2$  function on  $\Omega$  such that the metric  $g = u^{\frac{4}{n-2}} g_0$  on  $\Omega$  has scalar curvature  $R(g) = n(n-1)$ , then  $u$  must satisfy the equation

$$\Delta u + \frac{n(n-2)}{4} u^{\frac{n+2}{n-2}} = 0, \quad u > 0$$

on  $\Omega$ , where  $g_0$  is the Euclidean metric on  $\mathbb{R}^n$ .

Caffarelli, Gidas, and Spruck in their classical paper [11] studied non-negative smooth solutions of the conformally invariant equation

$$-\Delta u = u^{\frac{n+2}{n-2}}, \quad u \geq 0, \quad n \geq 3,$$

in a punctured ball  $B_1(0) \setminus \{0\} \subseteq \mathbb{R}^n$ , with an isolated singularity at the origin. In this paper, the authors introduced a heuristic idea of asymptotic symmetry technique which can roughly be described as follows: After an inversion, the function  $u$  becomes defined in the complement of  $B_1$ , is strictly positive of  $\partial B_1$ , and in some sense ‘goes to zero’ at infinity. If the function  $u$  can be extended to  $B_1$  as a super solution of our problem, then the reflection process at infinity can start and move all the way to  $\partial B_1$ . This would imply asymptotic radial symmetry at infinity. With this comprehensive report on this issue we would like, on the one hand, to emphasize the important contribution of this great article of Caffarelli, Gidas, and Spruck on the study on the direction of finding the radial solutions of our problem and on the other hand, we wish to make clear that in our procedural paper we do not care about the radial solutions but we do care about the existence of non-radial solutions.

Schoen in [57] built solutions of (P) with prescribed isolated singularities. Schoen, also, in [58], have used the geometrical meaning of problem (P) in order to derive, through ideas of conformal geometry, the existence of weak solutions having a singular set whose Hausdorff dimension is less than or equal to  $\frac{n-2}{2}$ . Let us notice that in this paper the authors explain how to build solutions of (P) with a singular set whose Hausdorff dimension is not necessarily an integer.

Bartsch and Schneider in [8] proved that for  $N > 2m$  the equation

$$(-\Delta)^m = |u|^{\frac{4m}{N-2m}} u$$

on  $\mathbb{R}^N$  has a sequence of nodal, finite energy solutions which is unbounded in  $\mathcal{D}^{m,2}(\mathbb{R}^N)$ , the completion of  $\mathcal{D}(\mathbb{R}^N)$  with respect to the scalar product:

$$(u, v) = \begin{cases} \int_{\mathbb{R}^N} \Delta^{\frac{m}{2}} u \cdot \Delta^{\frac{m}{2}} v, & m \text{ even} \\ \int_{\mathbb{R}^N} \nabla \Delta^{\frac{m-1}{2}} u \cdot \nabla \Delta^{\frac{m-1}{2}} v, & m \text{ odd.} \end{cases}$$

This result generalizes the result of Ding for  $m = 1$ , and provides interesting information concerning the number and the kind of the solutions of the equation.

Wang in [66] studied the following nonlinear Neumann elliptic problem:

$$(P_N) \quad \begin{cases} -\Delta u = u^{\frac{N+2}{N-2}}, & u > 0 \text{ in } \mathbb{R}^N \setminus \Omega, \\ u(x) \rightarrow 0 & \text{as } |x| \rightarrow +\infty, \\ \frac{\partial u}{\partial n} = 0 & \text{on } \partial\Omega, \end{cases}$$

where  $n$  denotes interior unit normal vector and  $\Omega$  is a smooth bounded domain in  $\mathbb{R}^N$ ,  $N \geq 4$ . In this paper, it is proved that if  $N \geq 4$ , (Wang believes that the results will also hold in the case of  $N = 3$ ), and  $\Omega$  is a smooth and bounded domain then the problem  $(P_N)$  has infinity many non-radial positive solutions, whose energy can be made arbitrarily large when  $\Omega$  is convex, as seen from inside (with some symmetries). We refer to the Wang’s problem  $(P_N)$  due to its close relationship with our problem and as we will see later, if we choose suitable  $\Omega$  we can have a result on this problem in almost all the space. In particular, in both problems we have to solve the same non-linear differential equation with critical exponent with boundary conditions Dirichlet and Neumann, respectively. In addition, in both cases the domain  $\Omega$  presents some symmetries. However, a subsequent process in each case is completely different from that of another. In our case, our goal is to solve the problem *almost in the whole space*, starting from an open symmetric domain  $\Omega$  of  $n$ -dimensional space and we extend the  $\Omega$  so that it remains symmetrical to fill almost all the space. In the other case is considered the corresponding Neumann problem in  $\mathbb{R}^N \setminus \Omega$  where  $\Omega$  is convex as seen from inside with some symmetries. If we choose appropriate a such  $\Omega$  with a small volume as much as we can say that the solutions of Wang satisfy the conditions of the problem *almost in the whole space*. Finally, in both problems we take an infinity number of non-radial solutions,



whose energy can be made arbitrary large, however in the first problem we find nodal solutions while in the second are founded positive solutions.

Concerning to the progress of the study of the problem (P) for  $p = 2$  a number of important articles are available (cf. [1, 3, 4, 9–11, 20, 21, 23, 27, 29, 30, 32, 37, 42, 43, 45, 47, 48, 52, 57, 58, 63, 66]).

The  $p$ -Laplace operator (or  $p$ -harmonic operator) occupies a similar position to the standard Laplace operator when it comes to nonlinear phenomena. In fact, many of the things that apply to the usual Laplace operator and consequently to the equations that relate to it also apply to the  $p$ -Laplace as well as his equations, except that the Principle of Superposition which is of course lost. A very detailed and complete study is provided by Lindqvist [44]. Also, a Morse theoretic study of a very general class of homogeneous operators that includes the  $p$ -Laplacian as a special case is presented by Perera, Agarwal, and O'Regan in [53]. However, the  $p$ -Laplacian operator also appears in many areas of physics, such as non-Newtonian fluid flows, turbulent filtration in porous media, plasticity theory, rheology, glaciology, radiation of heat (cf. [24, 35, 49]).

The  $p$ -Laplace operator is a particularly interesting and remarkable case and this fact is confirmed not only by the large number of articles dedicated to it but also by the multifaceted study of the problems related to it (cf. [13, 18, 22, 26, 28, 31, 36, 40, 46, 54, 55, 60, 67, 69]).

In the problem (P), considered for any  $1 < p < n - k$ , a main difficulty comes from the double lack of compactness. By lack of compactness, we mean that the functional that we consider do not satisfy the Palais-Smale condition (cf. [50, 51, 59, 61, 62, 70]), (i.e. there exists a sequence along which the functional remains bounded, its gradient goes to zero, and does not converge). However, for  $p \neq 2$ , a second difficulty arises from the fact that the  $p$ -Laplace operator is not conformal invariant operator so the methods used in the case of the Laplace operator cannot be applied.

Concerning the lack of compactness, the first difficulty comes from the fact that the exponent

$$p^*(k) = \frac{(n - k)p}{n - k - p}$$

is supercritical (in fact the critical of the supercritical), and the second one is some extra difficulty because of the lack of compactness in unbounded domains. But, it is well known (see in [15, 16, 25, 32]) that the symmetry property of the domain allows us to improve the Sobolev embedding in higher  $L^p$  spaces and we overcome the obstruction of the exponent. Regarding the problem of lack of compactness in unbounded domains we avoid solving problems in such domains by remaining in bounded domains and then we pass to unbounded with limit procedures. In addition, this ensures us the ability to overcome the problems due to the non-conformality of the  $p$ -Laplace operation.

To overcome all the above obstacles we consider the following corresponding problem

$$(P_\varepsilon) \quad \begin{cases} \Delta_p u_\varepsilon + \varepsilon a(x)|u_\varepsilon|^{p-2}u_\varepsilon = f(x)|u_\varepsilon|^{p^*(k)-2}u_\varepsilon \\ u_\varepsilon \not\equiv 0 \text{ in } \Omega_\varepsilon, u_\varepsilon = 0 \text{ on } \partial\Omega_\varepsilon \\ 1 < p < n - k, p^*(k) = \frac{(n-k)p}{n-k-p}, \end{cases}$$

where  $\Omega_\varepsilon$ ,  $\varepsilon > 0$ , is an expanding domain in  $\mathbb{R}^n$ ,  $n \geq 3$ , invariant under the action of a subgroup  $G$  of the isometry group  $O(n)$  and  $a, f \in C^\infty(\overline{\Omega_\varepsilon})$  are two smooth  $G$ -invariant functions on  $\overline{\Omega_\varepsilon}$ .

The problem  $(P_\varepsilon)$  has been studied by many authors. We refer to [3, 4, 10, 20, 23, 27, 29, 32] and the references therein for a further discussion of both the problem itself and several variants of it. Some special cases have been also studied. For example, no solution can exist if  $\Omega$  is starshaped, as a consequence of the Pohozaev identity (see in [56]). Furthermore, if  $\Omega$  is an annulus, there are infinite solutions (see in [43]). Also, a general result of Bahri and Coron guarantees the existence of positive solutions in domains  $\Omega$  having nontrivial topology (i.e. certain homology groups of  $\Omega$  are non trivial) (see in [7]). The existence and multiplicity of positive or nodal solutions of critical equations on bounded domains or in some contractible domains have been determined by other authors (see for example in [21, 27, 32, 52, 63]). Some more nonexistence results in this case are available, (see in [1, 4, 12, 37]).

Our proof is via approximation by an infinite sequence of problems defined on a sequence of expanding symmetric bounded domains. Firstly, we solve the problem  $(P_\varepsilon)$ . (see in [14, 42] for the case of the Laplacian and for the case of the  $p$ -Laplacian see in [15, 16], for  $n = 3$ , and for  $n \geq 3$ , respectively). Then we consider a sequence of problems  $(P_{\varepsilon_j})$ ,  $j = 1, 2, \dots$ , defined in a sequence of expanding domains  $\Omega_{\varepsilon_j}$ ,  $j = 1, 2, \dots$ , and henceforth, sending  $\varepsilon \rightarrow 0$ , we obtain the solution of the limit problem  $(P)$  as the limit of the sequence of the solutions of the problems  $(P_{\varepsilon_j})$ . This method is a generalization of the method we used in [42] and thus a uniform treatment of both cases  $p = 2$  and  $p \in (1, 2) \cup (2, n - k)$  is achieved. In addition, the used method is a different from previous ones and can be used to solve poly-harmonic equations with supercritical exponent and even in the critical of supercritical case, as in our case, providing an alternative way of utilizing the best constants of the appearing Sobolev inequalities. Furthermore, this method enables us to determine the kind and the number of solutions of the problem in both cases, i.e. for  $p = 2$  and for  $p \in (1, 2) \cup (2, n - k)$ .

This article is organized as follows: Section 2 is devoted to notations and in some necessary background material. In Section 3, we introduce our main tool, meaning the process through which an open symmetric domain of  $n$ -dimensional space can be extended in an appropriate manner to ‘fill’ eventually the entire space ‘almost everywhere’, remaining symmetric, and subsequently we solve the auxiliary problem  $(P_\varepsilon)$ . Section 4 is devoted to some basic definitions and to the proof of the main theorem.

## 2 Notations and Some Background Material

As referred in the beginning of this article, our main objective is to prove the existence of an infinite sequence of distinct non-radial nodal  $G$ -invariant solutions defined in ‘almost the whole’ Euclidean space for the supercritical nonlinear elliptic problem (P). However, before dealing with problem (P) let us consider the following basic problem which will play an important role in solving the problem (P).

$$(P_0) \quad \begin{cases} \Delta_p u + a(x)|u|^{p-2}u = f(x)|u|^{q-2}u \\ u \neq 0 \text{ in } \Omega, \quad u = 0 \text{ on } \partial\Omega, \end{cases}$$

where  $\Omega$  is a bounded, smooth, domain of  $\mathbb{R}^n$ ,  $n \geq 3$ .

If we consider

$$p^* = \frac{np}{n-p}$$

it is well known by Sobolev’s embedding theorem (cf. [6]) the embedding

$$H_0^{1,p}(\Omega) \hookrightarrow L^p(\Omega)$$

is compact for any  $p \in [1, p^*)$  but the embedding

$$H_0^{1,p}(\Omega) \hookrightarrow L^{p^*}(\Omega)$$

is only continuous.

We say that the exponent

$$p^* = \frac{np}{n-p}$$

for the Sobolev embedding

$$H_0^{1,p}(\Omega) \hookrightarrow L^p(\Omega)$$

is the *critical exponent* for this embedding and that the problem (P) is *supercritical*, *critical* or *subcritical* if  $q - 1 < p^*$ ,  $q - 1 = p^*$  or  $q - 1 > p^*$  respectively. If  $p > n$  the problem (P) is always sub-critical.

In order to make this article self-contained we will open at this point a parenthesis where we will introduce some useful background material from the geometry. (More details see in [9] or [38]).

Consider a group  $G$  acting on a set  $X$ . The *orbit* of a point  $x$  in  $X$  is the set of elements of  $X$  to which  $x$  can be moved by the elements of  $G$ . (Just as gravity moves a planet around in its orbit, the group action moves an element around in its orbit.)

The  $G$  – orbit of  $x$  is denoted by

$$O_G(x) = \{\tau(x), \tau \in G\},$$

and for any  $Y \subseteq X$ , we write

$$G(Y) = \{\tau(y) : y \in Y \text{ and } \tau \in G\}.$$

If for some subset  $Y \subseteq X$  is valid

$$G(Y) = Y,$$

then, we say that  $Y$  is *invariant* under the action of  $G$  and in this case we denote it by  $Y_G$ .

For every  $x \in X$ , we define the *stabilizer subgroup* of  $G$  with respect to  $x$  (also called the *isotropy group*) as the set of all elements in  $G$  that fix  $x$ :

$$S_G(x) = \{\tau \in G : \tau(x) = x\}.$$

Moreover, if the set  $X$  is equipped with a metric, then the *isometry group* of this metric space is the set of all isometries (i.e. distance-preserving maps) from the metric space onto itself, with the function composition as group operation. Its identity element is the identity function (i.e. the isometry group of a two-dimensional sphere is the orthogonal group  $O(3)$ ).

Given  $(M, g)$  a Riemannian manifold (complete or not, but connected), we define by  $I(M, g)$  its group of isometries. It is well known (see for instance [38]) that  $I(M, g)$  is a Lie group with respect to the compact open topology, and that  $I(M, g)$  acts differentiably on  $M$ . Since (this is actually due to E. Cartan) any closed subgroup of a compact Lie group is a Lie group, we get that any compact subgroup of  $I(M, g)$  is a sub-Lie group of  $I(M, g)$ . It is now classical (see [9] and [19]), that for any  $x \in M$ ,  $O_G(x)$  is a smooth compact sub-manifold of  $M$ .

We denote by  $|O_G(x)|$  the volume of  $O_G(x)$  for the Riemannian metric induced on  $O_G(x)$ . In the special case where  $O_G(x)$  has finite cardinal, then,

$$|O_G(x)| = \text{card } O_G(x).$$

Let  $G$  be a closed subgroup of  $I(M, g)$ . Assume that for any  $x \in M$ ,

$$\text{card } O_G(x) = +\infty,$$

and set

$$k = \min_{x \in M} \dim O_G(x).$$

Then  $k \geq 1$  (see [32]), and is called *minimum orbit dimension*.

We consider a bounded, smooth domain  $\Omega$  of  $\mathbb{R}^n = \mathbb{R}^k \times \mathbb{R}^{n-k}$ ,  $k \geq 2$ ,  $n-k \geq 1$  such that

$$\overline{\Omega} \subset (\mathbb{R}^k \setminus \{0\}) \times \mathbb{R}^{n-k}.$$

Suppose that  $\overline{\Omega}$  is invariant under the action of  $G_{k,n-k}$ , that is

$$\tau(\overline{\Omega}) = \overline{\Omega}, \text{ for all } \tau \in G_{k,n-k},$$

where  $G_{k,n-k} = O(k) \times Id_{n-k}$  (then denoted by  $G$ ), is the subgroup of the isometry group  $O(n)$  of the type

$$(x_1, x_2) \longrightarrow (\sigma(x_1), x_2), \quad \sigma \in O(k), \quad x_1 \in \mathbb{R}^k, \quad x_2 \in \mathbb{R}^{n-k}.$$

For example, a such  $\Omega$  in  $\mathbb{R}^3$  is the solid torus

$$\overline{T} = \left\{ (x, y, z) \in \mathbb{R}^3 : \left( \sqrt{x^2 + y^2} - R \right)^2 + z^2 \leq r^2, \quad R > r > 0 \right\}.$$

Also, as such  $\Omega$  we can see the part of the  $n$ -dimensional ball  $B_n$  from which we have removed a part of it in such a way that the rest is invariant under the action of the group  $G$  and its cover belongs to  $\overline{B}_n \subset (\mathbb{R}^k \setminus \{0\}) \times \mathbb{R}^{n-k}$ . This is because the balls enjoy a large number of symmetries in addition to the radial symmetry.

We define

$$C_G^\infty(\Omega) = \{u \in C^\infty(\Omega) : u \circ \tau = u, \forall \tau \in G\},$$

and

$$C_{0,G}^\infty(\Omega) = \{u \in C_0^\infty(\Omega) : u \circ \tau = u, \forall \tau \in G\},$$

where  $C^\infty(\Omega)$  denotes the space of smooth functions on  $\Omega$  and where  $C_0^\infty(\Omega)$  denotes the space of smooth functions with compact support on  $\Omega$ .

We define, also, the Sobolev space  $H^{1,p}(\Omega)$  as the completion of  $C^\infty(\Omega)$  with respect to the norm

$$\|u\|_{H^{1,p}(\Omega)} = \left( \|\nabla u\|_{L^p(\Omega)}^p + \|u\|_{L^p(\Omega)}^p \right)^{1/p}, \quad p \geq 1,$$

and the Sobolev space  $H_0^{1,p}(\Omega)$  as the closure of  $C_0^\infty(\Omega)$  in  $H^{1,p}(\Omega)$ .

Finally, we denote by  $H_G^{1,p}(\Omega)$  and  $H_{0,G}^{1,p}(\Omega)$  the subspaces of  $H^{1,p}(\Omega)$  and  $H_0^{1,p}(\Omega)$ , respectively, of all  $G$ -invariant functions defined on  $\Omega$ .

It is well known that the symmetry property of the domain allows us to improve the Sobolev embedding in higher  $L^p$  spaces. More precisely, let us consider a smooth compact  $n$ -dimensional,  $n \geq 3$ , Riemannian manifold  $(M, g)$  invariant under the action of an arbitrary compact subgroup  $G$  of  $\text{Isom}_g(M)$ . Let us also assume that

$$\text{Card } O_G^x = +\infty$$

for any orbit  $O_G^x$  of  $G$  and  $k \geq 1$ . It is well known that the Sobolev embedding

$$H_G^{1,p}(M) \hookrightarrow L^q(M)$$

is compact for any

$$1 \leq q < \frac{(n-k)p}{n-k-p}$$

but if

$$1 \leq q \leq \frac{(n-k)p}{n-k-p}$$

is only continuous (cf. in [14, 16, 20, 25, 33, 34]).

### 3 Preliminary Results

Let  $\Omega$  be a domain such that  $\overline{\Omega} \subset (\mathbb{R}^k \setminus \{0\}) \times \mathbb{R}^{n-k}$  and also invariant under the action of the group  $G$  defined above. For any small  $\varepsilon > 0$  and some  $m > 0$  (which will be determined later) we consider the family of expanding domains

$$\Omega_\varepsilon = \varepsilon^{-m} \Omega = \{\varepsilon^{-m} x : x \in \Omega\}$$

Then, it is very simple to be confirmed that  $\Omega_\varepsilon$ s inherit the symmetry properties of  $\Omega$  for any  $\varepsilon$ .

At this point we need to comment on the term ‘almost the whole’ space and specify the impact of this term on solutions to the problem. To do this we must describe the process by which we ‘fill’ the space by properly expanding the domain  $\Omega$  and then see how the method of solving the problem works. In fact we consider a sequence consisting of  $\Omega_{\varepsilon_j}$ , where the sequence of  $\varepsilon_j$ s (for the time being) is a sequence that tends to 0 in such a way that  $\Omega_{\varepsilon_j}$ s extend continuously and as  $\varepsilon \rightarrow 0$  they cover “almost everywhere” the entire space. This is because this extension also entails the inside boundary of the  $\Omega_{\varepsilon_j}$ s (i.e. the one that is on the zero side) and of course increases the volume of the orbit with the minimum dimension. This does not

pose a problem for us in the solution because it depends only on the volume of this orbit (apart from the other parameters) (see Theorem 3) so we can extend these  $\Omega_{\varepsilon_j}$ s to the inside as much as we want by extending with zero values the functions defined in them. The outer boundary of  $\Omega_{\varepsilon_j}$ s does not impose any restrictions and this is because any orbits close to it do not play any role since as mentioned above only the orbit with the minimum volume affects the solutions and it is on the opposite side, the side of zero. Finally, the fact that the domain is also expanding does not affect either the Sobolev inequalities associated with the problem or the solutions because we can, for example, normalize the functions  $u_j$  i.e. so that their norms are equal to 1.

We consider now the transformation

$$\phi : \Omega \rightarrow \Omega_\varepsilon : X = \varepsilon^{-m}x, \quad x \in \Omega, \quad X \in \Omega_\varepsilon \tag{1}$$

and for  $\ell > 0$  we set

$$u_\varepsilon(X) = \varepsilon^{-\ell} u(\varepsilon^m X).$$

In particular we obtain

$$|\nabla u| = \varepsilon^{-m} |\nabla u_\varepsilon| \tag{2}$$

and

$$\Delta_p u = -\varepsilon^{-mp} \operatorname{div} \left( |\nabla u_\varepsilon|^{p-2} \nabla u_\varepsilon \right). \tag{3}$$

Note the equality (3) remains valid for  $p = 2$ , i.e. for  $\Delta_2 = \Delta$ , the Laplace–Beltrami operator.

In the following, we will suppose that  $p \neq 2$ , since the case where  $p = 2$  was studied in [42].

Applying the transformation (1) in the equation of the problem (P<sub>0</sub>), because of (2) and (3), we obtain the following equation

$$\Delta_p u_\varepsilon + \varepsilon^{mp+\ell(2-p)} a(x) |u_\varepsilon|^{p-2} u_\varepsilon = \varepsilon^{mp+\ell(2-q)} f(x) |u_\varepsilon|^{q-2} u_\varepsilon.$$

Since  $\ell$  is an arbitrary positive real, we can choose  $\ell = \frac{mp}{q-2}$  and thus we obtain the following equation:

$$\Delta_p u_\varepsilon + \varepsilon^{mp+mp(2-p)/(q-2)} a(x) |u_\varepsilon|^{p-2} u_\varepsilon = f(x) |u_\varepsilon|^{q-2} u_\varepsilon. \tag{4}$$

Finally, replacing the  $\varepsilon^{mp+mp(2-p)/(q-2)}$  by  $\varepsilon$ , we can write the equation (4) in the following form

$$\Delta_p u_\varepsilon + \varepsilon a(x) |u_\varepsilon|^{p-2} u_\varepsilon = f(x) |u_\varepsilon|^{q-2} u_\varepsilon. \tag{5}$$

Let  $\Omega$  be a smooth bounded domain in  $\mathbb{R}^n$ ,  $G$ -invariant and  $k$  be the minimum of the dimensions of all orbits of  $G$  with infinite cardinal. Let, also,  $\Omega_\varepsilon$  as defined above. A such  $\Omega$  is the above defined solid torus  $T$  and a such  $\Omega_\varepsilon$ , in this case, is an expanding torus  $T_\varepsilon$ .

Now for any  $\varepsilon > 0$  consider the following auxiliary problem:

$$(P_\varepsilon) \quad \begin{cases} \Delta_p u_\varepsilon + \varepsilon a(x)|u_\varepsilon|^{p-2}u_\varepsilon = f(x)|u_\varepsilon|^{p^*(k)-2}u_\varepsilon \\ u_\varepsilon \not\equiv 0 \text{ in } \Omega_\varepsilon, \quad u_\varepsilon = 0 \text{ on } \partial\Omega_\varepsilon \\ 1 < p < n - k, \quad p^*(k) = \frac{(n-k)p}{n-k-p}, \end{cases}$$

where  $a, f$  are two smooth  $\mathcal{H}_\sigma^p$ -invariant functions (defined bellow).

Before, we solve the problem  $(P_\varepsilon)$ , we must compute the best constant  $K_G^p(\Omega_\varepsilon)$  in the following Sobolev inequality, which appears in this problem:

$$\left( \int_{\Omega_\varepsilon} |u|^{p^*(k)} dx \right)^{\frac{p}{p^*(k)}} \leq (K_G^p(\Omega_\varepsilon) + \epsilon) \int_{\Omega_\varepsilon} |\nabla u|^p dx + B_\epsilon \int_{\Omega_\varepsilon} |u|^p dx, \quad (6)$$

where  $\epsilon$  is a positive constant no matter how small, but it cannot disappear and  $B_\epsilon$  a positive constant.

In fact we will express the best constant  $K_G(\Omega_\varepsilon)$  of inequality (6) as a function of the optimal constant the best constant  $K_G^p(\Omega)$  and  $\varepsilon$ .

Concerning this best constant the following theorem holds:

**Theorem 1**

$$K_G(\Omega_\varepsilon) = \varepsilon^m K_G(\Omega) = \frac{K(n - k, p)}{\varepsilon^{-m} \mathcal{V}^{\frac{1}{n-k}}}$$

where  $K(n - k, p)$  is the best constant in the classical Sobolev inequality of  $\mathbb{R}^{n-k}$  and  $\mathcal{V}$  denotes the minimum of the volume of the  $k$ -dimensional orbits in  $\Omega$ .

**Proof** According to the Theorem 2.1 in [16], (see, also, Theorem 3.1 in [15]) we have

$$K_G^p(\Omega_\varepsilon) = \frac{K(n - k, p)}{\mathcal{V}_\varepsilon^{\frac{1}{n-k}}}$$

and since

$$\mathcal{V}_\varepsilon = \varepsilon^{-m(n-k)} \mathcal{V}$$



we obtain

$$K_G^p(\Omega_\varepsilon) = \frac{K(n-k, p)}{\varepsilon^{-m} \gamma^{\frac{1}{n-k}}}$$

□

Now, for the problem  $(P_\varepsilon)$ , consider the functional

$$J(u_\varepsilon) = \int_{\Omega_\varepsilon} (|\nabla u_\varepsilon|^p + \varepsilon a(x)|u_\varepsilon|^p) dx$$

and suppose that the operator

$$L_p(u_\varepsilon) = \Delta_p u_\varepsilon + \varepsilon a(x)|u_\varepsilon|^{p-2} u_\varepsilon$$

is coercive.

Denote

$$\mathcal{H}^p = \left\{ u_\varepsilon \in H_{0,G}^{1,p}(\Omega_\varepsilon) : \int_{\Omega_\varepsilon} f(x)|u_\varepsilon|^q dx = 1 \right\},$$

$$\mu_\varepsilon = \inf J(u_\varepsilon),$$

for all  $u_\varepsilon \in \mathcal{H}^p$ , and suppose that exists an isometry  $\sigma$  such that  $\sigma(\Omega_\varepsilon) = \Omega_\varepsilon$ . Moreover we suppose that the functions  $a(x)$  and  $f(x)$  are invariant under the action of  $\sigma$ , and

$$\mathcal{H}_\sigma^p = \mathcal{H}^p \cap \{ u_\varepsilon \in H_{0,G}^{1,p}(\Omega_\varepsilon) : u_\varepsilon \circ \sigma = -u_\varepsilon \} \neq \emptyset.$$

Then, we have the following theorems.

**Theorem 2** For  $p = 2$  and  $n \geq 3$ , the problem  $(P_\varepsilon)$ , always, has a non-radial nodal solution  $u$ . Moreover, if  $f(x) > 0$  for all  $x \in \overline{\Omega_\varepsilon}$ ,  $(P_0)$  has an infinity sequence  $\{u_{\varepsilon_i}\}$  of non-radial nodal solutions, such that

$$\lim_{i \rightarrow \infty} \int_{\Omega_\varepsilon} (|\nabla u_{\varepsilon_i}|^2 + u_{\varepsilon_i}^2) dx = +\infty.$$

In addition,  $u$  and  $\{u_{\varepsilon_i}\}_{i=1,2,\dots}$  are  $G$ -invariant and  $\sigma$ -antisymmetrical.

**Theorem 3** Let  $a$  and  $f$  be two smooth functions,  $\mathcal{H}_\sigma^p$ -invariant and  $p, q$  be two real numbers defined as in  $(P_\varepsilon)$ . Suppose that  $\sup_{x \in \Omega_\varepsilon} f(x) > 0$  and the operator  $L_p$  is coercive. Then the problem  $(P_\varepsilon)$  has a non-radial nodal  $\mathcal{H}_\sigma^p$ -invariant solution, that belongs to  $C^{1,\alpha}(\Omega_\varepsilon)$  for some  $\alpha \in (0, 1)$ , if

$$\mu_\varepsilon < K_G^p(\Omega_\varepsilon)^{-p} \left( \sup_{x \in \Omega_\varepsilon} f(x) \right)^{-p/q}.$$

The proofs of Theorems 2 and 3 use standard variational methods, under the assumptions of Lemma 3.6 in [16], (cf. [14, 15, 17, 25]).

### 4 Solution of the Problem (P)

We return to our main problem

$$(P) \begin{cases} \Delta_p u = |u|^{p^*(k)} u, & u \in C^2(\mathbb{R}^n), \quad n \geq 3 \\ 1 < p < n - k, \quad p^*(k) = \frac{(n-k)p}{n-k-p}. \end{cases}$$

In the problem (P) direct variational methods are not applicable because of the double lack of compactness. To overcome this problem we will use an approximate method. That is, we consider a sequence of expanding  $\Omega_{\varepsilon_j}$  (where  $\varepsilon_j \rightarrow 0$  as  $j \rightarrow \infty$ ) as well as the sequence of problems

$$(P_{\varepsilon_j}) \begin{cases} \Delta_p u_{\varepsilon_j} + \varepsilon_j a(x) |u_{\varepsilon_j}|^{p-2} u_{\varepsilon_j} = f(x) |u_{\varepsilon_j}|^{p^*(k)-2} u_{\varepsilon_j} \\ u_{\varepsilon_j} \neq 0 \text{ in } \Omega, \quad u_{\varepsilon_j} = 0 \text{ on } \partial\Omega \\ 1 < p < n - k, \quad p^*(k) = \frac{(n-k)p}{n-k-p}, \end{cases}$$

where  $a, f$  are as in the problem  $(P_{\varepsilon_j})$ .

According to the Theorems 2 and 3, every problem  $(P_{\varepsilon_j})$  has a non-radial nodal  $\mathcal{H}_\sigma^p$ -invariant solution. Thus, a solution to the problem (P) may be then obtained by the limit procedure as  $\varepsilon_j \rightarrow 0$ .

Before we will approximate the solutions in  $\mathbb{R}^n$  by solutions in bounded domains  $\Omega_{\varepsilon_j} \in \mathbb{R}^n$ , we note that, in the generalized setting of the problems in  $\Omega_{\varepsilon_j}$ s, the Dirichlet condition  $u_{\varepsilon_j}(x) = 0$  on  $\partial\Omega_{\varepsilon_j}$  may actually be included in the condition  $u_{\varepsilon_j} \in H_{0,G}^{1,p}(\Omega_{\varepsilon_j})$ .

Moreover, since any function  $u_{\varepsilon_j} \in H_{0,G}^{1,p}(\Omega_{\varepsilon_j})$  can be extended onto  $\mathbb{R}^n$  by

$$\tilde{u}_\varepsilon(x) = \begin{cases} u_{\varepsilon_j}(x), & x \in \Omega_{\varepsilon_j} \\ 0, & x \in \mathbb{R}^n \setminus \Omega_{\varepsilon_j}, \end{cases}$$

generalized solutions may be defined in  $\Omega_{\varepsilon_j}$ s analogously to the case in  $\mathbb{R}^n$ .

We need now the following two definitions:

**Definition 1** A function  $u_{\varepsilon_j} \in H_{0,G}^{1,p}(\Omega_{\varepsilon_j})$  is a *generalized solution* of  $(P_{\varepsilon_j})$  if the function

$$g(x, u_{\varepsilon_j}) = \varepsilon_j a(x)u_{\varepsilon_j} - f(x)|u_{\varepsilon_j}|^{p^*(k)-2}u_{\varepsilon_j}$$

is locally integrable and for all  $\varphi \in C_0^\infty(\Omega_{\varepsilon_j})$ , the following holds:

$$\int_{\Omega_{\varepsilon_j}} |\nabla u_{\varepsilon_j}|^{p-2}(\nabla u_{\varepsilon_j}, \nabla \varphi)dx + \int_{\Omega_{\varepsilon_j}} f(x, u_{\varepsilon_j})\varphi dx = 0.$$

**Definition 2** A function  $u_\varepsilon \in C^2(\Omega_\varepsilon) \cap C(\overline{\Omega_\varepsilon})$  is a *classical solution* to  $(P_\varepsilon)$  if after substituting it into equation of  $(P_\varepsilon)$ , this equation becomes the identity at each  $x \in \Omega_\varepsilon$  and  $u_\varepsilon(x) = 0$  provided  $x \in \partial\Omega_\varepsilon$ .

Provided that all the conditions of the Theorem 3 are satisfied, we apply it to the sequence of the problems  $(P_{\varepsilon_j})$  and denote by  $\{u_j\}_{j=1}^\infty$  the sequence of the corresponding solutions.

Under the above considerations to following theorem holds.

**Theorem 4** *The problem*

$$\Delta_p u = f(x)|u|^{p^*(k)-2}u \quad \text{in } \mathbb{R}^n, \quad n \geq 3$$

has a generalized non-radial nodal  $\mathcal{H}_\sigma^p$ -invariant solution  $u$  and there is a subsequence of  $\{u_j\}$  (again denoted by  $\{u_j\}$ ) such that

$$u_j \rightharpoonup u \quad \text{in } H_{0,G}^{1,p} \quad \text{as } j \rightarrow \infty.$$

*In addition*

$$\lim_{j \rightarrow \infty} \int_{\mathbb{R}^n} |\nabla u_j|^p dx = +\infty.$$

**Proof** The case  $p = 2$  is presented in [42], thus, we will prove the case  $p \in (1, 2) \cup (2, n - k)$ . However, we present a unified proof for both cases. For the proof we borrow ideas from [42] and carried out in 5 steps.

*Step 1.* According to the above Theorem 3, every problem  $(P_{\varepsilon_j})$  has at least one non-radial nodal  $G$ -invariant and  $\sigma$  antisymmetrical solution  $u_j$ . Let  $u_j, j = 1, 2, \dots$ , an arbitrary sequence of such solutions. Since the problem  $(P_{\varepsilon_j})$  has a nontrivial solution belonging to one of the spaces considered earlier, then for any  $\lambda > 0$  the function

$$v_j = \lambda^{\frac{1}{p^*(k)-p}} u_j \in H_{0,G}^{1,p}(\Omega_{\varepsilon_j})$$

is a non trivial solution to the problem:

$$(P_{\varepsilon_j}^\lambda) \quad \begin{cases} \Delta_p v_j + \varepsilon_j a(x) |v_j|^{p-2} v_j = \lambda f(x) |v_j|^{p^*(k)-2} v_j \\ v_j \not\equiv 0 \text{ in } \Omega, \quad v_j = 0 \text{ on } \partial\Omega \\ 1 < p < n - k, \quad p^*(k) = \frac{(n-k)p}{n-k-p}. \end{cases}$$

In this first step of the proof we prove that there exists a sub-sequence of the sequence of the solutions to the problems  $(P_{\varepsilon_j}^\lambda)$  which converges weakly in

$$H_0^{1,p}(\mathbb{R}^n).$$

For

$$\lambda = \|u_j\|_{H^{1,p}(\Omega_{\varepsilon_j})}^{-(p^*(k)-p)}$$

we obtain that

$$v_j = \frac{u_j}{\|u_j\|_{H^{1,p}(\Omega_{\varepsilon_j})}},$$

which means that the sequence  $\{v_{\varepsilon_j}\}$  is bounded in  $H^{1,p}(\Omega_{\varepsilon_j})$  for all  $j = 1, 2, \dots$

Therefore, there is a positive constant  $C$  not dependent on  $j$  and such that:

$$\|v_j\|_{H^{1,p}(\Omega_{\varepsilon_j})} \leq C, \quad \forall j = 1, 2, \dots \tag{7}$$

Because of the reflexivity of  $H_0^{1,p}(\mathbb{R}^n)$  and condition (7) we may choose a sub-sequence of  $\{v_j\}$  (again denoted by  $\{v_j\}$ ) such that:

$$v_j \rightharpoonup v \text{ in } H_0^{1,p}(\mathbb{R}^n) \text{ as } j \rightarrow +\infty. \tag{8}$$

*Step 2.* In this step we prove that the function  $v$  is a nontrivial  $G$ -invariant generalized solution of the limit problem obtained from the sequence of problems  $(P_{\varepsilon_j}^\lambda)$  as  $j \rightarrow \infty$ .

We choose an arbitrary  $\varphi \in C_0^\infty(\mathbb{R}^n)$ . Then, according to the definition of  $C_0^\infty(\mathbb{R}^n)$ , the support of  $\varphi$  is bounded in  $\mathbb{R}^n$ , which means that there is an  $\Omega_{\varepsilon_0}$  such that  $\text{supp}\varphi \subset \Omega_{\varepsilon_0}$ . Since, by definition, the  $\Omega_{\varepsilon_j}$ s constitute a family of expanding domains, we can choose the  $\Omega_{\varepsilon_0}$  such that  $\Omega_{\varepsilon_0} \subset \Omega_{\varepsilon_1}$  and so  $\Omega_{\varepsilon_0} \subset \Omega_{\varepsilon_j}$  for all  $j = 1, 2, \dots$

Let

$$g(x, v_j) = \varepsilon_j a(x)v_j - \lambda f(x)|v_j|^{p^*(k)-2}v_j.$$

Then, because the  $v_j$  is a generalized solution to  $(P_{\varepsilon_j}^\lambda)$ , it holds

$$\int_{\mathbb{R}^n} |\nabla v_j|^{p-2} (\nabla v_j, \nabla \varphi) dx = - \int_{\Omega_{\varepsilon_j}} g(x, v_j) \varphi dx = - \int_{\Omega_{\varepsilon_0}} g(x, v_j) \varphi dx. \tag{9}$$

for all  $\Omega_{\varepsilon_j}$ .

By the weak convergence (8), we obtain the following limit relation for the left-hand side of (9):

$$\lim_{j \rightarrow \infty} \int_{\mathbb{R}^n} |\nabla v_j|^{p-2} (\nabla v_j, \nabla \varphi) dx = \int_{\mathbb{R}^n} |\nabla v|^{p-2} (\nabla v, \nabla \varphi) dx. \tag{10}$$

In addition, the critical exponent of the embedding

$$H_G^{1,p}(\Omega_{\varepsilon_0}) \hookrightarrow L^p(\Omega_{\varepsilon_0})$$

is equal to

$$p^*(k) = \frac{(n-k)p}{n-k-k} > \frac{np}{n-p} = p^*.$$

Let some  $p_0$  such that

$$p^* < p_0 < p^*(k).$$

Then the embedding is compact and thus from the Sobolev and Kondrashov theorems together and (9) arises that

$$v_j \rightarrow v \text{ in } L^{p_0-1}(\Omega_{\varepsilon_0}), \text{ as } j \rightarrow +\infty. \tag{11}$$

Furthermore, by definition of  $a(x)$  and  $f(x)$ , there exists a positive constant  $C$  such that:

$$|g(x, t)| \leq C(|t| + |t|^{p_0-1}), \quad p^* < p_0 < p^*(k),$$

for almost all  $x \in \Omega_{\varepsilon_j}$ ,  $j = 1, 2, \dots$  and for all  $t \in \mathbb{R}$ .

Therefore, by Vainberg-Krasnoselskii Theorem (cf. [39, 65] or [41]) gives that:

$$\varphi g(\cdot, v_j(\cdot)) \rightarrow \varphi g(\cdot, v(\cdot)) \text{ in } L^{\frac{p_0}{p^*}}(\Omega_{\varepsilon_0}) \text{ as } j \rightarrow +\infty \tag{12}$$

and the Hölder inequality from (12) follows that:

$$\varphi g(\cdot, v_j(\cdot)) \rightarrow \varphi g(\cdot, v(\cdot)) \text{ in } L^r(\Omega_{\varepsilon_0}) \text{ as } j \rightarrow +\infty, \tag{13}$$

for all  $1 \leq r \leq \frac{p_0}{p^*}$ .

By (13) the limit relation from the right hand-side of (10) yields:

$$\lim_{j \rightarrow \infty} \int_{\Omega_{\varepsilon_0}} g(x, v_j) \varphi dx = \int_{\Omega_{\varepsilon_0}} g(x, v) \varphi dx. \tag{14}$$

Finally, passing to the limit in (9) because of (8) and (14), we obtain:

$$\int_{\mathbb{R}^n} |\nabla v|^{p-2} (\nabla v, \nabla \varphi) dx = - \int_{\Omega_{\varepsilon_0}} g(x, v) \varphi dx = - \int_{\mathbb{R}^n} g(x, v) \varphi dx,$$

which corresponds to the definition of a weak solution. This is a generalized solution by the force of (9) and since the function  $f$  is regular enough it is a classical solution, (see §§ 1.2 and 3.1 in [41]). As convergence in  $L^p$  spaces implies a.e. convergence by (11) follows that the function  $v$  will be  $G$ -invariant.

*Step 3.* In this step we prove that the solution  $v$  is non trivial, that is  $v \not\equiv 0$ . Suppose, by contradiction, that  $v \equiv 0$ . Then, for any  $\varepsilon > 0$  we have

$$|v| < \frac{\varepsilon}{2}. \tag{15}$$

On the other hand, from (13) arises that

$$v_j \rightarrow v \text{ in } L^1(\Omega_{\varepsilon_0}),$$

which means that for any  $\varepsilon > 0$  there exists a positive integer  $j_0$  such that:

$$|v_j - v| < \frac{\varepsilon}{2} \text{ for all } j > j_0. \tag{16}$$

Therefore, by the standard inequality

$$|v_j| \leq |v_j - v| + |v|$$

due to (15) and (16) we obtain that:

$$|v_j| < \varepsilon \text{ for any } j \geq j_0. \tag{17}$$

We recall now that every solution to the problem  $(P_{\varepsilon_j})$  belongs to the set

$$\mathcal{H}_{\varepsilon}^{\sigma} = \left\{ u_{\varepsilon} \in H_{0,G}^{1,p}(\Omega_{\varepsilon_j}) : u_{\varepsilon_j} \circ \sigma = -u_{\varepsilon_j} \text{ and } \int_{\Omega_{\varepsilon_j}} f(x) |u_{\varepsilon_j}|^{p^*(k)} dx = 1 \right\}.$$

Since every  $v_j$  corresponds to an  $u_{\varepsilon_j} \in \mathcal{H}_\varepsilon^\sigma$ , and  $v_{\varepsilon_j} = \lambda^{\frac{1}{p^*(k)-p}} u_{\varepsilon_j}$ , by definition, we have the following:

$$1 = \int_{\Omega_{\varepsilon_j}} f(x) \lambda^{-\frac{p^*(k)}{p^*(k)-p}} |v_j|^{p^*(k)} dx < \int_{\Omega_{\varepsilon_j}} f(x) \lambda^{-\frac{p^*(k)}{p^*(k)-p}} \varepsilon^{p^*(k)} dx,$$

which is false due to (17) as the  $\varepsilon > 0$  can be chosen as small as we want.

*Step 4.* We have proved that the limit problem

$$(P^\lambda) \quad \Delta_p v = \lambda f(x) |v|^{p^*(k)-2} v \quad \text{in } \mathbb{R}^n, \quad n \geq 3$$

has a generalized non-radial nodal  $G$ -invariant and  $\sigma$ -anti-symmetrical solution  $v$ , which means that the function  $u = \lambda^{\frac{1}{p^*(k)-p}} v$  is a generalized non-radial nodal  $G$ -invariant and  $\sigma$ -anti-symmetrical solution to the limit problem:

$$(P) \quad \Delta_p u = f(x) |u|^{p^*(k)-2} u \quad \text{in } \mathbb{R}^n, \quad n \geq 3.$$

*Step 5.* It remains to prove that

$$\lim_{j \rightarrow \infty} \int_{\mathbb{R}^n} |\nabla u_j|^p dx = +\infty.$$

The Sobolev inequality (6) after a normalization of the sequence  $u_j$ s so that  $\|u_j\|_{L^{p^*(k)}(\Omega_{\varepsilon_j})} = 1$  and provided that the constants  $B_\varepsilon$  are positive give us that

$$1 \leq (K_G^p(\Omega_{\varepsilon_j}) + \epsilon) \int_{\Omega_{\varepsilon_j}} |\nabla u_j|^p dx. \tag{18}$$

From (18) after a replacement of the constant  $K_G^p(\Omega_{\varepsilon_j})$  from the one calculated in Theorem 1 we obtain the inequality

$$\frac{1}{\varepsilon_j^m V^{-\frac{1}{n-k}} K(n-k, p) + \epsilon} < \int_{\Omega_{\varepsilon_j}} |\nabla u_j|^p dx. \tag{19}$$

By inequality (19) taking the limits for  $j \rightarrow \infty$  we have that  $\varepsilon_j \rightarrow 0$  and then

$$\int_{\mathbb{R}^n} |\nabla u_j|^p dx \rightarrow \infty.$$

This completes the proof of the theorem. □

**Corollary 1** *The problem:*

$$(P) \begin{cases} \Delta_p u = |u|^{p^*(k)} u, & u \in C^2(\mathbb{R}^n), \quad n \geq 3 \\ 1 < p < n - k, \quad p^*(k) = \frac{(n-k)p}{n-k-p}, \end{cases}$$

has a sequence  $\{u_j\}$  of non-radial nodal  $G$ -invariant and  $\sigma$ -anti-symmetrical solutions, such that:

$$\lim_{j \rightarrow +\infty} \int_{\mathbb{R}^n} |\nabla u_j|^p dx = +\infty.$$

**Proof** The result is obtained if we put

$$f(x) = \frac{1}{|\Omega_{\varepsilon_j}|} - \varepsilon_j |x|^\alpha, \quad \alpha > -n$$

and follows the spirit of the approach in Theorem 4. □

*Remark 1* The number of the sequences of non-radial nodal  $G$ -invariant and  $\sigma$ -anti-symmetrical solutions to the problem (P), depends on the number of all subgroups of  $O(n)$  of which the cardinal of orbits with minimum volume is infinite, that are on the dimension  $n$  of the domain.

## References

1. K. Adimurthi, S.L. Yadava, Elementary proof of the nonexistence of nodal solutions for the semilinear elliptic equations with critical Sobolev exponent. *Nonlinear Anal.* **14**(9), 785–787 (1990)
2. A. Ambrosetti, P.H. Rabinowitz, Dual variational methods in critical point theory and applications. *J. Funct. Anal.* **14**, 349–381 (1973)
3. A. Ambrosetti, M. Struwe, A note on the problem  $-\Delta u = \lambda u + |u|^{2^*-2}u$ . *Manuscripta Math.* **54**, 373–379 (1986)
4. F.V. Atkinson, H. Brezis, L.A. Peletier, Nodal solutions of elliptic equations with critical Sobolev exponents. *J. Diff. Eq.* **85**, 151–170 (1990)
5. T. Aubin, Equations différentielles non linéaires et problème de Yamabe concernant la courbure scalaire. *J. Math. Pures Appl.* **55**, 269–296 (1976)
6. T. Aubin, *Some Non Linear Problems in Riemannian Geometry* (Springer, Berlin, 1998)
7. A. Bahri, J.M. Coron, On a nonlinear elliptic equation involving the limiting Sobolev exponent. *Comm. Pure Appl. Math.* **41**, 253–294 (1988)
8. T. Bartsch, M. Schneider, Multiple solutions of a critical polyharmonic equation. *J. Reine Angew. Math.* **571**, 131–143 (2004)
9. H. Brezis, Elliptic equations with limiting Sobolev exponent—the impact of topology, in *Proceeding of the 50th Anniversary Courant Institute Communications on Pure Applied Mathematics*, vol. 39 (1986), pp. 517–539
10. H. Brezis, L. Nirenberg, Positive solutions of nonlinear elliptic equations involving critical Sobolev exponents. *Comm. Pure Appl. Math.* **36**(4), 437–477 (1983)



11. L. Caffarelli, B. Gidas, J. Spruck, Asymptotic symmetry and local behavior of semilinear elliptic equations with critical Sobolev growth. *Comm. Pure Appl. Math.*, **42**, 271–297 (1989)
12. A. Carpio Rodriguez, M. Comte, R. Lewandoski, A nonexistence result for nonlinear equation involving critical Sobolev exponent. *Ann. Inst. Henri Poincaré Anal. Non Linéaire* **9**, 243–261 (1992)
13. Y.-H. Cheng, Eigenvalue problems with p-Laplacian operators. *Electron. J. Differ. Equations* **2014**(139), 1–11 (2014)
14. A. Cotioliis, D. Iliopoulos, Equations elliptiques non lineaires a croissance de Sobolev sur-critique. *Bull. Sci. Math.* **119**, 419–431 (1995)
15. A. Cotioliis, N. Labropoulos, Dirichlet problem on a solid torus in the critical of supercritical case. *Bull. Greek Math. Soc.* **53**, 39–57 (2007)
16. A. Cotioliis, N. Labropoulos, Best constants in Sobolev inequalities on manifolds with boundary in the presence of symmetries and applications. *Bull. Sci. Math.* **132**, 562–574 (2008)
17. F. Demengel, E. Hebey, On some nonlinear equations involving the p-Laplacian with critical Sobolev growth. *Adv. Diff. Equations* **3**, 533–574 (1998)
18. P. De N'apoli1, M.C. Mariani, Mountain pass solutions to equations of p-Laplacian type. *Nonlinear Anal.* **54**, 1205–1219 (2003)
19. J. Dieudonné, *Eléments d'Analyse*, tome 3. Gauthier-Villars **1974**, 562–574 (2008)
20. W. Ding, On a conformally invariant elliptic equation on  $\mathbb{R}^n$ . *Commun. Math. Phys.* **107**, 331–335 (1986)
21. W. Ding, Positive solutions of  $\Delta u + u^{\frac{n+2}{n-2}} = 0$  on contactible domains. *J. Partial Diff. Equation* **2**, 83–88 (1989)
22. H. Ding, J. Zhou, Global existence and blow-up for a mixed pseudo-parabolic p-Laplacian type equation with logarithmic nonlinearity. *J. Math. Anal. Appl.* **478**, 393–420 (2019)
23. L. Dupaigne, M. Ghergu, V. Rădulescu, Lane-Emden-Fowler equations with convection and singular potential. *J. Math. Pures Appl.* **87**(6), 563–581 (2007)
24. J. Esteban, J. Vazquez, On the equation of turbulent filtration in one dimensional porous media. *Nonlinear Anal.* **10**, 1303–1325 (1986)
25. Z. Faget, Best constants in Sobolev inequalities on Riemannian manifolds in the presence of symmetries. *Potential Anal.* **17**, 105–124 (2002)
26. R. Filippucci, P. Pucci, F. Robert, On a p-Laplace equation with multiple critical nonlinearities. *J. Math. Pures Appl.* **91**, 156–177 (2009)
27. D. Fortunato, E. Jannelli, Infinity many nodal solutions for some nonlinear elliptic problems in symmetrical domains. *Proc. R. Soc. Edinb.* **105**(A), 205–213 (1987)
28. M. Fraas, Y. Pinchover, Positive liouville theorems and asymptotic behavior for p-Laplacian type elliptic equations with a fuchsian potential. *Confluentes Mathematici* **3**(2), 291–323 (2011)
29. M. Ghergu, V. Rădulescu, Ground state solutions for the singular Lane-Emden-Fowler equation with sublinear convection term. *J. Math. Anal. Appl.* **333**, 265–273 (2007)
30. B. Gidas, W.-M. Ni, L. Nirenberg, Symmetry and related properties via the maximum principle. *Commun. Math. Phys.* **68**, 209–243 (1979)
31. C.-Y. Guo, M. Kar, Quantitative uniqueness estimates for p-Laplace type equations in the plane. *Nonlinear Anal.* **143**, 19–44 (2016)
32. E. Hebey, M. Vaugon, Existence and multiplicity of nodal solutions for nonlinear elliptic equations with critical Sobolev growth. *J. Funct. Anal.* **119**, 298–318 (1994)
33. E. Hebey, M. Vaugon, Meilleures constantes dans le théorème d'inclusion de Sobolev. *Ann. Inst. Henri Poincaré, Analyse Non Linéaire* **13**, 57–93 (1996)
34. E. Hebey, M. Vaugon, Sobolev spaces in the presence of symmetries. *J. Math. Pures Appl.* **76**, 859–881 (1997)
35. S. Kamin, J. Vazquez, Fundamental solutions and asymptotic behaviour for the p-Laplacian Equation. *Revista Matematica Iberoamericana* **4**(2), 339–354 (1988)
36. B. Kawohl, J. Horák, On the geometry of the p-Laplacian operator. *Discrete Contin. Dyn. Syst. Ser. S* **10**, 799–813 (2017)

37. J. Kim, M. Zhu, Non-existence results about  $-\Delta u = u^{\frac{n+2}{n-2}}$  on non-starshaped domains. *J. Differ. Equations* **225**, 737–753 (2006)
38. S. Kobayashi, Transformation groups in differential geometry, in *Ergebnisse der Mathematik und ihrer Grenzgebiete*, vol. 70 (1972)
39. A. Krasnoselskii, *Topological Methods in the Theory of Nonlinear Integral Equations* (Macmillan, New York and London, 1964)
40. A. Kristály, H. Liseib, C. Vargab, Multiple solutions for p-Laplacian type equations. *Nonlinear Anal.* **68**, 1375–1381 (2008)
41. I. Kuzin, S. Pohozaev, Entire solutions of semilinear elliptic equations, in *Progress in Nonlinear Differential Equations and their Applications*, vol. 33 (Birkhauser, Basel, 1997)
42. N. Labropoulos, An alternative approach to critical PDEs. *Electron. J. Differ. Equations* **2017**(150), 1–22 (2017)
43. Y.Y. Li, Existence of many positive solutions of semilinear elliptic equations on annulus. *J. Diff. Eq.* **83**, 348–367 (1990)
44. P. Lindqvist, *Notes on the p-Laplacian Equation, University of Jyväskylä—Lecture notes*, 2nd edn. (2017)
45. C. Loewner, L. Nirenberg, Partial differential equations invariant under conformal and projective transformations, in *Contributions to Analysis* (Academic Press, New York, 1974), pp. 245–272
46. I. Ly, The first eigenvalue for the p-Laplacian operator. *J. Inequal. Pure Appl. Math.* **6**(3/91), 1–12 (2005)
47. R. Mazzeo, N. Smale, Conformally flat metrics of constant positive scalar curvature on subdomains of the sphere. *J. Differ. Geom.* **34**, 581–621 (1991)
48. L. Nirenberg, Variational and topological methods in nonlinear problems. *Bull. Am. Math. Soc.* **4**, 267–302 (1981)
49. J. Padiál, P. Takáč, L. Tello, An antimaximum principle for a degenerate parabolic problem. *Adv. Differ. Equations*, **15**(7–8), 601–648 (2010)
50. R.S. Palais, Morse theory on Hilbert manifolds. *Topology* **2**, 299–340 (1963)
51. R.S. Palais, Lusternik-Schnirelman theory on Banach manifolds. *Topology* **5**, 115–132 (1966)
52. D. Passaseo, Multiplicity of positive solutions of nonlinear elliptic equations with critical Sobolev exponent in some contractible domains. *Manuscripta Math.* **65**, 147–165 (1989)
53. K. Perera, R.P. Agarwal, D. O'Regan, Morse Theoretic Aspects of p-Laplacian Type Operators, in *Mathematical Surveys and Monographs*, vol. 161 (American Mathematical Society, Providence, 2010)
54. M. Pereira, R. Silva, Remarks on the p-Laplacian on thin domains. *Prog. Nonlin. Differ. Equ. Appl.* **86**, 389–403 (2015)
55. Y. Pinchover, K. Tintarev, On positive solutions of p-Laplacian-type equations in Analysis, in *Partial Differential Equations and Applications—The Vladimir Maz'ya Anniversary Volume. Operator Theory: Advances and Applications*, vol. 193, eds. by A. Cialdea et al. (Birkhäuser, Basel, 2009), pp. 245–268
56. S.I. Pohozaev, Eigenfunctions of the equation  $\Delta u + \lambda f(u) = 0$ . *Sov. Math. Dokl.* **6**, 1408–1411 (1965)
57. R. Schoen, Conformal deformation of a Riemannian metric to constant scalar curvature. *J. Diff. Geometry* **20**, 479–495 (1984)
58. R. Schoen, The existence of weak solutions with prescribed singular behavior for a conformally invariant scalar equation. *Commun. Pure Appl. Math.* **41**, 317–392 (1988)
59. J. Schwartz, *Nonlinear Functional Analysis* (Gordon and Breach, New York, 1968)
60. L. Shi, X. Chang, Multiple solutions to p-Laplacian problems with concave nonlinearities. *J. Math. Anal. Appl.* **363**, 155–160 (2010)
61. W. Smith, Morse theory without critical points. *Rocky Mountain J. Math.* **19**(2), 531–543 (1989)
62. M. Struwe, *Variational Methods Applications to Nonlinear Partial Differential Equations and Hamiltonian Systems* (Springer, Berlin, 2008)

63. G. Tarantello, Nodal solutions of semilinear elliptic equations with critical exponent. *Differ. Integral Equation* **5**, 25–42 (1992)
64. N.S. Trüdinger, On imbeddings into Orlicz spaces and some applications. *J. Math. Mech.* **17**(5), 473–483 (1967)
65. M. Vainberg, *The Variational Method for the Study of Non-linear Operators* (Holden-Day, San Francisco, 1964)
66. L. Wang, Infinitely many solutions for an elliptic problem with critical exponent in exterior domain. *J. Part. Differ. Equation* **23**(1), 80–104 (2010)
67. E.-M. Wang, Y. Zheng, Regularity of the first eigenvalue of the p-laplacian and Yamabe invariant along geometric flows. *Pacific J. Math.* **254**(1), 239–255 (2011)
68. H. Yamabe, On a deformation of Riemannian structures on compact manifolds. *Osaka. Math. J.* **12**, 21–37 (1960)
69. Z. Yang, D. Geng, H. Yan, Three solutions for singular p-Laplacian type equations. *Electron. J. Differ. Equations* **2008**(61), 1–12 (2008)
70. E. Zeidler, *Nonlinear Analysis and its Applications III: Variational Methods and Optimization* (Springer, Berlin, 1985)

# Financial Contagion in Interbank Networks: The Case of Erdős–Rényi Network Model



K. Loukaki, P. Boufounou, and J. Leventides

**Abstract** In this study, we extend the model developed in Leventides et al. (J Econ Behav Organ 158:500–525, 2019) to include a wide variety of network topologies and provide a better understanding of the relation between network structure, banks’ characteristics and interbank contagion. While the focus of this paper is on the various factors that affect interbank contagion such as bank capital ratios, leverage, interconnectedness and homogeneity across banks’ sizes, the model lacks flexibility as far as the variability of the networks links is concerned. In order to circumvent this problem, we introduce the Erdős–Rényi probabilistic network model in our study to provide a wider vicinity of scenarios concerning the network structure of the interbank system and study how homogeneity within the interbank network affects the propagation of financial distress from one institution to the other parts of the system through bilateral exposures.

## 1 Introduction

Meeting the SDGs has currently secured prior importance for both businesses’ and nations’ socio-economic-environmental transformation to achieving sustainable development. More than 10,000 companies around the world have already signed up to the principles of sustainable business behavior and an adequate number of special toolkits has been developed to assist them towards this transformation. As explicitly stated in [1] “Achieving the Global Goals would create a world that is comprehensively sustainable: socially fair; environmentally secure; economically prosperous; inclusive; and more predictable”. According to Oxford Analytica Foundation [2], “companies that see the business case – as well as the moral imperative – for achieving all the Global Goals will take a ‘Global Goals lens’ to every aspect of their business strategy to change the way they operate and put

---

K. Loukaki · P. Boufounou · J. Leventides (✉)

Department of Economics, National and Kapodistrian University of Athens, Athens, Greece  
e-mail: [ylevant@econ.uoa.gr](mailto:ylevant@econ.uoa.gr)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,  
Springer Optimization and Its Applications 173,  
[https://doi.org/10.1007/978-3-030-72563-1\\_13](https://doi.org/10.1007/978-3-030-72563-1_13)

277

more focus on inclusion". Accordingly, Weber [3] underlined that "the banking sector became aware of the opportunity to finance the change to more sustainable development instead of just focusing on risks for their lending business". Alves et al. [4] debated that "EU aid policy evolved over the last fifteen years in accordance to the notion of financial liberalization and to the importance of private initiative and combined with the SDGs increased the promotion of the financial sector as engine of growth and development in the developing countries' and that "the New International Financial architecture assigns new roles for developing nations in the global financial markets". Achieving the Agenda 2030 depends on aligning the entire global chain of the financial and the banking system with sustainability and long-term outcomes therefore delineating of interbank linkages network structure becomes of outmost importance.

Furthermore, in the wake of the aftermaths of global financial crisis of 2007–2009 and the European sovereign debt crisis, there is a lot of attention of systemic risk, interconnectedness and contagious effects. Thus, there is a critical need for a better understanding of the fragility of the financial systems, their inner interconnections, their interaction with real economy and the conditions that can drive them from stability to instability and complete breakdown. In recent years, both academics and regulators has started to study various architectures of the financial system in order to assess certain risks within the system that potentially lead to huge losses for the overall economy.

The global financial system can be represented as a large complex network in which banks, hedge funds and other financial institutions are interconnected to each other through various forms of financial linkages. For example, in the banking sector, banks can be interconnected through direct and indirect links. Direct interconnectedness arises from bilateral transactions; borrowing or lending relationships between banks. A default by one bank, for example, can impose distress on other entities that hold significant liabilities of the defaulting bank. Thus, the failure of a bank can jeopardize the ability of its creditors banks to meet their obligations to their interbank creditors which may lead to a domino effect. There are also indirect ways that banks can be interconnected, since they invest in common securities, namely portfolio overlap. If, for example, a bank holds identical assets with other banks the correlation between their portfolios can cause fire sales in the market during a crisis period depressing thus overall prices in the market, ultimately leading to downward spirals for asset sales and inducing significant losses for all the participants in the market. The complexity of the financial system led many academics to utilize the network theory to study the effects of the interconnectedness and network topology on financial stability. Studying the financial system as a network is one of the methods to investigate the emergence of systemic risk through the connections of banks. In such a network structure every node represents a bank, the connections between banks are represented by edges where edge's weight represents the magnitude of exposure between the two parties and edge directionality allows one to determine who is the creditor and who is the lender. A robust interbank market plays an important role on the stability of the financial system. Through the interbank market, banks which suffer a liquidity shortage can

borrow from banks with liquidity surpluses. The interbank market can stabilize the financial system, by redistributing the funds in a effective way among the banks but at the same time can make the system prone to contagion of financial trouble from one bank to another (linkages).

In this paper, we focus our attention on the direct contagion channel and aim to identify the main drivers that affect interbank contagion. The flourishing literature which ensued in recent years has developed both theoretical models and empirical applications aimed at addressing the various issues concerning systemic risk. Counterfactual simulations on data have been extensively employed to study interbank contagion under different scenarios related to the topology of the interbank network, the size of interbank exposures and the degree of heterogeneity and interconnectedness within the network. In our assessment of the various drivers that affect interbank contagion, we extend the model developed in Leventides et al. [5] to include a wide variety of network topologies and provide a better understanding of the relation between network structure, banks' characteristics and interbank contagion. While the focus of this paper is on the various factors that affect interbank contagion such as bank capital ratios, leverage, interconnectedness and homogeneity across banks' sizes, the model lacks flexibility as far as the variability of the networks links is concerned. In this effort, interbank exposure and capital equity among banks displayed a stochasticness and the ability to construct a wide range of scenarios regarding connective links among banks is limited.

The introduction of the Erdős–Rényi probabilistic network model provides us with a wider vicinity of scenarios concerning the network structure of the interbank system. Under this framework, we build up multiple scenarios of various network structures that include a satisfactory number of cases via Monte Carlo simulations. In every single network that we construct, we investigate the dynamics of cascading defaults from an initial random shock that hits the system. Erdős–Rényi random graph model is one of the earliest theoretical network models and was introduced in the early 1960s by the Hungarian mathematicians Paul Erdős and Alfréd Rényi. In this random graph, each possible link between any two nodes can occur with a certain independent and identical probability, the Erdős and Rényi probability.

The Erdős and Rényi random graph model is a model in which has been extensively applied for the study of contagion in financial networks, e.g. [23, 24]. However, a number of alternatives have been recently developed that differ in the probability law governing the distribution of links between nodes. Using the Erdős–Rényi network structure, the degree distribution or the connectivity among banks can vary with respect to the chosen probability  $p$ . Thus, each random network generated with the same parameters  $N, p$  looks slightly different. Not only the detailed wiring network graph changes between realizations, but so does the number of links. Random graphs or Erdős–Rényi graphs are useful for modeling, analysis, and solving of structural and algorithmic problems arising in mathematics, theoretical computer science, statistical mechanics, natural sciences, and even in social sciences. However, the utility of an Erdős–Rényi model lies mainly in its mathematical simplicity, not in its realism. Virtually, the comparison with real-

world networks indicates that the random network model does not capture the degree distribution of real networks but it provides a useful baseline for more complicated network models.

The remainder of the paper is organized as follows. The following section discusses briefly the recent literature that has addressed on the topic of interbank contagion. Various aspects of systemic risk and network structures that are either found in real-world data or used in some theoretical studies of interbank contagion are addressed before we introduce the model investigated in Section 3. In Section 4 we describe the variables, considered in our subsequent analysis, provide in full detail the computer experiments conducted and discuss our simulations results. Summary and concluding remarks are drawn in the final section.

## 2 Related Literature

According to Upper [11], the channels through which a shock spreads can be broken down into two groups: indirect and direct contagion channels. A direct contagion channel results from the direct interbank linkages among banks and can take effect when an idiosyncratic shock travels through the network. This shock can be due to the inability of some banks to meet their financial obligations or due to interbank exposures that are quite large relative to the lender's capital. The possibility of the occurrence and transmission of direct contagion depends mainly on the structure and size of the interbank market. On the other hand, indirect contagion is created by indirect linkages among banks such as identical assets, portfolio returns and overlapping portfolios. If, for example, a bank holds identical assets with other banks, the correlation between their portfolios can cause fire sales in the market during a crisis period, thus depressing overall prices in the market and inducing significant losses for all participants [12]. Distinguishing among the various contagion channels is crucial for understanding financial contagion and the mechanisms through which it spreads and evolves.

There are a number of recent studies that have dealt with the issue of interbank contagion. Memmel and Sachs [13] simulate interbank contagion effects for the German banking sector and find that bank capital ratios, the share of interbank assets in the system and the degree of equality in the distribution of interbank exposures are the most important determinants for financial stability. Georgescu [14] compares the contagion potential of accounting induced regulatory constraints to that of funding constraints in a bank network and concludes that the interplay between illiquidity and solvency can lead to bank failures which are manifested by the vulnerable funding structure of banks during a crisis. Tonzer [15] examines the relationship between cross-border bank linkages and financial stability and show that larger cross-border exposures increase bank risks, however, when bilateral interbank linkages exist there is a shift toward a more stable banking system. Fink et al. [16] model contagion in the German interbank market via the credit quality channel and propose a novel metric which estimates the potential regulatory capital loss to a banking system due to contagion via interbank loans. They show that

contagion effects can be reduced if banks alter their lending and borrowing habits in response to policy interventions.

Our analysis also relates to the role of heterogeneity in the structure of interbank networks and how this characteristic affect systemic risk. Iori et al. [6] use an Erdős–Rényi network model of 400 banks comprising the interbank market in which the lending and borrowing functions are endogenously generated. The authors find that the likelihood of contagion is lowered in case the interconnected institutions are homogeneous, i.e. they have similar characteristics such as size or investment opportunities and thus, no institution becomes significant for either borrowing or lending. The authors also suggest, in line with Allen and Gale [17], that as connectivity increases the system becomes more stable. In a related study, Caccioli et al. [18] study the role of heterogeneity in degree distributions (the number of incoming and outgoing links), balance sheet size and degree correlations between banks. They find that networks with heterogeneous degree distributions are shown to be more resilient to contagion triggered by the failure of a random bank, but more fragile with respect to contagion triggered by the failure of highly connected nodes. The authors also provide evidence that when the average degree of connectivity is low, the probability of contagion due to failure of highly connected banks is higher than that due to the failure of large banks. However, when the average degree of connectivity is high, the opposite holds. Since the second scenario seems to be more realistic (networks with high connectivity), having “too big to fail” banks is more effective in eliminating a shock. Ladley [19] develops a partial equilibrium model of a closed economy in which heterogeneous banks interact with borrowers and depositors through the interbank market. Banks in the model are subject to regulation and the aim of the model is to qualitatively show how regulation and network structure can constrain or enhance the risk of contagion. The results show that for high levels of connectivity the system is more stable when the shock is small, while the contagion effects are amplified in case of larger initial shocks. Chinazzi et al. [20] explore the interplay between heterogeneity, network structure and balance sheet composition in the transmission of contagion. They argue that heterogeneity in connectivity provides additional resiliency to the system when the initial default is random and also show that ‘too-connected-to-fail’ banks are more dangerous than ‘too-big-to-fail’ ones and should be the primary concern of policy makers since their failure can trigger systemic breakdowns. Amini et al. [10] focus on bank heterogeneity in terms of the number of banks included in the network and the magnitude of their interconnections with other banks. They conclude that the more heterogeneity is introduced, the less resilient the network becomes. Contrary to these findings, the study of Georg and Poschmann [21] finds no significant evidence that the heterogeneity of the financial system has a negative impact on financial stability.

Finally, as far as the structure of an interbank system is concerned, the most common network structures that are either found in real-world data or used in some theoretical studies of interbank contagion are the Erdős–Rényi random network structure, introduced in Erdős and Rényi (1960), the small-world structure,



introduced in Watts and Strogatz (1998) and the scale-free structure, introduced in Barabasi and Albert (1999).

The Erdős–Rényi network structure, which is applied in our study, can be obtained by connecting any two nodes with a fixed and independent probability  $p$ . Thus, in an Erdős–Rényi network structure the degree or the number of links of a node is  $p(n-1)$ . The expected degree distribution for such networks is Binomial, converging to Poisson for large  $n$ . The Erdős and Rényi (1960) random graph model is a model in which has been extensively applied for the study of contagion in financial networks, e.g. in the contributions from Iori et al. [6], Nier et al. [7], Gai and Kapadia [8], May and Arinaminpathy [9] and Amini et al. [10]. A number of alternatives models have been recently developed that differ in the probability law governing the distribution of links between nodes. Nier et al. [7] study the extent to which the resilience of an interbank network depends on a combination of variables characterizing the network topology, banks' characteristics in terms of net worth and interbank exposures, and market concentration. Using Monte Carlo simulation experiments in Erdős–Rényi random graphs, they find that the effect of the degree of connectivity is non-monotonic. Specifically, a small initial increase in connectivity increases the chance of contagion defaults. However, after a certain threshold value, connectivity improves the capacity of a banking system to withstand shocks. In addition, the authors find that the banking system is more resilient to contagious defaults if its banks are better capitalized and this effect is non-linear. Finally, the size of interbank liabilities tends to increase the risk of default cascades, even if banks hold capital against such exposures and more concentrated banking systems are shown to be prone to larger systemic risk. Gai and Kapadia [8] using a network model of a banking system study how the probability and potential impact of contagion is influenced by aggregate and idiosyncratic shocks, network structure and liquidity. The authors agree with Haldane (2009) concerning the “robust-yet-fragile” property that the financial system exhibit. Even when the probability of contagion is very low, its effects can have tremendous consequences to the financial system. Higher connectivity may reduce the probability of default when contagion has not started yet but it may also increase the probability of having large default cascades when contagion begins. May and Arinaminpathy [9] apply an Erdos–Renyi network structure of which they build on the models of Nier et al. [7] and Gai and Kapadia [8] and study the interplay between the characteristics of individual banks and the overall behavior of the network. The authors consider that banks interact through different asset classes and study contagion between those asset classes. May and Arinaminpathy [9] find that increasing the level of connectivity is beneficial only when the initial shock has been caused by a default on interbank loans. However, by contrast, the opposite holds in case of liquidity shocks since they do not experience attenuation and for a given asset class, they tend to grow as more and more banks hold the failing asset. Finally, the authors emphasize the importance of having large capital buffers that will make for greater robustness both of individual banks and of the system as a whole. Finally, Amini et al. [10] test the impact of heterogeneity in an interbank network structure and the relation between resilience and connectivity using three different network models; a scale-

free network with equal and heterogeneous weights and an Erdős–Rényi network with equal weights. The main result of this study is that the most heterogeneity is introduced, the least the resilience of the network.

### 3 Erdős–Rényi Random Graph Model

The random graph model which is one of the earliest theoretical network models was introduced by Erdős and Rényi (1960). In this random graph, each possible link between any two nodes can occur with a certain independent and identical probability,  $p$ . This model is typically denoted  $G(n, p)$  and has two parameters:  $n$  the number of vertices and  $p$ , the probability that each simple edge  $(i, j)$  exists, which is constant for each pair nodes.

The adjacency matrix of a random graph is given by

$$\forall i > j, A_{ij} = A_{ji} = \begin{cases} 1, & \text{edge } (i, j) \text{ exists; prob } (p) \\ 0, & \text{edge } (i, j) \text{ does not exist; prob } (1 - p) \end{cases}$$

In other words, each edge is included in the graph with probability  $p$ , independent from every other edge. The probability to create randomly a graph with  $n$  nodes and  $m$  edges is given by  $p^m(1 - p)^{\binom{n}{2} - m}$ . Furthermore, the probability  $p$  serves as the parameter of our model and as  $p$  increases, the graph is more likely to have more edges.

The restriction of  $i > j$  appears because edges are undirected or to put it differently, the adjacency matrix is symmetric across the diagonal, and there are no self loops. In the network there are  $n(n - 1)$  possible links to be created, resulting in an expected number of edges in the network equal to  $pn(n - 1)$ , so that the (expected) average degree is  $p(n - 1)$ . Thus, the degree distribution of such a graph is given by

$$p(k) = \binom{n - 1}{k} p^k (1 - p)^{n - 1 - k} \tag{1}$$

The mean degree,  $c$ , in the  $G(n,p)$  graph model is given by

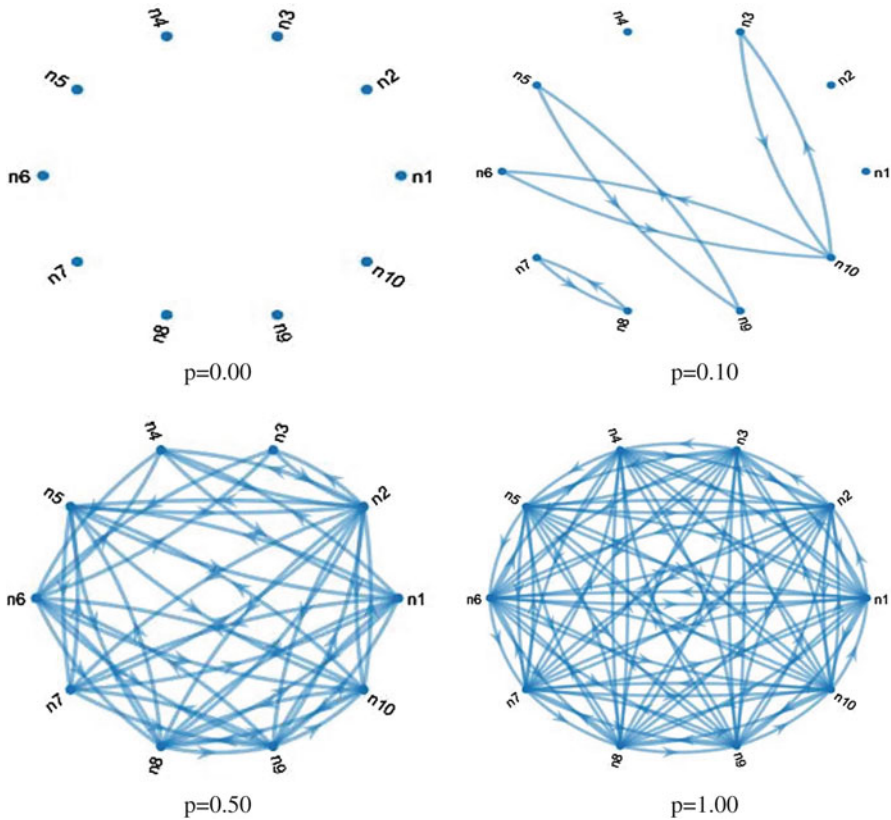
$$c = (n - 1) p \tag{2}$$

In other words, each vertex has  $(n-1)$  possible partners and each of these exist with the same independent probability  $p$ . Asymptotically, as  $n \rightarrow \infty$ , the degree distribution of a random graph converges to a Poisson ( $c$ ) distribution

$$p(k) = \frac{e^{-c} c^k}{k!} \tag{3}$$

Due to the above property, the Erdős–Rényi random graph model is sometimes referred as Poisson random graph or random graph. The Erdős and Rényi (1960) graph model results in networks with small diameters and short average path lengths, capturing very well the “small-world” property, observed in many real networks. The clustering coefficient of an Erdős–Rényi graph model is equal to the probability of an edge’s existence between two nodes,  $p$ . The Erdős and Rényi (1960) random graph model is a model in which has been extensively applied for the study of contagion in financial networks, e.g. in the contributions from Iori et al. [6], Nier et al. [7], Gai and Kapadia [8] and Montagna and Kok (2013).

In an Erdős –Rényi model we begin with  $n$  isolated nodes as presented in the first snapshot in Figure 1. Then, with probability  $p > 0$  each pair of nodes is connected by a link. Therefore, in this model the network is determined only by the number

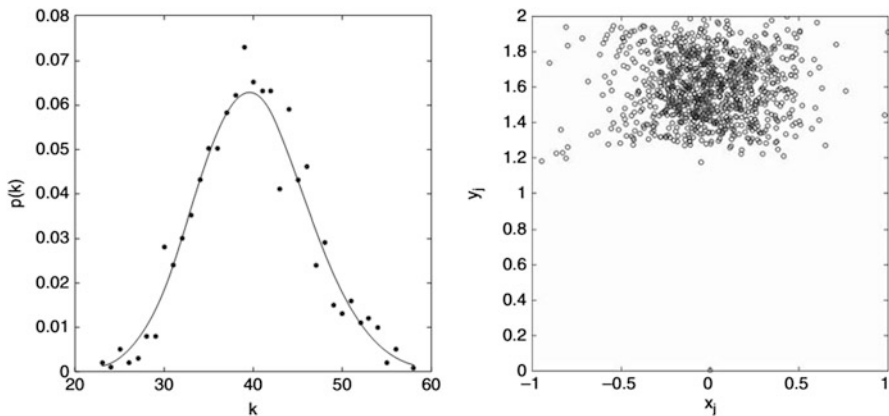


**Fig. 1** Erdős–Rényi random networks: Erdős–Rényi random networks with ten nodes and different probabilities of connecting a pair of nodes

of nodes,  $n$ , and edges,  $m$ , and usually an Erdős–Rényi random graph is written as  $G(n, m)$  or  $G(n, p)$ . In Figure 1 we present some examples of Erdős–Rényi random graphs with the same number of nodes and different linking probabilities. It is easy to understand that if we repeat the process for the same number of nodes and the same probability, we will not necessarily get the same network.

However, a number of alternatives models have been recently developed that differ in the probability law governing the distribution of links between nodes. Since, the Erdős–Rényi probability,  $p$ , is assumed to be equal and constant across all pairs of nodes, the resulting network structure does not present marked heterogeneity. Thus, modeling interbank networks using the Erdős–Rényi structure fails to mimic the heterogeneity observed in real interbank network systems.

In order to fully understand the heterogeneity of an Erdős–Rényi random network, we now consider one particular random realization of an Erdős–Rényi random network with 1000 nodes and  $p = 0.04$ , that is  $G(n = 1000, p = 0.04)$  and plot the probability  $p(k)$  of finding a node of degree  $k$ , versus the degree, we obtain Figure 2, where it can be seen that the maximum of the distribution is about the value  $k = (n - 1)p = 39$ . Obviously, the probability  $p(k)$  follows a binomial distribution of the form represented in Equation (1). As we explained above, for large values of  $n$ , the degree distribution of a random graph converges to a Poisson ( $c$ ) distribution. Figure 2 displays the heterogeneity plot for  $G(1000, 0.04)$ , where two characteristic features of the Erdős–Rényi networks are observed. The first is a typical dispersion of the points around the value  $x = 0$ , and the second is the very small value of  $\rho(G)$ , which in this case is 0.0066.



**Fig. 2** Heterogeneity of Erdős–Rényi random networks. A typical Poisson degree distribution of an Erdős–Rényi random network with 1000 nodes and  $p = 0.04$  (left), and the characteristic heterogeneity plot for the same network. (Source: Estrada [22])

### 3.1 The Mathematical Description of the Contagion Model

In this section we study the case of an Erdős–Rényi network model in which, as we stated earlier, all nodes have the same probability of being connected to another node in the network. Our model is tailored to simulate default cascades triggered by an exogenous shock in an interbank network as in Leventides et al. [5]. We first introduce the interbank network model, describe the default cascades initiated by a random negative shock on this network and analyze the parameters that affect interbank contagion.

### 3.2 The Interbank Network

As in Leventides et al. [5], we assume that the banking system contains  $i = 1, \dots, N$  banks. Every bank has its own balance sheet and the accounting equation holds at all times. Total assets are divided in three categories: interbank assets  $A_i^{IB}$ , other assets  $A_i^{OT}$  and cash reserves  $C_i$ . On the liabilities side of the balance sheet we have included: interbank liabilities  $L_i^{IB}$ , other liabilities  $L_i^{OT}$  and equity capital  $E_i$ . A schematic overview of the balance sheet is given in Table 1. Although the proposed balance sheet structure does not capture all elements of a bank balance sheet, it includes all those positions that are relevant to our study.

We introduce a standard notation for our model and we define a simple interbank network as  $G = (V, E)$ , where  $V$  represents the nodes of the graph while  $E$  represents the edges. We further consider  $A$ , the adjacency matrix of the graph, defined as

$$\forall i > j, A_{ij} = A_{ji} = \begin{cases} 1, & \text{edge } (i, j) \text{ exists} \\ 0, & \text{edge } (i, j) \text{ does not exist} \end{cases}$$

The  $u$ th row or column of  $A$  has  $k_u$  entries, where  $k_u$  is the *degree* of the node  $u$ , which is simply the number of nearest neighbours that  $u$  has. Denoting by  $\mathbf{1}$  a  $|V| \times 1$  vector, the column vector of node degrees  $\kappa$  is given by

**Table 1** Stylized Balance sheet structure

| Assets $A_i$                    | Liabilities $L_i$                    |
|---------------------------------|--------------------------------------|
| Interbank Assets ( $A_i^{IB}$ ) | Interbank Liabilities ( $L_i^{IB}$ ) |
| Other Assets ( $A_i^{OT}$ )     | Other Liabilities ( $L_i^{OT}$ )     |
| Cash ( $C_i$ )                  | Equity Capital ( $E_i$ )             |

The table presents a stylized balance sheet structure in the interbank network. Total assets are divided in three categories: Interbank assets ( $A_i^{IB}$ ), other assets ( $A_i^{OT}$ ), and cash reserves ( $C_i$ ). Total liabilities include: Interbank liabilities ( $L_i^{IB}$ ), other liabilities ( $L_i^{OT}$ ), and equity capital ( $E_i$ ). It is assumed that the accounting equation holds at all times

$$\kappa = \left(1^T A\right)^T = A^T 1 \tag{4}$$

We define the *indegree* as the number of links pointing toward a given node, and the *outdegree* as the number of links departing from the corresponding node. Specifically:

$$\kappa^{in} = \left(1^T A\right)^T = A^T 1 \tag{5}$$

$$\kappa^{out} = A 1 \tag{6}$$

Thus, our interbank network of credit exposures between  $n$  banks can be visualized by a graph  $G = (V, E)$  where  $V$  represents the set of financial institutions—nodes, and  $E$  is the set of the edges linking the banks, that is, the set of ordered couples  $(i, j) \in V \times V$  indicating the presence of a loan made by bank  $i$  to bank  $j$ . The number of nodes defines the size of the interbank network. Every edge  $(i, j)$  is weighted by the face value of the interbank claim and the representation of interbank claims is made by a single weighted  $N \times N$  matrix  $X$ :

$$X = \begin{bmatrix} 0 & \cdots & x_{1j} & \cdots & x_{1N} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{i1} & \cdots & 0 & \cdots & x_{iN} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{N1} & \cdots & x_{Nj} & \cdots & 0 \end{bmatrix}$$

where  $x_{ij}$  is the credit exposure of bank  $i$  vis-à-vis bank  $j$  and  $N$  is the number of banks in the network. Interbank assets are represented along the rows while columns represent interbank liabilities. Once  $X$  is in place, the interbank entries of each bank are given according to the following rules:

- (i)  $A_i = \sum_{j=1}^N x_{ij}$  (horizontal summation), where  $A_i$  is the total interbank assets of bank  $i$ .
- (ii)  $L_j = \sum_{i=1}^N x_{ij}$  (vertical summation), where  $L_j$  is the summation of the total interbank liabilities of bank  $j$ .

One can observe that the diagonal line contains zeros due to the fact that banks do not lend to themselves. In this framework, a random network is generated based on two parameters, the size of the network (number of nodes/banks) and the probability  $p_{ij}$  that there is a lending/borrowing link between two nodes/banks. Thus, each possible link between two nodes exists with an independent and identical probability, which is often called the Erdős–Rényi probability.

Although, we have undirected edges in this framework, we cannot really speak of undirected links, since the two directions of the same link are given different weights.

### 3.3 Shock Propagation and Contagion Dynamics

The failure of a bank can affect other banks through their interbank connections. Below, we describe the mechanism through which an initial shock affecting a bank propagates onto its counterparties along the network. Contrary to the recent literature, the term contagion here translates into total capital losses due to multiple default cascades. The cascade dynamics we use in this study are straightforward to implement and enable us to run a great number of simulations on a variety of different scenarios (Table 2).

The default procedure starts with an exogenous shock being simulated, typically by setting to zero the equity of one randomly chosen bank  $i$  and the cascade of defaults proceeds on a timestep-by-timestep basis, assuming zero recovery for shock transmissions. The zero recovery assumption, which is a realistic one in the short run, is often used in the literature to analyze worst case scenarios and

**Table 2** OLS regression analysis for Scenario 1 (Heterogeneous banks with homogeneous exposures)

|                         | N = 20                 | N = 50                 | N = 80                 | N = 100                |
|-------------------------|------------------------|------------------------|------------------------|------------------------|
| CATIN 1                 | 0.051<br>(16.198)***   | -0.002<br>(-0.459)     | -0.001<br>(-0.347)     | -0.007<br>(-2.044)**   |
| CATIN2                  | 0.098<br>(4.195)***    | 0.004<br>(0.170)       | 0.179<br>(8.073)***    | 0.104<br>(5.059)***    |
| LEVIN                   | 0.389<br>(17.018)***   | 0.413<br>(19.043)***   | 0.260<br>(12.205)***   | 0.315<br>(15.935)***   |
| NOUTGOING               | -0.080<br>(-3.915)***  | 0.097<br>(2.773)***    | -0.170<br>(-4.933)***  | -0.053<br>(-1.534)     |
| COUNT                   | 0.602<br>(138.571)***  | 0.572<br>(134.093)***  | 0.576<br>(136.735)***  | 0.540<br>(124.326)***  |
| VARCAP                  | -0.088<br>(-53.348)*** | -0.075<br>(-61.005)*** | -0.053<br>(-53.890)*** | -0.054<br>(-57.165)*** |
| P                       | -0.101<br>(-5.089)***  | -0.080<br>(-2.338)**   | 0.165<br>(4.885)***    | 0.107<br>(3.148)***    |
| Adjusted R <sup>2</sup> | 0.800                  | 0.763                  | 0.756                  | 0.749                  |

The table presents the regression results for Scenario 1. The dependent variable is CATEND measured as the total loss of capital due to contagion as percentage of total capital in the network. Explanatory variables are, CATIN1, CATIN2, LEVIN, NOUTGOING, COUNT, VARCAP and P, the probability for a link to exist between two nodes. Each cell displays the OLS standardized coefficients along with the corresponding  $t$ -statistics (shown in parentheses). The sample comprises of 18,000 realizations (simulated banking crises). \*, \*\* and \*\*\* denote significance at the 10, 5 and 1 percent level, respectively

refers to a situation where creditor banks lose all of their interbank assets held against a defaulting bank [8, 20]. A bank’s default implies that it is no longer able to meet its interbank liabilities to its counterparties. Since these liabilities constitute other banks’ assets, the banks that get into trouble affect simultaneously their counterparties, leading to write-downs in their balance sheets. The interbank asset loss due to failure of bank  $i$  is subtracted from the bank’s  $j$  capital. Bank  $j$  will fail if its exposure against bank  $i$  exceeds its equity. A second round of bank failure occurs if bank creditors cannot withstand the losses realized due to its default and eventually, contagion stops if no additional bank goes bankrupt, otherwise a third round of contagion takes place. An initial shock can be amplified through banks’ interconnections and further transmitted to other institutions, such that the overall effect on the system goes largely beyond the original shock. As Upper and Worms (2004) demonstrate, in response to a liquidity shock banks prefer to withdraw their deposits at other banks instead of liquidating their long-term assets, creating further instability and liquidity dry-ups in the financial system.

A general mathematical description of the dynamical system expressing the shock propagation mechanism is presented hereafter. We consider a network consisting of  $N$  banks numbered from 1 to  $N$ . We define  $b_i$  as the capital possessed by bank  $i$  in the network and

$$b_0 = (b_1, b_2, \dots, b_N) \tag{7}$$

stands for the initial vector of bank capital.  $X$  is defined as a  $N \times N$  matrix with entries:

$x_{ij}$  = the credit exposure of bank  $i$  vis-à-vis bank  $j$  in the network

$$x_{ii} = b_i \tag{8}$$

We consider the case where some of the banks (one or more) collapse. We wish to study how the crisis travels through the bank network and when exactly it comes to a fixed point. The collapse of banks  $i_1, i_2, \dots, i_k$  (where  $k \leq N$ ), can be described in the following way. Consider the element  $x_0 \in Z_N^2 = \{0, 1\}^N$  which has zero entries everywhere except the positions  $i_1, i_2, \dots, i_k$  where  $x_0$  takes on the value 1. Then,

$$b_1 = b_0 - X \cdot x_0 \tag{9}$$

is the new vector of capital of the  $N$  banks. We now take

$$x_1(i) = \begin{cases} 1, & b_1(i) \leq 0; \\ 0, & b_1(i) > 0. \end{cases} \tag{10}$$



Then  $x_1 \in Z_2^N$  and  $x_1$  indicates the banks that have collapsed after the bankruptcy of the first  $k$  banks. The vector  $x_1$  takes on the value 1 in the positions  $i_1, i_2, \dots, i_k$ . If  $x_1 \neq x_0$ , this indicates that the collapse of the first  $k$  banks has adversely affected other banks leading them to bankruptcy. Similarly, from  $x_1$  we take:

$$b_1 = b_0 - X \cdot x_1 \tag{11}$$

and then

$$x_2(i) = \begin{cases} 1, & b_2(i) \leq 0; \\ 0, & b_2(i) > 0. \end{cases} \tag{12}$$

The vector  $x_2$  indicates the banks that collapse after the bankruptcy of the banks of  $x_1$ . Therefore, we have a map:

$$F : Z_2^N \rightarrow Z_2^N \tag{13}$$

$$x \rightarrow F(x) = f(b_0 - X \cdot x) \tag{14}$$

The map  $F(x)$  defines a dynamical system  $x_{n+1} = F(x_n)$  which describes the evolution of contagion in the interbank network.

### 3.4 Monte Carlo simulations

In this section we apply Monte Carlo simulations in four different stages. As in Leventides et al. [5], we introduce randomness in three areas: amount of capital, interbank claims and network structure. The stochasticness introduced in our model provides us with a wide vicinity of scenarios that may come across in real world. Using the Erdős–Rényi network structure, the degree distribution or the connectivity among banks can vary with respect to the chosen probability  $p$ . Thus, each random network generated with the same parameters  $N, p$  looks slightly different.

The second stage involves estimating the parameters of interest, i.e. the value of the coefficients in the regression model. In the third stage the test statistics of interest are saved, while in the fourth stage we go back to the first stage and repeat  $N$  times. The quantity  $N$  is the number of replications which should be as large as is feasible. As Monte Carlo is based on random sampling from a given distribution (with results equal to their analytical counterparts asymptotically), setting a small number of replications will yield results that are sensitive to odd combinations of random number draws. Generally speaking, the sampling variation is measured by the standard error estimate, denoted  $S_x = \sqrt{\text{var}(x)/N}$ , where  $x$  denotes the value of

the parameter of interest and  $\text{var}(x)$  is the variance of the estimates of the quantity of interest over the  $N$  replications.

Similar to Leventides et al. [5], we consider four different scenarios, in line with Chinazzi et al. [20], where we let vary the balance sheet composition, the size of the network and the link probability among banks which is held constant for each pair of nodes. The four scenarios tested are as follows:

|                    |   |
|--------------------|---|
| <i>Scenario 1:</i> | <ul style="list-style-type: none"> <li>• <b>Heterogeneous banks with homogeneous exposures.</b> In this scenario, we construct interbank networks where banks have different equity size and their interbank claims are evenly distributed across the outgoing links</li> </ul> |
| <i>Scenario 2:</i> | <ul style="list-style-type: none"> <li>• <b>Heterogeneous banks with heterogeneous exposures.</b> In this scenario, the interbank networks allow for heterogeneous bank sizes and heterogeneous interbank claims among banks.</li> </ul>  |
| <i>Scenario 3:</i> | <ul style="list-style-type: none"> <li>• <b>Homogeneous banks with heterogeneous exposures.</b> In this scenario, we construct interbank networks where banks have the same equity size and unevenly distribute their exposures across creditor banks</li> </ul>                |
| <i>Scenario 4:</i> | <ul style="list-style-type: none"> <li>• <b>Homogeneous banks with homogeneous exposures.</b> In this last scenario, we construct interbank networks where banks have the same equity size and interbank claims are evenly distributed across creditor banks</li> </ul>         |

In each case, we do not control the number of outgoing links as in Leventides et al. [5] but for each network that is generated a random probability, which is constant for each pair of nodes, defines the lending/borrowing relation of each bank. The probability  $p_{ij}$  is assumed to be equal and constant across all pairs  $(i,j)$ . For simplicity, we denote the probability, termed as the Erdős–Rényi probability, by  $p$ . Since the probability of forming a link is homogeneous, the resulting network structure does not present marked heterogeneity.

We examine banking systems consisting of small banks with low, medium and large interbank exposures, as well as systems of large banks with corresponding exposure levels. We consider a basic model that uses only two components from a bank’s balance sheet, that is, equity and interbank loans—in the words of May and Arinaminpathy [9] ‘*a caricature for banking ecosystems*’. We generate our model in two separate steps. First, we construct a model structure of  $N$  nodes representing the banks in our system and randomly choose the probability  $p$  of forming a link between each of the  $\binom{N}{2}$  possible links.

For all the possible couples of nodes, a link is created with probability  $p$  which represent lending/borrowing relationship, while in a second step, we assign each node to a stylized balance sheet structure. Once the banking networks are created, the default propagation dynamics are implemented to examine the effects of an idiosyncratic shock hitting one bank. The effect of a shock is simulated, typically by setting to zero the equity of the affected bank. We estimate the initial loss of capital by letting the first bank default and subsequently record the loss as percentage of the total capital in the system. Consequently, the defaulted bank will be unable to

repay its creditors and the interbank loans that were granted will be written-off, as we have selected to work under a zero recovery assumption. This bad debt will be recorded and expressed as percentage of the total capital in the system. Moreover, the creditors of the defaulted bank will experience a shock on their balance sheets and the recorded losses will be subtracted from their equity.

If at any time the total losses realized by a bank exceed its net worth, the bank is deemed in default and is removed from the network. Note that time steps are modeled as being discrete and there is the possibility that many banks default simultaneously in each timestep. These shocks propagate to their creditors and take effect in the next timestep. When no further failures are observed, the default procedure terminates and various contagion indicators<sup>1</sup> are calculated based on the contagion map as described in Subsection 3.3.

## 4 Main Findings

This section discusses the main findings of this study. Subsection 4.1 describes in full detail the computer experiments conducted while Subsection 4.2 discusses the simulation results of all four scenarios considered.

### 4.1 Computer Experiments

Having generated banking systems via an Erdős–Rényi network structure framework and balance sheet allocation, the dynamics of an initial shock affecting a bank within the interbank network can be investigated. Given the complexity of the interbank network outlined above, it is extremely difficult to derive analytical solutions. In order to obtain data to describe the variables that affect contagion, we employ several Monte Carlo simulations. In each realization, we construct an interbank network with  $N \in [20, 50, 80, 100]$  nodes under the rewiring process of the Erdős–Rényi methodology. In a second step, we test the four scenarios mentioned before by varying the equity size of banks and the interbank exposure structure across creditor banks. For each scenario tested we check a wide range of link probabilities, such that we can observe dense or sparse interbank network systems. Since the probability of forming a link is homogeneous, the resulting network structure does not present marked heterogeneity.

When homogeneity across bank sizes is considered, all banks are assumed to have the same equity size and thus, each bank is endowed with a balance sheet that

---

<sup>1</sup>We refer the interested reader to Appendix in Leventides et al. [5] for a formalization of the aforementioned mechanism in a pseudocode which simulates the default cascade in the interbank network.

consists of 100 units of equity. On the other hand, when homogeneity is present with respect to interbank exposures, interbank claims are randomly allocated within the interbank network and are categorized as follows: small loans granted (4 units), medium loans (8 units) and large loans (14 units). With respect to scenarios tested where heterogeneity of bank size is introduced, the amount of equity of each bank is drawn from a uniform distribution in the range:  $b_i \in [0, 100]$ , whereas when heterogeneity is introduced with respect to interbank claims, credit is allocated in the following ranges:  $a_{ij} \in [0, 4]$ ,  $a_{ij} \in [0, 8]$ ,  $a_{ij} \in [0, 14]$ .<sup>2</sup> Interbank exposures are set 60% lower than these in Leventides et al. [5]. This is due to the fact that we cannot control the connectivity across banks since the link probability is randomly selected. The interbank exposure decrease was set by trial and error in order not to observe enormous high leveraged systems. In addition, we control the leverage of the system by setting the rule that the maximum leverage ratio of each network system cannot exceed five. Then, balance sheets are assigned to each node, consistent with each specific scenario tested. We randomly choose a single bank in the system to default due to an exogenous shock and the default cascades proceed sequentially, assuming zero recovery. When no further failures are observed results are recorded before another realization begins. We then impose another shock on the second bank in the network and this procedure continues until all banks in the interbank network are hit by an exogenous shock.

For each scenario tested and for each network size we have three cases in which we allow the weight of outgoing links (small, medium and large interbank claims) to vary among banks. Each case gives us 6000 realizations or, to put it differently, 6000 banking crises. We deem that 6000 realizations provide a satisfactory number of runs and robustness to our analysis. Thus, for each scenario tested and each network size we employ  $6000 \times 3 = 18,000$  realizations using the following variables in each realization:

- Total loss of capital due to contagion as percentage of total capital in the system (**CATEND**)
- Initial loss of capital by defaulting bank  $i$  as percentage of total capital in the system (**CATINI**), i.e. bank's  $i$  depleted equity divided by the total equity in the network
- Loss of capital at the first stage (interbank loans that cannot be repaid) by defaulting bank  $i$  as percentage of total capital in the system (**CATIN2**), i.e. total amount of loans granted to bank  $i$  that cannot be repaid divided by the total equity in the network
- Leverage of the interbank network (**LEVIN**), i.e. total interbank exposures as measured by the sum of matrix's  $A$  elements, divided by the total capital in the network

---

<sup>2</sup>Although those ranges have been selected arbitrarily, they are not sensitive to any regression model employed in the following analysis and thus, our regression results will be unaffected from a qualitative point of view if different ranges are used.

- Number of outgoing links of bank  $i$  (**NOUTGOING**), i.e. the outdegree of a bank  $i$  which corresponds to the number of creditors in the network. It is defined as the summation of the  $i$ th column of the adjacency matrix  $A$ .
- Shock propagation variable (**COUNT**) which measures the number of rounds needed until no further bank defaults
- Variance of capital (equity) (**VARCAP**) used in those scenarios tested where only heterogeneous bank sizes are considered
- Variance of interbank loans (**VARLOANS**) used in those scenarios tested where only heterogeneous interbank loan exposures are considered
- Erdős–Rényi probability  $p_{ij}$  (**p**) that there is a lending/borrowing link between two nodes/banks.

Our selection of variables is motivated by economic intuition and by the findings of previous studies on the dynamics of systemic risks [7] and Leventides et al. [5]. Table 3 presents summary statistics on these variables. In order to study the effect the aforementioned variables have on contagion risk, we estimate the following ordinary least squares (OLS) models:

$$CATEND = \beta_1 CATIN1 + \beta_2 CATIN2 + \beta_3 LEVIN + \beta_4 NOUTGOING + \beta_5 COUNT + \beta_6 VARCAP + \beta_7 p \quad (15)$$

$$CATEND = \beta_1 CATIN1 + \beta_2 CATIN2 + \beta_3 LEVIN + \beta_4 NOUTGOING + \beta_5 COUNT + \beta_6 VARCAP + \beta_7 VARLOANS + \beta_8 p \quad (16)$$

$$CATEND = \beta_1 CATIN2 + \beta_2 LEVIN + \beta_3 NOUTGOING + \beta_4 COUNT + \beta_5 VARLOANS + \beta_6 p \quad (17)$$

$$CATEND = \beta_1 CATIN2 + \beta_2 LEVIN + \beta_3 NOUTGOING + \beta_4 COUNT + \beta_5 p \quad (18)$$

The model described in Equation (15) is applied to scenarios involving heterogeneous bank sizes with homogeneous exposures in the network structure, Equation (15) refers to a situation where emphasis is placed on heterogeneous interbank loan exposures combined with heterogeneous bank sizes, Equation (17) takes into account homogeneous banks with heterogeneous exposures while Equation (18) considers only homogeneous bank sizes and interbank claims. The variable  $CATIN1$  has been omitted from Eqs. (17) and (18) due to the fact that banks in the interbank system are homogeneous, i.e. we keep constant the equity of each bank and thus  $CATIN1$  remains stable during our simulation runs. There is an explanation in the next subsection concerning the fact that in our experiments we have selected to

**Table 3** Summary statistics

| Variable            | Heterogeneous banks-homogeneous exposures |          |           | Heterogeneous banks-heterogeneous exposures |         |           | Homogeneous banks-heterogeneous exposures |        |           | Homogeneous banks-homogeneous exposures |        |           |
|---------------------|---|----------|-----------|---|---------|-----------|---|--------|-----------|---|--------|-----------|
|                     | Mean                                      | Median   | Std. Dev. | Mean  | Median  | Std. Dev. | Mean                                      | Median | Std. Dev. | Mean                                    | Median | Std. Dev. |
| <b>n = 20 banks</b> |   |          |           |   |         |           |   |        |           |   |        |           |
| CATEND              | 0.241                                     | 0.067    | 0.368     | 0.185                                       | 0.063   | 0.313     | 0.220                                     | 0.050  | 0.350     | 0.251                                   | 0.050  | 0.367     |
| CATIN1              | 0.050                                     | 0.050    | 0.028     | 0.050                                       | 0.050   | 0.029     | 0.050                                     | 0.050  | 0.000     | 0.050                                   | 0.050  | 0.000     |
| CATIN2              | 0.079                                     | 0.063    | 0.062     | 0.070                                       | 0.052   | 0.067     | 0.147                                     | 0.130  | 0.108     | 0.175                                   | 0.160  | 0.124     |
| COUNT               | 2.624                                     | 1.000    | 2.342     | 2.278                                       | 1.000   | 1.999     | 1.916                                     | 1.000  | 1.893     | 2.029                                   | 1.000  | 1.876     |
| LEVIN               | 1.572                                     | 1.273    | 1.193     | 1.406                                       | 1.081   | 1.241     | 2.935                                     | 2.679  | 1.888     | 3.510                                   | 3.510  | 2.046     |
| P                   | 0.491                                     | 0.486    | 0.286     | 0.459                                       | 0.438   | 0.278     | 0.402                                     | 0.371  | 0.253     | 0.276                                   | 0.214  | 0.226     |
| NOUTGOING           | 9.303                                     | 9.000    | 5.729     | 8.701                                       | 8.000   | 5.595     | 7.601                                     | 7.000  | 5.139     | 5.228                                   | 4.000  | 4.603     |
| VARCAP              | 829.371                                   | 824.415  | 171.611   | 836.673                                     | 832.512 | 172.930   | —   | —      | —         | —                                       | —      | —         |
| VARLOANS            | —   | —        | —         | 38.407                                      | 8.230   | 45.561    | 47.449                                    | 33.195 | 40.430    | —                                       | —      | —         |
| <b>n = 50 banks</b> |   |          |           |   |         |           |   |        |           |   |        |           |
| CATEND              | 0.364                                     | 0.034    | 0.459     | 0.352                                       | 0.033   | 0.453     | 0.177                                     | 0.020  | 0.348     | 0.215                                   | 0.020  | 0.367     |
| CATIN1              | 0.020                                     | 0.020    | 0.011     | 0.020                                       | 0.020   | 0.011     | 0.020                                     | 0.020  | 0.000     | 0.020                                   | 0.020  | 0.000     |
| CATIN2              | 0.047                                     | 0.044963 | 0.029     | 0.046                                       | 0.045   | 0.030     | 0.068                                     | 0.063  | 0.047     | 0.068                                   | 0.064  | 0.051     |
| COUNT               | 4.078                                     | 2.000    | 3.473     | 4.026                                       | 2.000   | 3.532     | 2.209                                     | 1.000  | 2.618     | 2.451                                   | 1.000  | 2.767     |
| LEVIN               | 2.330                                     | 2.320    | 1.365     | 2.310                                       | 2.336   | 1.347     | 3.380                                     | 3.392  | 1.996     | 3.401                                   | 3.240  | 2.002     |
| NOUTGOING           | 16.754                                    | 14.000   | 12.604    | 18.478                                      | 14.000  | 14.272    | 10.092                                    | 8.000  | 9.003     | 5.025                                   | 4.000  | 4.562     |
| P                   | 0.342                                     | 0.290    | 0.250     | 0.377                                       | 0.270   | 0.286     | 0.206                                     | 0.159  | 0.176     | 0.103                                   | 0.079  | 0.084     |
| VARCAP              | 834.415                                   | 832.740  | 106.833   | 830.829                                     | 829.254 | 115.273   | —   | —      | —         | —                                       | —      | —         |
| VARLOANS            | —   | —        | —         | 38.207                                      | 8.299   | 44.715    | 47.688                                    | 33.015 | 39.929    | —                                       | —      | —         |
| <b>n = 80 banks</b> |   |          |           |   |         |           |   |        |           |   |        |           |
| CATEND              | 0.383                                     | 0.022    | 0.469     | 0.359                                       | 0.021   | 0.460     | 0.180                                     | 0.012  | 0.362     | 0.198                                   | 0.012  | 0.3676    |
| CATIN1              | 0.012                                     | 0.012    | 0.007     | 0.012                                       | 0.012   | 0.007     | 0.012                                     | 0.012  | 0.000     | 0.012                                   | 0.012  | 0.000     |
| CATIN2              | 0.031                                     | 0.030    | 0.020     | 0.029                                       | 0.028   | 0.020     | 0.045                                     | 0.042  | 0.031     | 0.040                                   | 0.035  | 0.032     |

(continued)

**Table 3** (continued)

| Variable  | Heterogeneous banks-homogeneous exposures |         |           | Heterogeneous banks-heterogeneous exposures |          |           | Homogeneous banks-heterogeneous exposures |        |           | Homogeneous banks-homogeneous exposures |        |           |
|-----------|---|---------|-----------|---|----------|-----------|---|--------|-----------|---|--------|-----------|
|           | Mean                                      | Median  | Std. Dev. | Mean  | Median   | Std. Dev. | Mean                                      | Median | Std. Dev. | Mean                                    | Median | Std. Dev. |
| COUNT     | 4.638                                     | 3.000   | 4.259     | 4.662                                       | 2.000    | 4.442     | 2.359                                     | 1.000  | 2.906     | 2.455                                   | 1.000  | 2.882     |
| LEVIN     | 2.461                                     | 2.429   | 1.464     | 2.365                                       | 2.297    | 1.430     | 3.623                                     | 3.595  | 2.099     | 3.209                                   | 3.120  | 2.065     |
| NOUTGOING | 17.752                                    | 15.000  | 14.422    | 19.403                                      | 14.000   | 16.368    | 10.829                                    | 8.000  | 9.556     | 4.860                                   | 3.000  | 4.777     |
| P         | 0.225                                     | 0.196   | 0.178     | 0.246                                       | 0.164    | 0.203     | 0.137                                     | 0.104  | 0.115     | 0.061                                   | 0.046  | 0.055     |
| VARCAP    | 820.667                                   | 820.990 | 85.218    | 822.902                                     | 815.162  | 86.217    | —   | —      | —         | —                                       | —      | —         |
| VARLOANS  | —   | —       | —         | 38.955                                      | 8.376    | 46.160    | 47.854                                    | 33.478 | 39.566    | —                                       | —      | —         |
| CATEND    | 0.370                                     | 0.017   | 0.468     | 0.382                                       | 0.018    | 0.468     | 0.149                                     | 0.010  | 0.332     | 0.220                                   | 0.010  | 0.385     |
| CATIN1    | 0.010                                     | 0.010   | 0.006     | 0.010                                       | 0.010    | 0.006     | 0.010                                     | 0.010  | 0.000     | 0.010                                   | 0.010  | 0.000     |
| CATIN2    | 0.025                                     | 0.024   | 0.015     | 0.025                                       | 0.025    | 0.017     | 0.035                                     | 0.033  | 0.024     | 0.036                                   | 0.032  | 0.026     |
| COUNT     | 4.649                                     | 3.000   | 4.261     | 4.823                                       | 3.000    | 4.291     | 2.383                                     | 1.000  | 3.212     | 2.830                                   | 1.000  | 3.469     |
| LEVIN     | 2.474                                     | 2.458   | 1.423     | 2.555                                       | 2.629389 | 1.515     | 3.539                                     | 3.615  | 2.006     | 3.600                                   | 3.552  | 1.967     |
| NOUTGOING | 18.732                                    | 14.000  | 15.964    | 21.815                                      | 15.00000 | 18.283    | 10.898                                    | 8.000  | 9.488     | 5.267                                   | 4.000  | 4.565     |
| P         | 0.189                                     | 0.142   | 0.157     | 0.220                                       | 0.141615 | 0.181     | 0.110                                     | 0.079  | 0.092     | 0.053                                   | 0.041  | 0.040     |
| VARCAP    | 822.626                                   | 821.214 | 80.128    | 827.379                                     | 826.1881 | 76.812    | —   | —      | —         | —                                       | —      | —         |
| VARLOANS  | —   | —       | —         | 37.755                                      | 8.359417 | 44.330    | 48.180                                    | 33.479 | 39.959    | —                                       | —      | —         |

The mean, median, and standard deviation are depicted for interbank networks consisting of 20, 50, 80, and 100 banks, respectively. Four scenarios are included: (a) Heterogeneous banks-homogeneous exposures; (b) Heterogeneous banks-heterogeneous exposures; (c) Homogeneous banks-heterogeneous exposures; (d) Homogeneous banks-homogeneous exposures. The variables are: CATEND, defined as total loss of capital due to contagion as percentage of total capital in the system; CATIN1, defined as bank *i*'s depleted equity divided by the total equity in the network; CATIN2, defined as the total amount of loans granted to bank *i* that cannot be repaid, divided by the total equity in the network; LEVIN, defined as the leverage of the interbank network; NOUTGOING, defined as the number of outgoing links of bank *i*, which corresponds to the number of its creditors in the network; COUNT, defined as the number of rounds needed until no further bank defaults; VARCAP, defined as the variance of bank capital; VARLOANS, defined as the variance of interbank loans and  $p_i$ , the Erdős-Rényi probability  $p_{ij}$  that there is a lending/borrowing link between two nodes/banks

work with standardized variables—both dependent and independent variables—and have not included the intercept term in the regression models as it will be zero. Our concern is to measure effects not in terms of the original units of the dependent variable or the independent variables, but in standard deviation units.<sup>3</sup>

## 4.2 Simulation Results

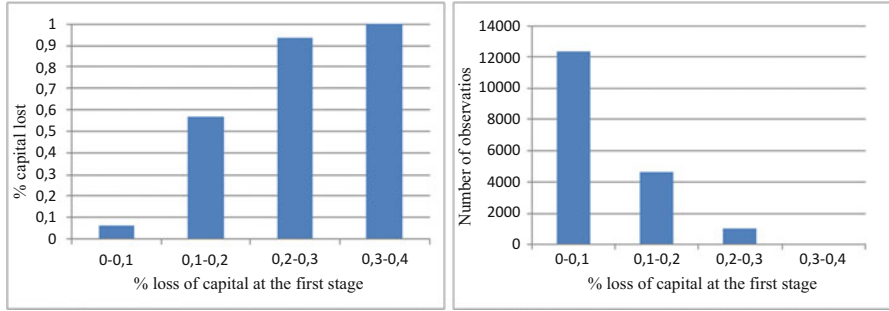
In this section, we discuss the regression results of all four scenarios. Since our variables are measured on different scales, we cannot directly infer which independent variable has the largest effect on the dependent variable. In order to circumvent this problem we standardize our series to have zero mean and unit variance. Table 2 presents the regression results of the **first scenario** using the OLS model described in Equation (15), where heterogeneous banks distribute evenly their interbank claims across the outgoing links of a network consisting of  $N = 20, 50, 80$  and 100 banks. Almost all regressor coefficients are found to be statistically significant for all the sizes of the network. We discern only two cases where regressor coefficients are found to be statistically insignificant and has to do with CATIN1 variable and one case that has to do with CATIN2.  $R$ -squared coefficients take on large values ranging from 74.9 to 80% and highlight the ability of our selected variables to explain financial distress in interbank networks.

The variable CATIN1 captures the initial effect defaulting bank  $i$  exerts on the network, whereas the magnitude of interconnectedness across the banks that comprise the interbank network is measured through parameter CATIN2. As we observe from Table 1, variable CATIN1 does not seem to affect much the dependent variable, whereas two regressor coefficients are found to be insignificant. Financial shocks will propagate into the defaulting bank's counterparties along the network, erode their capital and make them more vulnerable to subsequent shocks. The magnitude of the positive relationship between CATIN2 and CATEND – the dependent variable – seems to increase as the size of the interbank network increases with the only exception being the  $N = 50$  bank network segment which follows an autonomous path (although statistically insignificant). The increasing magnitude of the above relationship seems to cease as we move from the case of  $n = 80$  banks to the case of  $n = 100$  banks. This finding implies that as we move from smaller to larger network settings, systemic risk and the likelihood of contagion increases. However, when we move from the case of  $n = 80$  banks to the case of  $n = 100$  banks the likelihood of contagion seems to decrease. Figure 3 visually illustrates the extent of contagion as a function of the percentage loss of capital due to bank's  $i$  default. It is shown that as the network size increases from small to medium sized networks, we observe that capital losses rises, confirming the findings from the

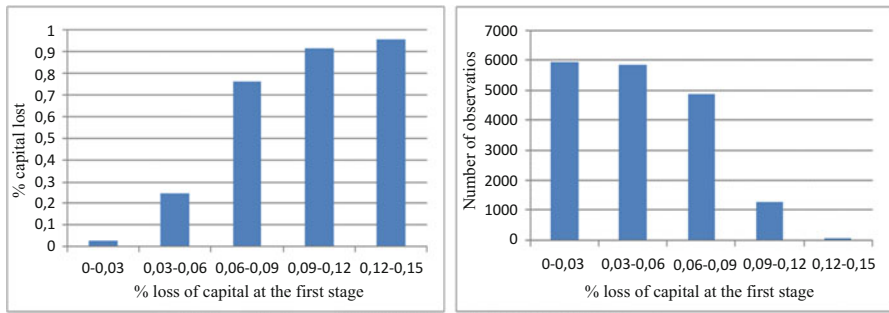
---

<sup>3</sup>See Wooldridge (2003) for an interesting discussion on standardization and explanation of the absence of the standardized intercept.

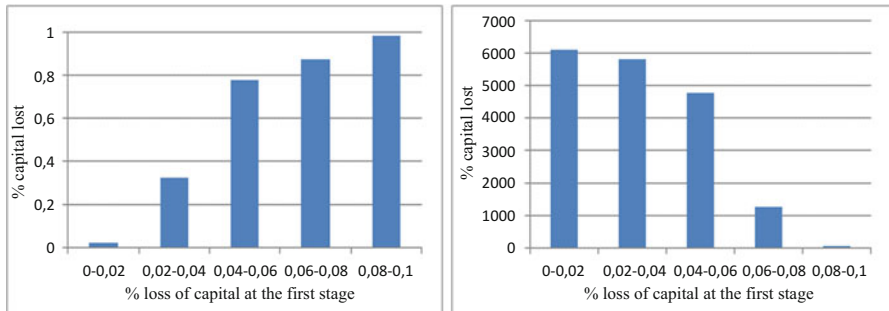




(a) N=20 banks

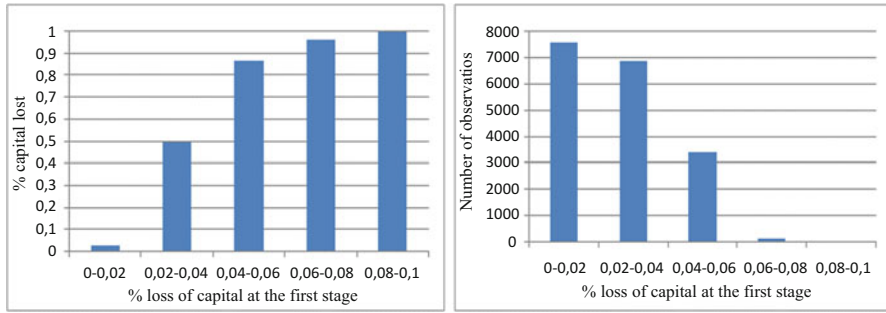


(b) N=50 banks



(c) N=80 banks

**Fig. 3 Scenario 1: Heterogeneous Banks with homogeneous exposures. Extent of contagion** (expressed as the total capital lost from the banking system due to the failure of at least one bank) **as a function of the % initial loss of capital** due to default of the first bank. Panels (a–d) show the relation between the % initial loss of capital due to default of the first bank and the extent of contagion across interbank networks with different number of banks



(d) N=100 banks

Fig. 3 (continued)

regression model. As we can observe from Figure 3, as we move from the  $n = 80$  interbank network scheme to  $n = 100$  the likelihood of contagion seems to decrease since we have very few cases that cause systemic break downs and defaults.

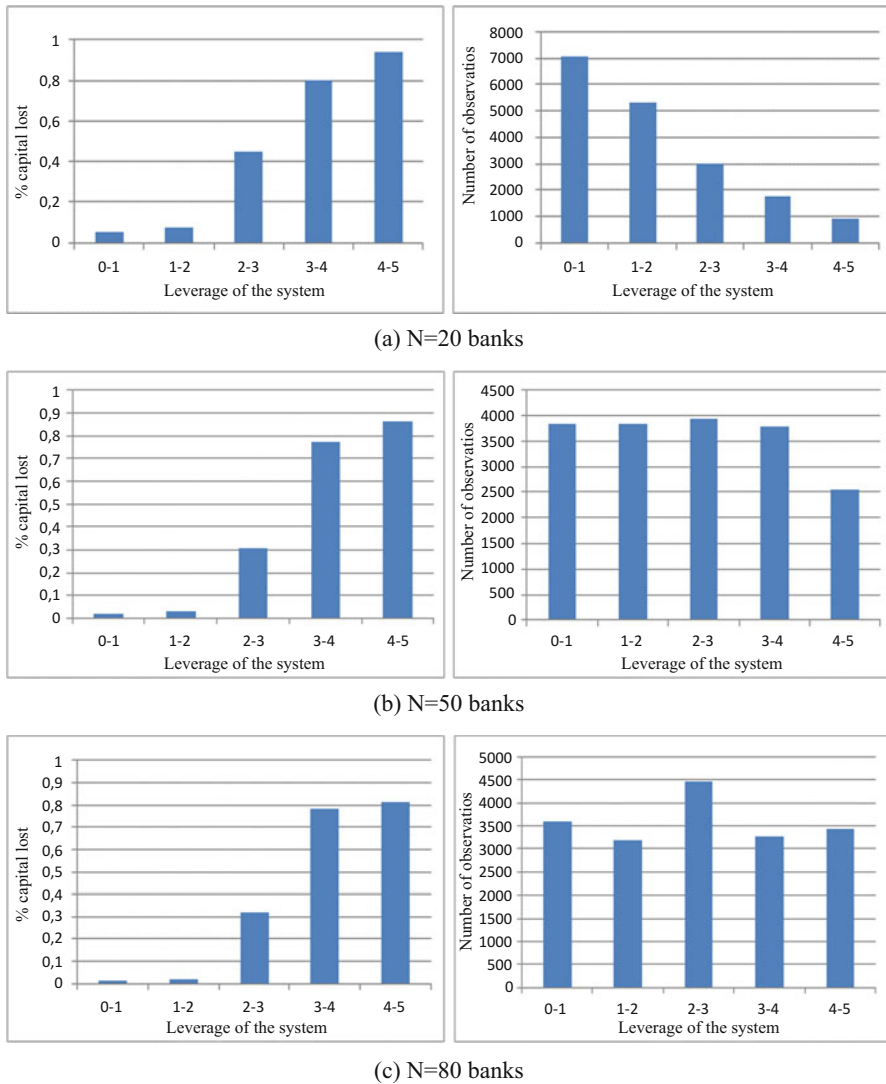
As expected, we also find that there is a positive relationship between the leverage of the network and the capital losses due to contagion, which is depicted by Figure 4. This result is in line with the findings of Nier et al. [7] who provide evidence that systemic risk increases when system-wide leverage increases. Highly leveraged banks in the interbank network are clearly more exposed to the risk of default on interbank loans. The greater the size of default on debt is, the larger the losses are that banks transmit to their neighbors, other things being equal. Thus, highly leveraged banks contribute more to systemic risk as they become a vehicle for transmitting shocks within the network. Moreover, it is shown that the magnitude of the positive relationship between the network’s leverage and contagion risk increases as we move from smaller to larger interbank networks (illustrated in Table 2) with the only exception being the  $n = 80$  bank network scheme where the magnitude of the standardized coefficients seems to decrease.

Our results also suggest that connectivity, expressed in our experiments as the outgoing<sup>4</sup> of the first bank that defaults, has a negative effect on interbank contagion with the only exception being the case of  $n = 50$  banks where we can observe a positive relationship between contagious defaults and connectivity.

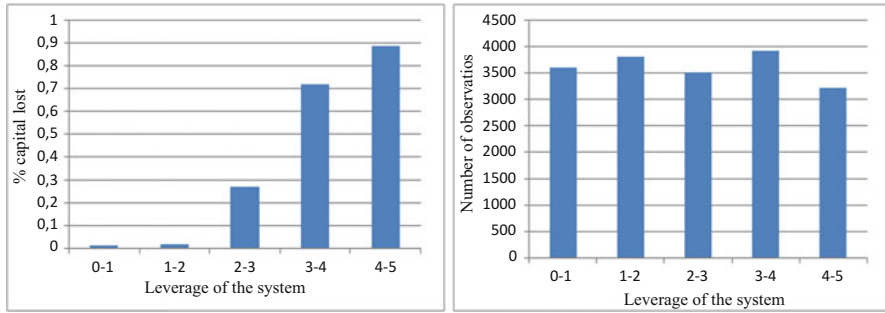
Interestingly, as we move from small networks consisted of twenty banks to networks consisted of 50 banks the effect of connectivity to interbank contagion turns from negative to positive and after then connectivity keeps affect negatively the systemic risk of the network. Thus, as we move from network systems consisted of fifty banks to networks consisted of 100 banks this negative relationship seems to decrease. In relatively small interbank networks, a high level of connectivity

<sup>4</sup>It should be highlighted that in the Erdős–Rényi network structure the outdegree equals the indegree since we have an undirected network structure. However, in our framework, the two directions of the same link are given different weights.

will allow an efficient absorption of shocks, whereas in medium size networks the increased connectivity will spread the shock throughout the system, potentially leading to many default cascades. The link probability, that is assumed to be equal across all pairs, seems to contribute to the resilience of the system for small and medium size networks. However, as we move from medium to large size networks



**Fig. 4** Scenario 1: **Heterogeneous Banks with homogeneous exposures. Extent of contagion** (expressed as the total capital lost from the banking system due to the failure of at least one bank) **as a function of the leverage of the system.** Panels (a–d) show the relation between the leverage of the system and the extent of contagion across interbank networks with different number of banks



(d) N=100 banks

Fig. 4 (continued)

this effect turns negative to the resilience of the system as it seems to contribute positively to systemic risk.

Our regression analysis also shows that the COUNT variable which measures the number of rounds until no further bank defaults, has a positive impact on interbank contagion. Heterogeneity expressed as the variance of capital exhibits a negative and statistically significant relationship with interbank contagion, showing that size heterogeneity can have positive effects on the stability of an interbank network.

However, the positive magnitude seems to decrease as we move from small to large interbank networks. An interbank network consisting of banks of various sizes can more easily withstand a negative shock, therefore no institution becomes significant for either borrowing or lending. Furthermore, in such network both smaller and larger banks can act as shock absorbers when an initial shock hits the banking system, making contagion a less likely phenomenon. This finding is in line with the results of Iori et al. [6] concerning bank size heterogeneity.

Table 4 presents the regression results of the **second scenario** using the model described in Equation (16), where banking institutions with heterogeneous bank sizes are linked to one another via heterogeneous interbank claims. The regressor coefficients are statistically significant in almost all cases and the *R*-squared values are quite high and lie in the vicinity of 75–83%, highlighting the good explanatory power of the model.

CATIN1 does not seem to impact much the dependent variable in all network segments and the regressor coefficients in the relatively large interbank networks becomes statistically insignificant. The magnitude of standardized coefficients is almost the same with the corresponding magnitude of those derived from the first scenario. In other words, an initial shock from defaulting bank *i* will spill over more easily in the network. Thus, the first bank defaulting has the dynamics to spread the initial shock and contaminate the entire interbank network. CATIN2 has a large positive impact on contagion risk, however, its magnitude fades away as we move from smaller to larger networks. It should also be highlighted that the CATIN2 coefficients are much larger than those recorded in the first scenario in all network

**Table 4** OLS regression analysis for Scenario 2 (Heterogeneous banks with heterogeneous exposures)

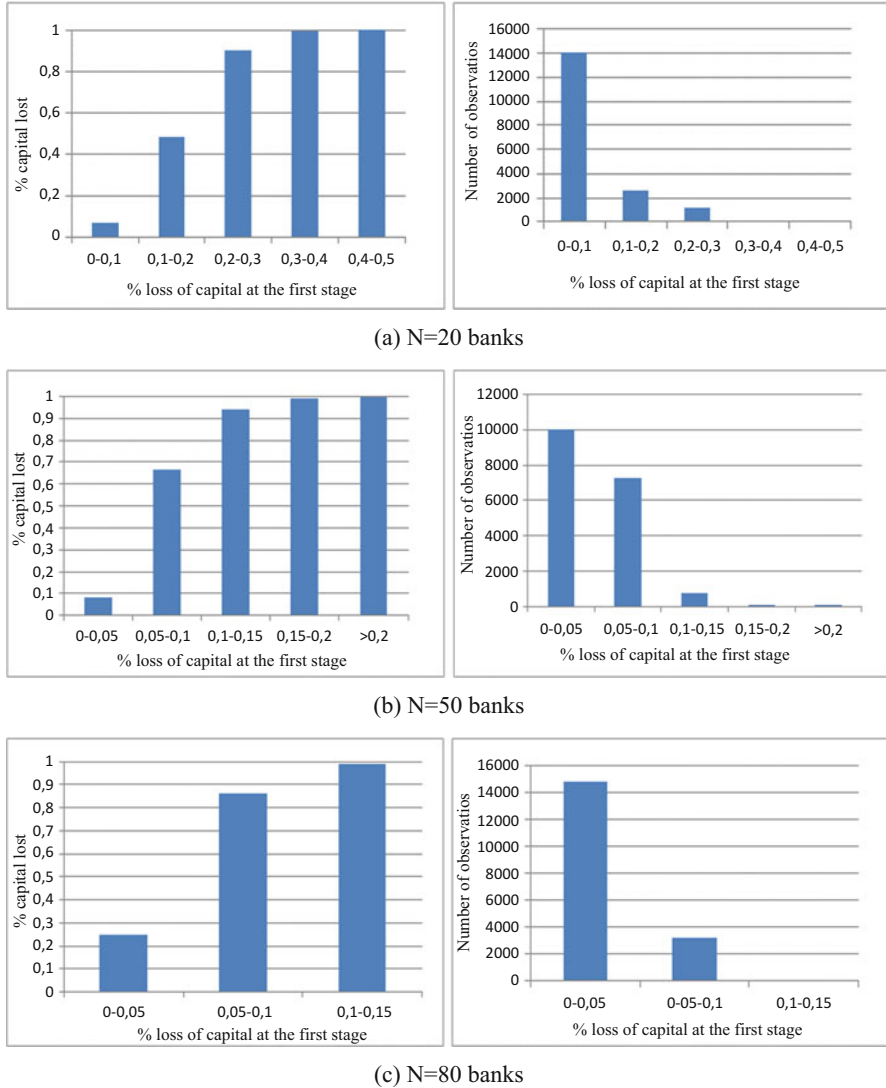
|                         | N = 20                 | N = 50                 | N = 80                 | N = 100                |
|-------------------------|------------------------|------------------------|------------------------|------------------------|
| CATIN 1                 | 0.070<br>(23.660)***   | 0.007<br>(2.024)**     | 0.000<br>(0.047)       | -0.001<br>(-0.283)     |
| CATIN2                  | 0.201<br>(19.541)***   | 0.113<br>(11.183)***   | 0.106<br>(9.669)***    | 0.071<br>(6.015)***    |
| LEVIN                   | 0.653<br>(58.484)***   | 0.346<br>(30.847)***   | 0.321<br>(26.320)***   | 0.399<br>(30.132)***   |
| NOUTGOING               | -0.136<br>(-11.540)*** | -0.150<br>(-6.539)***  | -0.052<br>(-2.253)**   | 0.038<br>(1.575)       |
| COUNT                   | 0.456<br>(111.687)***  | 0.630<br>(156.274)***  | 0.577<br>(141.939)***  | 0.573<br>(131.397)***  |
| VARCAP                  | -0.032<br>(-18.897)*** | -0.067<br>(-50.848)*** | -0.053<br>(-52.027)*** | -0.041<br>(-40.597)*** |
| VARLOANS                | -0.246<br>(-45.472)*** | -0.091<br>(-14.882)*** | -0.018<br>(-3.113)***  | -0.082<br>(-12.307)*** |
| P                       | -0.254<br>(-21.462)*** | 0.038<br>(1.620)       | 0.064<br>(2.678)***    | -0.110<br>(-4.311)***  |
| Adjusted R <sup>2</sup> | 0.830                  | 0.796                  | 0.776                  | 0.751                  |

The table presents the regression results for Scenario 2. The dependent variable is CATEND measured as the total loss of capital due to contagion as percentage of total capital in the network. Explanatory variables are the constant term CATIN1, CATIN2, LEVIN, NOUTGOING, COUNT, VARCAP, VARLOANS and P, the probability for a link to exist between two nodes. Each cell displays the OLS standardized coefficients along with the corresponding *t*-statistics (shown in parentheses). The sample comprises of 18,000 realizations (simulated banking crises). \*, \*\* and \*\*\* denote significance at the 10, 5 and 1 percent level, respectively

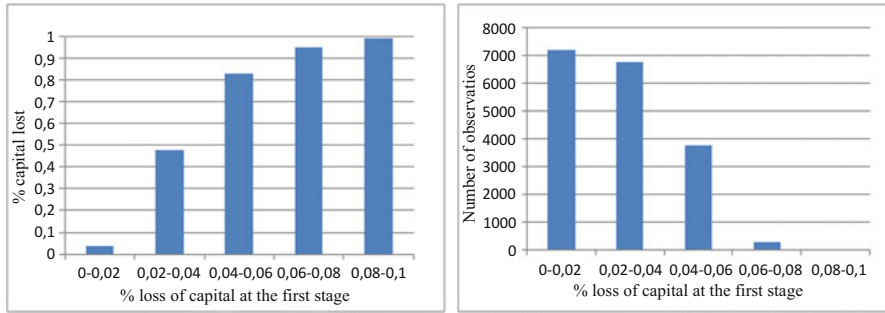
sizes. An initial shock following the default of bank *i* seems to contribute much to a banking crisis scenario within small and medium-sized networks and the size of total capital losses is smaller than that documented in the first scenario. Figure 5 depicts the extent of contagion as a function of the percentage loss of capital due to default of the first bank and confirms the results recorded in Table 5.

The results also show that there still exists a positive relationship between leverage and contagion (illustrated in Figure 6); however, the coefficient estimates are larger in almost all cases than those recorded in the previous scenario. Moreover, the magnitude of the leverage coefficients decreases as the number of banks in the interbank network increases, with the only exception being the 100 bank network segment where one can observe a slight increase compared to the 80 bank network segment.

Results on connectivity are qualitatively similar to those of the first scenario, showing that connectivity negatively impacts contagion risk especially in small and medium interbank networks with the only exception being the 100 bank network segment which follows an autonomous path and is positively related to contagion (although statistically insignificant).



**Fig. 5** Scenario 2: **Heterogeneous Banks with heterogeneous exposures. Extent of contagion** (expressed as the total capital lost from the banking system due to the failure of at least one bank) **as a function of the % initial loss of capital** due to default of the first bank. Panels (a–d) show the relation between the % initial loss of capital due to default of the first bank and the extent of contagion across interbank networks with different number of banks



(d) N=100 banks

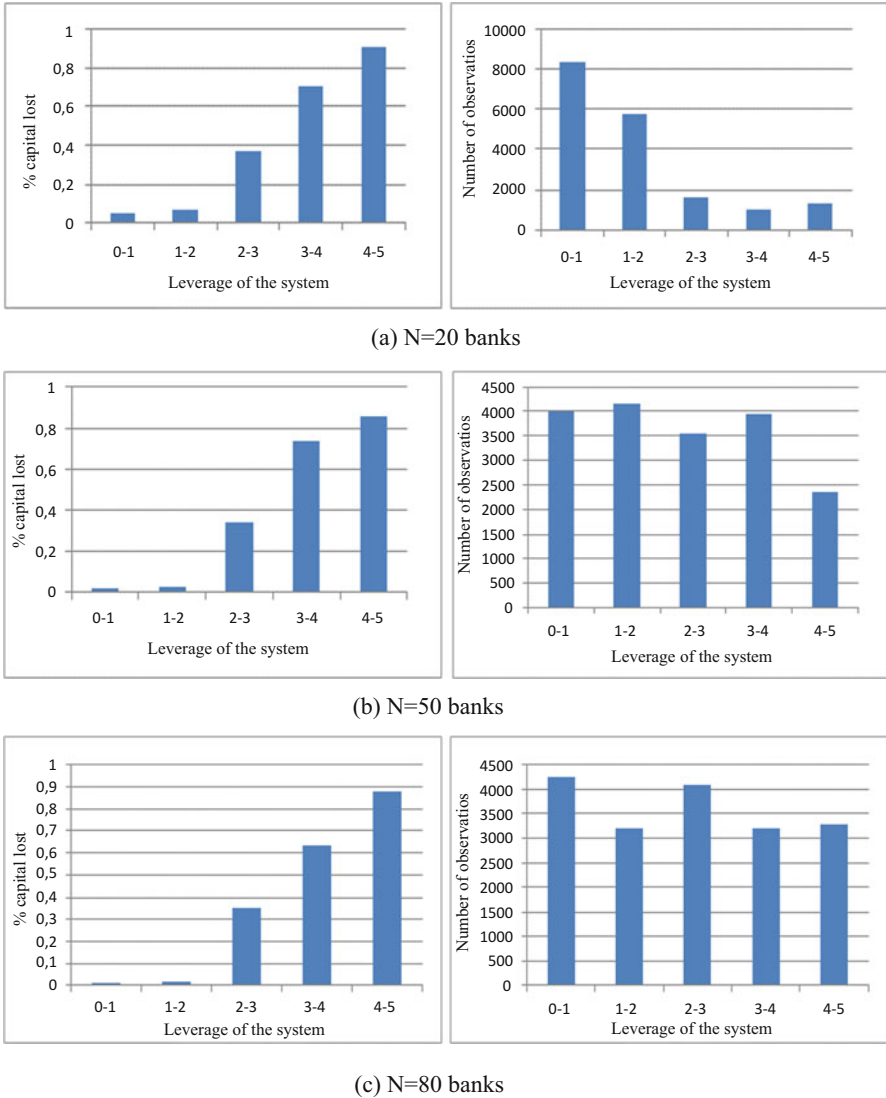
Fig. 5 (continued)

Table 5 OLS regression analysis for Scenario 3 (Homogeneous banks with heterogeneous exposures)

|                         | N = 20                 | N = 50                 | N = 80                 | N = 100                |
|-------------------------|------------------------|------------------------|------------------------|------------------------|
| CATIN2                  | 0.196<br>(25.178)***   | 0.143<br>(16.422)***   | 0.125<br>(15.232)***   | 0.088<br>(9.806)***    |
| LEVIN                   | 0.324<br>(39.268)***   | 0.298<br>(32.475)***   | 0.275<br>(31.619)***   | 0.279<br>(30.578)***   |
| NOUTGOING               | -0.163<br>(-15.308)*** | -0.168<br>(-10.438)*** | -0.126<br>(-8.707)***  | -0.087<br>(-5.561)***  |
| COUNT                   | 0.736<br>(191.690)***  | 0.761<br>(175.841)***  | 0.790<br>(195.383)***  | 0.793<br>(186.247)***  |
| VARLOANS                | -0.175<br>(-43.977)*** | -0.190<br>(-44.390)*** | -0.180<br>(-46.723)*** | -0.167<br>(-41.937)*** |
| P                       | -0.253<br>(-24.270)*** | -0.313<br>(-19.153)*** | -0.322<br>(-21.833)*** | -0.339<br>(-21.575)*** |
| Adjusted R <sup>2</sup> | 0.860                  | 0.823                  | 0.845                  | 0.809                  |

The table presents the regression results for Scenario 3. The dependent variable is CATEND measured as the total loss of capital due to contagion as percentage of total capital in the network. Explanatory variables are the constant term CATIN2, LEVIN, NOUTGOING, COUNT, VARLOANS and P, the probability for a link to exist between two nodes. Each cell displays the OLS standardized coefficients along with the corresponding *t*-statistics (shown in parentheses). The sample comprises of 18,000 realizations (simulated banking crises). \*, \*\* and \*\*\* denote significance at the 10, 5 and 1 percent level, respectively

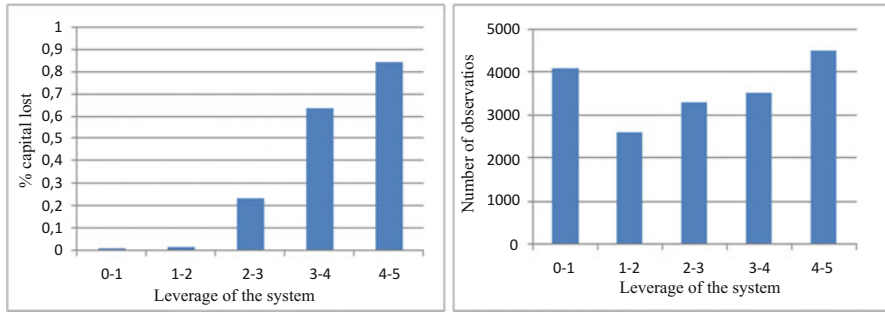
As far as the link probability is concerned, we can observe a different pattern from that of the first scenario. For small and large sized networks, link probability seems to contribute negatively to systemic risk while for medium sized networks there is a positive relationship between link probability and contagion. The number of rounds until no further bank defaults positively impacts contagion risk and contributes the most to total capital losses in the banking system when medium and large interbank networks are formed. Under this scenario, the heterogeneity allowed on both bank sizes and interbank exposures has had a great impact on the resilience



**Fig. 6** Scenario 2: **Heterogeneous Banks with heterogeneous exposures. Extent of contagion** (expressed as the total capital lost from the banking system due to the failure of at least one bank) **as a function of the leverage of the system.** Panels (a–d) show the relation between the leverage of the system and the extent of contagion across interbank networks with different number of banks

of the network system. Heterogeneity impacts negatively on interbank contagion although its intensity decreases as the size of the network increases. Moreover, as we can see from the Table 4 heterogeneity of bank size contributes less to the resilience





(d) N=100 banks

Fig. 6 (continued)

of the interbank network than heterogeneity of interbank exposures when it comes to small and medium sized networks.

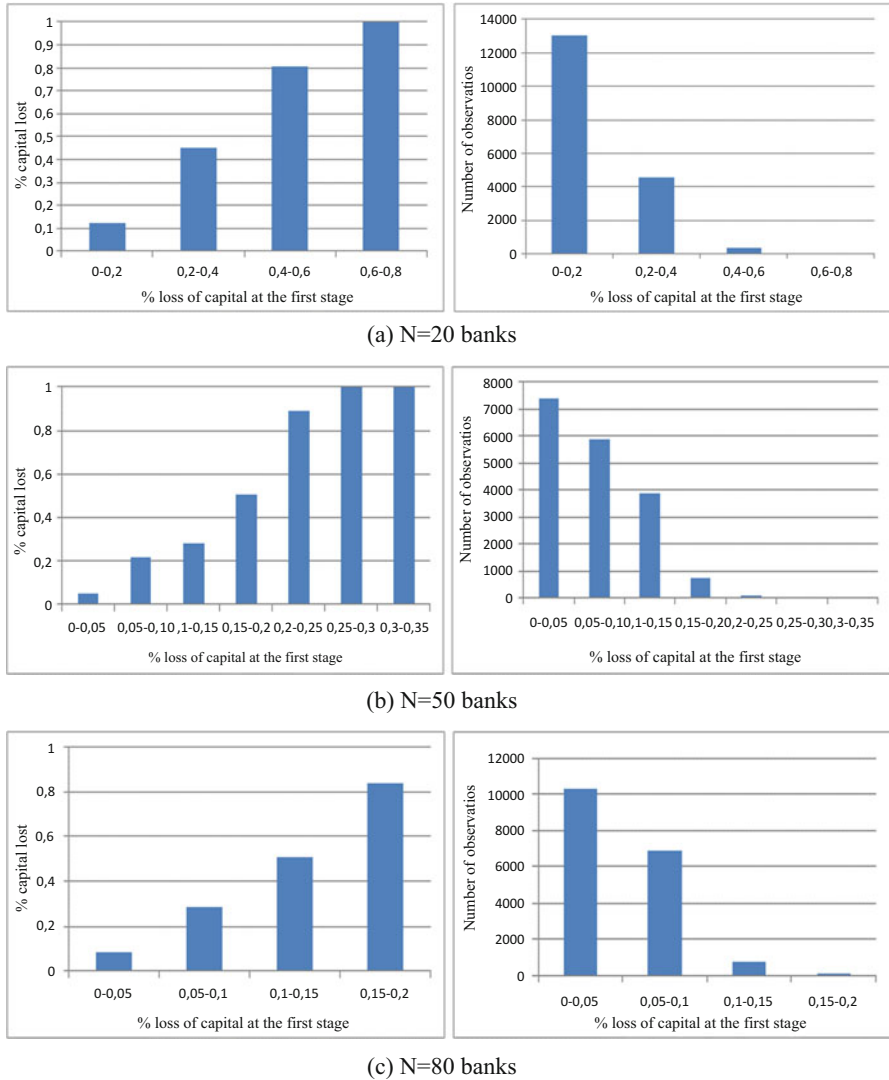
The heterogeneity of interbank exposures acts as a diversification tool and contributes to a smaller extent to an unfolding crisis compared to the scenario where homogeneous banks are interconnected via heterogeneous exposures (shown in Table 4).

Table 5 depicts the results of the **third scenario** as described in Equation (17). In this scenario, we construct network systems where banks have the same equity size and unevenly distribute their exposures across creditor banks. We note that an initial shock fades away with the failure of the first bank and does not spillover to other banks within the network. This is mainly due to our choice of parameters regarding the equity of each bank, the links among banks and the interbank claims among creditor banks. In order to observe default cascades we relax our initial assumptions and lower the equity of each bank in the network system.

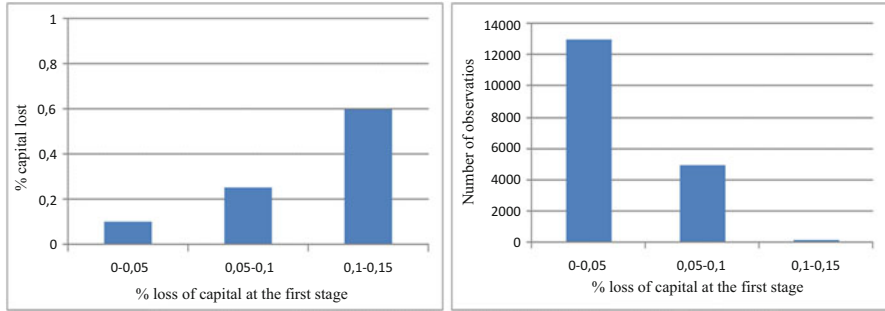
Specifically, each bank is now endowed with a balance sheet that consists of 25 units of equity and interbank claims among creditor banks are distributed in the following ranges:  $a_{ij} \in [0, 10]$ ,  $a_{ij} \in [0, 20]$ ,  $a_{ij} \in [0, 35]$ . Interbank exposures levels were kept the same as in Leventides et al. (2018). Moreover, we control the leverage of the system by setting the rule that the maximum leverage ratio of each network system cannot exceed seven. Similar to the previous scenarios, the regressor coefficients are statistically significant in all cases and the *R*-squared values are still large, in fact the largest of all three scenarios tested. Variable CATIN2 has a highly significant positive impact on systemic risk that fades away as the network system gets larger. The same observation holds for the level of connectivity in the banking system i.e. a strong negative impact on contagion risk that dissipates as *N* increases.

The leverage of the system has a positive impact on systemic risk and its magnitude decreases as the size of the network increases. Figures 7 and 8 illustrate the third scenario as a function of the percentage loss of capital due to default of the first bank in the network and as a function of leverage in the banking system, respectively. As in the previous cases, we find the number of rounds until no further

bank defaults to affect contagion risk positively and statistically significantly, and such impact is magnified in relatively larger interbank networks. The heterogeneity of interbank exposures plays a significant role in the stability of the financial network especially in the medium sized interbank networks.



**Fig. 7** Scenario 3: **Homogeneous banks with heterogeneous exposures** (expressed as the total capital lost from the banking system due to the failure of at least one bank) **as a function of the % initial loss of capital** due to default of the first bank. Panels (a–d) show the relation between the % initial loss of capital due to default of the first bank and the extent of contagion across interbank networks with different number of banks



(d) N=100 banks

Fig. 7 (continued)

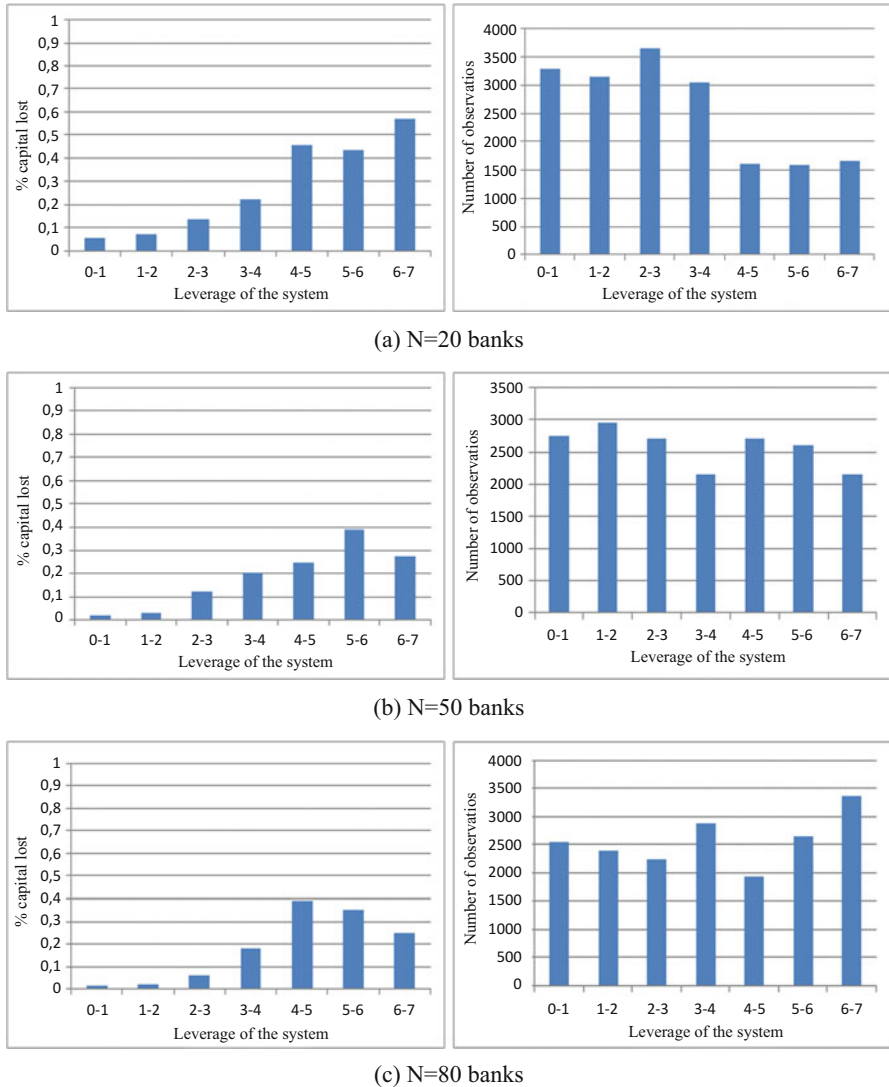
Table 6 OLS regression analysis for Scenario 4 (Homogeneous banks with homogeneous exposures)

|                         | N = 20                 | N = 50                 | N = 80                 | N = 100                |
|-------------------------|------------------------|------------------------|------------------------|------------------------|
| CATIN2                  | 0.228<br>(21.978)***   | 0.153<br>(14.098)***   | 0.137<br>(12.902)***   | 0.105<br>(9.426)***    |
| LEVIN                   | 0.137<br>(14.890)***   | 0.268<br>(28.512)***   | 0.352<br>(37.106)***   | 0.352<br>(37.707)***   |
| NOUTGOING               | -0.257<br>(-15.906)*** | -0.146<br>(-9.715)***  | -0.130<br>(-8.719)***  | -0.095<br>(-6.262)***  |
| COUNT                   | 0.645<br>(198.356)***  | 0.617<br>(172.925)***  | 0.568<br>(150.736)***  | 0.573<br>(148.381)***  |
| P                       | -0.156<br>(-10.231)*** | -0.304<br>(-21.593)*** | -0.378<br>(-26.723)*** | -0.379<br>(-27.197)*** |
| Adjusted R <sup>2</sup> | 0.834                  | 0.806                  | 0.817                  | 0.779                  |

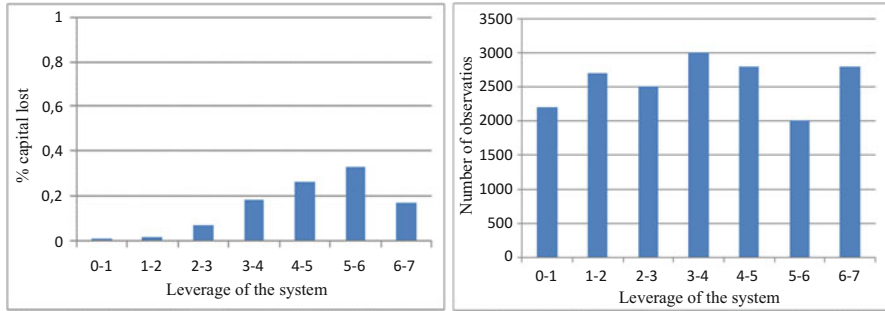
The table presents the regression results for Scenario 4. The dependent variable is CATEND measured as the total loss of capital due to contagion as percentage of total capital in the network. Explanatory variables are the constant term CATIN2, LEVIN, NOUTGOING, COUNT and P, the probability for a link to exist between two nodes.. Each cell displays the OLS standardized coefficients along with the corresponding *t*-statistics (shown in parentheses). The sample comprises of 18,000 realizations (simulated banking crises). \*, \*\* and \*\*\* denote significance at the 10, 5 and 1 percent level, respectively

Finally, Table 6 depicts the results of the **fourth scenario** as described in Equation (18). In this scenario, we construct network systems where banks have the same equity size and interbank claims are evenly distributed across creditor banks. We acknowledge the fact that this scenario is a bit unrealistic as banks in real-world interbank networks do not have the same equity size and do not necessarily distribute their interbank claims evenly across their creditors. However, by testing a wide range of link probabilities between any two nodes, we are in a position to effectively examine the effect of different calibrations on systemic risk. Thus, although this scenario can be regarded as a special case with magnifying effects, it provides useful insights on interbank market resiliency during periods of stress.

The variable CATIN2 has a strong positive impact on systemic risk that dissipates as the network system gets larger. Simulations show that this scenario yields qualitatively similar results with the previous three scenarios in relation to the leverage of the network, that is, leverage positively and significantly affects



**Fig. 8** Scenario 3: Homogeneous banks with heterogeneous exposures. Extent of contagion (expressed as the total capital lost from the banking system due to the failure of at least one bank) as a function of the leverage of the system. Panels (a–d) show the relation between the leverage of the system and the extent of contagion across interbank networks with different number of banks



(d) N=100 banks

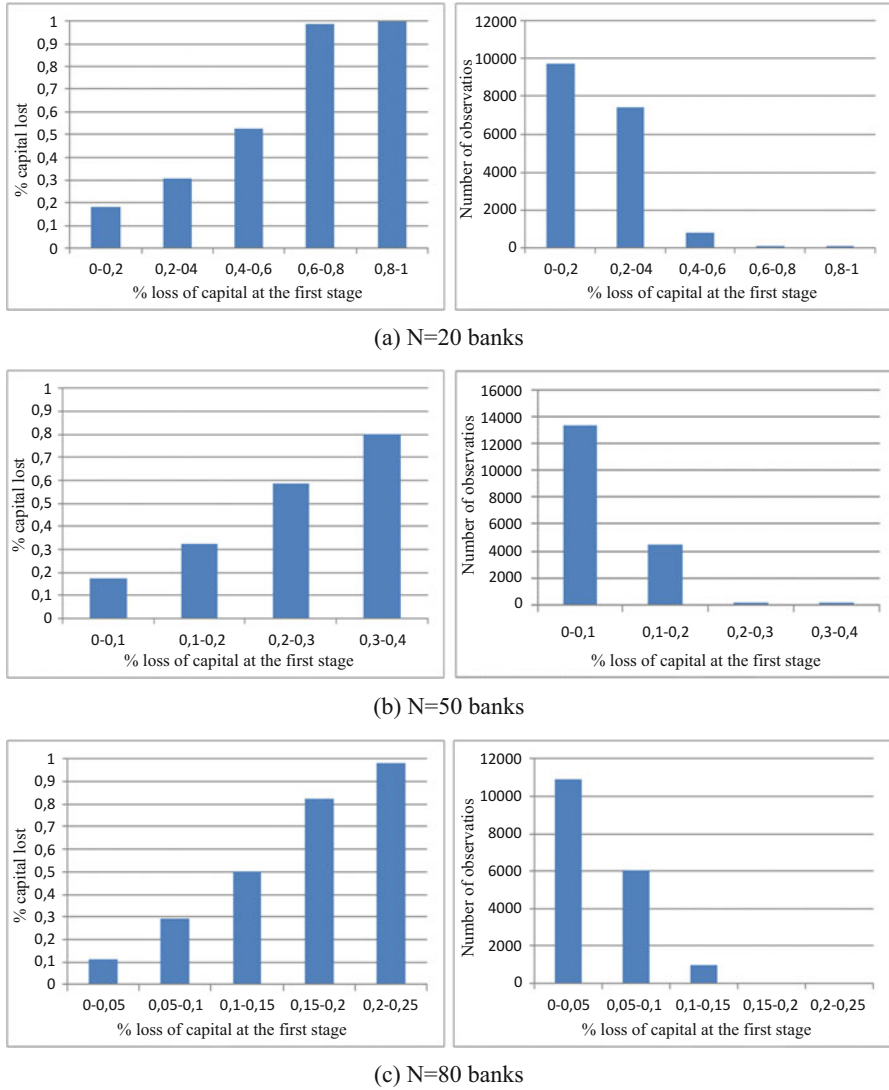
Fig. 8 (continued)

contagion risk (illustrated in Figure 10). However, in this scenario, we observe that this effect becomes stronger progressively when the number of constituent banks in the network increases. Figure 9 confirms the results recorded in the Table 6 concerning the relationship between the extent of contagion and the percentage loss of capital in the network. For instance, the likelihood of systemic breakdowns seems to decrease as we move from smaller to larger network systems since we have very few cases that cause large capital losses. Connectivity impacts negatively on interbank contagion, although this negative impact dissipates as the number of banks in the interbank networks increases. As expected, the link probability has the same negative impact as connectivity on the interbank contagion. Contrary to the previous findings concerning connectivity, the negative impact of the link probability on interbank contagion seems to scale up as we move from smaller to larger interbank networks (Figure 10).

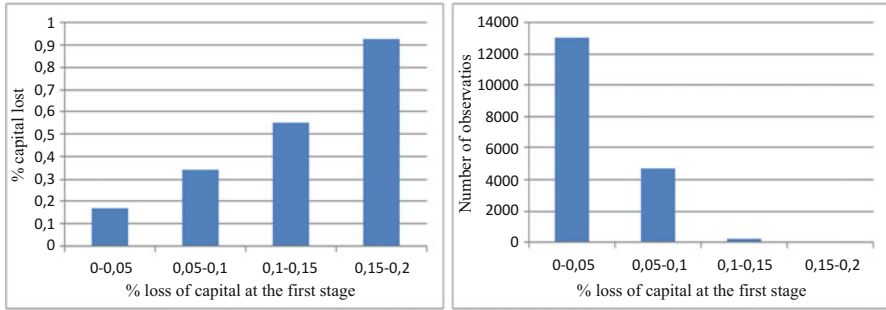
Finally, the number of rounds until no further bank defaults affects contagion risk in a statistically significant manner especially when small interbank networks are considered.

The main intuition behind these results is that increasing connectivity on a homogeneous interbank network can reduce the frequency of contagion in case the first bank that defaults is less leveraged as the interbank network has the dynamics to absorb more easily the shock and thus the initial shock is dissipated at a faster rate. This is the case for small network systems. As the size of the network increase and the system gets more leveraged, the stabilizing force of connectivity weakens and default cascades prevail.

Tables 7, 8, 9, and 10 depict robustness tests on all four scenarios based on random sampling. We have performed second run Monte Carlo simulations in order to examine whether the new results differ from the previous ones, thus checking how random sampling affects our main conclusions. We observe qualitatively similar results in all four cases to those from the first run providing evidence that our findings are stable across different simulation scenarios.



**Fig. 9** Scenario 4: **Homogeneous banks with homogeneous exposures** (expressed as the total capital lost from the banking system due to the failure of at least one bank) **as a function of the % initial loss of capital** due to default of the first bank. Panels (a–d) show the relation between the % initial loss of capital due to default of the first bank and the extent of contagion across interbank networks with different number of banks



(d) N=100 banks

Fig. 9 (continued)

## 5 Conclusions

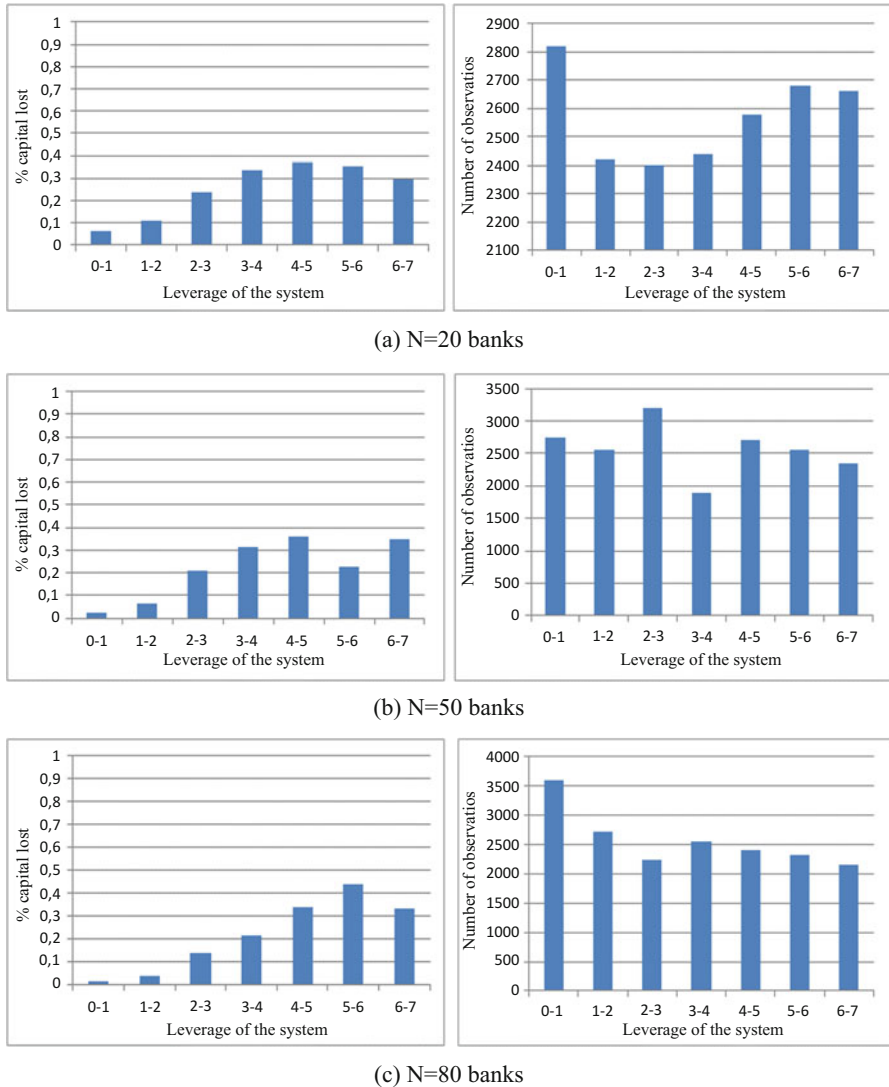
This paper investigates how complexity under a specific network structure, that has been extensively applied for the study of contagion in financial networks, affects interbank contagion. Similar to Leventides et al. [5], we explore the interplay between heterogeneity, balance sheet composition in the spreading of contagion using four basic scenarios, under an Erdős–Rényi network structure using a wide range of link probabilities between any two banks.

Our findings indicate a non-monotonic relation between diversification and interbank contagion across the different sizes of interbank networks for all scenarios tested. While for small or medium interbank networks, connectivity can act as an absorbing barrier, such that interbank systems of these sizes can withstand the initial shock, for large network systems connectivity does not seem to provide an effective shield against capital losses. Our results, for the four scenarios tested, indicate that small and thus more concentrated interbank network systems are more prone to contagion. In these cases, we observe that the size of total capital losses is, in most cases, larger than that documented in medium and large sized systems, which is in line with the findings of Nier et al. [7].

As far as heterogeneity is concerned, this enters in our experiments in the form of interbank claims and bank sizes. Our results clearly suggests that heterogeneity plays a significant role in the stability of the financial system. Similar to Leventides et al. [5], we still find that when heterogeneity is introduced with respect to the size of each bank, the system’s shock absorption capacity is enhanced. In addition, when we allow for heterogeneity on interbank exposures in our model, we observe additional resilience to the interbank network as an initial shock dissipates more easily than in the case of homogeneous interbank claims.

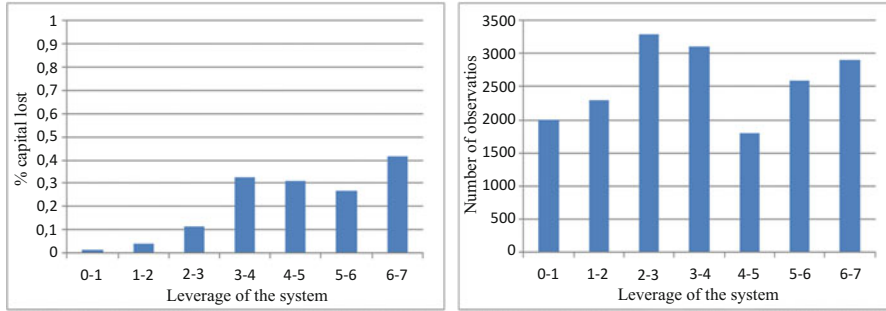
Finally, we should also justify the fact that we choose to work under an Erdős–Rényi network structure even if this network framework is not very realistic. In an such a network framework, where the probability of forming a link is homogeneous, the resulting network structure does not present marked heterogeneity. As we

observed from all the four scenarios tested, the initial shock that hits the system seems to propagate into the system jeopardizing thus the stability of the entire system. This strengthens even more our arguments concerning the critical role that heterogeneity plays in the resilience of the financial system.



**Fig. 10** Scenario 4: Homogeneous banks with homogeneous exposures. Extent of contagion (expressed as the total capital lost from the banking system due to the failure of at least one bank) as a function of the leverage of the system. Panels (a–d) show the relation between the leverage of the system and the extent of contagion across interbank networks with different number of banks





(d) N=100 banks

Fig. 10 (continued)

Table 7 Robustness tests: OLS regression analysis for Scenario 1 (Heterogeneous banks with homogeneous exposures)

|                         | N = 20                 | N = 50                 | N = 80                 | N = 100                |
|-------------------------|------------------------|------------------------|------------------------|------------------------|
| CATIN 1                 | 0.044<br>(13.363)***   | 0.001<br>(0.389)       | -0.004<br>(-1.135)     | 0.002<br>(0.579)       |
| CATIN2                  | 0.259<br>(10.095)***   | 0.163<br>(7.044)***    | 0.073<br>(3.301)***    | 0.050<br>(2.614)***    |
| LEVIN                   | 0.272<br>(10.899)***   | 0.255<br>(11.408)***   | 0.412<br>(19.261)***   | 0.402<br>(22.155)***   |
| NOUTGOING               | -0.213<br>(-10.339)**  | -0.162<br>(-4.518)***  | 0.014<br>(-4.933)***   | 0.042<br>(1.328)       |
| COUNT                   | 0.569<br>(128.765)***  | 0.604<br>(147.789)***  | 0.523<br>(126.895)***  | 0.539<br>(131.678)***  |
| VARCAP                  | -0.085<br>(-50.816)*** | -0.075<br>(-57.790)*** | -0.057<br>(-58.108)*** | -0.054<br>(-64.019)*** |
| P                       | 0.021<br>(-5.089)***   | 0.141<br>(3.998)**     | -0.006<br>(-0.171)     | -0.012<br>(-0.405)     |
| Adjusted R <sup>2</sup> | 0.786                  | 0.785                  | 0.768                  | 0.789                  |

The table presents the regression results for Scenario1 applied on a second run of Monte Carlo simulations based on random sampling as robustness test. The dependent variable is CATEND measured as the total loss of capital due to contagion as percentage of total capital in the network. Explanatory variables are, CATIN1, CATIN2, LEVIN, NOUTGOING, COUNT, VARCAP and P, the probability for a link to exist between two nodes. Each cell displays the OLS standardized coefficients along with the corresponding *t*-statistics (shown in parentheses). The sample comprises of 18,000 realizations (simulated banking crises). \*, \*\* and \*\*\* denote significance at the 10, 5 and 1 percent level, respectively

**Table 8** Robustness tests: OLS regression analysis for Scenario 2 (Heterogeneous banks with heterogeneous exposures)

|                         | N = 20                 | N = 50                 | N = 80                 | N = 100                |
|-------------------------|------------------------|------------------------|------------------------|------------------------|
| CATIN 1                 | 0.071<br>(23.366)***   | 0.001<br>(0.273)       | 0.008<br>(2.368)***    | 0.006<br>(1.730)*      |
| CATIN2                  | 0.207<br>(20.109)***   | 0.098<br>(8.941)***    | 0.101<br>(9.336)***    | 0.068<br>(6.721)***    |
| LEVIN                   | 0.602<br>(53.999)***   | 0.469<br>(38.008)***   | 0.313<br>(26.051)***   | 0.304<br>(26.494)***   |
| NOUTGOING               | -0.154<br>(-12.833)*** | -0.096<br>(-4.023)***  | -0.064<br>(-2.747)***  | 0.008<br>(0.391)       |
| COUNT                   | 0.459<br>(107.602)***  | 0.567<br>(131.004)***  | 0.590<br>(144.084)***  | 0.609<br>(156.107)***  |
| VARCAP                  | -0.038<br>(-21.431)*** | -0.067<br>(-41.713)*** | -0.053<br>(-54.105)*** | -0.051<br>(-40.597)*** |
| VARLOANS                | -0.220<br>(-42.365)*** | -0.091<br>(-24.628)*** | -0.018<br>(-2.107)**   | -0.009<br>(-12.307)*** |
| P                       | -0.223<br>(-18.398)*** | -0.080<br>(-3.217)***  | 0.092<br>(3.779)***    | 0.061<br>(2.751)***    |
| Adjusted R <sup>2</sup> | 0.817                  | 0.770                  | 0.772                  | 0.800                  |

The table presents the regression results for Scenario2 applied on a second run of Monte Carlo simulations based on random sampling as robustness test. The dependent variable is CATEND measured as the total loss of capital due to contagion as percentage of total capital in the network. Explanatory variables are the constant term CATIN1, CATIN2, LEVIN, NOUTGOING, COUNT, VARCAP, VARLOANS and P, the probability for a link to exist between two nodes. Each cell displays the OLS standardized coefficients along with the corresponding *t*-statistics (shown in parentheses). The sample comprises of 18,000 realizations (simulated banking crises). \*, \*\* and \*\*\* denote significance at the 10, 5 and 1 percent level, respectively

**Table 9** Robustness tests: OLS regression analysis for Scenario 3 (Homogeneous banks with heterogeneous exposures)

|                         | N = 20                 | N = 50                 | N = 80                 | N = 100                |
|-------------------------|------------------------|------------------------|------------------------|------------------------|
| CATIN2                  | 0.200<br>(25.728)***   | 0.153<br>(18.410)***   | 0.127<br>(13.078)***   | 0.157<br>(21.173)***   |
| LEVIN                   | 0.282<br>(34.036)***   | 0.189<br>(22.195)***   | 0.329<br>(32.672)***   | 0.308<br>(37.777)***   |
| NOUTGOING               | -0.187<br>(-16.831)*** | -0.168<br>(-11.145)*** | -0.114<br>(-6.923)***  | -0.145<br>(-11.721)*** |
| COUNT                   | 0.745<br>(190.987)***  | 0.773<br>(184.138)***  | 0.736<br>(164.477)***  | 0.765<br>(190.795)***  |
| VARLOANS                | -0.167<br>(-41.084)*** | -0.137<br>(-33.586)*** | -0.164<br>(-38.890)*** | -0.196<br>(-46.316)*** |
| P                       | -0.217<br>(-19.612)*** | -0.226<br>(-14.785)*** | -0.371<br>(-22.288)*** | -0.323<br>(-25.206)*** |
| Adjusted R <sup>2</sup> | 0.862                  | 0.824                  | 0.789                  | 0.864                  |

The table presents the regression results for Scenario3 applied on a second run of Monte Carlo simulations based on random sampling as robustness test. The dependent variable is CATEND measured as the total loss of capital due to contagion as percentage of total capital in the network. Explanatory variables are the constant term CATIN2, LEVIN, NOUTGOING, COUNT, VARLOANS and P, the probability for a link to exist between two nodes. Each cell displays the OLS standardized coefficients along with the corresponding *t*-statistics (shown in parentheses). The sample comprises of 18,000 realizations (simulated banking crises). \*, \*\* and \*\*\* denote significance at the 10, 5 and 1 percent level, respectively

**Table 10** Robustness tests: OLS regression analysis for Scenario4 (Homogeneous banks with homogeneous exposures)

|                         | N = 20                 | N = 50                 | N = 80                 | N = 100                |
|-------------------------|------------------------|------------------------|------------------------|------------------------|
| CATIN2                  | 0.266<br>(25.631)***   | 0.196<br>(18.103)***   | 0.153<br>(13.994)***   | 0.126<br>(11.749)***   |
| LEVIN                   | 0.163<br>(17.098)***   | 0.247<br>(25.943)***   | 0.283<br>(29.111)***   | 0.357<br>(37.852)***   |
| NOUTGOING               | -0.306<br>(-19.180)*** | -0.220<br>(-14.254)*** | -0.164<br>(-10.529)*** | -0.118<br>(-8.020)***  |
| COUNT                   | 0.616<br>(188.365)***  | 0.600<br>(161.885)***  | 0.609<br>(160.323)***  | 0.565<br>(145.072)***  |
| P                       | -0.150<br>(-9.783)***  | -0.256<br>(-17.542)*** | -0.309<br>(-20.730)*** | -0.371<br>(-26.650)*** |
| Adjusted R <sup>2</sup> | 0.834                  | 0.804                  | 0.790                  | 0.798                  |

The table presents the regression results for Scenario4 applied on a second run of Monte Carlo simulations based on random sampling as robustness test. The dependent variable is CATEND measured as the total loss of capital due to contagion as percentage of total capital in the network. Explanatory variables are the constant term CATIN2, LEVIN, NOUTGOING, COUNT and P, the probability for a link to exist between two nodes.. Each cell displays the OLS standardized coefficients along with the corresponding *t*-statistics (shown in parentheses). The sample comprises of 18,000 realizations (simulated banking crises). \*, \*\* and \*\*\* denote significance at the 10, 5 and 1 percent level, respectively

## References

1. Business and Sustainable Development Commission Report, Better business better world (2017)
2. Oxford Analytica Foundation, The corporate stewardship compass-Guiding values for sustainable development, Project Report for the Caux Round Table of Moral Capitalism, December 2017
3. O. Weber, The financial sector and the SDGs: Interconnections and future directions, CIGI (Center for International Governance) Innovation Papers, No. 201, November 2018
4. C. Alves, V. Boufounou, A. Dellis, C. Pitelis, Toporowski, J., Synthesis Report: empirical analysis for new ways of global engagement, FESSUD (Financialisation, Economy, Society and Sustainable Development) Working Paper Series, No.163, August 2016
5. J. Leventides, K. Loukaki, V.G. Papavassiliou, Simulating financial contagion dynamics in random interbank networks. *J. Econ. Behav. Organ.* **158**, 500–525 (2019)
6. G. Iori, S. Jafarey, F. Padilla, Systemic risk on the interbank market. *J. Econ. Behav. Organ.* **61**, 525–542 (2006)
7. E. Nier, J. Yang, T. Yorulmazer, A. Alentorn, Network models and financial stability. *J. Econ. Dyn. Control.* **31**, 2033–2060 (2007)
8. P. Gai, S. Kapadia, Contagion in financial networks. *Proc. R. Soc. A Math. Phys. Eng. Sci.* **466**, 2401–2423 (2010)
9. R.M. May, N. Arinaminpathy, Systemic risk: the dynamics of model banking systems. *J. R. Soc. Interface* **7**, 823–838 (2010)
10. H. Amini, R. Cont, A. Minca, Resilience to contagion in financial networks. *Math. Financ.* **26**, 329–365 (2016)
11. C. Upper, Simulation methods to assess the danger of contagion in interbank markets. *J. Financ. Stab.* **7**, 111–125 (2011)

12. F. Allen, A. Babus, Networks in finance, Working Paper No. 08-07, Wharton Financial Institutions Center, 2008
13. C. Memmel, A. Sachs, Contagion in the interbank market and its determinants. *J. Financ. Stab.* **9**, 46–54 (2013)
14. O.-M. Georgescu, Contagion in the interbank market: funding versus regulatory constraints. *J. Financ. Stab.* **18**, 1–18 (2015)
15. L. Tonzer, Cross-border interbank networks, banking risk and contagion. *J. Finan. Stabil.* **18**, 19–32 (2015)
16. K. Fink, U. Krüger, B. Meller, L.-H. Wong, The credit quality channel: modeling contagion in the interbank market. *J. Financ. Stab.* **25**, 83–97 (2016)
17. F. Allen, D. Gale, Financial contagion. *J. Polit. Econ.* **108**, 1–33 (2000)
18. F. Caccioli, T.A. Catanach, J.D. Farmer, Heterogeneity, correlations and financial contagion. *Adv. Comp. Syst.* **15**(Suppl 02), 1250058 (2012)
19. D. Ladley, Contagion and risk-sharing on the inter-bank market. *J. Econ. Dyn. Control.* **37**, 1384–1400 (2013)
20. M. Chinazzi, S. Pegoraro, G. Fagiolo, Defuse the bomb: rewiring interbank networks. LEM Working Paper Series, Institute of Economics, Scuola Superiore Sant' Anna (2015)
21. G. Georg, J. Poschmann, Systemic risk in a network model of interbank markets with central bank activity. *Jena Economic Research Papers*, No. 2010, 033, 2010
22. E. Estrada, *The Structure of Complex Networks: Theory and Applications* (Oxford University Press, 2011)
23. G. Iori, G. de Masi, O. Precup, G. Gabbi, G. Caldarelli, A network analysis of the Italian overnight money market. *J. Econ. Dyn. Control.* **32**, 259–278 (2008)
24. A. Krause, S. Giansante, Interbank lending and the spread of bank failures: a network model of systemic risk. *J. Econ. Behav. Organ.* **83**, 583–608 (2012)

# Higher Order Strongly $m$ -convex Functions



Muhammad Aslam Noor and Khalida Inayat Noor

**Abstract** Some new concepts of the  $m$ -convex functions, where  $m \in (0, 1]$  are introduced and studied. Basic properties of  $m$ -convex functions are discussed. New modified Regula Falsi methods are suggested for solving nonlinear equations. Characterizations of the higher order strongly  $m$ -convex functions are investigated under suitable conditions. It is shown that the parallelogram laws for Banach spaces can be obtained as applications of higher order strongly  $m$ -convex functions. Results obtained in this paper can be viewed as refinement and significant improvement of previously known results.

## 1 Introduction

Lin and Fukushima [11] introduced the concept of higher order strongly convex functions and used it in the study of mathematical program with equilibrium constraints. These mathematical programs with equilibrium constraints are defined by a parametric variational inequality or complementarity system and play an important role in many fields such as engineering design, economic equilibrium, and multilevel game. Mishra and Sharma [12] derived the Hermite–Hadamard type inequalities for higher order strongly convex functions. Characterizations of the higher order strongly convex functions discussed in Lin and Fukushima [11] are not correct. These facts and observations motivated Noor and Noor [16] to consider higher order strongly convex function. Several new characterizations of the higher order strongly convex functions were discussed. Parallelogram laws for uniformly

---

Nonlinear Analysis, Differential Equations and Applications (Edits: Themistocles M. Rassias), Springer Volume

---

M. A. Noor (✉) · K. I. Noor  
COMSATS University Islamabad, Islamabad, Pakistan

© Springer Nature Switzerland AG 2021  
Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,  
Springer Optimization and Its Applications 173,  
[https://doi.org/10.1007/978-3-030-72563-1\\_14](https://doi.org/10.1007/978-3-030-72563-1_14)

Banach spaces can be deduced from the definitions of the higher order strongly convex functions, which is itself a novel and interesting property.

Strongly convex functions were introduced and studied by Polyak [20], which play an important part in the optimization theory and related areas. Karmardian [9] used the strongly convex functions to discuss the unique existence of a solution of the nonlinear complementarity problems. Strongly convex functions also played important role in the convergence analysis of the iterative methods for solving variational inequalities and equilibrium problems, see Zu and Marcotte [24]. Nikodem and Pales [14] investigated the characterization of the inner product spaces using the strongly convex functions, which can be viewed as a novel and innovative application. It is also known that the minimum of the strongly convex functions is unique. Qu and Li [21] investigated the exponentially stability of primal-dual gradient dynamics using the concept of strongly convex functions. Awan et al. [3] have derived Hermite–Hadamard type inequalities for various classes of strongly convex functions, which provide upper and lower estimate for the integrand. For more applications and properties of the strongly convex functions, see [1–4, 9–19, 21] and the references therein.

Relevant to the convex set and convex functions, we also have the concept of  $m$ -convex sets and  $m$ -convex functions, which were introduced by Toader [22]. For the properties and other aspects of the  $m$ -convex functions, see [10] and the references therein. Lara et al. [10] introduced and investigated the properties of strongly  $m$ -convex functions. We would like to mention that the concepts of strongly  $m$ -convex [10] are not correct. To overcome these deficiencies of the higher order strongly convex functions and strongly  $m$ -convex functions, we consider some new classes of convex functions, which are called higher order strongly  $m$ -convex functions. Using the techniques and ideas of this paper, one can easily obtain the refine and correct versions of higher order strongly convex functions and strongly  $m$ -convex functions.

In Section 2, we introduce the new concepts of  $m$ -convex functions and discuss their characterizations. It is pointed out that these  $m$ -convex functions are distinctly different the concepts of  $m$ -functions of Toader [22], which are being investigated in recent years. Higher order strongly  $m$ -convex functions are introduced in Section 3. Some results are discussed in Section 4. In Section 5, we discuss some applications of the higher order strongly convex functions. It is shown that the weakly parallelogram laws can be deduced from the definitions, which characterize the uniformly reflex Banach spaces. As special cases, one can obtain various new and refined versions of known results. It is expected that the ideas and techniques of this paper may stimulate further research in this field.

## 2 $m$ -Convex Functions

Let  $K$  be a nonempty closed set in a real Hilbert space  $H$ . We denote by  $\langle \cdot, \cdot \rangle$  and  $\| \cdot \|$  be the inner product and norm, respectively.

**Definition 1 ([11, 19])** A set  $K$  in  $H$  is said to be a convex set, if

$$u + t(v - u) \in K, \quad \forall u, v \in K, t \in [0, 1].$$

We denote by  $m \in (0, 1]$ , unless otherwise specified.

**Definition 2 ([23])** The set  $K$  in  $H$  is said to be  $m$ -convex set, if

$$(1 - t)u + tmv \in K, \quad \forall u, v \in K, t \in [0, 1].$$

Also one can define the concept of  $m_1$ -convex set as:

**Definition 3 ([23])** The set  $K$  in  $H$  is said to be  $m_1$ -convex set, if

$$(1 - t)m_1u + tv \in K, \quad \forall u, v \in K, t \in [0, 1], m_1 \in (0, 1].$$

We remark that  $m$ -convex set and  $m_1$ -convex set are entirely two different generalizations of the so-called convex sets. that is  $m$ -convex set  $\neq$   $m_1$ -convex set.

For example,  $[m_1a, b] \neq [a, mb]$ ,  $m, m_1 \in (0, 1]$ . Consequently these sets have different properties.

In [10], the authors remarked that  $[m_1a, b] = [a, mb]$ , which is not true. In passing, we remark that the concept of strongly  $m$ -convex functions introduced and discussed in [10] is also wrong. Our results can be viewed as the refinement and improvement of the known results.

In this paper, we only consider the  $m$ -convex set  $K_m$ , unless otherwise specified.

We now introduce new concepts of  $m$ -convex functions.

**Definition 4 ([17])** A function  $F$  is said to be  $m$ -convex function, if

$$F((1 - t)u + tmv) \leq (1 - t)F(u) + tF(mv), \quad \forall u, v \in K_m, t \in [0, 1], m \in (0, 1].$$

We now consider the  $m$ -convex functions on the interval  $I = [a, mb]$ ,  $m \in (0, 1]$ .

**Definition 5** Let  $I = [a, mb]$ . Then  $F$  is  $m$ -convex function, if and only if,

$$\begin{vmatrix} 1 & 1 & 1 \\ a & x & mb \\ F(a) & F(x) & F(mb) \end{vmatrix} \geq 0; \quad a \leq x \leq mb.$$

One can easily show that the following are equivalent:

1.  $F$  is  $m$ -convex function.
2.  $F(x) \leq F(a) + \frac{F(mb)-F(a)}{mb-a}(x - a)$ .
3.  $\frac{F(x)-F(a)}{x-a} \leq \frac{F(mb)-F(a)}{mb-a}$ .
4.  $(mb - x)F(a) + (a - mb)F(x) + (x - a)F(mb) \geq 0$ .
5.  $\frac{F(a)}{(mb-a)(a-x)} + \frac{F(x)}{(x-mb)(a-x)} + \frac{F(mb)}{(mb-a)(x-mb)} \leq 0$ ,

where  $x = (1 - t)a + tmb \in [a, mb]$ .

We now suggest an iterative method for solving the nonlinear equations  $F(x) = 0$  in the interval  $[a, mb]$ ,  $m \in (0,1)$  and  $[ma, b]$ ,  $m \in (0,1)$ . From  $F(x) = F(a) + \frac{F(mb)-F(a)}{mb-a}(x - a)$ , and using the fact that  $F(x) = 0$ , we have

$$x = \frac{aF(mb) - mbF(a)}{mb - a}, x \in [a, mb], m \in (0, 1].$$

This enable to suggest the following iterative method for solving the nonlinear equations  $F(x) = .$

**Algorithm 1** For given  $u_0, u_1$ , find the approximate solution  $u_{n+1}$  by the iterative scheme

$$u_{n+2} = \frac{aF(mu_{n+1}) - mu_{n+1}F(u_n)}{mu_{n+1} - u_n}, \quad m \in (0, 1],$$

which is called the Modified Regula Falsi Method.

In a similar way, we can also suggest the following modified Regula Falsi Method.

**Algorithm 2** For given  $u_0, u_1$ , find the approximate solution  $u_{n+1}$  by the iterative scheme

$$u_{n+2} = \frac{mu_nF(u_{n+1}) - u_{n+1}F(mu_n)}{u_{n+1} - mu_n}, \quad m \in (0, 1].$$

It is worth mentioning that Algorithms 1 and 2 are quite different and provide different results.

Using the technique of Toader [22], we can introduce the following concept of  $m$ -convex functions:

**Definition 6 ([22])** A function  $F$  is said to be  $m$ -convex function in the Toader sense, if

$$F((1 - t)u + tmv) \leq (1 - t)F(u) + mtF(v), \quad \forall u, v \in K_m, \quad t \in [0, 1].$$

We would like to point out that these concepts defined in Definitions 2 and 6 are equivalent, if the functions  $F(mv) = mF(v)$ , that is, the function  $F$  is homogeneous. Consequently all the results proved in this paper can be extended for the  $m$ -convex functions in the Toader sense with suitable modifications.

In this section, we consider some basic properties of  $m$ -convex functions. Using the technique of Pavic and Ardic [18], we derive the following result.

**Theorem 1** Let  $I = [a, mb] \subset R$  be an interval containing the zero and let  $m \in (0, ]$  be a constant. Let  $a, b, c \in I$  be a point such that  $a \leq mc \leq b$ . Then the  $m$ -convex function satisfy the inequality



$$\int_a^b F(x)dx \leq \frac{mc - a}{2} F(a) + \frac{b - mc}{2} F(b) + \frac{b - a}{2} F(mc). \tag{1}$$

**Proof** Assume that  $a \leq x \leq mc$ . Then, from (1), we have

$$\begin{aligned} \int_a^{mc} F(x)dx &\leq \int_a^{mc} \frac{mc - x}{mc - a} F(a)dx + \int_a^{mc} \frac{x - a}{mc - a} F(mc)dx \\ &= \frac{mc - a}{2} (f(a) + f(mc)). \end{aligned} \tag{2}$$

In a similar way, for  $mc \leq x \leq b$ , we have

$$\begin{aligned} \int_{mc}^b F(x)dx &\leq \int_{mc}^b \frac{mc - x}{mc - a} F(mc)dx + \int_{mc}^b \frac{x - a}{mc - a} F(b)dx \\ &= \frac{mc - a}{2} (F(b) + F(mc)). \end{aligned} \tag{3}$$

From (2) and (3), we have

$$\begin{aligned} \int_a^b F(x)dx &= \int_a^{mc} F(x)dx + \int_{mc}^b F(x)dx \\ &= \frac{mc - a}{2} (F(a) + F(mc)) + \frac{mc - a}{2} (F(b) + F(mc)) \\ &= \frac{mc - a}{mc - a} F(a) + \frac{b - mc}{2} F(b) + \frac{b - a}{2} F(mc), \end{aligned}$$

the required (1).

*Remark 1* For the interval  $I = [a, b]$  containing the zero, one can choose an point  $c \in [a, b]$  in (3), since  $mc \in [a, b]$ . Using this information, we can obtain the following inequality for the  $m$ -convex functions for the case  $c = a$  or  $b = c$ .

$$\int_a^b F(x)dx \leq \frac{mb - a}{2} F(a) + \frac{b - ma}{2} F(b),$$

from which, we can have

$$\int_a^b F(x)dx \leq \frac{b - a}{2} \{F(a) + F(b)\}.$$

**Theorem 2** Let  $F$  be a strictly  $m$ -convex function. Then any local minimum of  $F$  is a global minimum.

**Proof** Let the  $m$ -convex function  $F$  have a local minimum at  $u \in K_m$ . Assume the contrary, that is,  $F(mv) < F(u)$  for some  $mv \in K_m$ . Since  $F$  is  $m$ -convex, so

$$F(u + t(mv - u)) < tF(mv) + (1 - t)F(u), \quad \text{for } 0 < t < 1.$$

Thus

$$F(u + t(mv - u)) - F(u) < t[F(mv) - F(u)] < 0,$$

from which it follows that

$$F(u + t(mv - u)) < F(u),$$

for arbitrary small  $t > 0$ , contradicting the local minimum.

**Theorem 3** *If the function  $F$  on the  $m$ -convex set  $K_m$  is  $m$ -convex, then the level set  $L_\alpha = \{u \in K_m : F(u) \leq \alpha, \alpha \in \mathbb{R}\}$  is a  $m$ -convex set.*

**Proof** Let  $u, mv \in L_\alpha$ . Then  $F(u) \leq \alpha$  and  $F(mv) \leq \alpha$ . Now,  $\forall t \in (0, 1)$ ,  $w = u + t(mv - u) \in K_m$ , since  $K_m$  is a  $m$ -convex set. Thus, by the  $m$ -convexity of  $F$ , we have

$$\begin{aligned} F(u + t(mv - u)) &\leq (1 - t)F(u) + tF(mv) \\ &\leq (1 - t)\alpha + t\alpha = \alpha, \end{aligned}$$

from which it follows that  $u + t(mv - u) \in L_\alpha$ . Hence  $L_\alpha$  is a  $m$ -convex set.

**Theorem 4** *The function  $F$  is a  $m$ -convex function, if and only if,*

$$\text{epi}(F) = \{(u, \alpha) : u \in K_m : F(u) \leq \alpha, \alpha \in \mathbb{R}\}$$

*is a  $m$ -convex set.*

**Proof** Assume that  $F$  is a  $m$ -convex function. Let  $(u, \alpha), (mv, \beta) \in \text{epi}(F)$ . Then it follows that  $F(u) \leq \alpha$  and  $F(mv) \leq \beta$ . Thus,  $\forall t \in [0, 1]$ ,  $u, mv \in K_m$ , we have

$$\begin{aligned} F(u + t(mv - u)) &\leq (1 - t)F(u) + tF(mv) \\ &\leq (1 - t)\alpha + t\beta, \end{aligned}$$

which implies that

$$(u + t(mv - u), (1 - t)\alpha + t\beta) \in \text{epi}(F).$$

Thus  $\text{epi}(F)$  is a  $m$ -convex set. Conversely, let  $\text{epi}(F)$  be a convex set. Let  $u, mv \in K$ . Then  $(u, F(u)) \in \text{epi}(F)$  and  $(mv, F(mv)) \in \text{epi}(F)$ . Since  $\text{epi}(F)$  is a  $m$ -convex set, we must have

$$(u + t(mv - u), (1 - t)F(u) + tF(mv)) \in \text{epi}(F),$$

which implies that

$$F(u + t(mv - u)) \leq (1 - t)F(u) + tF(mv).$$

This shows that  $F$  is a  $m$ -convex function.

**Theorem 5** *The function  $F$  is quasi  $m$ -convex, if and only if, the level set  $L_\alpha = \{u \in K_m, \alpha \in R : e^{F(u)} \leq \alpha\}$  is a  $m$ -convex set.*

**Proof** Let  $u, mv \in L_\alpha$ . Then  $u, mv \in K_m$  and  $\max(F(u), F(mv)) \leq \alpha$ . Now for  $t \in (0, 1)$ ,  $w = u + t(mv - u) \in K_m$ , We have to prove that  $u + t(mv - u) \in L_\alpha$ . By the quasi  $m$ -convexity of  $F$ , we have

$$F(u + t(mv - u)) \leq \max F(u), F(mv) \leq \alpha,$$

which implies that  $u + t(mv - u) \in L_\alpha$ , showing that the level set  $L_\alpha$  is indeed a  $m$ -convex set.

Conversely, assume that  $L_\alpha$  is a  $m$ -convex set. Then for any  $u, mv \in L_\alpha, t \in [0, 1]$ ,  $u + t(mv - u) \in L_\alpha$ . Let  $u, mv \in L_\alpha$  for

$$\alpha = \max(F(u), F(mv)) \quad \text{and} \quad F(mv) \leq F(u).$$

Then from the definition of the level set  $L_\alpha$ , it follows that

$$F(u + t(mv - u)) \leq \max F(u), F(mv) \leq \alpha.$$

Thus  $F$  is an quasi  $m$ -convex function. This completes the proof.

**Theorem 6** *Let  $F$  be a  $m$ -convex function. Let  $\mu = \inf_{u \in K_m} F(u)$ . Then the set  $E = \{u \in K_m : F(u) = \mu\}$  is a  $m$ -convex set of  $K$ . If  $F$  is a exponentially  $m$ -convex function, then  $E$  is a singleton.*

**Proof** Let  $u, mv \in E$ . For  $0 < t < 1$ , let  $w = u + t(mv - u)$ . Since  $F$  is a  $m$ -convex function, then

$$\begin{aligned} F(w) &= F(u + t(mv - u)) \leq (1 - t)F(u) + tF(mv) \\ &= t\mu + (1 - t)\mu = \mu, \end{aligned}$$

which implies that to  $w \in E$ . and hence  $E$  is a  $m$ -convex set. For the second part, assume to the contrary that  $F(u) = F(mv) = \mu$ . Since  $K$  is a  $m$ -convex set, then for  $0 < t < 1$ ,  $u + t(mv - u) \in K_m$ . Further, since  $F$  is strictly  $m$ -convex function,

$$F(u + t(mv - u)) < (1 - t)F(u) + tF(mv) = (1 - t)\mu + t\mu = \mu.$$

This contradicts the fact that  $\mu = \inf_{u \in K_m} F(u)$  and hence the result follows.

**Theorem 7** If  $F$  is a  $m$ -convex function such that  $F(mv) < F(u)$ ,  $\forall u, mv \in K_m$ , then  $F$  is a strictly quasi  $m$ -convex function.

**Proof** By the  $m$ -convexity of the function  $F$ ,  $\forall u, mv \in K_m$ ,  $m, t \in [0, 1]$ , we have

$$F(u + t(mv - u)) \leq (1 - t)F(u) + tF(mv) < F(u),$$

since  $F(mv) < F(u)$ , which shows that the function  $F$  is strictly quasi  $m$ -convex.

We now discuss some properties of the differentiable  $m$ -convex functions.

**Theorem 8** Let  $F$  be a differentiable function on the  $m$ -convex set  $K_m$ . Then the function  $F$  is  $m$ -convex function, if and only if,

$$F(mv) - F(u) \geq \langle F'(u), mv - u \rangle, \quad \forall mv, u \in K_m. \quad (4)$$

**Proof** Let  $F$  be a  $y$   $m$ -convex function. Then

$$F(u + t(mv - u)) \leq (1 - t)F(u) + tF(mv), \quad \forall u, mv \in K_m,$$

which can be written as

$$F(mv) - F(u) \geq \left\{ \frac{F(u + t(mv - u)) - F(u)}{t} \right\}.$$

Taking the limit in the above inequality as  $t \rightarrow 0$ , we have

$$F(mv) - F(u) \geq \langle F'(u), mv - u \rangle,$$

which is (4), the required result.

Conversely, let (4) hold. Then  $\forall u, mv \in K_m$ ,  $t \in [0, 1]$ ,  $v_t = u + t(mv - u) \in K_m$ , we have

$$F(mv) - F(v_t) \geq \langle F'(v_t), mv - v_t \rangle = (1 - t)\langle F'(v_t), mv - u \rangle. \quad (5)$$

In a similar way, we have

$$F(u) - F(v_t) \geq \langle F'(v_t), u - v_t \rangle = -t\langle F'(v_t), mv - u \rangle. \quad (6)$$

Multiplying (5) by  $t$  and (6) by  $(1 - t)$  and adding the resultant, we have

$$F(u + t(mv - u)) \leq (1 - t)F(u) + tF(mv),$$

showing that  $F$  is a  $m$ -convex function.

Theorem 8 enables us to introduce the concept of the  $m$ -monotone operators, which appears to be new ones.

**Definition 7** The differential  $F'(\cdot)$  is said to be  $m$ -monotone, if

$$\langle F'(u) - F'(mv), u - mv \rangle \geq 0, \quad \forall u, mv \in H.$$

**Definition 8** The differential  $F'(\cdot)$  is said to be pseudo  $m$ -monotone, if

$$\langle F'(u), mv - u \rangle \geq 0, \quad \Rightarrow \langle F'(mv), mv - u \rangle \geq 0, \quad \forall u, mv \in H.$$

From these definitions, it follows that  $m$ -monotonicity implies pseudo  $m$ -monotonicity, but the converse is not true.

**Theorem 9** Let  $F$  be differentiable  $m$ -convex function on the  $m$ convex set  $K_m$ . Then (19) holds, if and only if,  $F'(\cdot)$  satisfies

$$\langle F'(u) - F'(mv), u - mv \rangle \geq 0, \quad \forall u, mv \in K_m. \tag{7}$$

**Proof** Let  $F$  be a  $m$ -convex function on the  $m$ -convex set  $K_m$ . Then, from Theorem 8, we have

$$F(mv) - F(u) \geq \langle F'(u), mv - u \rangle, \quad \forall u, mv \in K_m. \tag{8}$$

Changing the role of  $u$  and  $mv$  in (8), we have

$$F(u) - F(mv) \geq \langle F'(mv), u - mv \rangle, \quad \forall u, mv \in K_m. \tag{9}$$

Adding (7) and (8), we have

$$\langle F'(u) - F'(mv), u - mv \rangle \geq 0,$$

which shows that  $F'$  is a  $m$ -monotone.

Conversely, from (7), we have

$$\langle F'(mv), u - mv \rangle \leq \langle F'(u), u - mv \rangle. \tag{10}$$

Since  $K$  is a  $m$ -convex set,  $\forall u, mv \in K_m, \quad t \in [0, 1] \quad v_t = u + t(mv - u) \in K_m$ .

Taking  $v = v_t$  in (10), we have

$$\langle F'(v_t), u - v_t \rangle \leq \langle F'(u), u - v_t \rangle = -t \langle F'(u), mv - u \rangle,$$

which implies that

$$\langle F'(v_t), mv - u \rangle \geq \langle F'(u), v - u \rangle. \tag{11}$$

Consider the auxiliary function

$$g(t) = e^{F(u+t(mv-u))},$$

from which, we have

$$g(1) = F(mv), \quad g(0) = F(u).$$

Then, from (11), we have

$$g'(t) = \langle F'(v_t), mv - u \rangle \geq \langle F'(u), mv - u \rangle. \tag{12}$$

Integrating (12) between 0 and 1, we have

$$g(1) - g(0) = \int_0^1 g'(t)dt \geq \langle F'(u), mv - u \rangle.$$

Thus it follows that

$$F(mv) - F(u) \geq \langle F'(u), mv - u \rangle,$$

which is the required (8).

### 3 Strongly $m$ -Convex Functions

We now introduce the concept of higher order strongly  $m$ -convex functions and its variant forms.

**Definition 9** A function  $F$  on the convex set  $K_m$  is said to be higher order strongly  $m$ -convex, if there exists a constant  $\mu > 0$ , such that

$$F(u + t(mv - u)) \leq (1-t)F(u) + tF(mv) - \mu\{t^p(1-t) + t(1-t)^p\}\|mv - u\|^p, \tag{13}$$

$$\forall u, mv \in K_m, t \in [0, 1], m \in (0, 1], p > 1.$$

A function  $F$  is said to higher order strongly  $m$ -concave, if and only if,  $-F$  is higher order strongly  $m$ -convex function.

If  $t = \frac{1}{2}$  and  $\mu = 1$ , then

$$F\left(\frac{u+mv}{2}\right) \leq \frac{F(u)+F(mv)}{2} - \mu \frac{1}{2^p} \|mv - u\|^p, \forall u, mv \in K_m, m \in (0, 1], p > 1. \tag{14}$$

The function  $F$  is said to be higher order strongly  $Jm$ -convex function.

For  $m = 1$ , Definition 9 reduces to:

**Definition 10** A function  $F$  on the convex set  $K$  is said to be higher order strongly convex, if there exists a constant  $\mu > 0$ , such that

$$F(u + t(v - u)) \leq (1 - t)F(u) + tF(v) - \mu\{t^p(1 - t) + t(1 - t)^p\}\|v - u\|^p, \quad \forall u, v \in K, t \in [0, 1]. \tag{15}$$

A function  $F$  is said to higher order strongly concave, if and only if,  $-F$  is higher order strongly convex.

If  $t = \frac{1}{2}$  and  $\mu = 1$ , then

$$F\left(\frac{u + v}{2}\right) \leq \frac{F(u) + F(v)}{2} - \mu \frac{1}{2^p} \|v - u\|^p, \quad \forall u, v \in K, p > 1. \tag{16}$$

The function  $F$  is said to be higher order strongly  $J$ -convex function.

**Definition 11** A function  $F$  on the  $m$ -convex set  $K_m$  is said to be a higher order strongly affine  $m$ -convex, if there exists a constant  $\mu > 0$ , such that

$$F(u + t(mv - u)) = (1 - t)F(u) + tF(mv) - \mu\{t^p(1 - t) + t(1 - t)^p\}\|mv - u\|^p, \tag{17}$$

$$\forall u, mv \in K_m, t \in [0, 1], P > 1.$$

Note that if a functions is both higher order strongly  $m$ -convex and higher order strongly  $m$ -concave, then it is higher order strongly affine  $m$ -convex function.

A function  $F$  is called higher order strongly quadratic equation, if there exists a constant  $\mu > 0$ , such that

$$F\left(\frac{u + mv}{2}\right) = \frac{F(u) + F(mv)}{2} - \mu \frac{1}{2^p} \|mv - u\|^p, \quad \forall u, mv \in K_m, t \in [0, 1]. \tag{18}$$

This function  $F$  is also called higher order strongly affine  $Jm$ -convex function.

We now discuss some special cases.

**I.** If  $p = 2$ , then the higher order strongly convex function becomes strongly convex functions, that is,

$$F(u + t(mv - u)) \leq (1 - t)F(u) + tF(mv) - \mu t(1 - t)\|mv - u\|^2, \quad \forall u, v \in K_m, t \in [0, 1].$$

For the properties of the strongly convex functions in variational inequalities and equilibrium problems, see Noor [15].

**Definition 12** A function  $F$  on the  $m$ -convex set  $K_m$  is said to be higher order strongly quasi  $m$ -convex, if there exists a constant  $\mu > 0$  such that

$$F(u + t(mv - u)) \leq \max\{F(u), F(mv)\} - \mu\{t^p(1 - t) + t(1 - t)^p\}\|mv - u\|^p, \quad \forall u, v \in K_m, t \in [0, 1], p > 1.$$

**Definition 13** A function  $F$  on the  $m$ -convex set  $K_m$  is said to be higher order strongly log  $m$ -convex, if there exists a constant  $\mu > 0$  such that

$$F(u + t(mv - u)) \leq (F(u))^{1-t}(F(mv))^t - \mu\{t^P(1 - t) + t(1 - t)^P\}\|mv - u\|^P, \\ \forall u, mv \in K_m, t \in [0, 1],$$

where  $F(\cdot) > 0$ .

From the above definitions, we have

$$F(u + t(mv - u)) \leq (F(u))^{1-t}(F(mv))^t - \mu\{t^P(1 - t) + t(1 - t)^P\}\|mv - u\|^P \\ \leq (1 - t)F(u) + tF(mv) - \mu\{t^P(1 - t) + t(1 - t)^P\}\|mv - u\|^P \\ \leq \max\{F(u), F(mv)\} - \mu\{t^P(1 - t) + t(1 - t)^P\}\|mv - u\|^P.$$

This shows that every higher order strongly log  $m$ -convex function is a higher order strongly  $m$ -convex function and every higher order strongly  $m$ -convex function is a higher order quasi  $m$ -convex function. However, the converse is not true.

**Definition 14** An operator  $T : K_m \rightarrow H$  is said to be:

1. Higher order strongly  $m$ -monotone, if and only if, there exists a constant  $\alpha > 0$  such that

$$\langle Tu - T(mv), u - mv \rangle \geq \alpha\|mv - u\|^P, \forall u, mv \in K_m.$$

2. Higher order strongly  $m$ -pseudomonotone, if and only if, there exists a constant  $v > 0$  such that

$$\langle Tu, mv - u \rangle + v\|mv - u\|^P \geq 0 \\ \Rightarrow \\ \langle T(mv), mv - u \rangle \geq 0, \forall u, mv \in K_m.$$

3. Higher order strongly relaxed  $m$ -pseudomonotone, if and only if, there exists a constant  $\mu > 0$  such that

$$\langle Tu, mv - u \rangle \geq 0 \\ \Rightarrow \\ -\langle T(mv), u - mv \rangle + \mu\|mv - u\|^P \geq 0, \forall u, v \in K_m.$$

**Definition 15** A differentiable function  $F$  on the convex set  $K_m$  is said to be higher order strongly pseudo  $m$ -convex function, if and only if, if there exists a constant  $\mu > 0$  such that

$$\langle F'(u), mv - u \rangle + \mu\|mv - u\|^P \geq 0 \Rightarrow F(mv) \geq F(u), \forall u, mv \in K_m.$$



### 4 Properties of Strongly $m$ -Convex Functions

In this section, we consider some basic properties of higher order strongly  $m$ -convex functions.

**Theorem 10** *Let  $F$  be a differentiable function on the  $m$ -convex set  $K_m$ . Then the function  $F$  is higher order strongly  $m$ -convex function, if and only if,*

$$F(mv) - F(u) \geq \langle F'(u), mv - u \rangle + \mu \|mv - u\|^p, \forall mv, u \in K_m. \tag{19}$$

**Proof** Let  $F$  be a higher order strongly  $m$ -convex function on the convex set  $K_m$ . Then

$$F(u+t(mv-u)) \leq (1-t)F(u)+tF(mv)-\mu\{t^p(1-t)+t(1-t)^p\}\|mv-u\|^p, \\ \forall u, mv \in K_m,$$

which can be written as

$$F(mv)-F(u) \geq \left\{ \frac{F(u+t(mv-u)) - F(u)}{t} \right\} + \{t^{p-1}(1-t) + (1-t)^p\} \|mv-u\|^p.$$

Taking the limit in the above inequality as  $t \rightarrow 0$ , we have

$$F(mv) - F(u) \geq \langle F'(u), mv - u \rangle + \mu \|mv - u\|^p, \forall u, mv \in K_m.$$

which is (19), the required result.

Conversely, let (19) hold. Then,  $\forall u, mv \in K, t \in [0, 1], v_t = u+t(mv-u) \in K$ , we have

$$F(mv) - F(v_t) \geq \langle F'(v_t), mv - v_t \rangle + \mu \|mv - v_t\|^p \\ = (1-t)F'(v_t), mv-u + \mu(1-t)^p \|mv-u\|^p, \forall u, mv \in K_m \tag{20}$$

In a similar way, we have

$$F(u) - F(v_t) \geq \langle F'(v_t), u - v_t \rangle + \mu \|u - v_t\|^p \\ = -tF'(v_t), mv-u + \mu t^p \|mv-u\|^p. \tag{21}$$

Multiplying (20) by  $t$  and (21) by  $(1-t)$  and adding the resultant, we have

$$F(u+t(mv-u)) \leq (1-t)F(u)+tF(mv)-\mu\{t^p(1-t)+t(1-t)^p\}\|mv-u\|^p, \\ \forall u, mv \in K_m,$$

showing that  $F$  is a higher order strongly  $m$ -convex function.

**Theorem 11** Let  $F$  be a differentiable higher order strongly  $m$ -convex function on the  $m$ -convex set  $K_m$ . Then  $F'(\cdot)$  is a higher order strongly  $m$ -monotone operator.

**Proof** Let  $F$  be a higher order strongly  $m$ -convex function on the  $m$ -convex set  $K_m$ . Then, from Theorem 10, we have

$$F(mv) - F(u) \geq \langle F'(u), mv - u \rangle + \mu \|mv - u\|^p, \quad \forall u, mv \in K_m. \quad (22)$$

Changing the role of  $u$  and  $mv$  in (22), we have

$$F(u) - F(mv) \geq \langle F'(mv), u - mv \rangle + \mu \|mv - u\|^p, \quad \forall u, mv \in K_m. \quad (23)$$

Adding (22) and (23), we have

$$\langle F'(u) - F'(mv), u - mv \rangle \geq 2\mu \|mv - u\|^p, \quad \forall u, mv \in K_m. \quad (24)$$

which shows that  $F'(\cdot)$  is a higher order strongly monotone operator.

We remark that the converse of Theorem 11 is not true. However, we have the following result.

**Theorem 12** If the differential operator  $F'(\cdot)$  of a differentiable higher order strongly  $m$ -convex function  $F$  is higher order strongly  $m$ -monotone operator, then

$$F(mv) - F(u) \geq \langle F'(u), mv - u \rangle + 2\mu \frac{1}{p} \|mv - u\|^p, \quad \forall u, mv \in K_m. \quad (25)$$

**Proof** Let  $F'(\cdot)$  be a higher order strongly  $m$ -monotone operator. Then, from (24), we have

$$\langle F'(v), u - v \rangle \geq \langle F'(u), u - v \rangle + 2\mu \|v - u\|^p. \quad \forall u, v \in K. \quad (26)$$

Since  $K$  is an convex set,  $\forall u, mv \in K_m, t \in [0, 1], v_t = u + t(mv - u) \in K_m$ . Taking  $v = v_t$  in (26), we have

$$\begin{aligned} \langle F'(v_t), u - v_t \rangle &\leq \langle F'(u), u - v_t \rangle - 2\mu \|mv - u\|^p \\ &= -t \langle F'(u), v - u \rangle - 2\mu t^p \|mv - u\|^p, \end{aligned}$$

which implies that

$$\langle F'(v_t), mv - u \rangle \geq \langle F'(u), mv - u \rangle + 2\mu t^{p-1} \|mv - u\|^p. \quad (27)$$

Consider the auxiliary function

$$g(t) = F(u + t(mv - u)), \quad \forall u, mv \in K_m,$$

from which, we have

$$g(1) = F(mv), \quad g(0) = F(u).$$

Then, from (27), we have

$$g'(t) = \langle F'(v_t), mv - u \rangle \geq \langle F'(u), mv - u \rangle + 2\mu t^{p-1} \|mv - u\|^p. \quad (28)$$

Integrating (28) between 0 and 1, we have

$$\begin{aligned} g(1) - g(0) &= \int_0^1 g'(t) dt \\ &\geq \langle F'(u), mv - u \rangle + 2\mu \frac{1}{p} \|mv - u\|^p. \end{aligned}$$

Thus it follows that

$$F(mv) - F(u) \geq \langle F'(u), mv - u \rangle + 2\mu \frac{1}{p} \|v - u\|^p, \quad \forall u, mv \in K_m,$$

which is the required (25).

We note that, if  $p = 2$ , then Theorem 12 can be viewed as the converse of Theorem 11.

We now give a necessary condition for higher order strongly pseudo  $m$ -convex function.

**Theorem 13** *Let  $F'(\cdot)$  be a higher order strongly relaxed pseudo  $m$ -monotone operator. Then  $F$  is a higher order strongly pseudo  $m$ -convex function.*

**Proof** Let  $F'$  be a higher order strongly relaxed pseudo  $m$ -monotone operator. Then,  $\forall u, mv \in K_m$ ,

$$\langle F'(u), mv - u \rangle \geq 0.$$

implies that

$$\langle F'(mv), mv - u \rangle \geq \mu \|mv - u\|^p, \quad \forall u, mv \in K_m. \quad (29)$$

Since  $K_m$  is an  $m$ -convex set,  $\forall u, mv \in K_m, \quad t \in [0, 1], v_t = u + t(mv - u) \in K$ .

Taking  $v = v_t$  in (29), we have

$$\langle F'(v_t), mv - u \rangle \geq \mu t^{p-1} \|mv - u\|^p. \quad (30)$$

Consider the auxiliary function

$$g(t) = F(u + t(mv - u)) = F(v_t), \quad \forall u, mv \in K_m, t \in [0, 1],$$

which is differentiable, since  $F$  is differentiable function. Then, using (30), we have

$$g'(t) = \langle F'(v_t), mv - u \rangle \geq \mu t^{p-1} \|mv - u\|^p.$$

Integrating the above relation between 0 to 1, we have

$$g(1) - g(0) = \int_0^1 g'(t)dt \geq \frac{\mu}{p} \|mv - u\|^p,$$

that is,

$$F(mv) - F(u) \geq \frac{\mu}{p} \|mv - u\|^p, \quad \forall u, mv \in K_m,$$

showing that  $F$  is a higher order strongly pseudo  $m$ -convex function.

**Definition 16** A function  $F$  is said to be sharply higher order strongly pseudo  $m$ -convex, if there exists a constant  $\mu > 0$  such that

$$\langle F'(u), mv - u \rangle \geq 0$$

$\Rightarrow$

$$F(mv) \geq F(mv + t(u - mv)) + \mu \{t^p(1 - t) + t(1 - t)^p\} \|mv - u\|^p, \quad \forall u, mv \in K_m.$$

**Theorem 14** Let  $F$  be a sharply higher order strongly pseudo  $m$ -convex function on  $K_m$  with a constant  $\mu > 0$ . Then

$$\langle F'(mv), mv - u \rangle \geq \mu \|mv - u\|^p, \quad \forall u, mv \in K_m.$$

**Proof** Let  $F$  be a sharply higher order strongly pseudo  $m$ -convex function on  $K_m$ . Then

$$F(mv) \geq F(mv + t(u - mv)) + \mu \{t^p(1 - t) + t(1 - t)^p\} \|mv - u\|^p, \\ \forall u, mv \in K_m, t \in [0, 1],$$

from which, we have

$$\left\{ \frac{F(mv + t(u - mv)) - F(mv)}{t} \right\} + \mu \{t^{p-1}(1 - t) + (1 - t)^p\} \|mv - u\|^p \geq 0.$$

Taking limit in the above inequality, as  $t \rightarrow 0$ , we have

$$\langle F'(mv), mv - u \rangle \geq \mu \|mv - u\|^p, \quad \forall u, mv \in K_m,$$

the required result.

**Definition 17** A function  $F$  is said to be a pseudo  $m$ convex function, if there exists a strictly positive bifunction  $b(., .)$ , such that

$$\begin{aligned}
 &F(mv) < F(u) \\
 &\Rightarrow \\
 &F(u + t(mv, u)) < F(u) + t(t - 1)b(mv, u), \forall u, mv \in K_m, t \in [0, 1].
 \end{aligned}$$

**Theorem 15** If the function  $F$  is higher order strongly  $m$ -convex function such that  $F(mv) < F(u)$ , then the function  $F$  is higher order strongly pseudo  $m$ -convex.

**Proof** Since  $F(mv) < F(u)$  and  $F$  is higher order strongly  $m$ -convex function, then  $\forall u, mv \in K_m, t \in [0, 1]$ , we have

$$\begin{aligned}
 F(u + t(mv - u)) &\leq F(u) + t(F(mv) - F(u)) - \mu\{t^p(1 - t) + t(1 - t)^p\}\|mv - u\|^p \\
 &< F(u) + t(1 - t)(F(mv) - F(u)) - \mu\{t^p(1 - t) + t(1 - t)^p\}\|mv - u\|^p \\
 &= F(u) + t(t - 1)(F(u) - F(mv)) - \mu\{t^p(1 - t) + t(1 - t)^p\}\|mv - u\|^p \\
 &< F(u) + t(t - 1)b(u, mv) - \mu\{t^p(1 - t) + t(1 - t)^p\}\|mv - u\|^p, \forall u, mv \in K_m,
 \end{aligned}$$

where  $b(u, mv) = F(u) - F(mv) > 0$ . Thus, it show that the function  $F$  is a higher order strongly  $m$ -convex function.

We now discuss the optimality for the differentiable generalized strongly convex functions, which is the main motivation of our next result.

**Theorem 16** Let  $F$  be a differentiable higher order strongly  $m$ -convex function with modulus  $\mu > 0$ . If  $u \in K_m$  is the minimum of the function  $F$ , then

$$F(mv) - F(u) \geq \mu\|mv - u\|^p, \quad \forall u, mv \in K_m. \tag{31}$$

**Proof** Let  $u \in K_m$  be a minimum of the function  $F$ . Then

$$F(u) \leq F(mv), \forall mv \in K_m. \tag{32}$$

Since  $K$  is a  $m$ -convex set, so,  $\forall u, mv \in K_m, t \in [0, 1]$ ,

$$v_t = (1 - t)u + tmv \in K_m.$$

Taking  $v = v_t$  in (32), we have

$$0 \leq \lim_{t \rightarrow 0} \left\{ \frac{F(u + t(mv - u)) - F(u)}{t} \right\} = \langle F'(u), mv - u \rangle. \tag{33}$$

Since  $F$  is differentiable higher order strongly  $m$ -convex function, so

$$F(u+t(mv-u)) \leq F(u)+t(F(mv)-F(u))-\mu\{t^p(1-t)+t(1-t)^p\}\|mv-u\|^p, \\ \forall u, mv \in K_m,$$

from which, using (33), we have

$$F(mv) - F(u) \geq \lim_{t \rightarrow 0} \left\{ \frac{F(u + t(mv - u)) - F(u)}{t} \right\} \\ + \mu\{t^{p-1}(1-t) + (1-t)^p\}\|mv-u\|^p \\ = \langle F'(u), mv-u \rangle + \mu\|mv-u\|^p,$$

the required result (31).

*Remark 2* We would like to mention that, if

$$\langle F'(u), mv-u \rangle + \mu\|mv-u\|^p \geq 0, \quad \forall u, mv \in K_m,$$

then  $u \in K_m$  is the minimum of the function  $F$ .

**Theorem 17** Let  $f$  be a higher order strongly affine  $m$ -convex function. Then  $F$  is a higher order strongly  $m$ -convex function, if and only if,  $g = F - f$  is a  $m$ -convex function.

**Proof** Let  $f$  be a higher order strongly affine  $m$ -convex function, Then

$$f((1-t)u+tmv)=(1-t)f(u)+tf(mv)-\mu\{t^p(1-t)+t(1-t)^p\}\|mv-u\|^p, \\ \forall u, mv \in K_m. \tag{34}$$

From the higher order strongly convexity of  $F$ , we have

$$F((1-t)u+tmv) \leq (1-t)F(u)+tF(mv)-\mu\{t^p(1-t)+t(1-t)^p\}\|mv-u\|^p, \\ \forall u, mv \in K_m. \tag{35}$$

From (34) and (35), we have

$$F((1-t)u+tmv) - f((1-t)u+tmv) \leq (1-t)(F(u) - f(u)) \\ +t(F(mv) - f(mv)), \tag{36}$$

from which it follows that

$$g((1-t)u+tmv) = F((1-t)u+tmv) - f((1-t)u+tmv) \\ \leq (1-t)F(u) + tF(mv) - (1-t)f(u) - tf(mv) \\ = (1-t)(F(u) - f(u)) + t(F(mv) - f(mv)),$$

which show that  $g = F - f$  is a  $m$ -convex function.

The inverse implication is obvious.

## 5 Applications

In this section, we show that the characterizations of uniformly Banach spaces involving the notion of higher order strong  $m$ -convexity.

Setting  $F(u) = \|u\|^p$  in Definition 9, we have

$$\|u+t(mv-u)\|^p \leq (1-t)\|u\|^p+t\|mv\|^p-\mu\{t^p(1-t)+t(1-t)^p\}\|mv-u\|^p, \quad (37)$$

$$\forall u, mv \in K_m, t \in [0, 1].$$

Taking  $t = \frac{1}{2}$  in (37), we have

$$\left\| \frac{u+mv}{2} \right\|^p + \mu \frac{1}{2^p} \|mv-u\|^p \leq \frac{1}{2} \|u\|^p + \frac{1}{2} \|mv\|^p, \quad \forall u, mv \in K_m, \quad (38)$$

which implies that

$$\|u+mv\|^p + \mu \|mv-u\|^p \leq 2^{p-1} \{ \|u\|^p + \|mv\|^p \}, \quad \forall u, mv \in K_m, \quad (39)$$

which is known as the lower parallelogram for the  $l^p$ -spaces. In a similar way, one can obtain the upper parallelogram law as

$$\|u+mv\|^p + \mu \|mv-u\|^p \geq 2^{p-1} \{ \|u\|^p + \|mv\|^p \}, \quad \forall u, mv \in K_m, \quad (40)$$

From Definition 11, we have

$$\|u+mv\|^p + \mu \|mv-u\|^p = 2^{p-1} \{ \|u\|^p + \|mv\|^p \}, \quad \forall u, mv \in K_m, \quad (41)$$

which is known as the parallelogram for the  $l^p$ -spaces.

Note that, for  $m = 1$ , we obtain the parallelogram laws for uniformly Banach spaces. To be more precise, Xi [21] obtained the characterizations of  $p$ -uniform convexity and  $q$ -uniform smoothness of a Banach space via the functionals  $\|\cdot\|^p$  and  $\|\cdot\|^q$ , respectively. Bynum [5] and Chen et al. [6–8] have studied the properties and applications of the parallelogram laws for the Banach spaces. It is interesting to note that these parallelogram laws follow from the concepts of the higher order strongly  $m$ -convex functions, which is surprising and novel applications of the higher order strongly convex functions. For the applications of the parallelogram laws in Banach spaces in prediction theory and applied sciences, see [5–8, 20] and the references therein.

## Conclusion

In this paper, we have introduced and studied a new class of convex functions, which is called higher order strongly  $m$ -convex function. It is shown that several new classes of strongly convex functions can be obtained as special cases of these higher order strongly  $m$ -convex functions. We have studied the basic properties of these functions. We have derived new parallelogram laws as applications of the higher order strongly  $m$ -convex functions. Some known results can be obtained for suitable and appropriate choices of  $m$ , which have been used to characterize the  $p$ -uniform convexity and  $q$ -uniform smoothness of a Banach spaces. The interested readers may explore the applications and other properties of the higher order strongly convex functions in various fields of pure and applied sciences. This is an interesting direction of future research.

**Acknowledgments** The authors would like to thank the Rector, COMSATS University Islamabad, Islamabad, Pakistan, for providing excellent research and academic environments. Authors are grateful to Prof. Dr. Themistocles M. Rassias for his kind invitation and support.

## References

1. M. Adamek, On a problem connected with strongly convex functions. *Math. Inequ. Appl.* **19**(4), 1287–1293 (2016)
2. H. Angulo, J. Gimenez, A.M. Moeos, K. Nikodem, On strongly  $h$ -convex functions. *Ann. Funct. Anal.* **2**(2), 85–91 (2011)
3. M.U. Awan, M.A. Noor, T.-S. Du, K.I. Noor, New refinements of fractional Hermite–Hadamard inequality. *RACSAM*, **113**, 21–29 (2019)
4. A. Azcar, J. Gimenez, K. Nikodem, J.L. Sanchez, On strongly midconvex functions. *Opuscula Math.* **31**(1), 15–26 (2011)
5. W.L. Bynum, Weak parallelogram laws for Banach spaces. *Can. Math. Bull.* **19**, 269–275 (1976)
6. R. Cheng, C.B. Harris, Duality of the weak parallelogram laws on Banach spaces. *J. Math. Anal. Appl.* **404**, 64–70 (2013)
7. R. Cheng, W.T. Ross, Weak parallelogram laws on Banach spaces and applications to prediction. *Period. Math. Hung.* **71**, 45–58 (2015)
8. R. Cheng, J. Mashreghi, W.T. Ross, Optimal weak parallelogram constants for  $L_p$  space. *Math. Inequal. Appl.* **21**(4), 1047–1058 (2018)
9. S. Karamardian, The nonlinear complementarity problems with applications, Part 2. *J. Optim. Theory Appl.* **4**(3), 167–181 (1969)
10. T. Lara, N. Merentes, R. Quintero, E. Rosales, Strongly  $m$ -convex functions. *Math. Aeterna*, **5**(3), 521–535 (2015)
11. G.H. Lin, M. Fukushima, Some exact penalty results for nonlinear programs and mathematical programs with equilibrium constraints. *J. Optim. Theory Appl.* **118**(1), 67–80 (2003)
12. S.M. Mishra, N. Sharma, On strongly generalized convex functions of higher order. *Math. Inequal. Appl.* **22**(1), 111–121 (2019)
13. C.P. Niculescu, L.E. Persson, *Convex Functions and their Applications* (Springer, New York, 2018)



14. K. Nikodem, Z.S. Pales, Characterizations of inner product spaces by strongly convex functions. *Banach J. Math. Anal.*, **1**, 83–87 (2011)
15. M.A. Noor, Fundamentals of equilibrium problems. *Math. Inequal. Appl.* **9**(3), 529–566 (2006)
16. M.A. Noor, K.I. Noor, On generalized strongly convex functions involving bifunction. *Appl. Math. Inf. Sci.* **13**(3), 411–416 (2019)
17. M.A. Noor, K.I. Noor, Properties of exponentially  $m$ -convex functions, in *Nonlinear Analysis and Global Optimization* eds. by P. Pardalos, T.M. Rassias (Springer, Berlin, 2019)
18. Z. Pavic, A. Ardic, The most important inequalities of  $m$ -convex functions. *Turkish J. Math.* **41**, 625–635 (2017)
19. J. Pecric, F. Proschan, Y. I. Tong, *Convex Functions, Partial Ordering and Statistical Applications* (Academic Press, New York, 1992)
20. B.T. Polyak, Existence theorems and convergence of minimizing sequences in extremum problems with restrictions. *Soviet Math. Dokl.* **7**, 2–75 (1966)
21. G. Qu, N. Li, On the exponentially stability of primal-dual gradient dynamics, *IEEE Control Syst. Letters* **3**(1), 43–48 (2019)
22. G. Toader, Some generalizations of the convexity. *Proc. Colloq. Approx. Optim. Cluj-Naploca (Romania)*, 329–338 (1984)
23. S. Toader, The order of a star-convex function. *Bull. Appl. Comp. Math. (Budapest)* **85-B** (BAM-1473), 347–350 (1998)
24. D.L. Zu, P. Marcotte, Co-coercivity and its role in the convergence of iterative schemes for solving variational inequalities. *SIAM J. Optim.* **6**(3), 714–726 (1996)

# Characterizations of Higher Order Strongly Generalized Convex Functions



Muhammad Aslam Noor, Khalida Inayat Noor, and Michael Th. Rassias

**Abstract** In this paper, we define and consider some new concepts of the higher order strongly generalized convex functions with respect to two arbitrary functions. Some properties of the higher order strongly generalized convex functions are investigated under suitable conditions. It is shown that the operator parallelogram laws for the characterization of uniformly Banach spaces can be obtained as a novel applications of higher order strongly affine functions. It is shown that the optimality conditions of the higher order strongly generalized convex functions are characterized by a new class of variational inequalities. Some special cases also discussed. Results obtained in this paper can be viewed as significant refinement and improvement of previously known results.

## 1 Introduction

Convexity theory is a branch of mathematical sciences with a wide range of applications in industry, physical, social, regional and engineering sciences. The general theory of the convexity started soon after the introduction of differential and integral calculus by Newton and Leibnitz, although some individual optimization problems had been investigated before that. To be more specific, Bernoulli's brothers (1697) were the first, who considered the variational problems in mathematical terms. It

---

In: *Nonlinear Analysis, Differential Equations and Applications* (Edit: Themistocles M. Rassias)

---

M. A. Noor (✉) · K. I. Noor  
COMSATS University Islamabad, Islamabad, Pakistan

M. T. Rassias  
Institute of Mathematics, University of Zurich, Zurich, Switzerland  
Moscow Institute of Physics and Technology, Dolgoprudny, Russia

Institute for Advanced Study, Program in Interdisciplinary Studies, Princeton, NJ, USA  
e-mail: [michail.rassias@math.uzh.ch](mailto:michail.rassias@math.uzh.ch)

is worth mentioning that the first phase of the development of calculus of variations was characterized by a combination of philosophical concepts, mathematical methods and physical problems. Euler (eighteenth century) created a new branch of mathematics known as the calculus of variations. Motivated by geometrical considerations, Euler deduced his first principle which is now referred to as Euler's differential equation for the determination of maximizing and minimizing arcs. By variational principles, we mean maximum and minimum problems arising in game theory, mechanics, geometrical optics, general relativity theory, economics, transportation, differential geometry and related areas. The Hamiltonian-Jacobi theory represents a general framework for the mathematical description of the propagation of actions in nature and optimal modeling of control processes in daily life. It is known that the gauge field theories are a continuation of Einstein's concept of describing physical effects mathematically in terms of differential geometry. These theories play a fundamental role in the modern theory of elementary particles and are right tool of building up a unified theory of elementary particles, which includes all kind of known interactions. For example, the Weinberg-Salam theory unifies weak and electromagnetic interactions. It is also known that the variational formulation of field theories allows for a degree of unification absent their versions in terms of differential equations. Convexity plays an important part in the existence and stability of soliton, which occur in almost every branch of physics.

Variational inequalities represent the optimality conditions for the differentiable convex functions on the convex sets in a normed space. which were introduced and considered in early 1960s by Stampacchia [41], Variational inequalities combine both theoretical and algorithmic advances with new and novel domain of applications. Analysis of these problems requires a blend of techniques from convex analysis, functional analysis and numerical analysis. In recent years, considerable interest has been shown in developing various extensions and generalizations of variational inequalities, both for their own sake and their applications.

Mohsen et al. [19] and Noor and Noor [30–34, 37] introduced the concept of higher order strongly convex functions and studied their properties. These results can be viewed as a significant refinement of the results of Adamek [1], Alabdali et al. [2] and Lin and Fukushima [18]. Higher order strongly convex functions include the strongly convex functions, which were introduced and studied by Polyak [39]. Karmardian [15] used the strongly convex functions to discuss the unique existence of a solution of the nonlinear complementarity problems. Awan et al. [4, 5] have derived Hermite–Hadamard type inequalities for various classes of strongly convex functions, which provide upper and lower estimate for the integrand. For the applications of strongly convex functions in optimization, variational inequalities and other branches of pure and applied sciences, see [1–8, 8–18, 20, 22–32, 34, 43–45] and the references therein.

It is known that the properties of the convex functions may not hold, in general, when the convex set is non-convex. Recently, the concept of convexity has been generalized in several directions. A significant generalization of the convex set is the introduction of the  $(h, g)$ -convex set and  $(h, g)$ -convex [14, 28, 29] function involving two arbitrary functions  $h$  and  $g$ , (say). These  $(h, g)$  convex functions are

called generalized convex functions. It has been shown that the generalized convex functions enjoy some nice properties which convex functions have. We would like to emphasize that the  $(h, g)$ -convex set and generalized convex functions may not be convex sets and convex functions. For  $h = g$ , generalized functions reduce to  $g$ -convex functions, which were introduced and studied by Youness [43]. It is known [24] that the minimum of a differentiable  $g$ -convex function on the  $g$ -convex set can be characterized by a class of variational inequalities, which is known as general (Noor) variational inequalities introduced and studied by Noor [22, 24, 25] in 1988. For suitable and appropriate choices of the arbitrary functions  $h$  and  $g$ , one can obtain a wide class of known and new classes of convex functions and related general variational inequalities. This clearly shows that the  $(h, g)$ -convex sets and the generalized convex functions are quite flexible and unifying ones. For the formulation, applications, numerical methods, sensitivity analysis and other aspects of general variational inequalities, see [11–14, 19–29, 34, 35, 37] and the references therein.

Noor and Noor [32] have considered and studied the strongly generalized functions and studied their properties. Motivated by the work of Noor et al. [30–34], Alabdali et al. [2] and Lin and Fukushima [18], we introduce and investigate a new class of higher order strongly generalized convex functions with respect to two arbitrary functions, which is the main motivation of this paper. Several new concepts of monotonicity are introduced. We establish the relationship between these classes and derive some new results under some mild conditions. It is shown that (operator) parallelogram laws, which characterize the uniform Banach spaces can be obtained from these definitions. We have shown that the minimum of the higher order strongly generalized convex functions on the generalized convex set can be characterized by higher order strongly generalized variational inequalities. Some new special cases are discussed, which can be viewed as novel and interesting applications of the higher order strongly generalized convex functions.

For the sake of readers' convenience, we include all the relevant details. In Section 2, we recall the basic concepts and results. Generalized convex functions and their properties are discussed in Section 3. In Section 4, higher order strongly generalized convex are introduced. It is shown that several important special cases are discussed, which can be viewed as significant refinement of the previous known results. Properties of the differential higher order strongly generalized convex are proved. It is shown that the parallelogram laws can be obtained from the generalized affine functions, which is itself a novel application. These concepts are discussed in Section 5. Connection with generalized variational inequalities is investigated with the optimality conditions of the differentiable generalized convex functions, which is the main motivation of Section "Conclusion". We would like to emphasize that the results obtained and discussed in this paper may motivate and bring a large number of novel, innovative and potential applications, extensions and interesting topics in these areas. We have given only a brief introduction of higher order strongly generalized convex functions and applications. The interested reader is advised to explore this field further and discover novel and fascinating applications of the generalized convex functions in other areas of sciences.

## 2 Formulations and Basic Facts

Let  $K$  be a nonempty closed set in a real Hilbert space  $H$ . We denote by  $\langle \cdot, \cdot \rangle$  and  $\| \cdot \|$  be the inner product and norm, respectively.

**Definition 1** ([10, 20]) A set  $K$  in  $H$  is said to be a convex set, if

$$u + t(v - u) \in K, \quad \forall u, v \in K, t \in [0, 1].$$

**Definition 2** A function  $F$  is said to be convex function, if

$$F((1 - t)u + tv) \leq (1 - t)F(u) + tF(v); \forall u, v \in K, t \in [0, 1]. \tag{1}$$

It is well known that  $u \in K$  of a differential convex functions  $F$  is equivalent to finding  $u \in K$  such that

$$\langle F'(u), v - u \rangle \geq 0, \forall v \in K, \tag{2}$$

which is called the variational inequality, introduced and studied by Stampacchia [41]. Variational inequalities can be regarded as a novel and significant extension of variational principles, the origin of which can be traced back to Euler, Lagrange, Newton, and Bernoulli brothers.

We would like to mention that the underlying the set may not be a convex set in many important applications. To overcome this drawback, the set can be made convex set with respect to an arbitrary function, which is called general convex set.

We would like to mention that the underlying the set may not be a convex set in many important applications. To overcome this drawback, the set can be made convex set with respect to two arbitrary functions, which is called a generalized or  $(h, g)$ -convex set [14, 27, 28].

**Definition 3** ([14, 28]) The set  $K \subseteq H$  is said to be a  $(h, g)$ -convex set, if there are two functions  $h$  and  $g$  such as

$$(1 - t)h(u) + tg(v) \in K; \quad \forall u, v \in K, t \in [0, 1]. \tag{3}$$

We now discuss some special cases of the  $(h, g)$ -convex set  $K \subseteq D$ .

- (I). If  $g(u) = I(u) = u = h(u)$ , the identity operator, then  $(h, g)$ -convex set reduces to the classical convex set. Clearly every convex set is a  $(h, g)$ -convex set, but the converse is not true.
- (II). If  $h(u) = I(u) = u$ , then the  $(h, g)$ -convex set becomes the  $g$ -convex set, that is,

**Definition 4** The set  $K$  is said to be  $g$ -convex set, if

$$(1 - t)u + tg(v) \in K \subseteq D, \quad \forall u, v \in K \subseteq D, t \in [0, 1],$$

which was introduced and studied by Noor [27]. Cristescu et al. [9] discussed various applications of the general convex sets related to the necessity of adjusting investment or development projects out of environmental or social reasons. For example, the easiest manner of constructing this kind of convex sets comes from the problem of modernizing the railway transport system. Shape properties of the general convex sets with respect to a projection are investigated.

(III). If  $g(u) = I(u) = u$ , then the  $(h, g)$ -convex set becomes the  $h$ -convex set, that is,

**Definition 5** The set  $K$  is said to be  $h$ -convex set, if

$$(1 - t)h(u) + tv \in K \subseteq D, \quad \forall u, v \in K \subseteq D, t \in [0, 1],$$

which is mainly due to Noor [28].

For the sake of simplicity, we always assume that function  $F : D \rightarrow R$  and  $K \cup h(K) \cup g(K) \subseteq D$ . If  $K$  is  $(g, h)$ -convex set, then this condition becomes  $K \subseteq D$  and  $(h, g)$ -convex is called generalized convex functions unless otherwise specified.

**Definition 6** A function  $F$  is said to be a generalized convex function on the  $(h, g)$ -convex set  $K \subseteq D$ , if there exist two arbitrary non-negative functions  $h, g$  such that

$$F((1-t)h(u)+tg(v)) \leq (1-t)F(h(u))+tF(g(v)), \quad \forall u, v \in K \subseteq D, t \in [0, 1](4)$$

The generalized convex functions were introduced by Noor [28]. Noor [28] proved that the minimum  $u \in K \subseteq D$  of a differentiable generalized convex functions  $F$  can be characterized by the class of variational inequalities of the type:

$$\langle F'(h(u)), g(v) - h(u) \rangle \geq 0, \quad \forall v \in K \subseteq D, \tag{5}$$

which is known as the extended general variational inequalities. For the applications of the general variational inequalities in various branches of pure and applied sciences, see [22–25, 27–29, 37] and the references therein.

We now define the generalized convex functions on the interval  $K = I_{hg} = [h(a), g(b)]$ .

**Definition 7** Let  $I_{hg} = [h(a), g(b)]$ . Then  $F$  is a generalized convex function, if and only if,

$$\left| \begin{array}{ccc} 1 & 1 & 1 \\ h(a) & x & g(b) \\ F(h(a)) & F(x) & F(g(b)) \end{array} \right| \geq 0; \quad h(a) \leq x \leq g(b).$$

One can easily show that the following are equivalent:

1.  $F$  is a generalized convex function.
2.  $F(x) \leq F(h(a)) + \frac{F(g(b))-F(h(a))}{g(b)-h(a)}(x - h(a))$ .
3.  $\frac{F(x)-F(h(a))}{x-h(a)} \leq \frac{F(g(b))-F(h(a))}{g(b)-h(a)}$ .
4.  $(g(b) - x)F(h(a)) + (h(a) - g(b))F(x) + (x - h(a))F(g(b)) \geq 0$ .
5.  $\frac{F(a)}{(g(b)-h(a))(h(a)-x)} + \frac{F(x)}{(x-g(b))(a-x)} + \frac{F(g(b))}{(g(b)-h(a))(x-g(b))} \leq 0$ ,

where  $x = (1 - t)h(a) + tg(b) \in [h(a), g(b)]$ .

**Definition 8** The function  $F$  on the  $(h, g)$ -convex set  $K$  is said to be a generalized quasi-convex, if

$$F(h(u) + t(g(v) - h(u))) \leq \max\{F(h(u)), F(g(v))\}, \quad \forall u, v \in K, t \in [0, 1].$$

**Definition 9** The function  $F$  on the  $(h, g)$ -convex set  $K$  is said to be a generalized log-convex, if

$$F(h(u) + t(g(v) - h(u))) \leq (F(h(u))^{1-t}(F(g(v))))^t, \quad \forall u, v \in K, t \in [0, 1],$$

where  $F(\cdot) > 0$ .

From the above definitions, we have

$$\begin{aligned} F(h(u) + t(g(v) - h(u))) &\leq (F(h(u))^{1-t}(F(g(v))))^t \\ &\leq (1 - t)F(h(u)) + tF(g(v)) \\ &\leq \max\{F(h(u)), F(g(v))\}, \quad \forall u, v \in K, t \in [0, 1]. \end{aligned}$$

This shows that every generalized log-convex function is a generalized convex function and every generalized convex function is a generalized quasi-convex function. However, the converse is not true.

### 3 Generalized Convex Functions

In section, we now consider some basic properties of generalized convex functions.

**Theorem 1** Let  $F$  be a strictly generalized convex function. Then any local minimum of  $F$  is a global minimum.

**Proof** Let the strictly generalized convex function  $F$  have a local minimum at  $u \in K$ . Assume the contrary, that is,  $F(g(v)) < F(h(u))$  for some  $g(v) \in K$ . Since  $F$  is strictly generalized convex function, so

$$F(h(u) + t(g(v) - h(u))) < tF(g(v)) + (1 - t)F(h(u)), \quad \text{for } 0 < t < 1.$$

Thus

$$F(h(u) + t(g(v) - h(u))) - F(h(u)) < -t[F(g(v)) - F(h(u))] < 0,$$

from which it follows that

$$F(h(u) + t(g(v) - h(u))) < F(h(u)),$$

for arbitrary small  $t > 0$ , contradicting the local minimum.

**Theorem 2** *If the function  $F$  on the  $(h, g)$ -convex set  $K$  is generalized convex, then the level set*

$$L_\alpha = \{u \in K : F(h(u)) \leq \alpha, \quad \alpha \in R\}$$

*is a  $(h, g)$ -convex set.*

**Proof** Let  $h(u), g(v) \in L_\alpha$ . Then  $F(h(u)) \leq \alpha$  and  $F(g(v)) \leq \alpha$ . Now,  $\forall t \in (0, 1)$ ,  $g(w) = u + t(g(v) - h(u)) \in K$ , since  $K$  is a  $(h, g)$ -convex set. Thus, by the generalized convexity of  $F$ , we have

$$F(h(u) + t(g(v) - h(u))) \leq (1 - t)F(h(u)) + tF(g(v)) \leq (1 - t)\alpha + t\alpha = \alpha,$$

from which it follows that  $h(u) + t(g(v) - h(u)) \in L_\alpha$ . Hence  $L_\alpha$  is a  $(h, g)$ -convex set.

**Theorem 3** *The function  $F$  is generalized convex function, if and only if,*

$$epi(F) = \{(h(u), \alpha) : h(u) \in K : F(h(u)) \leq \alpha, \alpha \in R\}$$

*is a  $(h, g)$ -convex set.*

**Proof** Assume that  $F$  is generalized convex function. Let

$$(h(u), \alpha), \quad (g(v), \beta) \in epi(F).$$

Then it follows that  $F(h(u)) \leq \alpha$  and  $F(g(v)) \leq \beta$ . Hence, we have

$$F(h(u) + t(g(v) - h(u))) \leq (1 - t)F(h(u)) + tF(g(v)) \leq (1 - t)\alpha + t\beta,$$

which implies that

$$((1 - t)h(u) + tg(v), (1 - t)\alpha + t\beta) \in epi(F).$$

Thus  $epi(F)$  is a  $(h, g)$ -convex set. Conversely, let  $epi(F)$  be a  $(h, g)$ -convex set. Let  $h(u), g(v) \in K$ . Then  $(h(u), F(h(u))) \in epi(F)$  and  $(g(v), F(g(v))) \in epi(F)$ . Since  $F$  is a  $(h, g)$ -convex set, we must have

$$(h(u) + t(g(v) - h(u)), (1 - t)F(h(u)) + tF(g(v))) \in epi(F),$$



which implies that

$$F((1 - t)h(u) + tg(v)) \leq (1 - t)F(h(u)) + tF(g(v)).$$

This shows that  $F$  is a generalized convex function.

**Theorem 4** *The function  $F$  is a generalized quasi-convex, if and only if, the level set*

$$L_\alpha = \{h(u) \in K, \alpha \in R : F(h(u)) \leq \alpha\}$$

*is a  $(h, g)$ -convex set.*

**Proof** Let  $h(u), g(v) \in L_\alpha$ . Then  $h(u), g(v) \in K$  and  $\max(F(h(u)), F(g(v))) \leq \alpha$ . Now for  $t \in (0, 1)$ ,  $g(w) = h(u) + t(g(v) - h(u)) \in K$ . We have to prove that  $h(u) + t(g(v) - h(u)) \in L_\alpha$ . By the generalized convexity of  $F$ , we have

$$F(h(u) + t(g(v) - h(u))) \leq \max(F(h(u)), F(g(v))) \leq \alpha,$$

which implies that  $h(u) + t(g(v) - h(u)) \in L_\alpha$ , showing that the level set  $L_\alpha$  is indeed a  $(h, g)$ -convex set.

Conversely, assume that  $L_\alpha$  is a  $(h, g)$ -convex set. Then,  $\forall h(u), g(v) \in L_\alpha, t \in [0, 1], h(u) + t(g(v) - h(u)) \in L_\alpha$ . Let  $h(u), g(v) \in L_\alpha$  for

$$\alpha = \max(F(h(u)), F(g(v))) \quad \text{and} \quad F(g(v)) \leq F(h(u)).$$

Then from the definition of the level set  $L_\alpha$ , it follows that

$$F(h(u) + t(g(v) - h(u))) \leq \max(F(h(u)), F(g(v))) \leq \alpha.$$

Thus  $F$  is an generalized quasi-convex function. This completes the proof.

**Theorem 5** *Let  $F$  be a generalized convex function. Let  $\mu = \inf_{h(u) \in K} F(u)$ . Then the set*

$$E = \{h(u) \in K : F(h(u)) = \mu\}$$

*is a  $(h, g)$ -convex set of  $K$ . If  $F$  is strictly generalized convex function, then  $E$  is a singleton.*

**Proof** Let  $h(u), g(v) \in E$ . For  $0 < t < 1$ , let  $g(w) = h(u) + t(g(v) - h(u))$ . Since  $F$  is a generalized convex function, then

$$\begin{aligned} F(g(w)) &= F(h(u) + t(g(v) - h(u))) \\ &\leq (1 - t)F(h(u)) + tF(g(v)) = t\mu + (1 - t)\mu = \mu, \end{aligned}$$

which implies  $g(w) \in E$  and hence  $E$  is a  $(h, g)$ -convex set. For the second part, assume to the contrary that  $F(h(u)) = F(g(v)) = \mu$ . Since  $K$  is a  $(h, g)$ -convex set, then for  $0 < t < 1$ ,  $h(u) + t(g(v) - h(u)) \in K$ . Further, since  $F$  is a strictly generalized convex function, so

$$F(h(u) + t(g(v) - h(u))) < (1 - t)F(h(u)) + tF(g(v)) = (1 - t)\mu + t\mu = \mu.$$

This contradicts the fact that  $\mu = \inf_{h(u) \in K} F(u)$  and hence the result follows.

**Theorem 6** *If the function  $F$  is a generalized convex such that*

$$F(g(v)) < F(h(u)), \forall h(u), g(v) \in K,$$

*then  $F$  is a strictly generalized quasi-convex function.*

**Proof** By the generalized convexity of the function  $F$ , we have

$$F(h(u) + t(g(v) - h(u))) \leq (1 - t)F(h(u)) + tF(g(v)), \forall h(u), g(v) \in K, t \in [0, 1] \\ < F(h(u)),$$

since  $F(g(v)) < F(h(u))$ , which shows that the function  $F$  is a strictly generalized quasi-convex.

## 4 Higher Order Strongly Generalized Convex Functions

In this section, some new classes of higher order strongly generalized convex functions and higher order strongly affine  $(h, g)$  functions on the  $(h, g)$ -convex set  $K \subseteq D$ .

**Definition 10** A function  $F$  on the  $(h, g)$ -convex set  $K \subseteq D$  is said to be higher order strongly generalized convex with respect to two functions  $h$  and  $g$ , if there exists a constant  $\mu > 0$ , such that

$$F(h(u) + t(g(v) - h(u))) \\ \leq (1 - t)F(h(u)) + tF(g(v)) - \mu\{t^p(1 - t) + t(1 - t)^p\}\|g(v) - h(u)\|^p, \\ \forall u, v \in K \subseteq D, t \in [0, 1], p > 1. \tag{6}$$

A function  $F$  is said to be higher order strongly generalized convex with respect to two functions  $h$  and  $g$ , if and only if,  $-F$  is a higher order strongly generalized convex function with respect to two functions  $h$  and  $g$ .

If  $t = \frac{1}{2}$  in Definition 10, then one gets the generalized Jensen-type property called higher order strongly generalized convex with respect to two functions  $h$  and  $g$ .

We now discuss some special cases.

(IV). If  $p = 2$ . then the Definition 10 reduces to:

**Definition 11** A function  $F$  on the  $(h, g)$ -convex set  $K \subseteq D$  is said to be strongly generalized convex with respect to two functions  $h$  and  $g$ , if there exists a constant  $\mu > 0$ , such that

$$\begin{aligned}
 F(h(u) + t(g(v) - h(u))) &\leq (1 - t)F(h(u)) + tF(g(v)) \\
 &\quad - \mu\{t(1 - t)\}\|g(v) - h(u)\|^2, \\
 \forall u, v \in K \subseteq D, t \in [0, 1], p > 1.
 \end{aligned}
 \tag{7}$$

For the characterizations and properties of the strongly generalized convex functions, see Noor and Noor [32].

(V). If  $h(u) = I(u) = u$ , then the higher order strongly generalized convex with respect to two functions  $h$  and  $g$ , becomes strongly  $g$ -convex functions, that is,

$$\begin{aligned}
 F(u + t(g(v) - u)) &\leq (1 - t)F(u) + tF(g(v)) - \mu\{t^p(1 - t) \\
 &\quad + t(1 - t)^p\}\|g(v) - u\|^p, \\
 \forall u, v \in K \subseteq D, t \in [0, 1].
 \end{aligned}$$

For the properties of the higher order strongly generalized convex functions in variational inequalities and equilibrium problems, see Noor [27].

(VI). If  $g(u) = I(u) = u$ , then the higher order strongly generalized convex with respect to two functions  $h$  and  $g$ , becomes higher order strongly  $h$ -convex functions, that is,

$$\begin{aligned}
 F(h(u) + t(v - h(u))) &\leq (1 - t)F(h(u)) + tF(v) - \mu\{t^p(1 - t) \\
 &\quad + t(1 - t)^p\}\|v - h(u)\|^p, \\
 \forall u, v \in K \subseteq D, t \in [0, 1].
 \end{aligned}$$

**Definition 12** A function  $F$  on the  $(h, g)$ -convex set  $K \subseteq D$  is said to be a higher order strongly generalized quasi-convex with respect to two functions  $h$  and  $g$ , if there exists a constant  $\mu > 0$  such that

$$\begin{aligned}
 F(h(u) + t(g(v) - h(u))) &\leq \max\{F(h(u)), F(g(v))\} - \mu\{t^p(1 - t) \\
 &\quad + t(1 - t)^p\}\|g(v) - h(u)\|^p, \\
 \forall u, v \in K \subseteq D, t \in [0, 1].
 \end{aligned}$$

**Definition 13** A function  $F$  on the  $(h, g)$ -convex set  $K \subseteq D$  is said to be higher order strongly generalized log-convex with respect to two functions  $h$  and  $g$ , if there exists a constant  $\mu > 0$  such that

$$\begin{aligned}
 F(h(u) + t(g(v) - h(u))) &\leq (F(h(u)))^{1-t} (F(g(v)))^t - \mu\{t^p(1-t) \\
 &\quad + t(1-t)^p\} \|g(v) - h(u)\|^p, \\
 \forall u, v \in K \subseteq D, t \in [0, 1],
 \end{aligned}$$

where  $F(\cdot) > 0$ .

From the above definitions, we have

$$\begin{aligned}
 F(h(u) + t(g(v) - h(u))) &\leq (F(h(u)))^{1-t} (F(g(v)))^t - \mu\{t^p(1-t) + t(1-t)^p\} \|g(v) - h(u)\|^p \\
 &\leq (1-t)F(h(u)) + tF(g(v)) - \mu\{t^p(1-t) + t(1-t)^p\} \|g(v) - h(u)\|^p \\
 &\leq \max\{F(h(u)), F(g(v))\} - \mu\{t^p(1-t) + t(1-t)^p\} \|g(v) - h(u)\|^p.
 \end{aligned}$$

This shows that every higher order strongly generalized log-convex function is a higher order strongly generalized convex function and every higher order strongly generalized convex function is a higher order strongly generalized quasi-convex function. However, the converse is not true.

**Definition 14** A function  $F$  on the  $(h, g)$ -convex set  $K \subseteq D$  is said to be a higher order strongly generalized affine with respect to two functions  $h$  and  $g$ , if there exists a constant  $\mu > 0$ , such that

$$\begin{aligned}
 F(h(u) + t(g(v) - h(u))) &= (1-t)F(h(u)) + t(F(g(v)) - \mu\{t^p(1-t) \\
 &\quad + t(1-t)^p\} \|g(v) - h(u)\|^p, \\
 \forall u, v \in K \subseteq D, t \in [0, 1].
 \end{aligned}$$

We would like to remark that, if  $t = 1/2$  in Definition 14, then one gets the higher order generalized Jensen type property called higher order strongly generalized  $J$ -affine function.

For appropriate and suitable choice of the arbitrary functions  $h, g$ , one can obtain several new and known classes of and their variant forms as special cases of higher order strongly generalized convex functions with respect to two functions  $h$  and  $g$ . This shows that the class of higher order strongly generalized convex functions with respect to two functions  $h$  and  $g$ , is quite broad and unifying one.

We now introduce some new concepts and definitions.

**Definition 15** Let  $K \subseteq D$  be a  $(h, g)$ -convex set. An operator  $T : K \rightarrow H$  is said to be:

1. Higher order strongly monotone, if and only if, there exists a constant  $\alpha > 0$  such that

$$\langle Tu - Tv, h(u) - g(v) \rangle \geq \alpha \|g(v) - h(u)\|^p, \forall u, v \in K \subseteq D.$$

2. Higher order strongly pseudomonotone, if and only if, there exists a constant  $\nu > 0$  such that

$$\begin{aligned} \langle Tu, g(v) - h(u) \rangle + \nu \|g(v) - h(u)\|^p &\geq 0 \\ \Rightarrow \\ \langle Tv, g(v) - h(u) \rangle &\geq 0, \forall u, v \in K \subseteq D. \end{aligned}$$

3. Higher order strongly relaxed pseudomonotone, if and only if, there exists a constant  $\mu > 0$  such that

$$\begin{aligned} \langle Tu, g(v) - h(u) \rangle &\geq 0 \\ \Rightarrow \\ -\langle Tv, h(u) - g(v) \rangle + \mu \|g(v) - h(u)\|^p &\geq 0, \forall u, v \in K \subseteq D. \end{aligned}$$

4. Generalized monotone with respect to two functions  $h$  and  $g$ , if

$$\langle T(h(u)) - T(g(v)), h(u) - g(v) \rangle \geq 0, \forall u, v \in K \subseteq D.$$

**Definition 16** A differentiable function  $F$  on the  $(h, g)$ -convex set  $K \subseteq D$  is said to be higher order strongly generalized pseudoconvex function, if and only if, if there exists a constant  $\mu > 0$ , such that

$$\begin{aligned} \langle F'(h(u)), g(v) - h(u) \rangle + \mu \|g(v) - h(u)\|^p \geq 0 &\Rightarrow F(g(v)) \geq F(h(u)), \\ \forall u, v \in K \subseteq D. \end{aligned}$$

For suitable and appropriate choices of the arbitrary functions  $h, g$  and the constant  $p$ , one can obtain various new and old concepts as special of the above definitions. Thus, it is obvious that these concepts are unifying ones and flexible.

We now consider some basic properties of higher order strongly generalized convex functions.

**Theorem 7** Let  $F$  be a differentiable function on the  $(h, g)$ -convex set  $K \subseteq D$ . Then the function  $F$  is higher order strongly generalized convex function, if and only if,

$$\begin{aligned} F(g(v)) - F(h(u)) &\geq \langle F'(h(u)), g(v) - h(u) \rangle + \mu \|g(v) - h(u)\|^p, \\ \forall u, v \in K \subseteq D, t \in [0, 1]. \end{aligned} \tag{8}$$

**Proof** Let  $F$  be a higher order strongly generalized convex function on the  $(h, g)$ -convex set  $K \subseteq D$ . Then

$$\begin{aligned}
 F(h(u) + t(g(v) - h(u))) &\leq (1 - t)F(h(u)) + tF(g(v)) \\
 &\quad - \mu\{t^p(1 - t) + t(1 - t)^p\}\|g(v) - h(u)\|^p, \\
 &\quad \forall u, v \in K, t \in [0,1].
 \end{aligned}$$

which can be written as

$$\begin{aligned}
 F(g(v)) - F(h(u)) &\geq \left\{ \frac{F(h(u) + t(g(v) - h(u)) - F(h(u))}{t} \right\} \\
 &\quad + \{t^{p-1}(1 - t) + (1 - t)^p\}\|g(v) - h(u)\|^p.
 \end{aligned}$$

Taking the limit in the above inequality as  $t \rightarrow 0$ , we have

$$\begin{aligned}
 F(g(v)) - F(h(u)) &\geq \langle F'(h(u)), g(v) - h(u) \rangle + \mu\|g(v) - h(u)\|^p, \\
 &\quad \forall u, v \in K \subseteq D.
 \end{aligned}$$

which is (8), the required result.

Conversely, let (8) hold. Then,  $\forall u, v \in K \subseteq D, t \in [0, 1]$ ,  $g(v_t) = h(u) + t(g(v) - h(u)) \in K \subseteq D$ , we have

$$\begin{aligned}
 F(g(v)) - F(g(v_t)) &\geq \langle F'(g(v_t)), g(v) - g(v_t) \rangle + \mu\|g(v) - g(v_t)\|^p \\
 &= (1-t)F'(g(v_t), g(v)-h(u)) + \mu(1-t)^p\|g(v)-h(u)\|^p \quad (9) \\
 &\quad \forall v, u \in K \subseteq D.
 \end{aligned}$$

In a similar way, we have

$$\begin{aligned}
 F(u) - F(v_t) &\geq \langle F'(v_t), u - v_t \rangle + \mu\|g(u) - g(v_t)\|^p \\
 &= -tF'(v_t, v - u) + \mu t^p\|g(v) - g(u)\|^p. \quad (10)
 \end{aligned}$$

Multiplying (9) by  $t$  and (10) by  $(1 - t)$  and adding the resultant, we have

$$\begin{aligned}
 F(h(u) + t(g(v) - h(u))) &\leq (1 - t)F(h(u)) + tF(g(v)) - \mu\{t^p(1 - t) \\
 &\quad + t(1 - t)^p\}\|g(v) - h(u)\|^p, \\
 &\quad \forall u, v \in K \subseteq D,
 \end{aligned}$$

showing that  $F$  is a higher order strongly generalized convex function.

**Theorem 8** *Let  $F$  be a differentiable higher order strongly generalized convex function on the  $(h, g)$ -convex set  $K \subseteq D$ . Then  $F'(\cdot)$  is a higher order strongly monotone operator.*

**Proof** Let  $F$  be a higher order strongly generalized convex function on the  $(h, g)$ -convex set  $K \subseteq D$ . Then, from Theorem 7, we have

$$F(g(v)) - F(h(u)) \geq \langle F'(h(u)), g(v) - h(u) \rangle + \mu \|g(v) - h(u)\|^p, \quad \forall u, v \in K \subseteq D. \tag{11}$$

Changing the role of  $h(u)$  and  $g(v)$  in (11), we have

$$F(h(u)) - F(g(v)) \geq \langle F'(g(v)), h(u) - g(v) \rangle + \mu \|g(v) - h(u)\|^p, \quad \forall u, v \in K \subseteq D. \tag{12}$$

Adding (11) and (12), we have

$$\langle F'(h(u)) - F'(g(v)), h(u) - g(v) \rangle \geq 2\mu \|g(v) - h(u)\|^p, \quad \forall u, v \in K \subseteq D. \tag{13}$$

which shows that  $F'(\cdot)$  is a higher order strongly monotone operator.

We remark that the converse of Theorem 8 is not true. However, we have the following result.

**Theorem 9** *If the differential operator  $F'(\cdot)$  of a differentiable higher order strongly generalized convex function  $F$  is higher order strongly monotone operator, then*

$$F(g(v)) - F(h(u)) \geq \langle F'(h(u)), g(v) - h(u) \rangle + 2\mu \frac{1}{p} \|g(v) - h(u)\|^p, \quad \forall u, v \in K \subseteq D. \tag{14}$$

**Proof** Let  $F'(\cdot)$  be a higher order strongly monotone operator. Then, from (13), we have

$$\langle F'(g(v)), h(u) - g(v) \rangle \geq \langle F'(h(u)), h(u) - g(v) \rangle + 2\mu \|g(v) - h(u)\|^p. \quad \forall u, v \in K \subseteq D. \tag{15}$$

Since  $K$  is the  $(h, g)$ -convex set,  $\forall u, v \in K \subseteq D, t \in [0, 1]$ ,  $g(v_t) = h(u) + t(g(v) - h(u)) \in K \subseteq D$ . Taking  $g(v) = g(v_t)$  in (15), we have

$$\begin{aligned} \langle F'(g(v_t)), h(u) - g(v_t) \rangle &\leq \langle F'(h(u)), h(u) - g(v_t) \rangle - 2\mu \|g(v) - h(u)\|^p \\ &= -t \langle F'(h(u)), g(v) - h(u) \rangle - 2\mu t^p \|g(v) - h(u)\|^p, \end{aligned}$$

which implies that

$$\langle F'(g(v_t)), g(v) - h(u) \rangle \geq \langle F'(h(u)), g(v) - h(u) \rangle + 2\mu t^{p-1} \|g(v) - h(u)\|^p. \tag{16}$$

Consider the auxiliary function

$$\xi(t) = F(h(u)) + t(g(v) - h(u)), \forall u, v \in K \subseteq D,$$

from which, we have

$$\xi(1) = F(g(v)), \quad \xi(0) = F(h(u)).$$

Then, from (16), we have

$$\begin{aligned} \xi'(t) = \langle F'(g(v_t)), g(v) - h(u) \rangle &\geq \langle F'(h(u)), g(v) - h(u) \rangle \\ &+ 2\mu t^{p-1} \|g(v) - h(u)\|^p. \end{aligned} \tag{17}$$

Integrating (17) between 0 and 1, we have

$$\xi(1) - \xi(0) = \int_0^1 \xi'(t) dt \geq \langle F'(h(u)), g(v) - h(u) \rangle + 2\mu \frac{1}{p} \|g(v) - h(u)\|^p.$$

Thus it follows that

$$\begin{aligned} F(g(v)) - F(h(u)) &\geq \langle F'(h(u)), g(v) - h(u) \rangle + 2\mu \frac{1}{p} \|g(v) - h(u)\|^p, \\ &\forall u, v \in K \subseteq D, \end{aligned}$$

which is the required (14).

We note that, if  $p = 2$ , then Theorem 9 can be viewed as the converse of Theorem 8.

We now give a necessary condition for higher order strongly generalized pseudo-convex function.

**Theorem 10** *Let  $F'(\cdot)$  be a higher order strongly relaxed pseudomonotone operator. Then  $F$  is a higher order strongly generalized pseudo-convex function.*

**Proof** Let  $F'$  be a higher order strongly relaxed pseudomonotone operator. Then,

$$\langle F'(h(u)), g(v) - h(u) \rangle \geq 0, \forall u, v \in K \subseteq D,$$

implies that

$$\langle F'(g(v)), g(v) - h(u) \rangle \geq \mu \|g(v) - h(u)\|^p, \forall u, v \in K \subseteq D. \tag{18}$$

Since  $K$  is an convex set,  $\forall u, v \in K \subseteq D, \quad t \in [0, 1], g(v_t) = h(u) + t(g(v) - h(u)) \in K \subseteq D.$



Taking  $g(v) = g(v_t)$  in (18), we have

$$\langle F'(g(v_t)), g(v) - h(u) \rangle \geq \mu t^{p-1} \|g(v) - h(u)\|^p. \tag{19}$$

Consider the auxiliary function

$$\xi(t) = F(h(u) + t(g(v) - h(u))) = F(g(v_t)), \quad \forall u, v \in K \subseteq D, t \in [0, 1],$$

which is differentiable, since  $F$  is differentiable function. Then, using (19), we have

$$\xi'(t) = \langle F'(g(v_t)), g(v) - h(u) \rangle \geq \mu t^{p-1} \|g(v) - h(u)\|^p.$$

Integrating the above relation between 0 to 1, we have

$$\xi(1) - \xi(0) = \int_0^1 \xi'(t) dt \geq \frac{\mu}{p} \|g(v) - h(u)\|^p,$$

that is,

$$F(g(v)) - F(h(u)) \geq \frac{\mu}{p} \|g(v) - h(u)\|^p, \quad \forall u, v \in K \subseteq D,$$

showing that  $F$  is a higher order strongly generalized pseudo-convex function.

**Definition 17** A function  $F$  is said to be sharply higher order strongly generalized pseudo convex, if there exists a constant  $\mu > 0$  such that

$$\langle F'(h(u)), g(v) - h(u) \rangle \geq 0$$

$\Rightarrow$

$$F(g(v)) \geq F(g(v) + t(h(u) - g(v))) + \mu \{t^p(1 - t) + t(1 - t)^p\} \|g(v) - h(u)\|^p, \\ \forall u, v \in K \subseteq D.$$

**Theorem 11** Let  $F$  be a sharply higher order strongly generalized pseudo convex function on the  $(h, g)$ -convex set  $K \subseteq D$  with a constant  $\mu > 0$ . Then

$$\langle F'(g(v)), g(v) - h(u) \rangle \geq \mu \|g(v) - h(u)\|^p, \quad \forall u, v \in K \subseteq D.$$

**Proof** Let  $F$  be a sharply higher order strongly generalized pseudo-convex function. Then

$$F(g(v)) \geq F(g(v) + t(h(u) - g(v))) + \mu \{t^p(1 - t) + t(1 - t)^p\} \|g(v) - h(u)\|^p, \\ \forall u, v \in K \subseteq D, t \in [0, 1],$$

from which, we have

$$\left\{ \frac{F(g(v)) + t(h(u) - g(v)) - F(g(v))}{t} \right\} + \mu \{t^{p-1}(1-t) + (1-t)^p\} \|g(v) - h(u)\|^p \geq 0.$$

Taking limit in the above inequality, as  $t \rightarrow 0$ , we have

$$\langle F'(g(v)), g(v) - g(u) \rangle \geq \mu \|g(v) - h(u)\|^p, \forall u, v \in K \subseteq D,$$

the required result.

**Theorem 12** *Let  $f$  be a higher order strongly affine generalized convex function. Then  $F$  is a higher order strongly generalized convex function, if and only if,  $H = F - f$  is a generalized convex function.*

**Proof** Let  $f$  be a higher order strongly affine generalized convex function, Then

$$\begin{aligned} f((1-t)h(u) + tg(v)) &= (1-t)f(h(u)) + tf(g(v)) - \mu \{t^p(1-t) \\ &\quad + t(1-t)^p\} \|g(v) - h(u)\|^p, \end{aligned} \tag{20}$$

$$\forall u, v \in K \subseteq D.$$

From the higher order strongly generalized convexity of  $F$ , we have

$$\begin{aligned} F((1-t)h(u) + tg(v)) &\leq (1-t)F(h(u)) + tF(g(v)) - \mu \{t^p(1-t) \\ &\quad + t(1-t)^p\} \|g(v) - h(u)\|^p, \end{aligned} \tag{21}$$

$$\forall u, v \in K \subseteq D.$$

From (20) and (21), we have

$$\begin{aligned} F((1-t)h(u) + tg(v)) - f((1-t)h(u) + tg(v)) &\leq (1-t)(F(h(u)) - f(h(u))) \\ &\quad + t(F(g(v)) - f(g(v))), \end{aligned} \tag{22}$$

from which it follows that

$$\begin{aligned} H((1-t)h(u) + tg(v)) &= F((1-t)h(u) + tg(v)) - f((1-t)h(u) + tg(v)) \\ &\leq (1-t)F(h(u)) + tF(g(v)) - (1-t)f(h(u)) - tf(g(v)) \\ &= (1-t)(F(h(u)) - f(h(u))) + t(F(g(v)) - f(g(v))), \end{aligned}$$

which show that  $H = F - f$  is a convex function.

The inverse implication is obvious.

**Definition 18** A function  $F$  is said to be a pseudoconvex function with respect to a strictly positive bifunction  $B(., .)$ , if

$$\begin{aligned}
 &F(g(v)) < F(h(u)) \\
 &\Rightarrow \\
 &F(h(u) + (1 - t)(g(v), h(u))) < F(h(u)) + t(t - 1)B(g(v), h(u)), \\
 &\forall u, v \in K \subseteq D, t \in [0, 1].
 \end{aligned}$$

**Theorem 13** *If the function  $F$  is higher order strongly generalized convex function such that  $F(g(v)) < F(h(u))$ , then the function  $F$  is higher order strongly generalized pseudo convex function*

**Proof** Since  $F(g(v)) < F(h(u))$  and  $F$  is higher order strongly generalized convex function, then

$\forall u, v \in K \subseteq D, \quad t \in [0, 1]$ , we have

$$\begin{aligned}
 &F(h(u) + t(g(v) - h(u))) \leq F(h(u)) + t(F(g(v)) - F(h(u))) \\
 &\quad - \mu\{t^p(1 - t) + t(1 - t)^p\}\|g(v) - hu\|^p \\
 &< F(h(u)) + t(1 - t)(F(g(v)) - F(h(u))) - \mu\{t^p(1 - t) + t(1 - t)^p\}\|g(v) - h(u)\|^p \\
 &= F(h(u)) + t(t - 1)(F(h(u)) - F(g(v))) - \mu\{t^p(1 - t) + t(1 - t)^p\}\|g(v) - h(u)\|^p \\
 &< F(h(u)) + t(t - 1)B(h(u), g(v)) - \mu\{t^p(1 - t) + t(1 - t)^p\}\|g(v) - h(u)\|^p, \\
 &\forall u, v \in K \subseteq D,
 \end{aligned}$$

where  $B(h(u), g(v)) = F(h(u)) - F(g(v)) > 0$ . Hence the function  $F$  is higher order strongly generalized convex function, the required result.

From Definition 10, we have

$$\begin{aligned}
 &F((1 - t)h(u) + tg(v)) + F(th(u) + (1 - t)g(v)) \\
 &\leq F(h(u)) + F(g(v)) - 2\mu\{t^p(1 - t) + t(1 - t)^p\}\|g(v) - h(u)\|^p, \\
 &\forall u, v \in K \subseteq D, t \in [0, 1], p > 1,
 \end{aligned}$$

which is called the higher order strongly Wright generalized convex functions. One can investigate the properties and applications of higher order strongly Wright generalized convex functions.

## 5 Applications

We now show that uniformly Banach spaces can be characterized by the parallelogram laws, which can be obtained from the higher order strongly generalized affine functions.

Setting  $F(u) = \|u\|^p$  in Definition 14, we have

$$\begin{aligned} & \|h(u) + t(g(v) - h(u))\|^p \tag{23} \\ &= (1 - t)\|h(u)\|^p + t\|g(v)\|^p - \mu\{t^p(1 - t) + t(1 - t)^p\}\|g(v) - h(u)\|^p, \\ & \forall u, v \in K \subseteq D, t \in [0, 1]. \end{aligned}$$

Taking  $t = \frac{1}{2}$  in (23), we have

$$\left\| \frac{h(u)+g(v)}{2} \right\|^p + \mu \frac{1}{2^p} \|g(v)-h(u)\|^p = \frac{1}{2} \|h(u)\|^p + \frac{1}{2} \|g(v)\|^p, \forall u, v \in K \subseteq D, \tag{24}$$

which implies that

$$\|h(u)+g(v)\|^p + \mu \|g(v)-h(u)\|^p = 2^{p-1} \{ \|h(u)\|^p + \|g(v)\|^p \}, \forall u, v \in K \subseteq D, \tag{25}$$

which is the generalized parallelogram law for the  $l^p$ -spaces.

For  $p = 2$ , the generalized parallelogram law (25) reduces to:

$$\|h(u) + g(v)\|^2 + \mu \|g(v) - h(u)\|^2 = 2\{\|h(u)\|^2 + \|g(v)\|^2\}, \forall u, v \in K \subseteq D, \tag{26}$$

which is a new parallelogram law involving two arbitrary functions characterizing the inner product spaces and can be viewed as a novel application of the strongly generalized affine functions.

For  $g(u) = h(u) = I$ , we obtain the well known parallelogram law, that is,

$$\|u + v\|^p + \mu \|v - u\|^p = 2^{p-1} \{ \|u\|^p + \|v\|^p \}, \forall u, v \in K \subseteq D, \tag{27}$$

which was derived by Xu [42] via the functionals  $\|\cdot\|^p$  and  $\|\cdot\|^q$ , respectively. These parallelogram laws characterize the  $p$ -uniform convexity and  $q$ -uniform smoothness of a Banach space. Bynum [6] and Chen et al. [7, 8] have studied the properties and applications of the parallelogram laws for the Banach spaces in the prediction theory. In brief, for suitable and appropriate choice of the arbitrary functions  $h$  and  $g$ , one can obtained a wide class of parallelogram laws, which can be used to characterize the inner product spaces and uniformly Banach spaces. It is an interesting problem to consider the applications in optimization and prediction theory.

Let  $B(H)$  be the space of all bounded linear operators on a separable complex Hilbert space  $H$ . The absolute value of an operator  $A \in B(H)$  is defined by  $\|A\|^p = \langle A, A \rangle$ .

By taking  $F(u) = \|A\|^p$  in (27), we have

$$\|A + B\|^p + \mu \|A - B\|^p = 2^{p-1} \{ \|A\|^p + \|B\|^p \}, \quad \forall A, B \in B(H) \subseteq D,$$

which is known as the Clarkson Inequalities for Operators and be viewed as Operator Parallelogram Laws characterizing the uniformly Banach spaces. For more details, see Hirzallah and Kittaneh [12].

### 5.1 Generalized Variational Inequalities

In this section, we show that the optimality conditions of the differentiable higher order strongly generalized convex functions can be characterized by a class of variational inequalities. This is the main motivation of our next result.

**Theorem 14** *Let  $F$  be a differentiable higher order strongly generalized convex function with modulus  $\mu > 0$ . If  $u \in K \subseteq D$  is the minimum of the function  $F$ , then*

$$F(g(v)) - F(h(u)) \geq \mu \|g(v) - h(u)\|^p, \quad \forall v \in K \subseteq D, \quad p \geq 1. \quad (28)$$

**Proof** *Let  $u \in K \subseteq D$  be a minimum of the function  $F$ . Then*

$$F(h(u)) \leq F(g(v)), \quad \forall v \in K \subseteq D. \quad (29)$$

*Since  $K \subseteq D$  is a  $(h, g)$ -convex set, so,  $\forall u, v \in K \subseteq D, \quad t \in [0, 1]$ ,*

$$g(v_t) = (1 - t)h(u) + tg(v) \in K \subseteq D.$$

*Taking  $g(v) = g(v_t)$  in (29), we have*

$$0 \leq \lim_{t \rightarrow 0} \left\{ \frac{F(h(u) + t(g(v) - h(u))) - F(h(u))}{t} \right\} = \langle F'(h(u)), g(v) - h(u) \rangle. \quad (30)$$

*Since  $F$  is differentiable higher order strongly convex generalized function, so*

$$F(h(u) + t(g(v) - h(u))) \leq F(h(u)) + t(F(g(v)) - F(h(u))) - \mu \{t^p(1 - t) + t(1 - t)^p\} \|g(v) - h(u)\|^p, \quad \forall u, v \in K \subseteq D,$$

*from which, using (30), we have*

$$\begin{aligned} F(g(v)) - F(h(u)) &\geq \lim_{t \rightarrow 0} \left\{ \frac{F(h(u) + t(g(v) - h(u))) - F(h(u))}{t} \right\} + \mu \|g(v) - h(u)\|^p \\ &= \langle F'(h(u)), g(v) - h(u) \rangle + \mu \|g(v) - h(u)\|^p, \\ &\geq \mu \|g(v) - h(u)\|^p, \end{aligned}$$

*the required result (28).*

Using the above technique and ideas of Theorem 7, one can obtain the following result.

**Theorem 15** *Let  $K \subseteq D$  be a  $(h, g)$ -convex set. If  $u \in K \subseteq D$  satisfies the inequality*

$$\langle F'(h(u)), g(v) - h(u) \rangle + \mu \|g(v) - h(u)\|^p \geq 0, \forall u, v \in K \subseteq D, p \geq 1, \tag{31}$$

*then  $u \in K \subseteq D$  is the minimum of the higher order strongly generalized convex function  $F$ .*

*Remark 1* The inequality (31) is called the higher order strongly extended strongly general variational inequality, which include strongly extended general variational inequalities, higher order strongly general variational inequalities and general variational inequalities as special cases. It is an interesting problem to study the both quantitative and qualitative properties of these variational inequalities.

**Theorem 16** *If the operator  $T$  is a generalized monotone with respect the arbitrary functions  $h, g$  and  $u \in K \subseteq D$  is the solution of the inequality (31), then  $u \in K \subseteq D$  satisfies the inequality*

$$\langle F'(g(v)), g(v) - h(u) \rangle + \mu \|g(v) - h(u)\|^p \geq 0, \forall u, v \in K \subseteq D, p \geq 1. \tag{32}$$

**Proof** Let  $u \in K \subseteq D$  satisfies the inequality (31). Then

$$\begin{aligned} 0 &\leq \langle F'(h(u)), g(v) - h(u) \rangle + \mu \|g(v) - h(u)\|^p \\ &= \langle F'(h(u)) - F'(g(v)) + F'(g(v)), g(v) - h(u) \rangle + \mu \|g(v) - h(u)\|^p \\ &\leq \langle F'(g(v)), g(v) - h(u) \rangle + \mu \|g(v) - h(u)\|^p, \end{aligned}$$

where we have used the generalized monotonicity of the operator  $T$ .

We note that the converse of Theorem 16 does not hold due the presence of the term  $\|g(v) - h(u)\|^p$ . However, for  $p = 1$ , one can prove the converse of Theorem 16 using the concept of hemicontinuity.

*Remark 2* The inequality (32) is called the Minty higher order strongly extended strongly general variational inequality. We remark that the projection method and its variant forms can not be used to study the higher order strongly general variational inequalities (31) due to its inherent structure. To overcome this drawback, one can consider the auxiliary principle technique for solving higher order strongly extended strongly general variational inequality, which is mainly due to Glowinski et al. [11] and Lions and Stampacchia [17] as developed by Noor [25, 26] and Noor et al. [37]. We have only given a glimpse of the higher order strongly extended general variational inequalities. We are going to consider these problems in future. These problems may applications in various branches of pure and applied sciences and need further efforts.

## Conclusion

In this paper, we have introduced and studied a new class of convex functions, which is called higher order strongly generalized convex function. It is shown that several new classes of strongly convex functions can be obtained as special cases of these higher order strongly generalized convex functions. We have studied the basic properties of these functions. We have shown that one can derive the parallelogram laws in Banach spaces, which have applications in prediction theory and stochastic analysis. It is shown that the minimum of the differentiable higher order strongly generalized convex functions can be characterized by a new class of variational inequalities. Several important special cases are also discussed, which can be obtained from our results. The interested readers may explore the applications and other properties of the higher order strongly convex functions in various fields of pure and applied sciences. This is an interesting direction of future research.

**Acknowledgments** We wish to express our deepest gratitude to our teachers, colleagues, collaborators and friends, who have direct or indirect contributions in the process of this paper. The authors would like to thank the Rector, COMSATS University Islamabad, Islamabad Pakistan, for providing excellent research and academic environments.

## References

1. M. Adamek, On a problem connected with strongly convex functions. *Math. Inequal. Appl.* **19**(4), 1287–1293 (2016)
2. O. Alabdali, A. Guessab, G. Schmeisser, Characterization of uniform convexity for differentiable functions. *Appl. Anal. Discrete Math.* **13**, 721–732 (2019)
3. H. Angulo, J. Gimenez, A.M. Moeos, K. Nikodem, On strongly  $h$ -convex functions. *Ann. Funct. Anal.* **2**(2), 85–91 (2011)
4. M.U. Awan, M.A. Noor, M.V. Mihai, K.I. Noor, N. Akhtar, On approximately harmonic  $h$ -convex functions depending on a given function. *Filomat* **33**(12), 3783–3793 (2019)
5. M.U. Awan, M.A. Noor, T.-S. Du, K.I. Noor, New refinements of fractional Hermite-Hadamard inequality. *RACSAM* **113**, 21–29 (2019)
6. W.L. Bynum, Weak parallelogram laws for Banach spaces. *Can. Math. Bull.* **19**, 269–275 (1976)
7. R. Cheng, C.B. Harris, Duality of the weak parallelogram laws on Banach spaces. *J. Math. Anal. Appl.* **404**, 64–70 (2013)
8. R. Cheng, W.T. Ross, Weak parallelogram laws on Banach spaces and applications to prediction. *Period. Math. Hung.* **71**, 45–58 (2015)
9. G. Cristescu, M. Găianu, Shape properties of Noors convex sets, in *Proceedings of the Twelfth Symposium of Mathematics and its Applications, Timisoara* (2009), pp. 1–13
10. G. Cristescu, L. Lupsa, *Non-Connected Convexities and Applications* (Kluwer Academic Publisher, Dordrecht, 2002)
11. F. Glowinski, J.L. Lions, R. Tremileres, *Numerical Analysis of Variational Inequalities* (North Holland, Amsterdam, 1981)
12. O. Hirzallah, F. Kittaneh, Non-commutative Clarkson inequalities for  $n$ -tuples of operators. *Integr. Equ. Oper. Theory* **60**, 369–379 (2008)

13. S. Jabeen, B.B. Mohsen, M.A. Noor, K.I. Noor, Inertial projection methods for solving general quasi-variational inequalities. *AIMS Math.* **6**(2), 1075–1086 (2021)
14. J.-B. Jian, On  $(E, F)$  generalized convexity. *Internat. J. Math. Sci.* **2**, 121–132 (2003)
15. S. Karamardian, The nonlinear complementarity problems with applications, Part 2. *J. Optim. Theory Appl.* **4**(3), 167–181 (1969)
16. D. Kinderlehrer, G. Stampacchia, *An Introduction to Variational Inequalities and their Applications* (Academic Press, London, 1980)
17. J.L. Lions, G. Stampacchia, Variational inequalities. *Commun. Pure Appl. Math.* **20**, 492–512 (1967)
18. G.H. Lin, M. Fukushima, Some exact penalty results for nonlinear programs and mathematical programs with equilibrium constraints. *J. Optim. Theory Appl.* **118**(1), 67–80 (2003)
19. B.B. Mohsen, M.A. Noor, K.I. Noor, M. Postolache, Strongly convex functions of higher order involving bifunction. *Mathematics* **7**(11), 1028 (2019). <https://doi.org/10.3390/math7111028>
20. C.P. Niculescu, L.E. Persson, *Convex Functions and their Applications* (Springer, New York, 2018)
21. K. Nikodem, Z.S. Pales, Characterizations of inner product spaces by strongly convex functions. *Banach J. Math. Anal.* **1**, 83–87 (2011)
22. M. A. Noor, General variational inequalities. *Appl. Math. Letts.* **1**, 119–121(1988)
23. M.A. Noor, Quasi variational inequalities. *Appl. Math. Letts.* **1**(4), 367–370 (1988)
24. M.A. Noor, New approximations schemes for general variational inequalities. *J. Math. Anal. Appl.* **251**, 217–230 (2000)
25. M.A. Noor, Some developments in general variational inequalities. *Appl. Math. Comput.* **152**, 199–277 (2004)
26. M.A. Noor, Fundamentals of equilibrium problems. *Math. Inequal. Appl.* **9**(3), 529–566 (2006)
27. M.A. Noor, Differentiable nonconvex functions and general variational inequalities. *Appl. Math. Comput.* **199**(2), 623–630 (2008)
28. M.A. Noor, Extended general variational inequalities. *Appl. Math. Lett.* **22**, 182–186 (2009)
29. M.A. Noor, Some aspects of extended general variational inequalities. *Abstract Appl. Anal.* **2012**, 16 (2012). ID 303569
30. M.A. Noor, K.I. Noor, On generalized strongly convex functions involving bifunction. *Appl. Math. Inform. Sci.* **13**(3), 411–416 (2019)
31. M.A. Noor, K.I. Noor, Higher order strongly general convex functions and variational inequalities. *AIMS Math.* **5**(4), 3646–3663 (2020)
32. M.A. Noor, K.I. Noor, Characterizations of strongly generalized convex functions. *TJMM* **11**(2), 117–127 (2020)
33. M.A. Noor, K.I. Noor, Higher order general convex functions and general variational inequalities. *Canada J. Appl. Math.* **3**(1), 1–17 (2021)
34. M.A. Noor, K.I. Noor, New classes of preinvex functions and variational-like inequalities. *Filomat* **35** (2021)
35. M.A. Noor, K.I. Noor, T.M. Rassias, Some aspects of variational inequalities. *J. Comput. Appl. Math.* **47**(3), 285–312 (1993)
36. M.A. Noor, K.I. Noor, A. Bnouhachem, Some new iterative methods for solving variational inequalities. *Canad. J. Appl. Math.* **2**(2), 1–7 (2020)
37. M.A. Noor, K.I. Noor, M.T. Rassias, New trends in general variational inequalities. *Acta Appl. Math.* **170**(1), 981–1064 (2020)
38. J. Pecric, F. Proschan, Y.I. Tong, *Convex Functions, Partial Ordering and Statistical Applications* (Academic Press, New York, 1992)
39. B.T. Polyak, Existence theorems and convergence of minimizing sequences in extremum problems with restrictions. *Soviet Math. Dokl.* **7**, 2–75 (1966)
40. G. Qu, N. Li, On the exponentially stability of primal-dual gradient dynamics. *IEEE Control Syst. Letters*, **3**(1), 43–48 (2019)
41. G. Stampacchia, Formes bilineaires coercitives sur les ensembles convexes. *C. R. Acad. Sci. Paris* **258**, 4423–4426 (1964)



42. H.-K. Xu, Inequalities in Banach spaces with applications. *Nonlinear Anal. Theory Meth. Appl.* **16**(12), 1127–1138 (1991)
43. E.A. Youness,  $E$ -convex sets,  $E$ -convex functions and  $E$ -convex programming. *J. Optim. Theory Appl.* **102**, 439–450 (1999)
44. Y. Zhao, D. Sun, Alternative theorems for nonlinear projection equations and applications to generalized complementarity problems. *Nonl. Anal.* **46**, 853–868 (2001)
45. D.L. Zu, P. Marcotte, Co-coercivity and its role in the convergence of iterative schemes for solving variational inequalities. *SIAM J. Optim.* **6**(3), 714–726 (1996)

# A Note on Generalized Nash Games Played on Networks



Mauro Passacantando and Fabio Raciti

**Abstract** We investigate a generalized Nash equilibrium problem where players are modeled as nodes of a network and the utility function of each player depends on his/her own action as well as on the actions of his/her neighbors in the network. In the case of a quadratic reference model with shared constraints we are able to derive the variational solution of the game as a series expansion which involves the powers of the adjacency matrix, thus extending a previous result. Our analysis is illustrated by means of some numerical examples.

## 1 Introduction

Economic and social sciences are probably the areas that have benefited the most from the mathematical development of Game Theory (GT), although in the last decades, specific game-theoretical models have also been applied to various problems from engineering, transportation and communication networks, biology, and other fields (see, e.g., [1, 19, 21]). The pervasive role of physical and virtual social interactions in the actions taken by individuals or groups, described through a network model, has led to consider Network Games as a powerful tool to describe the process of decision making. Indeed, in many social or economic environments our actions are influenced by the actions of our friends, acquaintances or colleagues. In network game models, each individual (player) is identified with the node of a graph and the players that can interact directly are connected through links of the graph. The specificity of these games is the central role played by the graph structure in the description of the patterns of interactions, and in the final social or economic

---

M. Passacantando

Department of Computer Science, University of Pisa, Pisa, Italy

e-mail: [mauro.passacantando@unipi.it](mailto:mauro.passacantando@unipi.it)

F. Raciti (✉)

Department of Mathematics and Computer Science, University of Catania, Catania, Italy

e-mail: [fabio.raciti@unict.it](mailto:fabio.raciti@unict.it)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,

Springer Optimization and Its Applications 173,

[https://doi.org/10.1007/978-3-030-72563-1\\_16](https://doi.org/10.1007/978-3-030-72563-1_16)

outcome. The mathematical consequence of this description is that some interesting results depend on quantities such as the spectral radius of the adjacency matrix, its minimum eigenvalue, and its powers. With the huge variety of possible networks and interactions, it is difficult to make progress in the analysis without positing some specific structure of the problems under consideration. In these regards it is interesting to investigate the two classes of *games with strategic complements and substitutes*. Roughly speaking, in the first case, the incentive for a player to take an action increases when the number of his/her social contacts who take the action increases while in the second case this monotonic relation is reversed. The linear-quadratic model, although its simplicity, deserves particular attention and has been investigated in detail because it can be solved exactly and the solution formula can be nicely interpreted from a graph-theoretical point of view. Among the numerous references we suggest that the reader who wishes to become familiar with the topic of network games see the beautiful survey by Jackson and Zenou [9], along with the seminal paper by Ballester et al. [3]. Very recently, a different approach, based on variational inequality theory, has been put forward to tackle this kind of problems and, in this respect, the interested reader can refer to the interesting and very detailed paper by Parise and Ozdaglar [20].

In this paper we investigate, within the simple frame of the linear-quadratic model, a generalized Nash equilibrium problem (GNEP) with shared constraints. This class of games was proposed a long time ago by Rosen [22], but the last decade has witnessed a renewed interest in the subject, due to its wide range of applications and to its connection with the theory of variational inequalities [5, 7, 16, 17]. By adding a shared constraint to the original quadratic problem we thus obtain a GNEP, and loose uniqueness of the solution. Among the solutions of the new problem, we select the so called *variational solution* and provide a closed form expression for it. Furthermore, the new formula can be written by a series expansion of the adjacency matrix, thus extending one of the results in [3] and allowing for a graph-theoretic (as well as socio-economic) interpretation. Namely, this expansion shows that although players only interact with their neighbors, the solution also depends, to a certain extent, on indirect contacts (i.e., neighbors of neighbors, neighbors of neighbors of neighbors, and so on).

The paper is organized as follows. In the following Section 2 we summarize the basic background material on network games and focus on the exactly solvable linear-quadratic model. Section 3 is devoted to a brief description of generalized Nash equilibrium problems with shared constraints, and to the solution of our specific problem, while Section 4 is dedicated to illustrate our result by means of two worked-out examples. The paper ends with a concluding section where we outline some promising avenues of research.

## 2 Network Games

### 2.1 Elements of Graph Theory and Game Classes

We begin this section by recalling a few concepts and definitions of graph theory that will be used in the sequel. We warn the reader that the terminology is not uniform in the related literature. Formally, a graph  $g$  is a pair of sets  $(V, E)$ , where  $V$  is the set of nodes (or vertexes) and  $E$  is the set of arcs (or edges), formed by pairs of nodes  $(v, w)$ . Arcs which have the same end nodes are called parallel, while arcs of the form  $(v, v)$  are called loops. We consider here *simple* graphs, that is graphs with no parallel arcs or loops. In our setting, the players will be represented by the  $n$  nodes in the graph. Moreover, we consider here indirect graphs: arcs  $(v, w)$  and  $(w, v)$  are the same. Two nodes  $v$  and  $w$  are adjacent if they are connected by an arc, i.e., if  $(v, w)$  is an arc. The information about the adjacency of nodes can be stored in the adjacency matrix  $G$  whose elements  $g_{ij}$  are equal to 1 if  $(v_i, v_j)$  is an arc, 0 otherwise.  $G$  is thus a symmetric and zero diagonal matrix. Given a node  $v$ , the nodes connected to  $v$  with an arc are called the *neighbors* of  $v$  and are grouped in the set  $N_v(g)$ . The number of elements of  $N_v(g)$  is the *degree* of  $v$ . A *walk* in the graph  $g$  is a finite sequence of the form

$$v_{i_0}, e_{j_1}, v_{i_1}, e_{j_2}, \dots, e_{j_k}, v_{j_k},$$

which consists of alternating nodes and arcs of the graph, such that  $v_{i_{t-1}}$  and  $v_{i_t}$  are end nodes of  $e_{j_t}$ . The *length* of a walk is the number of its arcs. Let us remark that it is allowed to visit a node or go through an arc more than once. A *path* is a walk with all different nodes (except possibly the initial and terminal ones if the walk is closed). The indirect connections between any two nodes in the graph are described by means of the powers of the adjacency matrix  $G$ . Indeed, it can be proved that the element  $g_{ij}^{[k]}$  of  $G^k$  gives the number of walks of length  $k$  between  $v_i$  and  $v_j$ .

We now proceed to specify the type of game that we will consider. For simplicity, the set of players will be denoted by  $\{1, 2, \dots, n\}$  instead of  $\{v_1, v_2, \dots, v_n\}$ . We denote with  $A_i \subset \mathbb{R}$  the action space of player  $i$ , while  $A = A_1 \times \dots \times A_n$  and the notation  $a = (a_i, a_{-i})$  will be used when we want to distinguish the action of player  $i$  from the action of all the other players. Each player  $i$  is endowed with a payoff function  $u_i : A \rightarrow \mathbb{R}$  that he/she wishes to maximize. The notation  $u_i(a, g)$  is often utilized when one wants to emphasize the influence of the graph structure (e.g., when studying perturbation with respect to the removal of an arc). The solution concept that we consider here is the Nash equilibrium of the game, that is, we seek an element  $a^* \in A$  such that for each  $i \in \{1, \dots, n\}$ :

$$u_i(a_i^*, a_{-i}^*) \geq u_i(a_i, a_{-i}^*), \quad \forall a_i \in A_i. \tag{1}$$

A peculiarity of network games is that the vector  $a_{-i}$  is only made up of components  $a_j$  such that  $j \in N_i(g)$ , that is,  $j$  is a neighbor of  $i$ .

We mentioned in the introduction that it is convenient to consider two specific classes of games which allow a deeper investigation of the patterns of interactions among players. For any given player  $i$  it is interesting to distinguish how variations of the actions of player's  $i$  neighbors affect his/her marginal utility. In the case where the utility functions are twice continuously differentiable the following definitions clarify this point.

**Definition 1** We say that the network game has the property of strategic substitutes if for each player  $i$  the following condition holds:

$$\frac{\partial^2 u_i(a_i, a_{-i})}{\partial a_j \partial a_i} < 0, \quad \forall (i, j) : g_{ij} = 1, \forall a \in A.$$

**Definition 2** We say that the network game has the property of strategic complements if for each player  $i$  the following condition holds:

$$\frac{\partial^2 u_i(a_i, a_{-i})}{\partial a_j \partial a_i} > 0, \quad \forall (i, j) : g_{ij} = 1, \forall a \in A.$$

For the subsequent development it is important to recall that if the  $u_i$  are continuously differentiable functions on  $A$ , the Nash equilibrium problem is equivalent to the variational inequality  $VI(F, A)$ : find  $a^* \in A$  such that

$$[F(a^*)]^\top (a - a^*) \geq 0, \quad \forall a \in A, \tag{2}$$

where

$$[F(a)]^\top := - \left( \frac{\partial u_1}{\partial a_1}(a), \dots, \frac{\partial u_n}{\partial a_n}(a) \right) \tag{3}$$

is also called the pseudo-gradient of the game, according to the terminology introduced by Rosen. For an account of variational inequalities the interested reader can refer to [6, 13, 18]. We recall here some useful monotonicity properties.

**Definition 3**  $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is said to be monotone on  $A$  iff:

$$[F(x) - F(y)]^\top (x - y) \geq 0, \quad \forall x, y \in A.$$

If the equality holds only when  $x = y$ ,  $F$  is said to be strictly monotone.

A stronger type of monotonicity is given by the following

**Definition 4** Let  $\beta > 0$ .  $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is said to be  $\beta$ -strongly monotone on  $A$  iff:

$$[F(x) - F(y)]^\top (x - y) \geq \beta \|x - y\|^2, \quad \forall x, y \in A.$$

For linear operators on  $\mathbb{R}^n$  the two concepts of strict and strong monotonicity coincide and are equivalent to the positive definiteness of the corresponding matrix.

Conditions that ensure the unique solvability of a variational inequality problem are given by the following theorem (see, e.g. [6, 13, 18]).

**Theorem 1** *If  $K \subset \mathbb{R}^n$  is a compact convex set and  $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is continuous on  $K$ , then the variational inequality problem  $VI(F, K)$  admits at least one solution. In the case that  $K$  is unbounded, existence of a solution may be established under the following coercivity condition:*

$$\lim_{\|x\| \rightarrow +\infty} \frac{[F(x) - F(x_0)]^\top (x - x_0)}{\|x - x_0\|} = +\infty,$$

for  $x \in K$  and some  $x_0 \in K$ .

Furthermore, if  $F$  is strictly monotone on  $K$  the solution is unique.

In the following subsection, we describe in detail the linear-quadratic reference model on which we will build our generalized Nash equilibrium problem.

## 2.2 The Linear-Quadratic Model

Let  $A_i = \mathbb{R}_+$  for any  $i \in \{1, \dots, n\}$ , hence  $A = \mathbb{R}_+^n$ . The payoff of player  $i$  is given by:

$$u_i(a, g) = -\frac{1}{2}a_i^2 + \alpha a_i + \phi \sum_{j=1}^n g_{ij} a_i a_j, \quad \alpha, \phi > 0. \tag{4}$$

In this simplified model  $\alpha$  and  $\phi$  take on the same value for all players, which then only differ according to their position in the network. The last term describes the interaction between neighbors and since  $\phi > 0$  this interaction falls in the class of strategic complements. The pseudo-gradient's components of this game are easily computed as:

$$F_i(a) = a_i - \alpha - \phi \sum_{j=1}^n g_{ij} a_j, \quad i \in \{1, \dots, n\},$$

which can be written in compact form as:

$$F(a) = (I - \phi G)a - \alpha \mathbf{1},$$

where  $\mathbf{1} = (1, \dots, 1)^\top \in \mathbb{R}^n$ . We will seek Nash equilibrium points by solving the variational inequality:

$$[F(a^*)]^\top (a - a^*) \geq 0, \quad \forall a \in \mathbb{R}_+^n. \tag{5}$$

Since the constraint set is unbounded, to ensure solvability we require that  $F$  be strongly monotone, which (implying coercivity, for linear operators) also guarantees the uniqueness of the solution.

**Lemma 1 (see e.g. [9])** *The matrix  $I - \phi G$  is positive definite iff  $\phi \rho(G) < 1$ , where  $\rho(G)$  is the spectral radius of  $G$ .*

**Proof** The symmetric matrix  $I - \phi G$  is positive definite if and only if  $\lambda_{\min}(I - \phi G) > 0$ . On the other hand,  $\lambda_{\min}(I - \phi G) = 1 - \phi \lambda_{\max}(G)$ . Since  $G$  is a symmetric non-negative matrix, the Perron-Frobenius Theorem guarantees that  $\lambda_{\max}(G) = \rho(G)$ , hence  $I - \phi G$  is positive definite if and only if  $\phi \rho(G) < 1$ .  $\square$

To be self consistent, in the next lemma we recall the following result about series of matrices.

**Lemma 2 (see e.g. [2])** *Let  $T$  be a square matrix and consider the series:*

$$I + T + T^2 + \dots + T^k + \dots$$

*The series converges provided that  $\lim_k T^k = 0$ , which is equivalent to  $\rho(T) < 1$ . In such case the matrix  $I - T$  is non singular and we have:*

$$(I - T)^{-1} = I + T + T^2 + \dots + T^k + \dots$$

**Theorem 2 (see e.g. [9])** *If  $\phi \rho(G) < 1$ , then the unique Nash equilibrium is*

$$a^* = \alpha (I - \phi G)^{-1} \mathbf{1} = \alpha \sum_{p=0}^{\infty} \phi^p G^p \mathbf{1}. \tag{6}$$

**Proof** Since  $\phi \rho(G) < 1$ , Lemma 1 guarantees that  $F$  is strongly monotone. Hence, Theorem 1 applies and we get a unique solution of (5). On the other hand, Lemma 2 implies that the matrix  $I - \phi G$  is non singular, thus the linear system  $F(a) = 0$ , which reads

$$(I - \phi G)a = \alpha \mathbf{1},$$

has a unique solution  $a^*$  given by (6). Moreover, looking at the expansion we get, by construction, that any component of  $a^*$  is strictly positive. Therefore,  $a^*$  is the unique solution of (5), thus it is the unique Nash equilibrium.  $\square$

*Remark 1* The expansion in (6) suggests an interesting interpretation. Indeed, it can be shown that the  $(i, j)$  entry,  $g_{ij}^{[p]}$ , of the matrix  $G^p$  gives the number of walks of length  $p$  between nodes  $i$  and  $j$ . Based on this observation, a measure of centrality on the network was proposed by Katz and Bonacich (see e.g. [4]). Specifically, for

any weight  $w \in \mathbb{R}_+^n$ , the weighted vector of Katz-Bonacich is given by:

$$b_w(G, \phi) = M(G, \phi) = (I - \phi G)^{-1} w = \sum_{p=0}^{\infty} \phi^p G^p w.$$

In the case where  $w = \mathbf{1}$ , the (non weighted) centrality measure of Katz-Bonacich of node  $i$  is given by:

$$b_{1,i}(G, \phi) = \sum_{j=1}^n M_{ij}(G, \phi)$$

and counts the total number of walks in the graph, which start at node  $i$ , exponentially damped by  $\phi$ .

*Remark 2* The game under consideration also falls in the class of potential games according to the definition introduced by Monderer and Shapley [15]. Indeed, a potential function is given by:

$$P(a, G, \phi) = \sum_{i=1}^n u_i(a, G) - \frac{\phi}{2} \sum_{i=1}^n \sum_{j=1}^n g_{ij} a_i a_j.$$

Monderer and Shapley have proved that, in general, the solutions of the problem

$$\max_{a \in A} P(a, G, \phi)$$

form a subset of the solution set of the Nash game. Because under the condition  $\phi \rho(G) < 1$  both problems have a unique solution, it follows that the two problems share the same solution.

### 3 Generalized Nash Equilibrium Problems on Networks

#### 3.1 An Overview of GNEPs and the Variational Inequality Approach to their Solution

In GNEPs each player’s strategy set may depend on the strategies of the other players. We consider here the simplified framework where  $A_i \subseteq \mathbb{R}_+$  and we are given a function  $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$  which describes the shared constraints. The strategy set of player  $i$  is then written as

$$K_i(a_{-i}) = \{a_i \in \mathbb{R}_+ : g(a) = g(a_i, a_{-i}) \leq 0\}.$$



Thus, players share a common constraint  $g$  and have an additional individual nonnegativity constraint. With these ingredients, the GNEP is the problem of finding  $a^* \in \mathbb{R}^n$  such that, for any  $i \in \{1, \dots, n\}$ ,  $a_i^* \in K_i(a_{-i}^*)$  and

$$u_i(a_i^*, a_{-i}^*) \geq u_i(a_i, a_{-i}^*), \quad \forall a_i \in K_i(a_{-i}^*). \tag{7}$$

We will work under the common (although not minimal) assumptions that for each fixed  $a_{-i} \in A_{-i}$ , the functions  $u_i(\cdot, a_{-i})$  are concave and continuously differentiable, and the components of  $g$  are convex and continuously differentiable. As a consequence, a necessary and sufficient condition for  $a_i^* \in K_i(a_{-i}^*)$  to satisfy (7) is

$$-\frac{\partial u_i(a_i^*, a_{-i}^*)}{\partial a_i}(a_i - a_i^*) \geq 0, \quad \forall a_i \in K_i(a_{-i}^*). \tag{8}$$

Thus, if we define  $F(a)$  as in (3), and

$$K(a) = K_1(a_{-1}) \times \dots \times K_n(a_{-n}),$$

it follows that  $a^*$  is a GNE if and only if  $a^* \in K(a^*)$  and

$$[F(a^*)]^\top (a - a^*) \geq 0, \quad \forall a \in K(a^*). \tag{9}$$

The problem above, where also the feasible set depends on the solution, is called a quasi-variational inequality and its solution is as difficult as the original GNEP.

Assume now that  $a^*$  is a solution of GNEP. Hence, for each  $i$ ,  $a_i^*$  solves the maximization problem

$$\max_{a_i} \{u_i(a_i, a_{-i}^*) : g(a_i, a_{-i}^*) \leq 0, a_i \geq 0\}.$$

Under some standard constraint qualification we can then write the KKT conditions for each maximization problem. We then introduce the Lagrange multiplier  $\lambda^i \in \mathbb{R}^m$  associated with the constraint  $g(a_i, a_{-i}^*) \leq 0$  and the multiplier  $\mu_i \in \mathbb{R}$  associated with the nonnegativity constraint  $a_i \geq 0$ . The Lagrangian function for each player  $i$  reads as:

$$L_i(a_i, a_{-i}^*, \lambda^i, \mu_i) = u_i(a_i, a_{-i}^*) - [g(a_i, a_{-i}^*)]^\top \lambda^i + \mu_i a_i$$

and the KKT conditions for all players are given by:

$$\nabla_{a_i} L_i(a_i^*, a_{-i}^*, \lambda^{i*}, \mu_i^*) = 0, \quad i = 1, \dots, n, \tag{10}$$

$$\lambda_j^{i*} g_j(a^*) = 0, \lambda_j^{i*} \geq 0, g_j(a^*) \leq 0, \quad i = 1, \dots, n, j = 1, \dots, m \tag{11}$$

$$\mu_i^* a_i^* = 0, \mu_i^* \geq 0, a_i^* \geq 0, \quad i = 1, \dots, n. \tag{12}$$

Conversely, under the assumptions made, if  $a^*, \lambda, \mu^*$ , where  $\lambda^* = (\lambda^{1*}, \dots, \lambda^{n*})$  and  $\mu^* = (\mu_1^*, \dots, \mu_n^*)$ , satisfy the KKT system (10)–(12), then  $a^*$  is a GNE.

**Definition 5** Let  $a^*$  be a GNE which together with the Lagrange multipliers  $\lambda^* = (\lambda^{1*}, \dots, \lambda^{n*})$  and  $\mu^* = (\mu_1^*, \dots, \mu_n^*)$  satisfies the KKT system of all players. We call  $a^*$  a *normalized equilibrium* if there exists a vector  $r \in \mathbb{R}_{++}^n$  and a vector  $\bar{\lambda} \in \mathbb{R}_+^m$  such that

$$\lambda^{i*} = \frac{\bar{\lambda}}{r_i}, \quad \forall i = 1, \dots, n,$$

which means that, for a normalized equilibrium, the multipliers of the constraints shared by all players are proportional to a common multiplier. In the special case  $r_i = 1$  for any  $i$ , i.e., the multipliers coincide for each player,  $a^*$  is called *variational equilibrium* (VE). Rosen [22] proved that if the feasible set, which in our case is:

$$K = \{a \in \mathbb{R}_+^n : g(a) \leq 0\}$$

is compact and convex, then there exists a normalized equilibrium for each  $r \in \mathbb{R}_{++}^n$ .

Now, let us define, for each  $r \in \mathbb{R}_{++}^n$ , the vector function  $F^r : \mathbb{R}^n \rightarrow \mathbb{R}^n$  as follows:

$$[F^r(a)]^\top := - \left( r_1 \frac{\partial u_1}{\partial a_1}(a), \dots, r_n \frac{\partial u_n}{\partial a_n}(a) \right).$$

The variational inequality approach for finding the normalized equilibria of the GNEP is expressed by the following theorem which can be viewed as a special case of Proposition 3.2 in [17] or of Theorem 6.1 in [14].

**Theorem 3**

1. Suppose that  $a^*$  is a solution of  $VI(F^r, K)$ , where  $r \in \mathbb{R}_{++}^n$ , a constraint qualification holds at  $a^*$  and  $(\bar{\lambda}, \bar{\mu}) \in \mathbb{R}^m \times \mathbb{R}^n$  are the multipliers associated to  $a^*$ . Then,  $a^*$  is a normalized equilibrium such that the multipliers  $(\lambda^{i*}, \mu_i^*)$  of each player  $i$  satisfy the following conditions:

$$\lambda^{i*} = \frac{\bar{\lambda}}{r_i}, \quad \mu_i^* = \frac{\bar{\mu}_i}{r_i}, \quad \forall i = 1, \dots, n.$$

2. If  $a^*$  is a normalized equilibrium such that the multipliers  $(\lambda^{i*}, \mu_i^*)$  of each player  $i$  satisfy the following conditions:

$$\lambda^{i*} = \frac{\bar{\lambda}}{r_i}, \quad \forall i = 1, \dots, n,$$

for some vector  $\bar{\lambda} \in \mathbb{R}_+^m$  and  $r \in \mathbb{R}_{++}^n$ , then  $a^*$  is a solution of  $VI(F^r, K)$  and  $(\bar{\lambda}, r_1\mu_1^*, \dots, r_n\mu_n^*)$  are the corresponding multipliers.

### 3.2 A Linear-Quadratic Network GNEP

In this section we investigate an extension to a GNEP of the linear-quadratic network game described in Section 2.2. Specifically, we assume the same network structure given by the adjacency matrix  $G$  and the same payoff functions defined as in (4), while the strategy set of player  $i$  is given by the usual individual constraint  $a_i \geq 0$  and an additional constraint, shared by all the players, on the total quantity of activities of all players, that is

$$K_i(a_{-i}) = \left\{ a_i \in \mathbb{R}_+ : \sum_{j=1}^n a_j \leq C \right\}, \quad i = 1, \dots, n,$$

where  $C > 0$  is a given parameter. Depending on the specific application, the additional constraint can have the meaning of a collective budget upper bound or of a limited availability of a certain commodity.

We know from Theorem 2 that if  $\phi\rho(G) < 1$ , then the linear-quadratic network game (without the new shared constraint) has a unique Nash equilibrium  $a^*$  given by (6). However, if  $a^*$  does not satisfy the shared constraint, i.e., it does not belong to the set

$$K = \left\{ a \in \mathbb{R}_+^n : \sum_{i=1}^n a_i \leq C \right\},$$

it cannot be a GNE for the new game. On the other hand, under the assumption  $\phi\rho(G) < 1$ , Theorem 1 guarantees that the pseudo-gradient  $F$  of the game, defined as

$$F(a) = (I - \phi G)a - \alpha \mathbf{1},$$

is strongly monotone, hence there exists a unique solution of  $VI(F, K)$ , i.e., there exists a unique variational equilibrium of the linear-quadratic network GNEP. The following result gives an explicit formula for such variational equilibrium and an expansion similar to (6).

**Theorem 4** *If  $\phi\rho(G) < 1$ , then the unique variational equilibrium  $\bar{a}$  of the linear-quadratic network GNEP, that is the unique solution of  $VI(F, K)$ , is given by the following formula:*

$$\bar{a} = \begin{cases} a^* = \alpha \sum_{p=0}^{\infty} \phi^p G^p \mathbf{1} & \text{if } \sum_{i=1}^n a_i^* \leq C, \\ \frac{Ca^*}{\sum_{i=1}^n a_i^*} = \frac{C \sum_{p=0}^{\infty} \phi^p G^p \mathbf{1}}{\sum_{p=0}^{\infty} \phi^p \mathbf{1}^\top G^p \mathbf{1}} & \text{if } \sum_{i=1}^n a_i^* > C, \end{cases} \quad (13)$$

where  $a^* = \alpha(I - \phi G)^{-1} \mathbf{1}$  is the Nash equilibrium of the linear-quadratic network game.

**Proof** Theorem 1 guarantees that the matrix  $I - \phi G$  is positive definite and the map  $F$  is strongly monotone. Therefore,  $VI(F, K)$  has a unique solution. If  $a^* \in K$ , then  $a^*$  solves  $VI(F, K)$  since  $F(a^*) = 0$ . Otherwise, if  $a^* \notin K$ , then  $\bar{a} = Ca^* / \sum_{i=1}^n a_i^* \in K$  since  $\bar{a}_i > 0$  for any  $i = 1, \dots, n$  and  $\sum_{i=1}^n \bar{a}_i = C$ . Moreover  $\bar{a}$  is a solution of the KKT system related to  $VI(F, K)$  with multipliers

$$\bar{\lambda} = \alpha \left( 1 - \frac{C}{\sum_{i=1}^n a_i^*} \right) > 0,$$

associated to the shared constraint, and  $\bar{\mu}_i = 0$  associated to  $a_i \geq 0$  for any  $i = 1, \dots, n$ . In fact, we have

$$F(\bar{a}) + \bar{\lambda} \mathbf{1} - \bar{\mu} = \frac{C\alpha(I - \phi G)(I - \phi G)^{-1} \mathbf{1}}{\sum_{i=1}^n a_i^*} - \alpha \mathbf{1} + \alpha \left( 1 - \frac{C}{\sum_{i=1}^n a_i^*} \right) \mathbf{1} = 0$$

$$\bar{\lambda} \geq 0, \quad \sum_{i=1}^n \bar{a}_i \leq C, \quad \bar{\lambda} \left( C - \sum_{i=1}^n \bar{a}_i \right) = 0$$

$$\bar{\mu}_i \geq 0, \quad \bar{a}_i \geq 0, \quad \bar{\mu}_i \bar{a}_i = 0, \quad \forall i = 1, \dots, n.$$

Therefore,  $\bar{a}$  solves  $VI(F, K)$ . Finally, the ratio between the expansions in (13) follows from the one for  $a^*$  given in (6):

$$\bar{a} = \frac{Ca^*}{\mathbf{1}^\top a^*} = \frac{C\alpha \sum_{p=0}^{\infty} \phi^p G^p \mathbf{1}}{\alpha \sum_{p=0}^{\infty} \phi^p \mathbf{1}^\top G^p \mathbf{1}} = \frac{C \sum_{p=0}^{\infty} \phi^p G^p \mathbf{1}}{\sum_{p=0}^{\infty} \phi^p \mathbf{1}^\top G^p \mathbf{1}}.$$

This concludes the proof. □

Notice that if  $a^* \in K$ , then the formula giving  $\bar{a}$  contains  $\alpha$  but does not  $C$ , while if  $a^* \notin K$ , it contains  $C$  but not  $\alpha$ .

### 4 Numerical Experiments

In this section, we show some preliminary numerical experiments for the linear-quadratic network GNEP described in Section 3.2 by means of two small-size test problems.

*Example 1* We consider the network shown in Figure 1 (see also [3]) with 11 nodes (players). The spectral radius of the adjacency matrix  $G$  is  $\rho(G) \simeq 4.4040$ . We set parameter  $\alpha = 1$  and chose five different values for  $\phi$ :

$$\phi = 0.3/\rho(G), \quad \phi = 0.5/\rho(G), \quad \phi = 0.7/\rho(G), \quad \phi = 0.9/\rho(G), \quad \phi = 0.95/\rho(G),$$

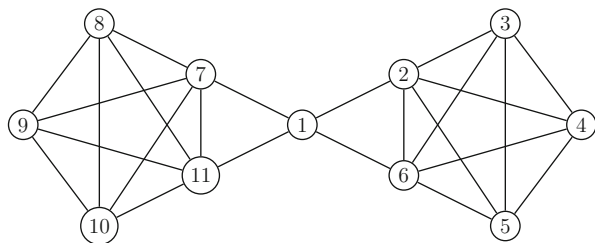
to guarantee the assumption of Theorem 4 holds. Moreover, we set  $C = 20$  in the shared constraint so that the variational equilibrium  $\bar{a}$  of the GNEP is different from the Nash equilibrium  $a^*$  of the classical network game. It follows from expansions in (13) that  $a^*$  and  $\bar{a}$  can be approximated by the sequences  $\{a_k^*\}$  and  $\{\bar{a}_k\}$ , respectively:

$$a_k^* = \alpha \sum_{p=0}^k \phi^p G^p \mathbf{1}, \quad \bar{a}_k = \frac{C \sum_{p=0}^k \phi^p G^p \mathbf{1}}{\sum_{p=0}^k \phi^p \mathbf{1}^\top G^p \mathbf{1}}.$$

Table 1 shows, for both sequences, the number of sums needed to get an approximation error less than  $10^{-t}$  for any  $t = 1, \dots, 10$ . Specifically, for any value of  $\phi$ , the numbers

$$\min \{k : \|a_k^* - a^*\|_\infty < 10^{-t}\} \quad \text{and} \quad \min \{k : \|\bar{a}_k - \bar{a}\|_\infty < 10^{-t}\}$$

**Fig. 1** Network topology of Example 1

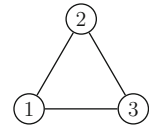


**Table 1** Speed of convergence of the sequences  $\{a_k^*\}$  and  $\{\bar{a}_k\}$  to  $a^*$  and  $\bar{a}$ , respectively

| Error      | $\phi = \frac{0.3}{\rho(G)}$ |    | $\phi = \frac{0.5}{\rho(G)}$ |    | $\phi = \frac{0.7}{\rho(G)}$ |    | $\phi = \frac{0.9}{\rho(G)}$ |     | $\phi = \frac{0.95}{\rho(G)}$ |     |
|------------|------------------------------|----|------------------------------|----|------------------------------|----|------------------------------|-----|-------------------------------|-----|
|            | NE                           | VE | NE                           | VE | NE                           | VE | NE                           | VE  | NE                            | VE  |
| $10^{-1}$  | 2                            | 1  | 4                            | 1  | 10                           | 1  | 44                           | 1   | 105                           | 1   |
| $10^{-2}$  | 4                            | 2  | 7                            | 3  | 16                           | 5  | 66                           | 9   | 150                           | 11  |
| $10^{-3}$  | 6                            | 4  | 11                           | 6  | 23                           | 11 | 88                           | 27  | 194                           | 43  |
| $10^{-4}$  | 8                            | 5  | 14                           | 9  | 29                           | 17 | 110                          | 48  | 239                           | 86  |
| $10^{-5}$  | 9                            | 7  | 17                           | 13 | 35                           | 23 | 132                          | 70  | 284                           | 131 |
| $10^{-6}$  | 11                           | 9  | 21                           | 16 | 42                           | 30 | 153                          | 92  | 329                           | 176 |
| $10^{-7}$  | 13                           | 11 | 24                           | 19 | 48                           | 36 | 175                          | 114 | 374                           | 221 |
| $10^{-8}$  | 15                           | 13 | 27                           | 23 | 55                           | 43 | 197                          | 136 | 419                           | 265 |
| $10^{-9}$  | 17                           | 15 | 31                           | 26 | 61                           | 49 | 219                          | 157 | 464                           | 310 |
| $10^{-10}$ | 19                           | 17 | 34                           | 29 | 68                           | 56 | 241                          | 179 | 509                           | 355 |

For any value of  $\phi$ , the first column (NE) reports  $\min \{k : \|a_k^* - a^*\|_\infty < 10^{-t}\}$ , while the second column (VE) reports  $\min \{k : \|\bar{a}_k - \bar{a}\|_\infty < 10^{-t}\}$ , for any  $t = 1, \dots, 10$

**Fig. 2** Network topology of Example 2



are reported in the first (NE) and second (VE) column, respectively, for any  $t = 1, \dots, 10$ .

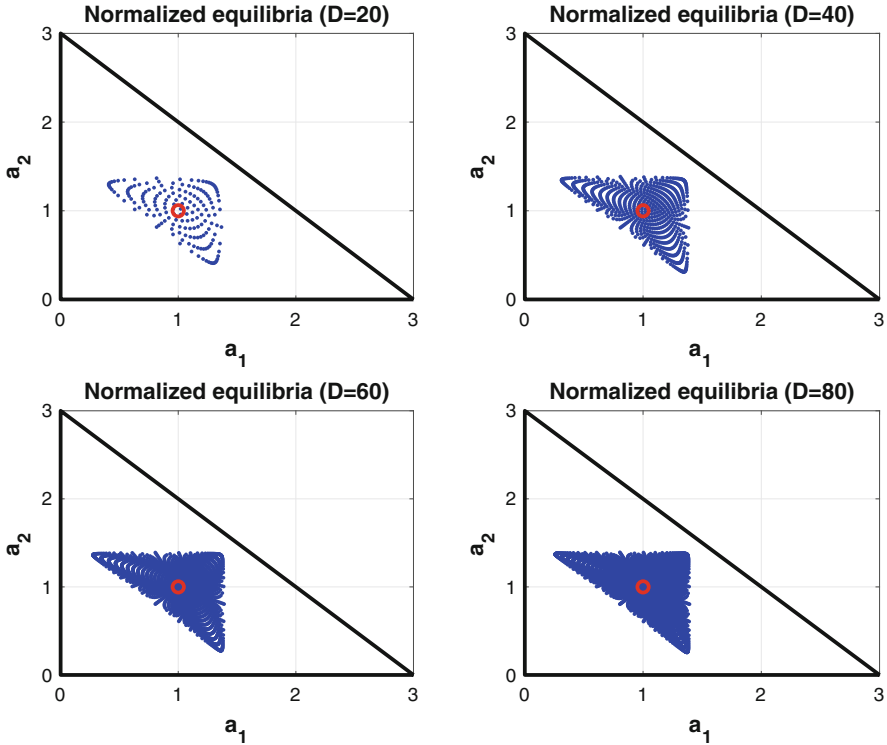
The results in Table 1 show that the convergence of  $\{\bar{a}_k\}$  to the variational equilibrium  $\bar{a}$  seems to be faster than the convergence of  $\{a_k^*\}$  to the Nash equilibrium  $a^*$ . Moreover, the more the value of  $\phi$  is close to  $1/\rho(G)$ , the more this gap is evident.

*Example 2* We now consider the network shown in Figure 2 with 3 nodes (players). The spectral radius of the adjacency matrix  $G$  is  $\rho(G) = 2$ . We set parameter  $\alpha = 1$ ,  $\phi = 0.25$  and  $C = 3$ , so that the Nash equilibrium is  $a^* = (2, 2, 2)^\top$  and the variational equilibrium is  $\bar{a} = (1, 1, 1)^\top$ . We exploit Theorem 3 to approximate the set of normalized equilibria. Consider the simplex

$$W = \left\{ r \in \mathbb{R}_{++}^3 : r_1 + r_2 + r_3 = 1 \right\}$$

of weights of the parametrized  $VI(F^r, K)$  and its discretization given by the finite set of vectors

$$\left( \frac{q_1}{D}, \frac{q_2}{D}, \frac{q_3}{D} \right)$$



**Fig. 3** Normalized equilibria (blue points) and variational equilibrium (red circle) of the linear-quadratic network GNEP

such that  $q_1, q_2, q_3$  and  $D$  are positive integers and  $q_1 + q_2 + q_3 = D$ . Figure 3 shows the set of normalized equilibria, projected on the plane  $(a_1, a_2)$ , for different values of  $D$ . Notice that all the found normalized equilibria belong to the plane  $a_1 + a_2 + a_3 = 3$ .

The results in Figure 3 suggest that the set of normalized equilibria is equal to

$$\left\{ a \in \mathbb{R}_+^3 : a_1 + a_2 + a_3 = 3, \quad a_i \leq 1.3924 \quad i = 1, 2, 3 \right\}$$

$$= \text{conv} \left\{ \begin{pmatrix} 1.3924 \\ 1.3924 \\ 0.2152 \end{pmatrix}, \begin{pmatrix} 1.3924 \\ 0.2152 \\ 1.3924 \end{pmatrix}, \begin{pmatrix} 0.2152 \\ 1.3924 \\ 1.3924 \end{pmatrix} \right\}.$$

Notice that, due to the symmetry of the considered network, the variational equilibrium  $\bar{a}$  is equal to the barycenter of the set of normalized equilibria.

## 5 Conclusions and Further Research Perspectives

In this note, we dealt with a network GNEP and derived a closed formula for its solution which involves the powers of the adjacency matrix, thus extending a previous result. To the best of our knowledge, this kind of formulas have been derived only in a few special cases and, because of their very interesting interpretation, it would be desirable to obtain similar results for more general problem classes. Another promising direction of research is the inclusion of random data in the model (see e.g. [10, 11]), which could be done by using tools from infinite-dimensional duality theory (see e.g. [8, 12, 14]).

**Acknowledgments** The authors are members of the Gruppo Nazionale per l'Analisi Matematica, la Probabilità e le loro Applicazioni (GNAMPA—National Group for Mathematical Analysis, Probability and their Applications) of the Istituto Nazionale di Alta Matematica (INdAM—National Institute of Higher Mathematics).

## References

1. T. Alpcan, T. Başar, A game-theoretic framework for congestion control in general topology networks, in *Proceedings of the IEEE 41st Conference on Decision and Control Las Vegas, December 10–13* (2002)
2. K. Atkinson, W. Han, *Theoretical Numerical Analysis* (Springer, Berlin, 2007)
3. C. Ballester, A. Calvo-Armengol, Y. Zenou, Who's who in networks. Wanted: the key player. *Econometrica* **74**, 1403–1417 (2006)
4. P. Bonacich, Power and centrality: a family of measures. *Am. J. Sociol.* **92**, 1170–1182 (1987)
5. F. Facchinei, C. Kanzov, Generalized Nash equilibrium problems. *Ann. Oper. Res.* **175**, 177–211 (2010)
6. F. Facchinei, J.-S. Pang, *Finite-Dimensional Variational Inequalities and Complementarity Problems* (Springer, Berlin, 2003)
7. F. Facchinei, A. Fischer, V. Piccialli, On generalized Nash games and variational inequalities. *Oper. Res. Lett.* **35**, 159–164 (2007)
8. F. Faraci, F. Raciti, On generalized Nash equilibrium problems in infinite dimension: the Lagrange multipliers approach. *Optimization* **64**, 321–338 (2015)
9. M.O. Jackson, Y. Zenou, Games on networks, in *Handbook of Game Theory with Economic Applications* (Elsevier, Amsterdam, 2015), pp. 95–163
10. B. Jadamba, F. Raciti, On the modelling of some environmental games with uncertain data. *J. Optim. Theory Appl.* **167**, 959–968 (2015)
11. B. Jadamba, A.A. Khan, F. Raciti, Regularization of stochastic variational inequalities and a comparison of an  $L_p$  and a sample-path approach. *Nonlinear Anal. Theory Methods Appl.* **94**, 65–83 (2014)
12. A.A. Khan, M. Sama, A new conical regularization for some optimization and optimal control problems: convergence analysis and finite element discretization. *Numer. Func. Anal. Opt.* **34**(8), 861–895 (2013)
13. D. Kinderlehrer, G. Stampacchia, *An Introduction to Variational Inequalities and their Applications* (Academic Press, New York, 1980)
14. G. Mastroeni, M. Pappalardo, F. Raciti, Generalized Nash equilibrium problems and variational inequalities in Lebesgue spaces. *Minimax Theory Appl.* **5**, 47–64 (2020)
15. D. Monderer, L.S. Shapley, Potential games. *Games Econ. Behav.* **14**, 124–143 (1996)



16. K. Nabetani, *Variational Inequality Approaches to Generalized Nash Equilibrium Problems*, Master Thesis. Department of Applied Mathematics and Physics Graduate School of Informatics, Kyoto University (2008)
17. K. Nabetani, P. Tseng, M. Fukushima, Parametrized variational inequality approaches to generalized Nash equilibrium problems with shared constraints. *Comput. Optim. Appl.* **48**, 423–452 (2011)
18. A. Nagurney, *Network Economics a Variational Inequality Approach* (Springer, Berlin, 1999)
19. A. Orda, R. Rom, N. Shimkin, Competitive routing in multiuser communication networks. *IEEE/ACM Trans. Netw.* **1**, 510–521 (1993)
20. F. Parise, A. Ozdaglar, A variational inequality framework for network games: Existence, uniqueness, convergence and sensitivity analysis. *Games Econ. Behav.* **114**, 47–82 (2019)
21. M. Passacantando, F. Raciti, Optimal road maintenance investment in traffic networks with random demands. *Optim. Lett.* (2019). <https://doi.org/10.1007/s11590-019-01493-y>
22. J.B. Rosen, Existence and uniqueness of equilibrium points for concave n-person games. *Econometrica* **33**, 520–534 (1965)

# Piecewise Polynomial Inversion of the Radon Transform in Three Space Dimensions via Plane Integration and Applications in Positron Emission Tomography



Nicholas E. Protonotarios, George A. Kastis, Nikolaos Dikaios,  
and Athanassios S. Fokas

**Abstract** The inversion of the celebrated Radon transform in three dimensions involves two-dimensional plane integration. This inversion provides the mathematical foundation of the important field of medical imaging, known as three-dimensional positron emission tomography (3D PET). In this chapter, we present an analytical expression for the inversion of the three-dimensional Radon transform, as well as a novel numerical implementation of this formula, based on piecewise polynomials of the third degree.

---

N. E. Protonotarios (✉)

Department of Applied Mathematics and Theoretical Physics, University of Cambridge, Cambridge, UK

Research Center of Mathematics, Academy of Athens, Athens, Greece  
e-mail: [np558@cam.ac.uk](mailto:np558@cam.ac.uk)

G. A. Kastis · N. Dikaios

Research Center of Mathematics, Academy of Athens, Athens, Greece  
e-mail: [gkastis@academyofathens.gr](mailto:gkastis@academyofathens.gr); [ndikaios@academyofathens.gr](mailto:ndikaios@academyofathens.gr)

A. S. Fokas

Research Center of Mathematics, Academy of Athens, Athens, Greece

Department of Applied Mathematics and Theoretical Physics, University of Cambridge, Cambridge, UK

Department of Biomedical Engineering, University of Southern California, Los Angeles, CA, USA

e-mail: [t.fokas@damtp.cam.ac.uk](mailto:t.fokas@damtp.cam.ac.uk)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*, Springer Optimization and Its Applications 173, [https://doi.org/10.1007/978-3-030-72563-1\\_17](https://doi.org/10.1007/978-3-030-72563-1_17)

## 1 Introduction

The celebrated Radon transform of a two-dimensional function is defined as the set of all its line integrals [1, 2]. The transition to three space dimensions yields a certain generalization of line integration, namely *plane integration*. Indeed, the Radon transform of a three-dimensional function is defined as the set of all its plane integrals.<sup>1</sup> The inversion of the three-dimensional Radon transform provides the mathematical foundation of the important field of medical imaging, known as 3D positron emission tomography (3D PET). The 3D Radon transform gives rise to an associated *inverse problem*, namely to “reconstruct” a function from its plane integrals. The main task in 3D PET imaging is the numerical implementation of the inversion of the 3D Radon transform.

In 3D PET, contrary to the conventional 2D PET, there is a certain generalization of the notion of image reconstruction: in the 2D case, the integration occurs in planes instead of lines. The difficulties arising in the 2D cases and their generalizations [3] are overcome in the 3D case. The inversion of the 3D Radon transform seems more straightforward than the one of the conventional Radon transform [4]. There are several numerical implementation methods in the literature, including: (i) the introduction of the concept of three-dimensional image reconstruction from “complete” projections [5]; (ii) the formulation of the 3D Radon transform for discrete 3D images (volumes), based on the summation over planes with small absolute slopes [6]; and (iii) the reconstruction of conductivities in the context of electric impedance tomography (EIT) [7]. The differences between 2D and 3D Radon transform inversion are emphasized in [8], and [9], where an analytic filter-backprojection method is introduced based on the spatially invariant detector point spread function. The authors of [10] proposed a spline-based inversion of the Radon transform in two and three dimensions; also the PET image reconstruction algorithms proposed in [11] show that analytic algorithms in 3D are linear and therefore allow easier control of the spatial resolution and noise correlations than in the case of the 2D reconstructions.

In this chapter, we present a novel formula for the inversion of the 3D Radon transform, as well as a novel numerical implementation of this formula, based on piecewise polynomial interpolation. We expect that our novel numerical implementation will enhance three-dimensional medical image reconstruction, especially in the case of 3D PET.

---

<sup>1</sup>Plane integrals are special cases of surface integrals, where the surface of integration is a plane.

## 2 The Radon Transform in Two Space Dimensions

In order to elucidate the properties of the three-dimensional Radon transform, it is essential to review the corresponding properties of the two-dimensional Radon transform.

A line  $L$  on the plane can be specified by the signed distance from the origin  $\rho$ , with  $-\infty < \rho < \infty$ , and the angle with the  $x_1$ -axis  $\theta$ , with  $0 \leq \theta < 2\pi$ , as in Figure 1. We denote the corresponding unit vectors perpendicular and parallel to  $L$  by  $\mathbf{n}$  and  $\mathbf{p}$ , respectively. These unit vectors are given by

$$\mathbf{n} = (-\sin \theta, \cos \theta)^T \quad \text{and} \quad \mathbf{p} = (\cos \theta, \sin \theta)^T, \tag{1}$$

with

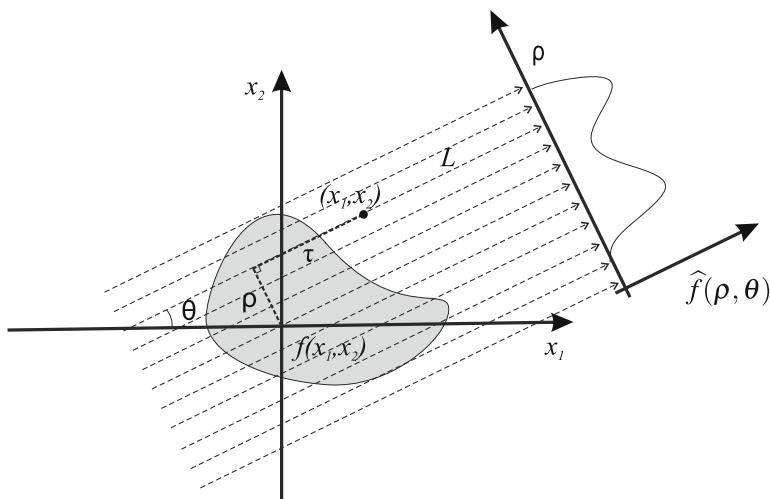
$$\mathbf{n} \cdot \mathbf{p} = 0. \tag{2}$$

Every point  $\mathbf{x} = (x_1, x_2)^T$  lying on the line  $L$  in Cartesian coordinates can be expressed in terms of the so-called *local coordinates*  $(\rho, \tau)$  via

$$\mathbf{x} = \rho \mathbf{n} + \tau \mathbf{p},$$

where  $\tau$  denotes the arc length. Therefore, we parameterize each point  $\mathbf{x}$  on the line  $L$  in the following manner:

$$\mathbf{x} := \mathbf{x}(\rho, \tau; \theta) = \begin{bmatrix} x_1(\rho, \tau; \theta) \\ x_2(\rho, \tau; \theta) \end{bmatrix} = \begin{bmatrix} \tau \cos \theta - \rho \sin \theta \\ \tau \sin \theta + \rho \cos \theta \end{bmatrix} \tag{3}$$



**Fig. 1** A two-dimensional function  $f(x_1, x_2)$  expressed in Cartesian coordinates, and its projections  $\hat{f}(\rho, \theta)$ , expressed in local coordinates

Through Equation (3), we can express the local coordinates  $(\rho, \tau)$  in terms of Cartesian coordinates  $(x_1, x_2)$  and the associated angle  $\theta$ :

$$\begin{bmatrix} \rho \\ \tau \end{bmatrix} := \begin{bmatrix} \rho(x_1, x_2; \theta) \\ \tau(x_1, x_2; \theta) \end{bmatrix} = \begin{bmatrix} x_2 \cos \theta - x_1 \sin \theta \\ x_2 \sin \theta + x_1 \cos \theta \end{bmatrix} \tag{4}$$

We define the line integral over all lines  $L$ , defined in Equation (1), of a two-dimensional Schwartz function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $f \in S(\mathbb{R}^2)$ , as its *two-dimensional Radon transform*,  $\mathcal{R}_2 f$ . In the context of 2D PET, the 2D Radon transform of the function  $f$  is usually stored in the form of the so-called *sinogram*, denoted by  $\widehat{f}(\rho, \theta)$

$$\mathcal{R}_2 f = \widehat{f}(\rho, \theta) = \int_L f ds, \tag{5}$$

where  $ds$  denotes an arc length differential, and  $S(\mathbb{R}^2)$  denotes the space of Schwartz functions in  $\mathbb{R}^2$ ,

$$S(\mathbb{R}^2) = \left\{ f \in C^\infty(\mathbb{R}^2) : \|f\|_{\alpha, \beta} < \infty \right\} \subset C^\infty(\mathbb{R}^2), \tag{6}$$

and

$$\|f\|_{\alpha, \beta} = \sup_{x \in \mathbb{R}^2} |x^\alpha D^\beta f(x)|, \quad \forall \text{ multi-index } \alpha, \beta, \quad |x^\alpha D^\beta f(x)| \rightarrow 0, \quad \text{as } |x| \rightarrow \infty. \tag{7}$$

Equation (5) may be rewritten via a parameterization  $\mathbf{x} := \mathbf{x}(\tau)$  of the line  $L$ , with  $\mathbf{x} : \mathbb{R}^2 \rightarrow L$ , as follows:

$$\widehat{f}(\rho, \theta) = \int_{-\infty}^{\infty} f(\mathbf{x}(\tau)) \|\mathbf{x}'(\tau)\|_2 d\tau, \tag{8}$$

where  $\|\cdot\|_2$  denotes the  $L^2$ -norm in  $\mathbb{R}^2$ . The parameterization provided by (3) will be proven to be very convenient and easy-to-manipulate, especially for the description of parallel lines. In this case, it is worth noting that

$$\|\mathbf{x}'(\tau)\|_2 = \sqrt{\left(\frac{dx_1}{d\tau}\right)^2 + \left(\frac{dx_2}{d\tau}\right)^2} = \cos^2 \theta + \sin^2 \theta = 1. \tag{9}$$

Hence, Equation (3) is a natural parameterization of the set of parallel lines  $L$ . Therefore, the 2D Radon transform defined in Equation (5) may be expressed as follows:

$$\mathcal{R}_2 f = \widehat{f}(\rho, \theta) = \int_{-\infty}^{\infty} f(\tau \cos \theta - \rho \sin \theta, \tau \sin \theta + \rho \cos \theta) d\tau, \tag{10}$$

with  $0 \leq \theta < 2\pi$  and  $-\infty < \rho < \infty$ . If we use a Dirac delta function, or a line impulse, then the Radon transform, denoted by  $\mathcal{R}_{2D}$  defined in (10) may be rewritten in the form:

$$\mathcal{R}_2 f = \widehat{f}(\rho, \theta) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x_1, x_2) \delta(\rho + x_1 \sin \theta - x_2 \cos \theta) dx_1 dx_2, \quad (11)$$

taking into account Equation (4).

The 2D Radon transform (10) gives rise to one of the most significant inverse problems in emission tomography. This specific inverse problem implies the “reconstruction” of the function  $f(x_1, x_2)$ , from its two-dimensional Radon transform, i.e. the function  $\widehat{f}(\rho, \theta)$ .

### 3 The Radon Transform in Three Space Dimensions

In the two-dimensional case, the Radon transform is considered on sets of parallel lines. This consideration implies the involvement of line integrals. However, in the three-dimensional case, the Radon transform is restricted on two-dimensional planes. In this direction, the transition to three space dimensions yields a certain generalization of line integration, namely *plane integration*.<sup>2</sup>

Therefore, we define the surface integral over all planes  $P$  of a three-dimensional Schwartz function  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ ,  $f \in S(\mathbb{R}^3)$ , as its *three-dimensional Radon transform*,  $\mathcal{R}_3 f$ . In the context of 3D PET, the 3D Radon transform of the function  $f$  is usually stored in the form of the so-called *3D sinogram*, denoted by  $\widehat{f}(\rho, \theta, \phi)$ :

$$\mathcal{R}_3 f = \widehat{f}(\rho, \theta, \phi) = \iint_P f ds, \quad (12)$$

where  $ds$  denotes an area differential and  $S(\mathbb{R}^3)$  denotes the space of Schwartz functions in  $\mathbb{R}^3$ :

$$S(\mathbb{R}^3) = \left\{ f \in C^\infty(\mathbb{R}^3) : \|f\|_{\alpha, \beta} < \infty \right\} \subset C^\infty(\mathbb{R}^3), \quad (13)$$

and

$$\|f\|_{\alpha, \beta} = \sup_{x \in \mathbb{R}^3} |x^\alpha D^\beta f(x)|, \quad \forall \text{ multi-index } \alpha, \beta, \quad |x^\alpha D^\beta f(x)| \rightarrow 0, \quad \text{as } |x| \rightarrow \infty. \quad (14)$$

Equation (12) may be rewritten via a parameterization  $\mathbf{x} := \mathbf{x}(u, v)$  of the plane  $P$ , with  $\mathbf{x} : \mathbb{R}^2 \rightarrow P$ , as follows:

---

<sup>2</sup>Plane integrals are special cases of surface integrals, where the surface of integration is a plane.

$$\widehat{f}(\rho, \theta, \phi) = \int_{-\infty}^{\infty} f(\mathbf{x}(u, v)) \left| \frac{\partial \mathbf{x}}{\partial u} \times \frac{\partial \mathbf{x}}{\partial v} \right| du dv, \tag{15}$$

where, for the area differential we employed

$$ds = \left| \frac{\partial \mathbf{x}}{\partial u} \times \frac{\partial \mathbf{x}}{\partial v} \right| du dv. \tag{16}$$

In the three-dimensional setting, for convenience, we characterize each two-dimensional plane  $P$  by a vector and a scalar, namely:

(i) the unit normal vector  $\mathbf{n}$

$$\mathbf{n} = \begin{bmatrix} n_1 \\ n_2 \\ n_3 \end{bmatrix}, \quad \text{with } \sqrt{n_1^2 + n_2^2 + n_3^2} = 1, \quad \text{and} \tag{17}$$

(ii) the signed distance from the origin  $\rho$ .

The normal from the origin to the plane intersects the plane at the point  $\rho \mathbf{n}$ . Thus, if  $\mathbf{x} = (x_1, x_2, x_3)^T$  is a point on the plane under investigation, then

$$(\rho \mathbf{n} - \mathbf{x}) \cdot \mathbf{n} = 0 \tag{18}$$

The above implies

$$\begin{bmatrix} \rho n_1 - x_1 \\ \rho n_2 - x_2 \\ \rho n_3 - x_3 \end{bmatrix} \cdot \begin{bmatrix} n_1 \\ n_2 \\ n_3 \end{bmatrix} = 0,$$

or

$$\rho(n_1^2 + n_2^2 + n_3^2) - (n_1 x_1 + n_2 x_2 + n_3 x_3) = 0.$$

Hence, the equation of the plane is

$$\rho - \mathbf{n} \cdot \mathbf{x} = 0, \quad \forall \mathbf{x} \in P. \tag{19}$$

We suppose that the plane of integration  $P$ , specified by its signed distance from the origin  $\rho$  and its unit normal vector  $\mathbf{n}$ , intersects the  $x_1 x_2$ -plane in an angle  $\theta$ , and the  $x_2 x_3$ -plane in an angle  $\phi$  (spherical angles). In this connection, the unit normal vector  $\mathbf{n}$  is uniquely specified by the two spherical angles, i.e.  $\mathbf{n} := \mathbf{n}(\theta, \phi)$ . Thus  $f(\rho, \mathbf{n})$  involves three variables, namely  $f(\rho, \mathbf{n}) = f(\rho, \theta, \phi)$ . In this direction, we characterize  $\mathbf{n}$  in terms of spherical angles, as follows:

$$\mathbf{n}(\theta, \phi) = \begin{bmatrix} \sin \theta \cos \phi \\ \sin \theta \sin \phi \\ \cos \theta \end{bmatrix}, \quad 0 \leq \theta \leq \pi \quad 0 \leq \phi \leq 2\pi. \tag{20}$$

If a point  $\mathbf{x} = (x_1, x_2, x_3)^T$  lies on the plane of integration  $P$ , i.e.  $\mathbf{x} \in P$ , then Equation (19) implies

$$\rho - x_1 \sin \theta \cos \phi - x_2 \sin \theta \sin \phi - x_3 \cos \theta = 0. \quad (21)$$

Equation (21) provides a convenient parameterization of the plane under investigation via  $x_1$  and  $x_2$ , treated hereafter as independent variables. In this setting,  $x_3$  will be considered as a dependent variable. Taking into account the parameterization induced by Equation (21), we rewrite the equation of a point  $\mathbf{x}$  lying on the plane of integration  $P$  as  $\mathbf{x} = \mathbf{x}(x_1, x_2)$ . More specifically, if  $\mathbf{x} \in P$ , then:

$$\mathbf{x} := \mathbf{x}(x_1, x_2; \rho, \theta, \phi) = (x_1, x_2, \csc \theta (\rho - x_1 \sin \theta \cos \phi - x_2 \sin \theta \sin \phi))^T, \quad (22)$$

where  $\csc \theta$  denotes the cosecant of the angle  $\theta$ , i.e.,

$$\csc \theta = \frac{1}{\cos \theta}. \quad (23)$$

Hence, the area differential  $ds$  is given by

$$ds = \left| \frac{\partial \mathbf{x}(x_1, x_2)}{\partial x_1} \times \frac{\partial \mathbf{x}(x_1, x_2)}{\partial x_2} \right| dx_1 dx_2, \quad (24)$$

where

$$\frac{\partial \mathbf{x}}{\partial x_1} \times \frac{\partial \mathbf{x}}{\partial x_2} = \begin{vmatrix} \hat{\mathbf{x}}_1 & \hat{\mathbf{x}}_2 & \hat{\mathbf{x}}_3 \\ \frac{\partial x_1}{\partial x_1} & \frac{\partial x_2}{\partial x_1} & \frac{\partial x_3}{\partial x_1} \\ \frac{\partial x_1}{\partial x_2} & \frac{\partial x_2}{\partial x_2} & \frac{\partial x_3}{\partial x_2} \end{vmatrix} = \begin{vmatrix} \hat{\mathbf{x}}_1 & \hat{\mathbf{x}}_2 & \hat{\mathbf{x}}_3 \\ 1 & 0 & -\tan \theta \cos \phi \\ 0 & 1 & \tan \theta \cos \phi \end{vmatrix} = \begin{bmatrix} \tan \theta \cos \phi \\ \tan \theta \cos \phi \\ 1 \end{bmatrix}, \quad (25)$$

and  $\hat{\mathbf{x}}_1$ ,  $\hat{\mathbf{x}}_2$ , and  $\hat{\mathbf{x}}_3$  are the corresponding unit vectors in the  $x_1$ ,  $x_2$ , and  $x_3$  directions, respectively. Taking into account Equation (25), the magnitude of the above ‘‘Jacobian’’ vector is,

$$\left| \frac{\partial \mathbf{x}}{\partial x_1} \times \frac{\partial \mathbf{x}}{\partial x_2} \right| = \sqrt{\tan^2 \theta + 1} = \csc \theta. \quad (26)$$

Hence Equation (24) yields

$$ds = \csc \theta dx_1 dx_2. \quad (27)$$

Thus, Equation (15) becomes

$$\hat{f}(\rho, \theta, \phi) = \int_{-\infty}^{\infty} dx_1 \int_{-\infty}^{\infty} \csc \theta f(x_1, x_2, \csc \theta (\rho - x_1 \sin \theta \cos \phi - x_2 \sin \theta \sin \phi)) dx_2. \quad (28)$$



An alternative way to express the three-dimensional Radon transform of a function  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$  involves a Dirac delta, or “plane impulse”, namely

$$\mathcal{R}_3 f = \widehat{f}(\rho, \theta, \phi) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x_1, x_2, x_3) \times \delta(\rho - x_1 \sin \theta \cos \phi - x_2 \sin \theta \sin \phi - x_3 \cos \theta) dx_1 dx_2 dx_3. \tag{29}$$

The alternative definition given by Equation (29) will be proven very useful (see Theorem 1) for the inversion and the numerical implementation of the three-dimensional Radon transform, as discussed in Sections 4 and 5.

### 4 The Inversion of the Radon Transform in Three Space Dimensions via Plane Integration

For the analytical inversion of the Radon transform in three space dimensions defined in Equation (28) we shall employ plane integration. In this direction, we will make use of the so-called *central slice theorem* (CST). This specific theorem, applied in the three-dimensional case, provides a fundamental tool for the Fourier-based inversion of the 3D Radon transform.

**Theorem 1 (Central Slice Theorem in 3D)** *The three-dimensional Fourier transform  $\mathcal{F}_3$  of a function  $f(x_1, x_2, x_3)$ , usually denoted by  $\tilde{f} = \mathcal{F}_3 f$ , equals the one-dimensional Fourier transform with respect to the signed distance from the origin  $\mathcal{F}_1^{(\rho)}$  of the three-dimensional Radon transform  $\mathcal{R}_3$  of the same function  $\widehat{f} = \mathcal{R}_3 f$ , i.e.*

$$\mathcal{F}_3 f = \mathcal{F}_1^{(\rho)} \mathcal{R}_3 f, \quad \text{or} \quad \tilde{f} = \mathcal{F}_1^{(\rho)} \widehat{f}, \tag{30}$$

where

$$\begin{aligned} \tilde{f}(k_1, k_2, k_3) &:= (\mathcal{F}_3 f)(k_1, k_2, k_3) \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x_1, x_2, x_3) e^{-2\pi i(k_1 x_1 + k_2 x_2 + k_3 x_3)} dx_1 dx_2 dx_3, \end{aligned} \tag{31}$$

and

$$\mathcal{F}_1^{(\rho)} \widehat{f} = \int_{-\infty}^{\infty} \widehat{f}(\rho, \theta, \phi) e^{-2\pi i k \rho} d\rho. \tag{32}$$

**Proof** For the proof of the central slice theorem in three dimensions, it is convenient to employ the alternative definition of the three-dimensional Radon transform as

provided via a delta function in Equation (29). In this case, we expand Equation (32) as follows:

$$\begin{aligned}
 \mathcal{F}_1^{(\rho)} \hat{f} &= \int_{-\infty}^{\infty} \hat{f}(\rho, \theta, \phi) e^{-2\pi i k \rho} d\rho & (33) \\
 &= \int_{-\infty}^{\infty} \left( \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x_1, x_2, x_3) \right. \\
 &\quad \left. \times \delta(\rho - x_1 \sin \theta \cos \phi - x_2 \sin \theta \sin \phi - x_3 \cos \theta) dx_1 dx_2 dx_3 \right) e^{-2\pi i k \rho} d\rho \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x_1, x_2, x_3) \left( \int_{-\infty}^{\infty} e^{-2\pi i k \rho} \right. \\
 &\quad \left. \times \delta(\rho - x_1 \sin \theta \cos \phi - x_2 \sin \theta \sin \phi - x_3 \cos \theta) d\rho \right) dx_1 dx_2 dx_3 \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x_1, x_2, x_3) e^{-2\pi i k (x_1 \sin \theta \cos \phi + x_2 \sin \theta \sin \phi + x_3 \cos \theta)} \\
 &\quad dx_1 dx_2 dx_3 \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x_1, x_2, x_3) e^{-2\pi i [(k \sin \theta \cos \phi)x_1 + (k \sin \theta \sin \phi)x_2 + (k \cos \theta)x_3]} \\
 &\quad dx_1 dx_2 dx_3 \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x_1, x_2, x_3) e^{-2\pi i (k_1 x_1 + k_2 x_2 + k_3 x_3)} dx_1 dx_2 dx_3 \\
 &= \tilde{f},
 \end{aligned}$$

where we introduced the new  $k$ -variables in the Fourier space vector  $\mathbf{k} := (k_1, k_2, k_3)^T$  of spatial frequencies as follows:

$$k_1 = k \sin \theta \cos \phi, \quad k_2 = k \sin \theta \sin \phi, \quad k_3 = k \cos \theta, \tag{34}$$

as the new  $k$ -variables in the Fourier space. □

Hence, the Fourier transform with respect to  $\rho$  of the “data” equals the three-dimensional Fourier transform of the function under investigation evaluated at the new set of variables. The inversion of Equation (30) yields

$$f = \mathcal{F}_3^{-1} \tilde{f}. \tag{35}$$

For the inversion of the 3D Radon transform we will utilize the following corollary.

**Corollary 1 (One-Dimensional Fourier Transform of the Second Derivative)**

For any twice differentiable function  $g$  with respect to  $\rho$ , the one-dimensional Fourier transform of the second derivative of  $g$ ,  $\mathcal{F}_1^{(\rho)} g''$ , is related with the one-dimensional Fourier transform of  $g$ ,  $\mathcal{F}_1^{(\rho)} g$ , by the following expression:

$$\left(\mathcal{F}_1^{(\rho)} g''\right)(\xi) = -4\pi^2 \xi^2 G(\xi), \tag{36}$$

where  $\mathcal{F}_1^{(\rho)}$  denotes the one-dimensional Fourier transform with respect to  $\rho$  defined in Equation (32),  $g''$  denotes the second derivative of  $g$  with respect to  $\rho$ , i.e.,

$$g'' = \frac{\partial^2 g}{\partial \rho^2}, \tag{37}$$

and  $G$  denotes the one-dimensional Fourier transform of  $g$ ,

$$G = \mathcal{F}_1^{(\rho)} g. \tag{38}$$

**Proof** Inverting Equation (32) and employing  $g$  instead of  $\tilde{f}$  yields

$$g = \left\{ \mathcal{F}_1^{(\rho)} \right\}^{-1} G = \int_{-\infty}^{\infty} G(\xi) e^{2\pi i \rho \xi} d\xi, \tag{39}$$

where  $G$  is defined in Equation (38). As in [12], we take the second derivative of both sides of Equation (39):

$$\begin{aligned} g'' &:= \frac{\partial^2 g}{\partial \rho^2} = \frac{\partial^2}{\partial \rho^2} \left( \int_{-\infty}^{\infty} G(\xi) e^{2\pi i \rho \xi} d\xi \right) \\ &= \int_{-\infty}^{\infty} G(\xi) \left[ \frac{\partial^2}{\partial \rho^2} \left( e^{2\pi i \rho \xi} \right) \right] d\xi \\ &= \int_{-\infty}^{\infty} G(\xi) \left[ (2\pi i \xi)^2 e^{2\pi i \rho \xi} \right] d\xi \\ &= \int_{-\infty}^{\infty} \left[ -4\pi^2 \xi^2 G(\xi) \right] e^{2\pi i \rho \xi} d\xi. \end{aligned} \tag{40}$$

From the above it is clear that the functions  $g''$  and  $-4\pi^2 \xi^2 G(\xi)$  form a Fourier transform pair. Hence, Equation (40), combined with Equation (39), imply Equation (36). □

**Theorem 2 (Inversion of the Three-Dimensional Fourier Transform)** The inverse of the three dimensional Radon transform  $\tilde{f} = \mathcal{R}_3 f$ , defined in Equations (28) and (29), of a Schwartz function  $f \in \mathcal{S}(\mathbb{R}^3)$  is given by

$$f(x_1, x_2, x_3) = -\frac{1}{4\pi^2} \int_0^\pi \sin \theta d\theta \int_0^{2\pi} \tilde{f}''(\rho^*, \theta, \phi) d\phi, \tag{41}$$

where, as in Equation (37), prime denotes differentiation with respect to  $\rho$ , i.e.

$$\tilde{f}''(\rho^*, \theta, \phi) = \left. \frac{\partial^2}{\partial \rho^2} \widehat{f}(\rho, \theta, \phi) \right|_{\rho=\rho^*}, \tag{42}$$

and  $\rho^*$  is given by

$$\rho^* = x_1 \sin \theta \cos \phi + x_2 \sin \theta \sin \phi + x_3 \cos \theta. \tag{43}$$

**Proof** Equation (35) implies

$$\left( \mathcal{F}_3^{-1} g \right) (x_1, x_2, x_3) = \int_{-\infty}^\infty \int_{-\infty}^\infty \int_{-\infty}^\infty g(x_1, x_2, x_3) e^{2\pi i(k_1 x_1 + k_2 x_2 + k_3 x_3)} dk_1 dk_2 dk_3. \tag{44}$$

In this direction, the inversion of the three-dimensional Fourier transform will reveal the unknown function  $f$  in the sense that:

$$f(x_1, x_2, x_3) = \int_{-\infty}^\infty \int_{-\infty}^\infty \int_{-\infty}^\infty \tilde{f}(k_1, k_2, k_3) e^{2\pi i(k_1 x_1 + k_2 x_2 + k_3 x_3)} dk_1 dk_2 dk_3. \tag{45}$$

We proceed by making a change of variables from  $(k_1, k_2, k_3)$  to  $(k, \theta, \phi)$  defined by Equation (34). The corresponding Jacobian is given by

$$J(k, \theta, \phi) = \begin{vmatrix} \frac{\partial k_1}{\partial k} & \frac{\partial k_1}{\partial \theta} & \frac{\partial k_1}{\partial \phi} \\ \frac{\partial k_2}{\partial k} & \frac{\partial k_2}{\partial \theta} & \frac{\partial k_2}{\partial \phi} \\ \frac{\partial k_3}{\partial k} & \frac{\partial k_3}{\partial \theta} & \frac{\partial k_3}{\partial \phi} \end{vmatrix} = \begin{vmatrix} \sin \theta \cos \phi & k \cos \theta \cos \phi & -k \sin \theta \sin \phi \\ \sin \theta \sin \phi & k \cos \theta \sin \phi & k \sin \theta \cos \phi \\ \cos \theta & -k \sin \theta & 0 \end{vmatrix} = k^2 \sin \theta. \tag{46}$$

Thus we modify Equation (45) in the following manner:

$$f(x_1, x_2, x_3) = \int_0^\pi \int_0^{2\pi} \int_{-\infty}^\infty \tilde{f}(k \sin \theta \cos \phi, k \sin \theta \sin \phi, k \cos \theta) \times e^{2\pi i k(\sin \theta \cos \phi x_1 + \sin \theta \sin \phi x_2 + \cos \theta x_3)} J(k, \theta, \phi) dk d\theta d\phi. \tag{47}$$

However, a point  $\mathbf{x} = (x_1, x_2, x_3)^T$  lying on the plane of integration  $P$  ( $\mathbf{x} \in P$ ), according to Equation (21), satisfies

$$\rho^* - x_1 \sin \theta \cos \phi - x_2 \sin \theta \sin \phi - x_3 \cos \theta = 0, \tag{48}$$

where  $\rho^*$  is the signed distance of the plane  $P$  from the origin, see Equation (43). Hence, taking into account Equations (46), (48), and (47) may be rewritten as

$$f(x_1, x_2, x_3) = \int_0^\pi \sin \theta d\theta \int_0^{2\pi} \left[ \int_{-\infty}^\infty \tilde{f}(k \sin \theta \cos \phi, k \sin \theta \sin \phi, k \cos \theta) \times e^{2\pi i k \rho^*} k^2 dk \right] d\phi. \tag{49}$$

We combine Equations (30) with (32) and (34) to obtain

$$\tilde{f}(k \sin \theta \cos \phi, k \sin \theta \sin \phi, k \cos \theta) = \int_{-\infty}^\infty \hat{f}(\rho, \theta, \phi) e^{-2\pi i k \rho} d\rho. \tag{50}$$

In Equation (49), we replace  $\tilde{f}$  by the right-hand side of Equation (50)

$$f(x_1, x_2, x_3) = \int_0^\pi \sin \theta d\theta \int_0^{2\pi} \left[ \int_{-\infty}^\infty \left( \int_{-\infty}^\infty \hat{f}(\rho, \theta, \phi) e^{-2\pi i k \rho} d\rho \right) e^{2\pi i k \rho^*} k^2 dk \right] d\phi. \tag{51}$$

We denote the one-dimensional Fourier transform of the three-dimensional Radon transform of  $f$  by  $\hat{F}$ ,

$$\hat{F}(k) := \int_{-\infty}^\infty \hat{f}(\rho, \theta, \phi) e^{-2\pi i k \rho} d\rho, \tag{52}$$

and insert Equation (52) into Equation (51):

$$f(x_1, x_2, x_3) = \int_0^\pi \sin \theta d\theta \int_0^{2\pi} \left[ \int_{-\infty}^\infty k^2 \hat{F}(k) e^{2\pi i k \rho^*} dk \right] d\phi. \tag{53}$$

The final step involves the rewriting of Equation (53) in the following manner:

$$f(x_1, x_2, x_3) = -\frac{1}{4\pi^2} \int_0^\pi \sin \theta d\theta \int_0^{2\pi} \left[ \int_{-\infty}^\infty \left( -4\pi^2 k^2 \hat{F}(k) \right) e^{2\pi i k \rho^*} dk \right] d\phi. \tag{54}$$

We employ Corollary 1, and replace the integral inside the brackets on the left-hand side of Equation (54), by the left-hand side of the first line of Equation (40) to obtain

$$f(x_1, x_2, x_3) = -\frac{1}{(4\pi)^2} \int_0^\pi \sin \theta d\theta \int_0^{2\pi} \frac{\partial^2 \tilde{f}(\rho, \theta, \phi)}{\partial \rho^2} \Big|_{\rho=\rho^*} d\phi, \tag{55}$$

which, via Equation (42), is Equation (41). □

## 5 Numerical Implementation of the Inversion of the Radon Transform in Three Space Dimensions via Piecewise Cubic Polynomials

For the numerical implementation of the inversion of the Radon transform in three space dimensions we will employ piecewise continuous cubic polynomials, namely cubic splines. It is important to note that all integrals involving the second derivative with respect to  $\rho$  of the 3D Radon transform will be evaluated at  $\rho = \rho^*$ , namely at

$$\rho^* = x_1 \sin \theta \cos \phi + x_2 \sin \theta \sin \phi + x_3 \cos \theta. \tag{56}$$

As shown in the previous section, the 3D inverse Radon transform can be expressed as

$$f(x_1, x_2, x_3) = -\frac{1}{4\pi^2} \int_0^\pi \sin \theta d\theta \int_0^{2\pi} \widehat{f}''(\rho, \theta, \phi) \Big|_{\rho=\rho^*} d\phi \tag{57}$$

where  $\rho^*$ ,  $\widehat{f}$  and  $\widehat{f}''$  are defined in Equations (56), (28), and (42), respectively.

We assume that the three-dimensional Radon transform,  $\widehat{f}$ , is given for every  $\theta$  and every  $\phi$  at the  $n$  knots  $\{\rho_i\}_1^n$ . We denote the value of  $\widehat{f}$  at  $\rho_i$  by  $\widehat{f}_i$ , namely

$$\widehat{f}_i = \widehat{f}(\rho_i, \theta, \phi), \quad \theta \in [0, \pi], \quad \phi \in [0, 2\pi], \quad i = 1, \dots, n - 1. \tag{58}$$

We also assume that both  $\widehat{f}(\rho, \theta, \phi)$  and  $\widehat{f}'(\rho, \theta, \phi)$ , where

$$\widehat{f}'(\rho, \theta, \phi) = \frac{\partial \widehat{f}(\rho, \theta, \phi)}{\partial \rho}, \tag{59}$$

vanish at the endpoints  $\rho_1 = -1$  and  $\rho_n = 1$ , i.e.

$$\widehat{f}(\rho_1, \theta, \phi) = \widehat{f}(\rho_n, \theta, \phi) = 0, \quad \theta \in [0, \pi], \quad \phi \in [0, 2\pi], \tag{60}$$

and

$$\frac{\partial}{\partial \rho} \widehat{f}(\rho_1, \theta, \phi) = \frac{\partial}{\partial \rho} \widehat{f}(\rho_n, \theta, \phi) = 0, \quad \theta \in [0, \pi], \quad \phi \in [0, 2\pi]. \tag{61}$$

In each interval  $[\rho_i, \rho_{i+1}]$ ,  $i = 1, \dots, n - 1$ , we approximate  $\widehat{f}(\rho, \theta, \phi)$  by the third-degree spline  $S_i^{(3)}$ , namely

$$\widehat{f}(\rho, \theta, \phi) \sim S_i^{(3)}(\rho, \theta, \phi), \quad \rho \in [\rho_i, \rho_{i+1}] \quad \theta \in [0, \pi], \quad \phi \in [0, 2\pi]. \tag{62}$$

The cubic spline  $S_i^{(3)}$  interpolates  $\widehat{f}$  at the points  $\{\rho_i\}_1^n$ :

$$S_i^{(3)}(\rho_i, \theta, \phi) = f_i, \quad i = 1, \dots, n. \tag{63}$$

Therefore, for  $\rho \in [\rho_i, \rho_{i+1}]$

$$S_i^{(3)}(\rho, \theta, \phi) = a_i(\theta, \phi) + b_i(\theta, \phi)\rho + c_i(\theta, \phi)\rho^2 + d_i(\theta, \phi)\rho^3. \tag{64}$$

Then, following Equation (62),

$$\frac{\partial}{\partial \rho} \widehat{f}(\rho, \theta, \phi) \sim \frac{\partial}{\partial \rho} S_i^{(3)}(\rho, \theta, \phi) =: S_i^{(2)}(\rho, \theta, \phi), \tag{65}$$

where

$$S_i^{(2)}(\rho, \theta, \phi) = b_i(\theta, \phi) + 2c_i(\theta, \phi)\rho + 3d_i(\theta, \phi)\rho^2. \tag{66}$$

Similarly,

$$\frac{\partial^2}{\partial \rho^2} \widehat{f}(\rho, \theta, \phi) \sim \frac{\partial^2}{\partial \rho^2} S_i^{(3)}(\rho, \theta, \phi) = \frac{\partial}{\partial \rho} S_i^{(2)}(\rho, \theta, \phi) =: S_i^{(1)}(\rho, \theta, \phi), \tag{67}$$

where

$$S_i^{(1)}(\rho, \theta, \phi) = 2c_i(\theta, \phi) + 6d_i(\theta, \phi)\rho \tag{68}$$

Hence, Equation (57) becomes

$$f(x, y, z) = -\frac{1}{4\pi^2} \int_0^\pi \int_0^\pi [2c_i(\theta, \phi) + 6d_i(\theta, \phi)\rho] \sin \theta d\phi d\theta \tag{69}$$

The constants  $c_i(\theta, \phi)$  and  $d_i(\theta, \phi)$  involved in the above inversion integral, are given by the following expressions, see [13]:

$$c_i(\theta) = \frac{1}{2\Delta_i} (\rho_{i+1} \widehat{f}_i'' - \rho_i \widehat{f}_{i+1}''), \tag{70a}$$

$$d_i(\theta) = \frac{\widehat{f}_{i+1}'' - \widehat{f}_i''}{6\Delta_i}, \tag{70b}$$

where

$$\Delta_i = \rho_{i+1} - \rho_i, \tag{70c}$$

and

$$\widehat{f}'_i := \frac{\partial^2}{\partial \rho^2} \widehat{f}(\rho_i, \theta, \phi). \tag{70d}$$

It is worth noting that the inversion formula (57) involves the known constants  $\{\widehat{f}\}_1^n$  and the unknown constants  $\{\widehat{f}''\}_1^n$ . For the computation of  $\{\widehat{f}''\}_1^n$ , we employ the continuity of the first derivative of the cubic spline, i.e.

$$S_i^{(2)}(\rho_{i+1}, \theta, \phi) = S_{i+1}^{(2)}(\rho_i, \theta, \phi), \quad i = 1, 2, \dots, n-2, \quad \theta \in [0, \pi], \quad \phi \in [0, 2\pi], \tag{71a}$$

and

$$S_1^{(2)}(\rho_1, \theta, \phi) = S_{n-1}^{(2)}(\rho_n, \theta, \phi) = 0, \quad \theta \in [0, \pi], \quad \phi \in [0, 2\pi]. \tag{71b}$$

The above consists of a system of  $n$  unknowns and of  $n$  equations, namely  $n - 2$  equations arising from Equation (71a) for  $i = 1, 2, \dots, n - 2$ , and 2 equations arising from Equation (71b). The continuity of the cubic spline itself, i.e.

$$S_i^{(3)}(\rho_{i+1}, \theta, \phi) = S_{i+1}^{(3)}(\rho_i, \theta, \phi) = 0, \quad i = 1, 2, \dots, n-2, \quad \theta \in [0, \pi], \quad \phi \in [0, 2\pi] \tag{72a}$$

and

$$S_1^{(3)}(\rho_1, \theta, \phi) = S_{n-1}^{(3)}(\rho_n, \theta, \phi) = 0, \quad \theta \in [0, \pi], \quad \phi \in [0, 2\pi]. \tag{72b}$$

The continuity of the cubic spline,  $S_i^{(3)}(\rho, \theta, \phi)$ , as expressed in Equations (72), implies that the knots  $\{\rho_i\}_1^n$  are *removable* logarithmic singularities.

**Acknowledgments** This work was partially supported by the research programme “Inverse Problems and Medical Imaging” (200/947) of the Research Committee of the Academy of Athens. A.S. Fokas has been supported by EPSRC, UK in the form of a senior fellowship.

## References

1. J. Radon, Über die bestimmung von funktionen durch ihre integralwerte längs gewisser mannigfaltigkeiten. *Akad. Wiss.* **69**, 262–277 (1917)
2. P. Kuchment, The Radon transform and medical imaging, in *CBMS-NSF Regional Conference Series in Applied Mathematics*. Society for Industrial and Applied Mathematics (2014)
3. N.E. Protonotarios, G.A. Kastis, A.S. Fokas, *A New Approach for the Inversion of the Attenuated Radon Transform* (Springer, Cham, 2019), pp. 433–457. [https://doi.org/10.1007/978-3-030-31339-5\\_16](https://doi.org/10.1007/978-3-030-31339-5_16)
4. M.-Y. Chiu, H.H. Barrett, R.G. Simpson, Three-dimensional reconstruction from planar projections. *JOSA* **70**(7), 755–762 (1980). <https://doi.org/10.1364/JOSA.70.000755>
5. M. Defrise, D. Townsend, R. Clack, Three-dimensional image reconstruction from complete projections. *Phys. Med. Biol.* **34**(5), 573 (1989). <https://doi.org/10.1088/0031-9155/34/5/002>



6. A. Averbuch, Y. Shkolnisky, 3d fourier based discrete radon transform. *Appl. Comput. Harmon. Anal.* **15**(1), 33–69 (2003). [https://doi.org/10.1016/S1063-5203\(03\)00030-7](https://doi.org/10.1016/S1063-5203(03)00030-7)
7. J. Bikowski, K. Knudsen, J.L. Mueller, Direct numerical reconstruction of conductivities in three dimensions using scattering transforms. *Inverse Probl.* **27**(1), 015002 (2010). <https://doi.org/10.1088/0266-5611/27/1/015002>
8. F.H. Fahey, Data acquisition in PET imaging. *J. Nucl. Med. Technol.* **30**(2), 39–49 (2002)
9. P. Kinahan, J. Rogers, Analytic 3D image reconstruction using all detected events. *IEEE Trans. Nucl. Sci.* **36**, 964–968 (1989). <https://doi.org/10.1109/23.34585>
10. P. La Rivière, X. Pan, Spline-based inverse radon transform in two and three dimensions. *IEEE Trans. Nucl. Sci.* **45**(4), 2224–2231 (1998). <https://doi.org/10.1109/23.708352>
11. M. Defrise, P.E. Kinahan, C.J. Michel, Image reconstruction algorithms in PET, in *Positron Emission Tomography* (Springer, Berlin, 2005), pp. 63–91. [https://doi.org/10.1007/1-84628-007-9\\_4](https://doi.org/10.1007/1-84628-007-9_4)
12. J. Gaskill, J.W. Sons, *Linear Systems, Fourier Transforms, and Optics*. *Wiley Series in Pure and Applied Optics* (Wiley, New York, 1978)
13. W.H. Press, S.A. Teukolsky, W.T. Vetterling, B.P. Flannery, *Numerical Recipes: The Art of Scientific Computing*, 3rd edn. (Cambridge University, New York, 2007)

# Factorization and Solution of Linear and Nonlinear Second Order Differential Equations with Variable Coefficients and Mixed Conditions



E. Providas

**Abstract** This chapter deals with the factorization and solution of initial and boundary value problems for a class of linear and nonlinear second order differential equations with variable coefficients subject to mixed conditions. The technique for nonlinear differential equations is based on their decomposition into linear components of the same or lower order and the factorization of the associated second order linear differential operators. The implementation and efficiency of the procedure is shown by solving several examples.

## 1 Introduction

One of the most important categories of ordinary differential equations is the second order differential equations with variable coefficients. Many problems from engineering and science are within this large class of differential equations. These equations, in addition to their natural significance, have also been used as a vehicle for the study of other higher order differential equations. Both exact and numerical methods have been developed for their solution [2]. Most of the explicit techniques rely on the knowledge of fundamental solutions. The factorization method does not require any fundamental solution of the given second order differential equation, but its applicability is limited to certain problems. For a review of the factorization of differential operators the interested reader can look at the selected articles [1, 3–9, 14–16].

Following the work in [10–13] and [17], this paper is concerned with the exact solution of a class of linear and nonlinear differential equations of second order with variable coefficients subject to nonlocal boundary conditions by direct factorization of the differential equation as well as the boundary conditions. Specifically, in Section 2, we recall some basic results and consider linear first order problems with

---

E. Providas (✉)

Department of Environmental Sciences, University of Thessaly, Larissa, Greece

e-mail: [providas@uth.gr](mailto:providas@uth.gr)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,

Springer Optimization and Its Applications 173,

[https://doi.org/10.1007/978-3-030-72563-1\\_18](https://doi.org/10.1007/978-3-030-72563-1_18)

397

a mixed boundary condition. In Section 3, we present the operator factorization method for solving, under certain conditions, the linear second order differential equation

$$u''(x) + p(x)u(x) + q(x)u(x) = f(x), \quad x \in (a, b), \quad (1)$$

where the coefficients  $p(x), q(x) \in C[a, b]$  and the forcing function  $f(x) \in C[a, b]$ , subject to general boundary conditions

$$\begin{aligned} \mu_{11}u(a) + \mu_{12}u(b) &= \beta_1, \\ \mu_{21}u'(a) + \mu_{22}u'(b) + \mu_{23}u(a) + \mu_{24}u(b) &= \beta_2, \end{aligned} \quad (2)$$

where  $\mu_{ij}, \beta_i \in \mathbb{R}, i = 1, 2, j = 1, 2, 3, 4$ . In Section 4, we deal with the construction of explicit solutions to two kinds of nonlinear differential equations of second order, which can be decomposed initially into linear second order differential equations. First, we consider the equation of the form

$$u''(x)u'(x) + [q(x)u'(x) + g(x)u''(x)]u(x) + q(x)g(x)u^2(x) = 0, \quad (3)$$

for  $x \in (a, b)$  and  $q(x), g(x) \in C[a, b]$ , along with the general boundary conditions (2). Also, we consider the nonlinear differential equation of the type

$$F\left(\frac{u''(x)}{u(x)}, x\right) = F(w(x), x) = w^2(x) + a(x)w(x) + b(x) = 0, \quad x \in (a, b), \quad (4)$$

where the nonlinear function  $F$  is a second degree polynomial of  $w(x) = u''(x)/u(x)$  and the coefficients  $a(x), b(x) \in C[a, b]$ , subject to general boundary conditions (2). Finally, some conclusions are quoted in Section 5.

## 2 Preliminaries

We first recall some basic results. A linear operator  $P : C[a, b] \rightarrow C[a, b]$  is said to be *correct* if  $P$  is injective,  $R(P) = C[a, b]$  and its inverse  $P^{-1}$  is bounded on  $C[a, b]$ . Let  $A : C[a, b] \rightarrow C[a, b]$  be the linear first order operator

$$Ay(x) = y'(x) + a(x)y(x), \quad D(A) = C^1[a, b], \quad (5)$$

where  $a(x) \in C[a, b]$ , and  $\widehat{A}$  be its restriction on

$$D(\widehat{A}) = \{y(x) \in D(A) : y(x_0) = y_0\}, \quad (6)$$

where  $x_0 \in [a, b]$  and  $y_0$  is an arbitrary real initial value. Then the following fundamental theorem holds.

**Theorem 1** *The linear operator  $\widehat{A}$  in (5) and (6) is correct and the unique solution of the initial value problem*

$$\widehat{A}y(x) = f(x), \quad \forall f(x) \in C[a, b], \tag{7}$$

is given by

$$y(x) = \widehat{A}^{-1}f(x) = e^{-\int_{x_0}^x a(t)dt} \left( y_0 + \int_{x_0}^x f(t)e^{\int_{t_0}^t a(\tau)d\tau} dt \right). \tag{8}$$

Accordingly, let  $B : C[a, b] \rightarrow C[a, b]$  be the linear second order operator

$$By(x) = y''(x) + b_1(x)y'(x) + b_2(x)y(x), \quad D(B) = C^2[a, b], \tag{9}$$

where  $b_1(x), b_2(x) \in C[a, b]$ , and  $\widehat{B}$  be its restriction on

$$D(\widehat{B}) = \{y(x) \in D(B) : y(x_0) = y_0, y'(x_0) = y'_0\}, \tag{10}$$

where  $x_0 \in [a, b]$  and  $y_0, y'_0$  is a couple of given real numbers. Then we have the next fundamental theorem.

**Theorem 2** *The linear operator  $\widehat{B}$  in (9) and (10) is correct and the initial value problem*

$$\widehat{B}y(x) = f(x), \quad \forall f(x) \in C[a, b], \tag{11}$$

has exactly one solution  $y(x) = \widehat{B}^{-1}f(x)$ .

We now consider a problem for a first order differential equation and a nonlocal boundary condition, which we will encounter below. For this, we prove the next theorem.

**Theorem 3** *Let the general linear first order problem with a nonlocal boundary condition*

$$\begin{aligned} \widehat{Q}y(x) &= y'(x) + q(x)y(x) = f(x), \\ D(\widehat{Q}) &= \left\{ y(x) \in C^1[a, b] : \mu_1y(a) + \mu_2y(b) = \beta \right\}, \end{aligned} \tag{12}$$

where the operator  $\widehat{Q} : C[a, b] \rightarrow C[a, b]$ , the given functions  $q(x), f(x) \in C[a, b]$ , and the constants  $\mu_1, \mu_2, \beta \in \mathbb{R}$ . If

$$\mu_1 + \mu_2e^{-\int_a^b q(t)dt} \neq 0, \tag{13}$$

then the operator  $\widehat{Q}$  is correct and the unique solution of problem (12) is given by

$$y(x) = \widehat{Q}^{-1} f(x) = e^{-\int_a^x q(t)dt} \left( C + \int_a^x f(t) e^{\int_a^t q(\tau)d\tau} dt \right), \quad (14)$$

where

$$C = \left( \mu_1 + \mu_2 e^{-\int_a^b q(t)dt} \right)^{-1} \left( \beta - \mu_2 e^{-\int_a^b q(t)dt} \int_a^b f(t) e^{\int_a^t q(\tau)d\tau} dt \right).$$

**Proof** It is known that the general solution of the first order differential equation in (12) is

$$y(x) = e^{-\int_{x_0}^x q(t)dt} \left( C + \int_{x_0}^x f(t) e^{\int_{x_0}^t q(\tau)d\tau} dt \right), \quad (15)$$

where  $x_0 \in [a, b]$ . For  $x_0 = a$ , we have

$$y(a) = C, \quad y(b) = e^{-\int_a^b q(t)dt} \left( C + \int_a^b f(t) e^{\int_a^t q(\tau)d\tau} dt \right).$$

Substituting these values into the boundary condition in (12), we obtain

$$\left( \mu_1 + \mu_2 e^{-\int_a^b q(t)dt} \right) C = \beta - \mu_2 e^{-\int_a^b q(t)dt} \int_a^b f(t) e^{\int_a^t q(\tau)d\tau} dt.$$

If relation (13) holds, then

$$C = \left( \mu_1 + \mu_2 e^{-\int_a^b q(t)dt} \right)^{-1} \left( \beta - \mu_2 e^{-\int_a^b q(t)dt} \int_a^b f(t) e^{\int_a^t q(\tau)d\tau} dt \right). \quad (16)$$

From (15) and (16) it is implied (14).  $\square$

### 3 Factorization Method for Linear Differential Equations

Let the linear differential operators of first order  $L_1 : C[a, b] \rightarrow C[a, b]$  and  $L_2 : C[a, b] \rightarrow C[a, b]$  be defined by

$$L_1 u(x) = [D + r(x)] u(x), \quad D(L_1) = C^1[a, b], \quad (17)$$

$$L_2 u(x) = [D + s(x)] u(x), \quad D(L_2) = C^1[a, b], \quad (18)$$

respectively, where  $D = \frac{d}{dx}$ , and the coefficients  $r(x) \in C[a, b]$  and  $s(x) \in C^1[a, b]$ . Consider the composition,

$$\begin{aligned} L_1 L_2 u(x) &= L_1 (L_2 u(x)) \\ &= [D + r(x)] ([D + s(x)] u(x)) \\ &= \left[ D^2 + (r(x) + s(x)) D + (s'(x) + r(x)s(x)) \right] u(x). \end{aligned} \quad (19)$$

This gives rise to the following proposition.

**Proposition 1** *Let the linear differential operator of second order  $L : C[a, b] \rightarrow C[a, b]$  be defined by*

$$Lu(x) = \left[ D^2 + p(x)D + q(x) \right] u(x), \quad D(L) = C^2[a, b], \quad (20)$$

where the coefficients  $p(x), q(x) \in C[a, b]$ . If there exist two functions  $r(x) \in C[a, b]$  and  $s(x) \in C^1[a, b]$  satisfying the relations

$$r(x) + s(x) = p(x), \quad (21)$$

$$s'(x) + r(x)s(x) = q(x), \quad (22)$$

then the operator  $L$  can be factorized into a product of the two linear differential operators of first order  $L_1, L_2$  in (17) and (18), respectively, such that

$$Lu(x) = L_1 L_2 u(x). \quad (23)$$

*Remark 1* By solving equation (21) with respect to  $r(x)$  and then substituting into (22), we get

$$r(x) = p(x) - s(x), \quad (24)$$

$$s'(x) + p(x)s(x) - s^2(x) = q(x), \quad (25)$$

where (25) is the nonlinear Riccati equation.

Consider the linear second order initial value problem

$$Lu(x) = f(x), \quad u(x_0) = \beta_1, \quad u'(x_0) = \beta_2, \quad (26)$$

where  $f(x) \in C[a, b]$  is a forcing function,  $x_0$  is a point in  $[a, b]$ ,  $\beta_i \in \mathbb{R}$ ,  $i = 1, 2$ , and  $u(x) \in C^2[a, b]$  is the unknown function describing the response of the system modeled by (26). If (21) and (22) hold true, then this problem can be factorized and solved in closed form as it is shown in the next theorem.

**Theorem 4** Let  $L$  be the linear differential operator of second order in (20) and  $\widehat{L}$  be its restriction on

$$D(\widehat{L}) = \{u(x) \in D(L) : u(x_0) = \beta_1, u'(x_0) = \beta_2\}, \tag{27}$$

where  $x_0 \in [a, b]$  and  $\beta_1, \beta_2 \in \mathbb{R}$ . If the prerequisites (21) and (22) are fulfilled then:

(i) The operator  $\widehat{L}$  can be factorized as

$$\widehat{L}u(x) = \widehat{L}_1\widehat{L}_2u(x), \tag{28}$$

where  $\widehat{L}_1, \widehat{L}_2$  are correct restrictions of the linear first order differential operators  $L_1, L_2$ , defined in (17) and (18), on

$$D(\widehat{L}_1) = \{z(x) \in D(L_1) : z(x_0) = \beta_2 + s(x_0)\beta_1\}, \tag{29}$$

$$D(\widehat{L}_2) = \{u(x) \in D(L_2) : u(x_0) = \beta_1\}, \tag{30}$$

respectively.

(ii) The operator  $\widehat{L}$  is correct and the unique solution of the initial value problem

$$\widehat{L}u(x) = f(x), \quad \forall f(x) \in C[a, b], \tag{31}$$

is given in closed form by

$$\begin{aligned} u(x) &= \widehat{L}^{-1}u(x) = \widehat{L}_2^{-1}\widehat{L}_1^{-1}f(x) = \widehat{L}_2^{-1}z(x) \\ &= e^{-\int_{x_0}^x s(t)dt} \left( \beta_1 + \int_{x_0}^x z(t)e^{\int_{t_0}^t s(\tau)d\tau} dt \right), \end{aligned} \tag{32}$$

where

$$z(x) = \widehat{L}_1^{-1}f(x) = e^{-\int_{x_0}^x r(t)dt} \left( \beta_2 + s(x_0)\beta_1 + \int_{x_0}^x f(t)e^{\int_{t_0}^t r(\tau)d\tau} dt \right). \tag{33}$$

**Proof**

(i) From the definition of  $\widehat{L}$  and Proposition 1, we have

$$Lu(x) = L_1L_2u(x) = f(x), \quad u(x_0) = \beta_1, \quad u'(x_0) = \beta_2. \tag{34}$$

By setting

$$L_2u(x) = u'(x) + s(x)u(x) = z(x), \tag{35}$$

we get

$$L_1 z(x) = z'(x) + r(x)z(x) = f(x). \quad (36)$$

From (35) it is implied that

$$u'(x_0) = z(x_0) - s(x_0)u(x_0),$$

which when is substituted into the second condition in (34) yields

$$z(x_0) = \beta_2 + s(x_0)\beta_1.$$

Whence, we have the two linear first order initial value problems

$$L_1 z(x) = f(x), \quad z(x_0) = \beta_2 + s(x_0)\beta_1, \quad (37)$$

$$L_2 u(x) = z(x), \quad u(x_0) = \beta_1. \quad (38)$$

That is  $\widehat{L}u(x) = \widehat{L}_1 \widehat{L}_2 u(x)$ . It remains to show that  $D(\widehat{L}) = D(\widehat{L}_1 \widehat{L}_2)$ . By using (29) and (30), we obtain

$$\begin{aligned} D(\widehat{L}_1 \widehat{L}_2) &= \{u(x) \in D(\widehat{L}_2) : \widehat{L}_2 u(x) \in D(\widehat{L}_1)\} \\ &= \{u(x) \in D(L_2) : u(x_0) = \beta_1, \quad u'(x) + s(x)u(x) \in D(\widehat{L}_1)\} \\ &= \left\{u(x) \in C^2[a, b] : u(x_0) = \beta_1, \quad u'(x_0) + s(x_0)u(x_0) = \beta_2 + s(x_0)\beta_1\right\} \\ &= \left\{u(x) \in C^2[a, b] : u(x_0) = \beta_1, \quad u'(x_0) = \beta_2\right\}. \end{aligned} \quad (39)$$

- (ii) The linear first order initial value problem (37) possesses exactly one solution  $z(x)$ , which can be found by using the standard means, such as the method of integrating factors [2], and is given in (33). Having obtained  $z(x)$ , we can solve the linear first order initial value problem (38) in like manner to obtain the solution  $u(x)$  in (32), which is the solution of the linear second order initial value problem (31). The operator  $\widehat{L} = \widehat{L}_1 \widehat{L}_2$  is correct because  $\widehat{L}_1$  and  $\widehat{L}_2$  are correct.

□

The factorization method also applies to some types of boundary value problems, although it is more complicated. Let the linear second order differential equation,

$$Lu(x) = u''(x) + p(x)u(x) + q(x)u(x) = f(x), \quad x \in (a, b), \quad (40)$$

where the coefficients  $p(x)$ ,  $q(x) \in C[a, b]$  and  $f(x) \in C[a, b]$ , and assume that the operator  $L : C[a, b] \rightarrow C[a, b]$  is factorable, i.e. there exist  $r(x) \in C[a, b]$  and



$s(x) \in C^1[a, b]$  such that  $r(x) + s(x) = p(x)$  and  $s'(x) + r(x)s(x) = q(x)$ . Let also the two boundary conditions

$$\begin{aligned} \mu_{11}u(a) + \mu_{12}u(b) &= \beta_1, \\ \mu_{21}[u'(a) + s(a)u(a)] + \mu_{22}[u'(b) + s(b)u(b)] &= \beta_2, \end{aligned} \tag{41}$$

where  $\mu_{ij}, \beta_i \in \mathbb{R}, i = 1, 2, j = 1, 2$ . Notice that (41) are the boundary conditions as in (2) when

$$\mu_{23} = s(a)\mu_{21}, \quad \mu_{24} = s(b)\mu_{22}. \tag{42}$$

We claim that the boundary value problem for the differential equation (40) and the boundary conditions (41) can be factorized and solved explicitly. We prove the following theorem.

**Theorem 5** *Let  $L$  be the linear second order differential operator in (40) and assume that there exist two functions  $r(x) \in C[a, b]$  and  $s(x) \in C^1[a, b]$  which satisfy (21) and (22). Let  $\bar{L}$  be a restriction of  $L$  on*

$$\begin{aligned} D(\bar{L}) = \{u(x) : u(x) \in D(L), \mu_{11}u(a) + \mu_{12}u(b) &= \beta_1, \\ \mu_{21}[u'(a) + s(a)u(a)] + \mu_{22}[u'(b) + s(b)u(b)] &= \beta_2\}, \end{aligned} \tag{43}$$

where  $\mu_{ij}, \beta_i \in \mathbb{R}, i = 1, 2, j = 1, 2$ . Then:

(i) *The operator  $\bar{L}$  can be factorized as follows*

$$\bar{L}u(x) = \bar{L}_1\bar{L}_2u(x), \tag{44}$$

where  $\bar{L}_1, \bar{L}_2$  are restrictions of the two first order linear differential operators  $L_1, L_2$ , defined in (17) and (18), on

$$D(\bar{L}_1) = \left\{z(x) \in C^1[a, b] : \mu_{21}z(a) + \mu_{22}z(b) = \beta_2\right\}, \tag{45}$$

$$D(\bar{L}_2) = \left\{u(x) \in C^1[a, b] : \mu_{11}u(a) + \mu_{12}u(b) = \beta_1\right\}, \tag{46}$$

respectively.

(ii) *If*

$$\mu_{21} + \mu_{22}e^{-\int_a^b r(t)dt} \neq 0, \quad \mu_{11} + \mu_{12}e^{-\int_a^b s(t)dt} \neq 0, \tag{47}$$

then the operator  $\bar{L}$  is correct and the unique solution of the boundary value problem

$$\bar{L}u(x) = f(x), \quad \forall f(x) \in C[a, b], \tag{48}$$

is given by

$$\begin{aligned}
 u(x) &= \bar{L}^{-1} f(x) = \bar{L}_2^{-1} \bar{L}_1^{-1} f(x) = \bar{L}_2^{-1} z(x) \\
 &= e^{-\int_a^x s(t)dt} \left( C_2 + \int_a^x z(t) e^{\int_a^t s(\tau)d\tau} dt \right), \tag{49}
 \end{aligned}$$

where

$$z(x) = \bar{L}_1^{-1} f(x) = e^{-\int_a^x r(t)dt} \left( C_1 + \int_a^x f(t) e^{\int_a^t r(\tau)d\tau} dt \right), \tag{50}$$

and

$$\begin{aligned}
 C_1 &= \left( \mu_{21} + \mu_{22} e^{-\int_a^b r(t)dt} \right)^{-1} \left( \beta_2 - \mu_{22} e^{-\int_a^b r(t)dt} \int_a^b f(t) e^{\int_a^t r(\tau)d\tau} dt \right), \\
 C_2 &= \left( \mu_{11} + \mu_{12} e^{-\int_a^b s(t)dt} \right)^{-1} \left( \beta_1 - \mu_{12} e^{-\int_a^b s(t)dt} \int_a^b z(t) e^{\int_a^t s(\tau)d\tau} dt \right).
 \end{aligned}$$

**Proof**

(i) From the definition of  $\bar{L}$  and Proposition 1, we have

$$Lu(x) = L_1 L_2 u(x) = f(x), \tag{51}$$

and

$$\begin{aligned}
 \mu_{11}u(a) + \mu_{12}u(b) &= \beta_1, \\
 \mu_{21}[u'(a) + s(a)u(a)] + \mu_{22}[u'(b) + s(b)u(b)] &= \beta_2. \tag{52}
 \end{aligned}$$

Let

$$L_2 u(x) = u'(x) + s(x)u(x) = z(x). \tag{53}$$

It follows that

$$u'(a) + s(a)u(a) = z(a), \quad u'(b) + s(b)z(b) = z(b),$$

and upon substitution into the second boundary condition in (52), we get

$$\mu_{21}z(a) + \mu_{22}z(b) = \beta_2.$$

Thus, we have

$$L_1 z(x) = z'(x) + r(x)z(x) = f(x), \quad \mu_{21}z(a) + \mu_{22}z(b) = \beta_2, \tag{54}$$

$$L_2u(x) = u'(x) + s(x)u(x) = z(x), \quad \mu_{11}u(a) + \mu_{12}u(b) = \beta_1. \quad (55)$$

That is  $\bar{L}u(x) = \bar{L}_1\bar{L}_2u(x)$ . It remains to show that  $D(\bar{L}) = D(\bar{L}_1\bar{L}_2)$ . By using the definition of  $D(\bar{L}_1\bar{L}_2)$  we obtain

$$\begin{aligned} D(\bar{L}_1\bar{L}_2) &= \{u(x) \in D(\bar{L}_2) : \bar{L}_2u(x) \in D(\bar{L}_1)\} \\ &= \left\{ u(x) \in C^1[a, b] : \mu_{11}u(a) + \mu_{12}u(b) = \beta_1, \right. \\ &\quad \left. u'(x) + s(x)u(x) \in D(\bar{L}_1) \right\} \\ &= \left\{ u(x) \in C^1[a, b] : \mu_{11}u(a) + \mu_{12}u(b) = \beta_1, \right. \\ &\quad z(x) = u'(x) + s(x)u(x) \in C^1[a, b], \\ &\quad \left. \mu_{21}[u'(a) + s(a)u(a)] + \mu_{22}[u'(b) + s(b)u(b)] = \beta_2 \right\} \\ &= \left\{ u(x) \in C^2[a, b] : \mu_{11}u(a) + \mu_{12}u(b) = \beta_1, \right. \\ &\quad \left. \mu_{21}[u'(a) + s(a)u(a)] + \mu_{22}[u'(b) + s(b)u(b)] = \beta_2 \right\} \\ &= D(\bar{L}). \end{aligned} \quad (56)$$

(ii) Application of Theorem 3 to solve boundary value problem (54) yields (50). Substituting this unique solution  $z(x) = \bar{L}_1^{-1}f(x)$  into (55) and applying Theorem 3 once more, we obtain (49), which is the solution to boundary value problem (48). The correctness of  $\bar{L} = \bar{L}_1\bar{L}_2$  follows from the correctness of  $\bar{L}_1$  and  $\bar{L}_2$ . □

To elucidate the implementation of the above procedure, we solve the following example problem.

*Example 1* Let the boundary value problem

$$\begin{aligned} u''(x) - \frac{x+2}{x+1}u'(x) + \frac{1}{x+1}u(x) &= 3(x+1), \quad 0 < x < 1, \\ u(0) - 5u(1) &= 0, \\ 3u'(0) - 4u'(1) - 3u(0) + 4u(1) &= 2. \end{aligned} \quad (57)$$

We take

$$p(x) = -\frac{x+2}{x+1}, \quad \text{and} \quad q(x) = \frac{1}{x+1},$$

which are continuous on  $[0, 1]$ . Notice that equations (21) and (22) are satisfied by

$$r(x) = -\frac{1}{x+1}, \quad \text{and} \quad s(x) = -1,$$

which are continuous on  $[0, 1]$  and  $s'(x) = 0$ . Lastly, the second of the boundary conditions (57) can be put in the form

$$3[u'(0) + (-1)u(0)] - 4[u'(1) + (-1)u(1)] = 2.$$

Thus (57) is carried to

$$\begin{aligned} \bar{L}u(x) &= u''(x) - \frac{x+2}{x+1}u'(x) + \frac{1}{x+1}u(x) = f(x), \\ D(\bar{L}) &= \left\{ u(x) : u(x) \in C^2[0, 1], u(0) - 5u(1) = 0, \right. \\ &\quad \left. 3[u'(0) + s(0)u(0)] - 4[u'(1) + s(1)u(1)] = 2 \right\}, \end{aligned} \tag{58}$$

where  $f(x) = 3(x+1)$ . By Theorem 5, the operator  $\bar{L}$  can be factorized as  $\bar{L}u(x) = \bar{L}_1\bar{L}_2u(x)$ , where

$$\begin{aligned} \bar{L}_1z(x) &= z'(x) - \frac{1}{x+1}z(x), \quad D(\bar{L}_1) = \left\{ z(x) \in C^1[0, 1] : 3z(0) - 4z(1) = 2 \right\}, \\ \bar{L}_2u(x) &= u'(x) - u(x), \quad D(\bar{L}_2) = \left\{ u(x) \in C^1[0, 1] : u(0) - 5u(1) = 0 \right\}. \end{aligned}$$

Furthermore,

$$\mu_{21} + \mu_{22}e^{-\int_0^1 r(t)dt} = -5 \neq 0, \quad \mu_{11} + \mu_{12}e^{-\int_0^1 s(t)dt} = 1 - 5e \neq 0, \tag{59}$$

and therefore (58) has only one solution. To construct the solution, we first solve the problem  $\bar{L}_1u(x) = f(x)$  by means of (50), which yields

$$z(x) = (x+1)\left(3x - \frac{26}{5}\right). \tag{60}$$

Then by utilizing (60) and solving  $\bar{L}_2u(x) = z(x)$  by (49), we get

$$u(x) = \frac{142e^x}{5(5e-1)} - \frac{15x^2 + 19x - 7}{5}. \tag{61}$$

This is the unique solution of the given boundary value problem (57).

## 4 Factorization Method for Nonlinear Differential Equations

In this section, we deal with the solution of a class of nonlinear boundary value problems for second order differential equations. Let the nonlinear differential equation of the form

$$u''(x)u'(x) + [q(x)u'(x) + g(x)u''(x)]u(x) + q(x)g(x)(u(x))^2 = 0, \quad (62)$$

for  $x \in (a, b)$ , and where  $q(x), g(x) \in C[a, b]$ , together with the boundary conditions

$$\begin{aligned} \mu_{11}u(a) + \mu_{12}u(b) &= \beta_1, \\ \mu_{21}u'(a) + \mu_{22}u'(b) + \mu_{23}u(a) + \mu_{24}u(b) &= \beta_2, \end{aligned} \quad (63)$$

where  $\mu_{ij}, \beta_i \in \mathbb{R}, i = 1, 2, j = 1, 2, 3, 4$ .

The nonlinear equation (62) can be decomposed as the product

$$[u''(x) + q(x)u(x)][u'(x) + g(x)u(x)] = 0,$$

and hence, either

$$u''(x) + q(x)u(x) = 0, \quad (64)$$

or

$$u'(x) + g(x)u(x) = 0. \quad (65)$$

As a consequence, the solutions of the nonlinear boundary value problem (62) and (63) may be obtained by solving the linear second order problem (64) and (63) and the linear first order problem (65) and (63).

For the solution of the linear second order problem (64) and (63), we may employ Theorem 5 provided that prerequisites (21) and (22) are satisfied, i.e. there exist  $r(x) \in C[a, b]$  and  $s(x) \in C^1[a, b]$  such that

$$r(x) = -s(x), \quad s'(x) - (s(x))^2 = q(x), \quad (66)$$

and if

$$\mu_{23} = s(a)\mu_{21}, \quad \mu_{24} = s(b)\mu_{22}. \quad (67)$$

In this case problem (64), (63) can be put in the form

$$\bar{L}u(x) = u''(x) + q(x)u(x) = 0,$$

$$D(\bar{L}) = \left\{ u(x) \in C^2[a, b] : \mu_{11}u(a) + \mu_{12}u(b) = \beta_1, \right. \\ \left. \mu_{21}[u'(a) + s(a)u(a)] + \mu_{22}[u'(b) + s(b)u(b)] = \beta_2 \right\}. \quad (68)$$

Problem (68) can be now solved by means of Theorem 5.

The linear first order problem (65) and (63) is subjected to more conditions than the order of the differential equation and it is most likely to possess no solution. Nevertheless, we can proceed as follows. By utilizing (65) evaluate  $u'(a)$  and  $u'(b)$  and substitute into the second of the boundary conditions in (63). Taking into account (67), we get

$$\mu_{21}[s(a) - g(a)]u(a) + \mu_{22}[s(b) - g(b)]u(b) = \beta_2. \quad (69)$$

Thus, problem (65) and (63) may be formulated as

$$Tu(x) = u'(x) + g(x)u(x) = 0, \\ D(T) = \{u(x) \in C^1[0, 1] : \mu_{11}u(a) + \mu_{12}u(b) = \beta_1, \\ \mu_{21}[s(a) - g(a)]u(a) + \mu_{22}[s(b) - g(b)]u(b) = \beta_2\}. \quad (70)$$

By employing Theorem 3, we find the unique solution of the problem

$$T_0u(x) = \overline{u'(x) + g(x)u(x)} = 0, \\ D(T_0) = \{u(x) \in C^1[0, 1] : \mu_{11}u(a) + \mu_{12}u(b) = \beta_1\}. \quad (71)$$

If the solution  $u(x)$  of this problem satisfies the second boundary condition in (70), then  $u(x)$  is a solution of (70); otherwise (70) has no solution.

*Example 2* Let us find the solutions of the nonlinear second order boundary value problem

$$u''(x)u'(x) - \left[ \frac{2}{(x+1)^2}u'(x) + \frac{1}{x+3}u''(x) \right] u(x) + \frac{2}{(x+3)(x+1)^2} (u(x))^2 = 0, \\ u(0) + 5u(1) = 0, \\ -u'(0) + 6u'(1) - u(0) + 3u(1) = 4, \quad (72)$$

where  $x \in [0, 1]$  and  $u(x) \in C^2[0, 1]$ .

The nonlinear second order differential equation (72) is of the type (62) with

$$q(x) = -\frac{2}{(x+1)^2}, \quad g(x) = -\frac{1}{x+3},$$

and it can be decomposed as

$$\left[ u''(x) - \frac{2}{(x+1)^2} u(x) \right] \left[ u'(x) - \frac{1}{x+3} u(x) \right] = 0.$$

Thus, we get the following two linear problems

$$\begin{aligned} \bar{L}u(x) &= u''(x) - \frac{2}{(x+1)^2} u(x) = 0, \\ D(\bar{L}) &= \{u(x) \in C^2[0, 1] : u(0) + 5u(1) = 0, \\ &\quad -u'(0) + 6u'(1) - u(0) + 3u(1) = 4\}, \end{aligned} \quad (73)$$

and

$$\begin{aligned} Tu(x) &= u'(x) - \frac{1}{x+3} u(x) = 0, \\ D(T) &= \{u(x) \in C^2[0, 1] : u(0) + 5u(1) = 0, \\ &\quad -u'(0) + 6u'(1) - u(0) + 3u(1) = 4\}. \end{aligned} \quad (74)$$

In solving the boundary value problem (73), notice that the functions

$$r(x) = -\frac{1}{x+1}, \quad s(x) = \frac{1}{x+1},$$

obey (66),  $r(x) \in C[0, 1]$ ,  $s(x) \in C^1[0, 1]$  and  $s(0) = 1$ ,  $s(1) = \frac{1}{2}$ , and that the preconditions (67) are met. Hence, problem (73) may be written in the form (68), namely

$$\begin{aligned} \bar{L}u(x) &= u''(x) - \frac{2}{(x+1)^2} u(x) = 0, \\ D(\bar{L}) &= \{u(x) \in C^2[0, 1] : u(0) + 5u(1) = 0, \\ &\quad -[u'(0) + s(0)u(0)] + 6[u'(1) + s(1)u(1)] = 4\}. \end{aligned} \quad (75)$$

By Theorem 5, the boundary value problem (75) is factorized into the following two first order problems

$$\begin{aligned} \bar{L}_1 z(x) &= z'(x) - \frac{1}{x+1} z(x) = 0, \\ D(\bar{L}_1) &= \left\{ z(x) \in C^1[0, 1] : -z(0) + 6z(1) = 4 \right\}, \end{aligned} \quad (76)$$

and

$$\begin{aligned} \bar{L}_2 u(x) &= u'(x) + \frac{1}{x+1} u(x) = z(x), \\ D(\bar{L}_2) &= \left\{ u(x) \in C^1[0, 1] : u(0) + 5u(1) = 0 \right\}. \end{aligned} \tag{77}$$

The first of the uniqueness requirements (47) is fulfilled, viz.

$$\mu_{21} + \mu_{22} e^{-\int_0^1 r(t) dt} = -1 + 6 \left( e^{\int_0^1 \frac{1}{t+1} dt} \right) = 11 \neq 0, \tag{78}$$

and therefore the operator  $\bar{L}_1$  is correct and the unique solution of (76) is derived through (50), which is

$$z(x) = \frac{4}{11}(x + 1). \tag{79}$$

By substituting (79) into (77) and verifying that the second of the uniqueness conditions (47) is also satisfied, viz.

$$\mu_{11} + \mu_{12} e^{-\int_0^1 s(t) dt} = 1 + 5 \left( e^{-\int_0^1 \frac{1}{t+1} dt} \right) = \frac{7}{2} \neq 0, \tag{80}$$

it follows that the operator  $\bar{L}_2$  is correct and the unique solution of (77), obtained via (49), is

$$u(x) = \frac{4(x^3 + 3x^2 + 3x - 5)}{33(x + 1)}. \tag{81}$$

The function  $u(x)$  in (81) is a solution to nonlinear second order boundary value problem (72).

We now examine the existence of a solution of the linear first order problem (74). By applying Theorem 3, we find that the problem

$$\begin{aligned} T_1 u(x) &= u'(x) - \frac{1}{x+3} u(x) = 0, \\ D(T_1) &= \{u(x) \in C^2[0, 1] : u(0) + 5u(1) = 0\} \end{aligned} \tag{82}$$

has no solution except the trivial  $u(x) = 0$ , which however does not satisfy the second of the boundary conditions in (74).

Summing up, the nonlinear second order boundary value problem (72) admits only the solution (81).

The technique presented above and explained in Example 2 can be extended to solve and other types of nonlinear boundary value problems. For example, consider the nonlinear differential equation of second order of the form



$$(u''(x))^2 + a(x)u''(x)u(x) + b(x)(u(x))^2 = 0, \quad x \in (a, b), \tag{83}$$

subject to two general boundary constraints

$$\begin{aligned} \mu_{11}u(a) + \mu_{12}u(b) &= \beta_1, \\ \mu_{21}u'(a) + \mu_{22}u'(b) + \mu_{23}u(a) + \mu_{24}u(b) &= \beta_2, \end{aligned} \tag{84}$$

where  $a(x), b(x) \in C[a, b]$  and  $\mu_{ij}, \beta_i \in \mathbb{R}, i = 1, 2, j = 1, 2, 3, 4$ .

The differential equation (83) can be put in the form

$$F\left(\frac{u''(x)}{u(x)}, x\right) = F(w(x), x) = (w(x))^2 + a(x)w(x) + b(x) = 0,$$

where  $w(x) = u''(x)/u(x)$  and the nonlinear function  $F$  is a second degree polynomial of  $w(x)$ . Hence, it can be decomposed as

$$[w(x) + q^-(x)][w(x) + q^+(x)] = 0,$$

where  $q^-(x), q^+(x) \in C[a, b]$  and  $a(x) = q^-(x) + q^+(x), b(x) = q^-(x)q^+(x)$ . By substituting back  $w(x) = u''(x)/u(x)$ , we get

$$[u''(x) + q^-(x)u(x)][u''(x) + q^+(x)u(x)] = 0, \tag{85}$$

from where follows that, either

$$u''(x) + q^-(x)u(x) = 0, \quad x \in (a, b), \tag{86}$$

or

$$u''(x) + q^+(x)u(x) = 0, \quad x \in (a, b). \tag{87}$$

Thus, the solution of the nonlinear boundary value problem (83) and (84) is reduced to the solution of the two linear second order boundary value problems (86), (84), and (87), (84). Whenever the conditions (21), (22) and (42) are met, Theorem 5 may be applied to acquire the solutions in closed form.

## 5 Conclusions

A practical technique has been presented for factorizing and solving linear initial and boundary value problems for second order differential equations with nonlocal boundary conditions. Two types of nonlinear boundary value problems for second order differential equations have also been considered where the factorization

method was used to construct their solutions in closed form. The main advantage of the factorization method is that no fundamental or particular solutions are required. Its main disadvantage is that it cannot be applied to all boundary value problems except to those where certain conditions are satisfied. The efficiency of the method encourages the pursuit of further research for the extension of the method to problems with fully mixed boundary conditions and multipoint conditions.

## References

1. L.M. Berkovich, Factorization and transformations of linear and nonlinear ordinary differential equations. *Nucl. Instrum. Methods Phys. Res. A* **502**, 646–648 (2003). [https://doi.org/10.1016/S0168-9002\(03\)00531-X](https://doi.org/10.1016/S0168-9002(03)00531-X)
2. W.E. Boyce, R.C. DiPrima, *Elementary Differential Equations and Boundary Value Problems* (Wiley, London, 2012)
3. D.I. Caruntu, Factorization of self-adjoint ordinary differential equations. *Appl. Math. Comput.* **219**, 7622–7631 (2013). <https://doi.org/10.1016/j.amc.2013.01.049>
4. J. Clegg, A new factorisation of a general second order differential equation. *Int. J. Math. Edu. Sci. Technol.* **37**, 51–64 (2006). <https://doi.org/10.1080/00207390500186339>
5. E. García, L. Littlejohn, J.L. López, E.P. Sinusía, Factorization of second-order linear differential equations and Liouville–Neumann expansions. *Math. Comput. Modell.* **57**, 1514–1530 (2013). <https://doi.org/10.1016/j.mcm.2012.12.012>
6. M.N. Hounkonnou, A. Ronveaux, Factorization of generalized Lamé and Heun’s differential equations. *Commun. Math. Anal.* **11**, 121–136 (2011). <https://projecteuclid.org/euclid.cma/1293054278>
7. L. Infeld, T.E. Hull, The factorization method. *Rev. Mod. Phys.* **23**, 21–68 (1951)
8. K. Janglajew, E. Schmeidel, The factorization of the linear differential operator. *Tatra Mt. Math. Publ.* **63**, 139–151 (2015). <https://doi.org/10.1515/tmmp-2015-0026>
9. J.J. Kovacic, An algorithm for solving second order linear homogeneous differential equations. *J. Symb. Comput.* **2**, 3–43 (1986). [https://doi.org/10.1016/S0747-7171\(86\)80010-4](https://doi.org/10.1016/S0747-7171(86)80010-4)
10. I.N. Parasidis, P. Hahamis, Factorization method for solving multipoint problems for second order difference equations with polynomial coefficients, in *Discrete Mathematics and Applications*. Springer Optimization and Its Applications, vol. 165, eds. by A.M. Raigorodskii, M.T. Rassias (Springer, Cham, 2020). [http://doi-org-443.webvpn.fjmu.edu.cn/10.1007/978-3-030-55857-4\\_17](http://doi-org-443.webvpn.fjmu.edu.cn/10.1007/978-3-030-55857-4_17)
11. I.N. Parasidis, E. Providas, Factorization method for the second order linear nonlocal difference equations, in *International Conference Polynomial Computer Algebra ’2018*, ed. by N.N. Vassiliev (Euler International Mathematical Institute, St. Petersburg, 2018), pp. 85–89
12. I.N. Parasidis, E. Providas, P.C. Tsekrekos, Factorization of linear operators and some eigenvalue problems of special operators. *Vestn. Bashkir. Univ.* **17**, 830–839 (2012)
13. E. Providas, Operator factorization and solution of second-order nonlinear difference equations with variable coefficients and multipoint constraints, in *Nonlinear Analysis and Global Optimization*. Springer Optimization and Its Applications, vol. 167, eds. by T.M. Rassias, P.M. Pardalos, pp. 427–443 (Springer, Cham, 2021). [https://doi.org/10.1007/978-3-030-61732-5\\_20](https://doi.org/10.1007/978-3-030-61732-5_20)
14. W. Robin, Operator factorization and the solution of second-order linear ordinary differential equations. *Int. J. Math. Educ. Sci. Technol.* **38**, 189–211 (2007). <https://doi.org/10.1080/00207390601002815>
15. F. Schwarz, Decomposition of ordinary differential equations. *Bull. Math. Sci.* **7**, 575–613 (2017). <https://doi.org/10.1007/s13373-017-0110-0>

16. S. Tsarev, Symbolic manipulation of integro differential expressions and factorization of linear ordinary differential operators over transcendental extensions of a differential field, in *ISSAC '97: Proceedings of the 1997 International Symposium on Symbolic and Algebraic Computation*, vol. 1997 (1997), pp. 310–315. <https://doi.org/10.1145/258726.258824>
17. N.N. Vassiliev, I.N. Parasidis, E. Providas, Exact solution method for Fredholm integro-differential equations with multipoint and integral boundary conditions. Part 2. Decomposition-extension method for squared operators. *Inf. Control Syst.* **2**, 2–9 (2019). <https://doi.org/10.31799/1684-8853-2019-2-2-9>

# A General Framework for Studying Certain Generalized Topologically Open Sets in Relator Spaces



Themistocles M. Rassias and Árpád Száz

**Abstract** A family  $\mathcal{R}$  of binary relations on a set  $X$  is called a relator on  $X$ , and the ordered pair  $X(\mathcal{R}) = (X, \mathcal{R})$  is called a relator space. Sometimes relators on  $X$  to  $Y$  are also considered.

By using an obvious definition of the generated open sets, each generalized topology  $\mathcal{T}$  on  $X$  can be easily derived from the family  $\mathcal{R}_{\mathcal{T}}$  of all Pervin's preorder relations  $R_V = V^2 \cup V^c \times X$  with  $V \in \mathcal{T}$ , where  $V^2 = V \times V$  and  $V^c = X \setminus V$ .

For a subset  $A$  of the relator space  $X(\mathcal{R})$ , we define

$$A^\circ = \text{int}_{\mathcal{R}}(A) = \{x \in X : \exists R \in \mathcal{R} : R(x) \subseteq A\}$$

and  $A^- = \text{cl}_{\mathcal{R}}(A) = \text{int}_{\mathcal{R}}(A^c)^c$ . And, for instance, we write  $A \in \mathcal{T}_{\mathcal{R}}$  if  $A \subseteq A^\circ$ .

Moreover, following some basic definitions in topological spaces, for a subset  $A$  of the relator space  $X(\mathcal{R})$  we write

- (1)  $A \in \mathcal{T}_{\mathcal{R}}^r$  if  $A = A^{-\circ}$ ;
- (2)  $A \in \mathcal{T}_{\mathcal{R}}^p$  if  $A \subseteq A^{-\circ}$ ;
- (3)  $A \in \mathcal{T}_{\mathcal{R}}^s$  if  $A \subseteq A^{\circ-}$ ;
- (4)  $A \in \mathcal{T}_{\mathcal{R}}^\alpha$  if  $A \subseteq A^{\circ--}$ ;
- (5)  $A \in \mathcal{T}_{\mathcal{R}}^\beta$  if  $A \subseteq A^{\circ--}$ ;
- (6)  $A \in \mathcal{T}_{\mathcal{R}}^a$  if  $A \subseteq A^{-\circ} \cap A^{\circ-}$ ;
- (7)  $A \in \mathcal{T}_{\mathcal{R}}^b$  if  $A \subseteq A^{-\circ} \cup A^{\circ-}$ .

The members of the above families will be called the topologically regular open, preopen, semi-open,  $\alpha$ -open,  $\beta$ -open,  $a$ -open and  $b$ -open subsets of the relator space  $X(\mathcal{R})$ , respectively.

---

Th. M. Rassias

Department of Mathematics Zografou Campus, National Technical University of Athens, Athens, Greece

e-mail: [trassias@math.ntua.gr](mailto:trassias@math.ntua.gr)

Á. Száz (✉)

Department of Mathematics, University of Debrecen, Debrecen, Hungary

e-mail: [szaz@science.unideb.hu](mailto:szaz@science.unideb.hu)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,

Springer Optimization and Its Applications 173,

[https://doi.org/10.1007/978-3-030-72563-1\\_19](https://doi.org/10.1007/978-3-030-72563-1_19)

In our former papers, having in mind the original definitions of N. Levine [49] and H. H. Corson and E. Michael [11], we have also investigated four further, closely related, families of generalized topologically open sets in  $X(\mathcal{R})$ .

Now, we shall offer a general framework for studying these families. Moreover, motivated by a definition of Á. Császár [15] and his predecessors, we shall also consider a further important class of generalized topologically open sets.

For the latter purpose, for a subset  $A$  of the relator space  $X(\mathcal{R})$ , we shall write

$$(8) \quad A \in \mathcal{A}_{\mathcal{R}} \quad \text{if} \quad A^{-\circ} \subseteq A^{\circ-}.$$

Thus, according to Császár’s terminology, the members of the family  $\mathcal{A}_{\mathcal{R}}$  should be called the topologically quasi-open subsets of the relator space  $X(\mathcal{R})$ . However, in the earlier literature, these sets have been studied under different names.

While, for the former purpose, for any two subsets  $A$  and  $B$  of the relator space  $X(\mathcal{R})$  we shall write

$$(9) \quad A \in \text{Ln}_{\mathcal{R}}(B) \text{ and } B \in \text{Un}_{\mathcal{R}}(A) \quad \text{if} \quad A \subseteq B \subseteq A^{-}.$$

Moreover, for a family  $\mathcal{A}$  of subsets of  $X(\mathcal{R})$  we shall define

$$(10) \quad \mathcal{A}^{\ell} = \text{cl}_{\text{Ln}_{\mathcal{R}}}(\mathcal{A}) = \text{Ln}_{\mathcal{R}}^{-1}[\mathcal{A}] \text{ and } \mathcal{A}^u = \text{cl}_{\text{Un}_{\mathcal{R}}}(\mathcal{A}) = \text{Un}_{\mathcal{R}}^{-1}[\mathcal{A}].$$

Thus,  $\mathcal{A}^{\ell}$  and  $\mathcal{A}^u$  may be called the lower and upper nearness closures of  $\mathcal{A}$ , respectively. Namely, if  $A \in \text{Ln}_{\mathcal{R}}(B)$ , then we may naturally say that  $A$  is near to  $B$  from below and  $B$  is near to  $A$  from above.

The most important particular cases are when  $\mathcal{A}$  is a minimal structure or a generalized topology on  $X$ . Or even more specially,  $\mathcal{A}$  is one of the families  $\mathcal{T}_{\mathcal{R}}$ ,  $\mathcal{T}_{\mathcal{R}}^{\ell}$  or  $\mathcal{T}_{\mathcal{R}}^u$ .

### 1 Motivations

If  $\mathcal{T}$  is a family of subsets of a set  $X$  such that  $\mathcal{T}$  is closed under finite intersections and arbitrary unions, then the family  $\mathcal{T}$  is called a *topology* on  $X$ , and the ordered pair  $X(\mathcal{T}) = (X, \mathcal{T})$  is called a *topological space*.

The members of  $\mathcal{T}$  are called the *open subsets* of  $X$ . While, the members of  $\mathcal{F} = \{A^c : A \in \mathcal{T}\}$ , where  $A^c = X \setminus A$ , are called the *closed subsets* of  $X$ . Moreover, the members of  $\mathcal{T} \cap \mathcal{F}$  are called the *clopen subsets* of  $X$ .

Since,  $\emptyset = \bigcup \emptyset$  and  $X = \bigcap \emptyset$ , we necessary have  $\{\emptyset, X\} \subseteq \mathcal{T} \cap \mathcal{F}$ . Therefore, if in particular  $\mathcal{T} = \{\emptyset, X\}$ , then  $\mathcal{T}$  is called *minimal* [69] instead of indiscrete. While, if  $\mathcal{T} \cap \mathcal{F} = \{\emptyset, X\}$ , then  $\mathcal{T}$  is called *connected* [97, p. 31].

For a subset  $A$  of  $X(\mathcal{T})$ , the sets  $A^{\circ} = \text{int}(A) = \bigcup \mathcal{T} \cap \mathcal{F}(A)$ ,

$$A^{-} = \text{cl}(A) = \text{int}(A^c)^c \quad \text{and} \quad A^{\dagger} = \text{res}(A) = \text{cl}(A) \setminus A$$

are called the *interior*, *closure* and *residue* of  $A$ , respectively.

Thus,  $-$  is a *Kuratowski closure operation* on  $X$ . That is,  $\emptyset^- = \emptyset$ , and  $-$  is *extensive, idempotent and additive* in the sense that, for any  $A, B \subseteq X$ , we have  $A \subseteq A^-$ ,  $A^{- -} = A^-$  and  $(A \cup B)^- = A^- \cup B^-$ .

In particular, the members of the families

$$\mathcal{D} = \{ A \subseteq X : A^- = X \} \quad \text{and} \quad \mathcal{N} = \{ A \subseteq X : A^{-\circ} = \emptyset \}$$

are called the *dense and rare (or nowhere dense) subsets* of  $X(\mathcal{T})$ , respectively.

In 1922, a subset  $A$  of a closure space  $X(-)$  was called *regular open* by Kuratowski [44] if  $A = A^{-\circ}$ . While, in 1937, a subset  $A$  of a topological space  $X(\mathcal{T})$  was called *regular open* by Stone [72] if  $A = B^\circ$  for some  $B \in \mathcal{T}$ .

The importance of regular open subsets of  $X(\mathcal{T})$  lies mainly in the fact that their family forms a complete *Boolean algebra* [32, p. 66] with respect to the operations defined by  $A' = A^{-c}$ ,  $A \wedge B = A \cap B$  and  $A \vee B = (A \cup B)''$ .

In 1982, a subset  $A$  of  $X(\mathcal{T})$  was called *preopen* by Mashhour et al. [56] if  $A \subseteq A^{-\circ}$ . However, by Dontchev [22], preopen sets, under different names, were much earlier studied by several mathematicians.

For instance, in 1964, Corson and Michael [11] called a subset  $A$  of  $X(\mathcal{T})$  *locally dense* if it is a dense subset of some  $V \in \mathcal{T}$  in the sense that  $A \subseteq V \subseteq A^-$ . Moreover, they noted that this property is equivalent to the inclusion  $A \subseteq A^{-\circ}$ .

This equivalence was later also stated by Jun at al. [38]. Moreover, Ganster [28] proved that  $A$  is preopen if and only if there exist  $V \in \mathcal{T}$  and  $B \in \mathcal{D}$  such that  $A = V \cap B$ . (See also Dontchev [22].)

In 1963, a subset  $A$  of  $X(\mathcal{T})$  was called *semi-open* by Levine [49] if there exists  $V \in \mathcal{T}$  such that  $V \subseteq A \subseteq V^-$ . First of all, he showed that the set  $A$  is semi-open if and only if  $A \subseteq A^{\circ-}$ .

Moreover, he also proved that if  $A$  is a semi-open subset of  $X(\mathcal{T})$ , then there exist  $V \in \mathcal{T}$  and  $B \in \mathcal{N}$  such that  $A = V \cup B$  and  $V \cap B = \emptyset$ . In addition, he also noted that the converse statement is false.

Levine's statement closely resembles to a famous stability theorem of Hyers [36] which says that an  $\varepsilon$ -approximately additive function of one Banach space to another is the sum of an additive function and an  $\varepsilon$ -small function.

Analogously to the paper of Hyers, Levine's paper has also attracted the interest of a surprisingly great number of mathematicians. For instance, by the Google Scholar, it has been cited by 2985 works.

Moreover, the above statement of Levine was improved by Dłaska et al. [21] who observed that a subset  $A$  of  $X(\mathcal{T})$  is semi-open if and only if there exist  $V \in \mathcal{T}$  and  $B \subseteq V^\dagger$  such that  $A = V \cup B$ .

The latter observation was later reformulated, in a more convenient form, by Duszyński and Noiri [23] who noted that a subset  $A$  of  $X(\mathcal{T})$  is semi-open if and only if there exists  $B \subseteq A^{\circ\dagger}$  such that  $A = A^\circ \cup B$ .

In particular, in 1965 and 1971, Njåstad [62] and Isomichi [37], being not aware of the paper of Levine, studied semi-open sets under the names  *$\beta$ -sets* and *subcondensed sets*, respectively.

Moreover, Njåstad called a subset  $A$  of  $X(\mathcal{T})$  an  $\alpha$ -set if  $A \subseteq A^{\circ\circ}$ . And, he proved that the set  $A$  is an  $\alpha$ -set if and only if there exist  $V \in \mathcal{T}$  and  $B \in \mathcal{N}$  such that  $A = V \setminus B$ .

In 1983, the subset  $A$  was called  $\beta$ -open by Abd El-Monsef et al. [1] if  $A \subseteq A^{-\circ}$ . Moreover, in 1986 Andrijević [3] used the term *semi-preopen* instead of  $\beta$ -open without knowing of [1].

Actually, Andrijević called a subset  $A$  of  $X(\mathcal{T})$  semi-preopen if there exists a preopen subset  $V$  of  $X(\mathcal{T})$  such that  $V \subseteq A \subseteq V^-$ . And, he showed that this is equivalent to the inclusion  $A \subseteq A^{-\circ}$ .

Moreover, in 1996, a subset  $A$  of  $X(\mathcal{T})$  was called  $b$ -open by Andrijević [4] if  $A \subseteq A^{\circ\circ} \cup A^{-\circ}$ . He proved that  $A$  is  $b$ -open if and only if there exist a preopen subset  $B$  and a semi-open subset  $C$  of  $X(\mathcal{T})$  such that  $A = B \cup C$ .

In 1961, a subset  $A$  of a topological space  $X(\mathcal{T})$  was said to have *property Q* by Levine [48] if  $A^{\circ\circ} = A^{-\circ}$ . He proved that  $A$  has property  $Q$  if and only if there exist  $V \in \mathcal{T} \cap \mathcal{F}$  and  $B \in \mathcal{N}$  such that  $A = V \Delta B$ . (See also [7, 10].)

While, in 1991, a subset  $A$  of  $X(\mathcal{T})$  was called a  $\delta$ -set by Chattopadhyay and Bandyopadhyay [8] if  $A^{-\circ} \subseteq A^{\circ\circ}$ . Moreover, in 2001,  $\delta$ -open sets, under the name *quasi-open sets*, were more systematically studied by Császár [15, 16].

In 1992, Ganster et al. [29] already proved that  $A$  is a  $\delta$ -set if and only if  $A = V \cup N$  for some  $V \in \mathcal{T}$  and  $B \in \mathcal{N}$ . Thus,  $\delta$ -sets coincide with the *simply open sets* of Biswas [5] and Neubrunnová [61]. (See also [43, 59, 60].)

Actually, such sets were also first studied by Kuratowski [45, p. 69] in a more general framework. By his definition, a subset  $A$  of  $X(\mathcal{T})$  has to be called *open modulo nowhere dense sets* if there exists  $V \in \mathcal{T}$  such that  $A \Delta V \in \mathcal{N}$ .

## 2 Preliminaries

In our former papers [67, 68], we have shown that the above definitions and several theorems on generalized open sets can be naturally extended not only to generalized topological and closure spaces, but also to relator spaces.

In the sequel, following a terminology introduced by the second author [73], a family  $\mathcal{R}$  of binary relations on a set  $X$  will be called a *relator* on  $X$ , and the ordered pair  $X(\mathcal{R}) = (X, \mathcal{R})$  will be called a *relator space*.

Thus, relator spaces are generalizations of not only *ordered sets* [19] and *uniform spaces* [27], but also *topological, closure and proximity spaces* [57]. However, to include *context spaces* [30] relators on  $X$  to  $Y$  are also needed [81, 82].

For instance, by [85], each *generalized topology*  $\mathcal{T}$  on  $X$  can be easily derived from the family  $\mathcal{R}_{\mathcal{T}}$  of all *Pervin's preorder relations*  $R_V = V^2 \cup V^c \times X$  with  $V \in \mathcal{T}$ . Thus, generalized topologies need not be studied separately.

For a subset  $A$  of the relator space  $X(\mathcal{R})$ , we define

$$A^{\circ} = \text{int}_{\mathcal{R}}(A) = \{ x \in X : \exists R \in \mathcal{R} : R(x) \subseteq A \}$$

and  $A^- = \text{cl}_{\mathcal{R}}(A) = \text{int}_{\mathcal{R}}(A^c)^c$ . And, for instance,  $A \in \mathcal{T}_{\mathcal{R}}$  if  $A \subseteq A^\circ$ .

Now, according to the former definitions on open-like subsets of topologically spaces mentioned in the Motivations, we may also naturally write

- (1)  $A \in \mathcal{T}_{\mathcal{R}}^r$  if  $A = A^{-\circ}$ ;
- (2)  $A \in \mathcal{T}_{\mathcal{R}}^p$  if  $A \subseteq A^{-\circ}$ ;
- (3)  $A \in \mathcal{T}_{\mathcal{R}}^s$  if  $A \subseteq A^{\circ-}$ ;
- (4)  $A \in \mathcal{T}_{\mathcal{R}}^\alpha$  if  $A \subseteq A^{\circ-\circ}$ ;
- (5)  $A \in \mathcal{T}_{\mathcal{R}}^\beta$  if  $A \subseteq A^{-\circ-}$ ;
- (6)  $A \in \mathcal{T}_{\mathcal{R}}^a$  if  $A \subseteq A^{-\circ} \cap A^{\circ-}$ ;
- (7)  $A \in \mathcal{T}_{\mathcal{R}}^b$  if  $A \subseteq A^{-\circ} \cup A^{\circ-}$ ;
- (8)  $A \in \mathcal{T}_{\mathcal{R}}^q$  if there exists  $V \in \mathcal{T}_{\mathcal{R}}$  such that  $V \subseteq A \subseteq V^-$ ;
- (9)  $A \in \mathcal{T}_{\mathcal{R}}^{ps}$  if there exists  $V \in \mathcal{T}_{\mathcal{R}}$  such that  $A \subseteq V \subseteq A^-$ ;
- (10)  $A \in \mathcal{T}_{\mathcal{R}}^\gamma$  if there exists  $V \in \mathcal{T}_{\mathcal{R}}^s$  such that  $A \subseteq V \subseteq A^-$ ;
- (11)  $A \in \mathcal{T}_{\mathcal{R}}^\delta$  if there exists  $V \in \mathcal{T}_{\mathcal{R}}^p$  such that  $V \subseteq A \subseteq V^-$ .

Moreover, the members of the above families may be called the *topologically regular open, preopen, semi-open,  $\alpha$ -open,  $\beta$ -open,  $a$ -open,  $b$ -open, quasi-open, pseudo-open,  $\gamma$ -open and  $\delta$ -open subsets* of the relator space  $X(\mathcal{R})$ , respectively.

Here, the use of the extra term “topologically” can only be motivated by the fact, for any two subsets  $A$  and  $B$  of the relator space  $X(\mathcal{R})$ , we may also naturally write  $B \in \text{Int}_{\mathcal{R}}(A)$  if  $R[B] \subseteq A$  for some  $R \in \mathcal{R}$ .

Thus, by using the plausible notations  $\text{Cl}_{\mathcal{R}}(A) = \text{Int}_{\mathcal{R}}(A^c)^c$ , and  $A \in \tau_{\mathcal{R}}$  if  $A \in \text{Int}_{\mathcal{R}}(A)$ , we may also naturally introduce some reasonable notions of certain *generalized proximally open sets*.

By using the *topological closure (refinement)*

$$\mathcal{R}^\wedge = \{ S \subseteq X^2 : \forall x \in X : x \in S(x)^\circ \}$$

of the relator  $\mathcal{R}$ , it can be shown that  $\text{Int}_{\mathcal{R}^\wedge}(A) = \mathcal{P}(\text{int}_{\mathcal{R}}(A))$  and  $\tau_{\mathcal{R}^\wedge} = \mathcal{T}_{\mathcal{R}}$ .

Therefore, the topological properties of  $\mathcal{R}$  can, in principle, be immediately derived from some proximal ones. For instance,  $\mathcal{R}$  may be called *topologically compact* if  $\mathcal{R}^\wedge$  is *properly compact* in the sense that for each  $S \in \mathcal{R}^\wedge$  there exists a finite subset  $A$  of  $X$  such that  $X = S[A]$ . (See [78] and [80].)

Unfortunately, up till now we have not been able to find the right proximal versions of definitions (1)–(11). However, we can now offer a general framework for introducing and studying the families given by (8)–(11).

Moreover, following the ideas of Chattopadhyay and Bandyopadhyay [8], Császár [15] and others mentioned in the Motivations, we shall now also consider a further important class of generalized topologically open sets.

For the latter purpose, for a subset  $A$  of the relator space  $X(\mathcal{R})$ , we shall write

$$(12) A \in \mathcal{A}_{\mathcal{R}} \quad \text{if} \quad A^{-\circ} \subseteq A^{\circ-}.$$

Thus, the members of the family  $\mathcal{A}_{\mathcal{R}}$  may be called the *topologically simply open subsets* of the relator space  $X(\mathcal{R})$ .

While, for the former purpose, for any two subsets  $A$  and  $B$  of the relator space  $X(\mathcal{R})$  we shall write



$$(13) \quad A \in \text{Ln}_{\mathcal{R}}(B) \quad \text{and} \quad B \in \text{Un}_{\mathcal{R}}(A) \quad \text{if} \quad A \subseteq B \subseteq A^-.$$

Moreover, for a family  $\mathcal{A}$  of subsets of  $X(\mathcal{R})$  we shall define

$$(14) \quad \mathcal{A}^{\ell} = \text{cl}_{\text{Ln}_{\mathcal{R}}}(\mathcal{A}) = \text{Ln}_{\mathcal{R}}^{-1}[\mathcal{A}] \quad \text{and} \quad \mathcal{A}^u = \text{cl}_{\text{Un}_{\mathcal{R}}}(\mathcal{A}) = \text{Un}_{\mathcal{R}}^{-1}[\mathcal{A}].$$

Thus,  $\mathcal{A}^{\ell}$  and  $\mathcal{A}^u$  may be called the *lower and upper nearness closures* of  $\mathcal{A}$ , respectively. Namely, if  $A \in \text{Ln}_{\mathcal{R}}(B)$ , then we may naturally say that  $A$  is near to  $B$  from below and  $B$  is near to  $A$  from above.

The most important particular cases are when  $\mathcal{A}$  is a minimal structure or a generalized topology on  $X$ . Or even more specially,  $\mathcal{A}$  is one of the families  $\mathcal{I}_{\mathcal{R}}$ ,  $\mathcal{I}_{\mathcal{R}}^{\ell}$  or  $\mathcal{I}_{\mathcal{R}}^u$ .

In a subsequent paper, we shall also try to work out a general framework for introducing and studying the families given by definitions (2)–(7). However, for this, in addition to ordinary and hyper relators we shall also need super relators.

The necessary prerequisites on relations and relators, which are certainly unfamiliar to the reader, will be briefly laid out in the subsequent preparatory sections which will also contain several new observations.

These sections may also be useful for all those readers who are not very much interested in the various generalizations of open sets having been studied recently by a surprisingly great number of topologists.

### 3 A Few Basic Facts on Relations

A subset  $F$  of a product set  $X \times Y$  is called a *relation on  $X$  to  $Y$* . In particular, a relation on  $X$  to itself is called a *relation on  $X$* . And,  $\Delta_X = \{(x, x) : x \in X\}$  is called the *identity relation of  $X$* .

If  $F$  is a relation on  $X$  to  $Y$ , then by the above definitions we can also state that  $F$  is a relation on  $X \cup Y$ . However, for several purposes, the latter view of the relation  $F$  would be quite unnatural.

If  $F$  is a relation on  $X$  to  $Y$ , then for any  $x \in X$  and  $A \subseteq X$  the sets  $F(x) = \{y \in Y : (x, y) \in F\}$  and  $F[A] = \bigcup\{F(x) : x \in A\}$  are called the *images or neighbourhoods of  $x$  and  $A$  under  $F$* , respectively.

If  $(x, y) \in F$ , then instead of  $y \in F(x)$ , we may also write  $x F y$ . However, instead of  $F[A]$ , we cannot write  $F(A)$ . Namely, it may occur that, in addition to  $A \subseteq X$ , we also have  $A \in X$ .

Now, the sets  $D_F = \{x \in X : F(x) \neq \emptyset\}$  and  $R_F = F[X]$  may be called the *domain* and *range* of  $F$ , respectively. If in particular  $D_F = X$ , then we may say that  $F$  is a *relation of  $X$  to  $Y$* , or that  $F$  is a *non-partial relation on  $X$  to  $Y$* .

In particular, a relation  $f$  on  $X$  to  $Y$  is called a *function* if for each  $x \in D_f$  there exists  $y \in Y$  such that  $f(x) = \{y\}$ . In this case, by identifying singletons with their elements, we may simply write  $f(x) = y$  instead of  $f(x) = \{y\}$ .

Moreover, a function  $\star$  of  $X$  to itself is called a *unary operation on  $X$* . While, a function  $*$  of  $X^2$  to  $X$  is called a *binary operation on  $X$* . And, for any  $x, y \in X$ , we usually write  $x^\star$  and  $x * y$  instead of  $\star(x)$  and  $*$ (( $x, y$ )), respectively.

If  $F$  is a relation on  $X$  to  $Y$ , then a function  $f$  of  $D_F$  to  $Y$  is called a *selection function* of  $F$  if  $f(x) \in F(x)$  for all  $x \in D_F$ . By using the Axiom of Choice, it can be shown that every relation is the union of its selection functions.

For a relation  $F$  on  $X$  to  $Y$ , we may naturally define two *set-valued functions*  $\varphi_F$  of  $X$  to  $\mathcal{P}(Y)$  and  $\Phi_F$  of  $\mathcal{P}(X)$  to  $\mathcal{P}(Y)$  such that  $\varphi_F(x) = F(x)$  for all  $x \in X$  and  $\Phi_F(A) = F[A]$  for all  $A \subseteq X$ .

Functions of  $X$  to  $\mathcal{P}(Y)$  can be naturally identified with relations on  $X$  to  $Y$ . While, functions of  $\mathcal{P}(X)$  to  $\mathcal{P}(Y)$  are more general objects than relations on  $X$  to  $Y$ . In [88, 93, 94], they were briefly called *corelations* on  $X$  to  $Y$ .

However, a relation on  $\mathcal{P}(X)$  to  $Y$  should be rather called a *super relation* on  $X$  to  $Y$ , and a relation on  $\mathcal{P}(X)$  to  $\mathcal{P}(Y)$  should be rather called a *hyper relation* on  $X$  to  $Y$ . Thus,  $\text{cl}_{\mathcal{R}}$  is a super relation and  $\text{Cl}_{\mathcal{R}}$  is a hyper relation on  $X$ .

If  $F$  is a relation on  $X$  to  $Y$ , then one can easily see that  $F = \bigcup_{x \in X} \{x\} \times F(x)$ . Therefore, the images  $F(x)$ , where  $x \in X$ , uniquely determine  $F$ . Thus, a relation  $F$  on  $X$  to  $Y$  can also be naturally defined by specifying  $F(x)$  for all  $x \in X$ .

For instance, the *complement*  $F^c$  and the *inverse*  $F^{-1}$  can be defined such that  $F^c(x) = F(x)^c = Y \setminus F(x)$  for all  $x \in X$  and  $F^{-1}(y) = \{x \in X : y \in F(x)\}$  for all  $y \in Y$ . Thus, it can be easily seen that  $F^c = X \times Y \setminus F$ .

Moreover, if in addition  $G$  is a relation on  $Y$  to  $Z$ , then the *composition*  $G \circ F$  can be defined such that  $(G \circ F)(x) = G[F(x)]$  for all  $x \in X$ . Thus, it can be easily seen that  $(G \circ F)[A] = G[F[A]] = \bigcup_{y \in F[A]} G(y)$  for all  $A \subseteq X$ .

While, if  $G$  is a relation on  $Z$  to  $W$ , then the *box product*  $F \boxtimes G$  can be defined such that  $(F \boxtimes G)(x, z) = F(x) \times G(z)$  for all  $x \in X$  and  $z \in Z$ . Thus, it can be shown that  $(F \boxtimes G)[A] = G \circ A \circ F^{-1}$  for all  $A \subseteq X \times Z$  [87].

Hence, by taking  $A = \{(x, z)\}$ , and  $A = \Delta_Y$  if  $Y = Z$ , one can at once see that the box and composition products are actually equivalent tools. However, the box product can be immediately defined for any family of relations.

Now, a relation  $R$  on  $X$  may be briefly defined to be *reflexive* on  $X$  if  $\Delta_X \subseteq R$ , and *transitive* if  $R \circ R \subseteq R$ . Moreover,  $R$  may be briefly defined to be *symmetric* if  $R^{-1} \subseteq R$ , and *antisymmetric* if  $R \cap R^{-1} \subseteq \Delta_X$ .

Thus, a reflexive and transitive (symmetric) relation may be called a *preorder (tolerance) relation*. And, a symmetric (antisymmetric) preorder relation may be called an *equivalence (partial order) relation*.

For any relation  $R$  on  $X$ , we may also naturally define  $R^0 = \Delta_X$  and  $R^n = R \circ R^{n-1}$  if  $n \in \mathbb{N}$ . Moreover, we may also naturally define  $R^\infty = \bigcup_{n=0}^\infty R^n$ . Thus,  $R^\infty$  is the smallest preorder relation on  $X$  containing  $R$  [33].

For  $A \subseteq X$ , the Pervin relation  $R_A = A^2 \cup A^c \times X$  is an important preorder on  $X$  [66]. While, for a *pseudometric*  $d$  on  $X$ , the *Weil surrounding*  $B_r = \{(x, y) \in X^2 : d(x, y) < r\}$ , with  $r > 0$ , is an important tolerance on  $X$  [99].

Note that  $S_A = R_A \cap R_A^{-1} = R_A \cap R_{A^c} = A^2 \cap (A^c)^2$  is already an equivalence relation on  $X$ . And, more generally if  $\mathcal{A}$  is a *cover (partition)* of  $X$ , then  $S_{\mathcal{A}} = \bigcup_{A \in \mathcal{A}} A^2$  is a tolerance (equivalence) relation on  $X$ .

As an important generalization of the Pervin relation  $R_A$ , for any  $A \subseteq X$  and  $B \subseteq Y$ , we may also naturally consider the *Hunsaker-Lindgren relation*  $R_{(A,B)} = A \times B \cap A^c \times Y$  [35]. Namely, thus we evidently have  $R_A = R_{(A,A)}$ .

The Pervin relations  $R_A$  and the Hunsaker-Lindgren relations  $R_{(A,B)}$  were actually first used by Davis [20] and Császár [12, pp. 42 and 351] in some less explicit and convenient forms, respectively.

## 4 A Few Basic Facts on Relators

A family  $\mathcal{R}$  of relations on one set  $X$  to another  $Y$  is called a *relator on  $X$  to  $Y$* , and the ordered pair  $(X, Y)(\mathcal{R}) = ((X, Y), \mathcal{R})$  is called a *relator space*. For the origins of this notion, see [73, 81], and the references in [73].

If in particular  $\mathcal{R}$  is a relator on  $X$  to itself, then  $\mathcal{R}$  is simply called a *relator on  $X$* . Thus, by identifying singletons with their elements, we may naturally write  $X(\mathcal{R})$  instead of  $(X, X)(\mathcal{R})$ . Namely,  $(X, X) = \{\{X\}, \{X, X\}\} = \{\{X\}\}$ .

Relator spaces of this simpler type are already substantial generalizations of the various *ordered sets* [19] and *uniform spaces* [27]. However, they are insufficient for some important purposes. (See, for instance, [30] and [81, 91].)

A relator  $\mathcal{R}$  on  $X$  to  $Y$ , or the relator space  $(X, Y)(\mathcal{R})$ , is called *simple* if  $\mathcal{R} = \{R\}$  for some relation  $R$  on  $X$  to  $Y$ . Simple relator spaces  $(X, Y)(R)$  and  $X(R)$  were called *formal contexts* and *gosets* in [30] and [90], respectively.

Moreover, a relator  $\mathcal{R}$  on  $X$ , or the relator space  $X(\mathcal{R})$ , may, for instance, be naturally called *reflexive* if each member of  $\mathcal{R}$  is reflexive on  $X$ . Thus, we may also naturally speak of *preorder, tolerance, and equivalence relators*.

For instance, for a family  $\mathcal{A}$  of subsets of  $X$ , the family  $\mathcal{R}_{\mathcal{A}} = \{R_A : A \in \mathcal{A}\}$ , where  $R_A = A^2 \cup A^c \times X$ , is an important preorder relator on  $X$ . Such relators were first used by Pervin [66] and Levine [52].

While, for a family  $\mathcal{D}$  of *pseudo-metrics* on  $X$ , the family  $\mathcal{R}_{\mathcal{D}} = \{B_r^d : r > 0, d \in \mathcal{D}\}$ , where  $B_r^d = \{(x, y) : d(x, y) < r\}$ , is an important tolerance relator on  $X$ . Such relators were first considered by Weil [99].

Moreover, if  $\mathfrak{S}$  is a family of *partitions* of  $X$ , then the family  $\mathcal{R}_{\mathfrak{S}} = \{S_{\mathcal{A}} : \mathcal{A} \in \mathfrak{S}\}$ , where  $S_{\mathcal{A}} = \bigcup_{A \in \mathcal{A}} A^2$ , is an equivalence relator on  $X$ . Such practically important relators were first investigated by Levine [51].

If  $\star$  is a unary operation for relations on  $X$  to  $Y$ , then for any relator  $\mathcal{R}$  on  $X$  to  $Y$  we may naturally define  $\mathcal{R}^\star = \{R^\star : R \in \mathcal{R}\}$ . However, this plausible notation may cause some confusions whenever, for instance,  $\star = c$ .

In particular, for any relator  $\mathcal{R}$  on  $X$ , we may naturally define  $\mathcal{R}^\infty = \{R^\infty : R \in \mathcal{R}\}$ . Moreover, we may also naturally define  $\mathcal{R}^\partial = \{S \subseteq X^2 : S^\infty \in \mathcal{R}\}$ . These operations were first introduced by Mala [53, 55] and Pataki [64, 65].

While, if  $*$  is a binary operation for relations, then for any two relators  $\mathcal{R}$  and  $\mathcal{S}$  we may naturally define  $\mathcal{R} * \mathcal{S} = \{R * S : R \in \mathcal{R}, S \in \mathcal{S}\}$ . However, this plausible notation may again cause some confusions whenever, for instance,  $*$  =  $\cap$ .

Therefore, in general we rather write  $\mathcal{R} \wedge \mathcal{S} = \{R \cap S : R \in \mathcal{R}, S \in \mathcal{S}\}$ . Moreover, for instance, we also write  $\mathcal{R} \Delta \mathcal{R}^{-1} = \{R \cap R^{-1} : R \in \mathcal{R}\}$ . Note that thus  $\mathcal{R} \Delta \mathcal{R}^{-1}$  is a symmetric relator such that  $\mathcal{R} \Delta \mathcal{R}^{-1} \subseteq \mathcal{R} \wedge \mathcal{R}^{-1}$ .

A function  $\square$  of the family of all relators on  $X$  to  $Y$  is called a *direct (indirect) unary operation for relators* if, for every relator  $\mathcal{R}$  on  $X$  to  $Y$ , the value  $\mathcal{R}^\square = \square(\mathcal{R})$  is a relator on  $X$  to  $Y$  (on  $Y$  to  $X$ ).

For instance,  $c$  and  $-1$  are *involution operations* for relators. While,  $\infty$  and  $\partial$  are *projection operations* for relators. Moreover, the operation  $\square = c, \infty$  or  $\partial$  is *inversion compatible* in the sense that  $\mathcal{R}^{\square^{-1}} = \mathcal{R}^{-1}\square$ .

More generally, a function  $\mathfrak{F}$  of the family of all relators on  $X$  to  $Y$  is called a *structure for relators* if, for every relator  $\mathcal{R}$  on  $X$  to  $Y$ , the value  $\mathfrak{F}_{\mathcal{R}} = \mathfrak{F}(\mathcal{R})$  is in a power set depending only on  $X$  and  $Y$ .

For instance, if  $\text{cl}_{\mathcal{R}}(B) = \bigcap \{R^{-1}[B] : R \in \mathcal{R}\}$  for every relator  $\mathcal{R}$  on  $X$  to  $Y$  and  $B \subseteq Y$ , then the function  $\mathfrak{F}$ , defined by  $\mathfrak{F}(\mathcal{R}) = \text{cl}_{\mathcal{R}}$ , is a structure for relators such that  $\mathfrak{F}(\mathcal{R}) \subseteq \mathcal{P}(Y) \times X$ , and thus  $\mathfrak{F}(\mathcal{R}) \in \mathcal{P}(\mathcal{P}(Y) \times X)$ .

A structure  $\mathfrak{F}$  for relators is called *increasing* if  $\mathcal{R} \subseteq \mathcal{S}$  implies  $\mathfrak{F}_{\mathcal{R}} \subseteq \mathfrak{F}_{\mathcal{S}}$  for any two relators  $\mathcal{R}$  and  $\mathcal{S}$  on  $X$  to  $Y$ . And,  $\mathfrak{F}$  is called *quasi-increasing* if  $R \in \mathcal{R}$  implies  $\mathfrak{F}_R \subseteq \mathfrak{F}_{\mathcal{R}}$  for any relator  $\mathcal{R}$  on  $X$  to  $Y$ . Note that here  $\mathfrak{F}_R = \mathfrak{F}_{\{R\}}$ .

Moreover, the structure  $\mathfrak{F}$  is called *union-preserving* if  $\mathfrak{F}_{\bigcup_{i \in I} \mathcal{R}_i} = \bigcup_{i \in I} \mathfrak{F}_{\mathcal{R}_i}$  for any family  $(\mathcal{R}_i)_{i \in I}$  of relators on  $X$  to  $Y$ . It can be shown that  $\mathfrak{F}$  is union-preserving if and only if  $\mathfrak{F}_{\mathcal{R}} = \bigcup_{R \in \mathcal{R}} \mathfrak{F}_R$  for every relator  $\mathcal{R}$  on  $X$  to  $Y$  [88].

In particular, an increasing operation  $\square$  for relators on  $X$  to  $Y$  is called a *projection or modification operation* for relators if it is idempotent in the sense that  $\mathcal{R}^{\square\square} = \mathcal{R}^\square$  holds for any relator  $\mathcal{R}$  on  $X$  to  $Y$ .

Moreover, a projection operation  $\square$  for relators on  $X$  to  $Y$  is called a *closure or refinement operation* for relators if it is extensive in the sense that  $\mathcal{R} \subseteq \mathcal{R}^\square$  holds for any relator  $\mathcal{R}$  on  $X$  to  $Y$ .

By using Pataki connections [64, 95], several closure operations can be derived from union-preserving structures. However, more generally, one can find first the Galois adjoint  $\mathfrak{G}$  of such a structure  $\mathfrak{F}$ , and then take  $\square_{\mathfrak{F}} = \mathfrak{G} \circ \mathfrak{F}$  [84].

Now, for an operation  $\square$  for relators, a relator  $\mathcal{R}$  on  $X$  to  $Y$  may be naturally called  $\square$ -*fine* if  $\mathcal{R}^\square = \mathcal{R}$ . And, for some structure  $\mathfrak{F}$  for relators, two relators  $\mathcal{R}$  and  $\mathcal{S}$  on  $X$  to  $Y$  may be naturally called  $\mathfrak{F}$ -*equivalent* if  $\mathfrak{F}_{\mathcal{R}} = \mathfrak{F}_{\mathcal{S}}$ .

Moreover, for a structure  $\mathfrak{F}$  for relators, a relator  $\mathcal{R}$  on  $X$  to  $Y$  may, for instance, be naturally called  $\mathfrak{F}$ -*simple* if  $\mathfrak{F}_{\mathcal{R}} = \mathfrak{F}_R$  for some relation  $R$  on  $X$  to  $Y$ . Thus, in particular singleton relators have to be actually called *properly simple*.

## 5 Structures Derived from Relators

**Definition 1** If  $\mathcal{R}$  is a relator on  $X$  to  $Y$ , then for any  $A \subseteq X, B \subseteq Y$  and  $x \in X, y \in Y$  we define:

- (1)  $A \in \text{Int}_{\mathcal{R}}(B)$  if  $R[A] \subseteq B$  for some  $R \in \mathcal{R}$ ;
- (2)  $A \in \text{Cl}_{\mathcal{R}}(B)$  if  $R[A] \cap B \neq \emptyset$  for all  $R \in \mathcal{R}$ ;

- (3)  $x \in \text{int}_{\mathcal{R}}(B)$  if  $\{x\} \in \text{Int}_{\mathcal{R}}(B)$ ;      (4)  $x \in \sigma_{\mathcal{R}}(y)$  if  $x \in \text{int}_{\mathcal{R}}(\{y\})$ ;
- (5)  $x \in \text{cl}_{\mathcal{R}}(B)$  if  $\{x\} \in \text{Cl}_{\mathcal{R}}(B)$ ;      (6)  $x \in \rho_{\mathcal{R}}(y)$  if  $x \in \text{cl}_{\mathcal{R}}(\{y\})$ ;
- (7)  $B \in \mathcal{E}_{\mathcal{R}}$  if  $\text{int}_{\mathcal{R}}(B) \neq \emptyset$ ;      (8)  $B \in \mathcal{D}_{\mathcal{R}}$  if  $\text{cl}_{\mathcal{R}}(B) = X$ .

*Remark 1* The relations  $\text{Int}_{\mathcal{R}}$ ,  $\text{int}_{\mathcal{R}}$  and  $\sigma_{\mathcal{R}}$  are called the proximal, topological and infinitesimal interiors generated by  $\mathcal{R}$ , respectively. While, the members of the families,  $\mathcal{E}_{\mathcal{R}}$  and  $\mathcal{D}_{\mathcal{R}}$  are called the fat and dense subsets of the relator space  $(X, Y)(\mathcal{R})$ , respectively.

The origins of the relations  $\text{Cl}_{\mathcal{R}}$  and  $\text{Int}_{\mathcal{R}}$  go back to Efremović’s proximity  $\delta$  [24] and Smirnov’s strong inclusion  $\Subset$  [71], respectively. While, the convenient notations  $\text{Cl}_{\mathcal{R}}$  and  $\text{Int}_{\mathcal{R}}$ , and family  $\mathcal{E}_{\mathcal{R}}$ , together with its dual  $\mathcal{D}_{\mathcal{R}}$ , was first explicitly used by the second author [73, 75, 76, 83].

The following simple, but important theorem shows that the big interior and closure are equivalent tools in a relator space.

**Theorem 1** *If  $\mathcal{R}$  is a relator on  $X$  to  $Y$ , then for any  $B \subseteq Y$  we have*

- (1)  $\text{Cl}_{\mathcal{R}}(B) = \text{Int}_{\mathcal{R}}(B^c)^c$ ;      (2)  $\text{Int}_{\mathcal{R}}(B) = \text{Cl}_{\mathcal{R}}(B^c)^c$ .

*Remark 2* By using the notation  $\mathcal{C}_Y(B) = B^c$ , assertion (1) can be expressed in the more concise form that  $\text{Cl}_{\mathcal{R}} = (\text{Int}_{\mathcal{R}} \circ \mathcal{C}_Y)^c = (\text{Int}_{\mathcal{R}})^c \circ \mathcal{C}_Y$ .

From Theorem 1, we can easily derive the following more familiar

**Theorem 2** *If  $\mathcal{R}$  is a relator on  $X$  to  $Y$ , then for any  $B \subseteq Y$  we have*

- (1)  $\text{cl}_{\mathcal{R}}(B) = \text{int}_{\mathcal{R}}(B^c)^c$ ;      (2)  $\text{int}_{\mathcal{R}}(B) = \text{cl}_{\mathcal{R}}(B^c)^c$ .

*Remark 3* By using the convenient notations  $B^- = \text{cl}_{\mathcal{R}}(B)$  and  $B^\circ = \text{int}_{\mathcal{R}}(B)$ , assertion (1) can be expressed in the more concise form that  $- = c \circ c$ , or equivalently  $-c = c \circ -$ .

The small closure and interior are usually much weaker tools than the big ones. Namely, in general, we can only prove the following

**Theorem 3** *If  $\mathcal{R}$  is a relator on  $X$  to  $Y$ , then for any  $A \subseteq X$  and  $B \subseteq Y$*

- (1)  $A \in \text{Int}_{\mathcal{R}}(B)$  implies  $A \subseteq \text{int}_{\mathcal{R}}(B)$ ;
- (2)  $A \cap \text{cl}_{\mathcal{R}}(B) \neq \emptyset$  implies  $A \in \text{Cl}_{\mathcal{R}}(B)$ .

Concerning closures and interiors, we can also prove the following two theorems which show that, despite their equivalences, closures are sometimes more convenient tools than interiors.

**Theorem 4** *For any relator  $\mathcal{R}$  on  $X$  to  $Y$ , we have*

- (1)  $\text{Cl}_{\mathcal{R}^{-1}} = \text{Cl}_{\mathcal{R}}^{-1}$ ;      (2)  $\text{Int}_{\mathcal{R}^{-1}} = \mathcal{C}_Y \circ \text{Int}_{\mathcal{R}}^{-1} \circ \mathcal{C}_X$ .

**Theorem 5** *If  $\mathcal{R}$  is a relator on  $X$  to  $Y$ , then for any  $B \subseteq Y$ , we have (1)*

$$\text{cl}_{\mathcal{R}}(B) = \bigcap_{R \in \mathcal{R}} R^{-1}[B]; \quad (2) \text{int}_{\mathcal{R}}(B) = \bigcup_{R \in \mathcal{R}} R^{-1}[B^c]^c.$$

From the  $B = \{y\}$  particular case of this theorem, we can easily derive

**Corollary 1** For any relator  $\mathcal{R}$  on  $X$  to  $Y$ , we have

$$\rho_{\mathcal{R}} = \bigcap \mathcal{R}^{-1} = \left( \bigcap \mathcal{R} \right)^{-1}.$$

Moreover, by using the  $\mathcal{R} = \{R\}$  particular case of Theorem 5, we can prove

**Theorem 6** If  $R$  is a relation on  $X$  to  $Y$ , then for any  $A \subseteq X$  and  $B \subseteq Y$

$$A \subseteq \text{int}_R(B) \iff \text{cl}_{R^{-1}}(A) \subseteq B.$$

*Remark 4* This shows that the mappings  $A \mapsto \text{cl}_{R^{-1}}(A)$  and  $B \mapsto \text{int}_R(B)$  establish a Galois connection between the posets  $\mathcal{P}(X)$  and  $\mathcal{P}(Y)$ .

The above important closure-interior Galois connection, used first in [92], is not independent from the well-known upper and lower bound one [86].

The following two closely related theorems show that the fat and dense sets are also equivalent tools in a relator space.

**Theorem 7** If  $\mathcal{R}$  is a relator on  $X$  to  $Y$ , then for any  $B \subseteq Y$  we have (1)  $B \in \mathcal{D}_{\mathcal{R}} \iff B^c \notin \mathcal{E}_{\mathcal{R}}$ ; (2)  $B \in \mathcal{E}_{\mathcal{R}} \iff B^c \notin \mathcal{D}_{\mathcal{R}}$ .

**Theorem 8** If  $\mathcal{R}$  is a relator on  $X$  to  $Y$ , then for any  $B \subseteq Y$  we have

- (1)  $B \in \mathcal{D}_{\mathcal{R}}$  if and only if  $B \cap E \neq \emptyset$  for all  $E \in \mathcal{E}_{\mathcal{R}}$ ;
- (2)  $B \in \mathcal{E}_{\mathcal{R}}$  if and only if  $B \cap D \neq \emptyset$  for all  $D \in \mathcal{D}_{\mathcal{R}}$ .

*Remark 5* By the corresponding definitions, we have  $R(x) \in \mathcal{E}_{\mathcal{R}}$  and thus also  $R(x)^c \notin \mathcal{D}_{\mathcal{R}}$  for all  $x \in X$  and  $R \in \mathcal{R}$ .

While, by using the notation  $\mathcal{U}_{\mathcal{R}}(x) = \text{int}_{\mathcal{R}}^{-1}(x) = \{B \subseteq Y : x \in \text{int}_{\mathcal{R}}(B)\}$ , we can note that  $\mathcal{E}_{\mathcal{R}} = \bigcup_{x \in X} \mathcal{U}_{\mathcal{R}}(x)$ .

By using Definition 1, we may easily introduce several further important definitions. For instance, we may also naturally have the following

**Definition 2** If  $\mathcal{R}$  is a relator on  $X$  to  $Y$ , then for any  $B \subseteq Y$ , we define (1)  $\text{bnd}_{\mathcal{R}}(B) = \text{cl}_{\mathcal{R}}(B) \setminus \text{int}_{\mathcal{R}}(B)$ .

Moreover, if in particular  $\mathcal{R}$  is a relator on  $X$ , then for any  $A \subseteq X$  we also define (2)  $\text{res}_{\mathcal{R}}(A) = \text{cl}_{\mathcal{R}}(A) \setminus A$ ; (3)  $\text{bor}_{\mathcal{R}}(A) = A \setminus \text{int}_{\mathcal{R}}(A)$ .

*Remark 6* Somewhat differently, the *border*, *boundary* and *residue* of a set in neighbourhood and closure spaces were also introduced by Hausdorff and Kuratowski [44, pp. 4–5]. (See also Elez and Papaz [26] for a recent treatment.)

If in particular  $\mathcal{R}$  is a reflexive relator on  $X$ , then by Definition 1, for any  $A \subseteq X$ , we have  $A^\circ \subseteq A \subseteq A^-$ . Therefore,

$$\text{bnd}_{\mathcal{R}}(A) = \text{res}_{\mathcal{R}}(A) \cup \text{bor}_{\mathcal{R}}(A) = \text{res}_{\mathcal{R}}(A) \cup \text{res}_{\mathcal{R}}(A^c).$$

Namely, by using Definition 2 and Theorem 2, we can easily see that

$$\text{res}_{\mathcal{R}}(A^c) = A^{c-} \setminus A^c = A^{c-} \cap A^{cc} = A^{\circ c} \cap A = A \setminus A^{\circ} = \text{bor}_{\mathcal{R}}(A).$$

Note that if in particular  $A \in \mathcal{I}_{\mathcal{R}}$  in the sense that  $A \subseteq A^{\circ}$ , then  $\text{bor}_{\mathcal{R}}(A) = \emptyset$ . Therefore, in this particular case, by the above equality, we can simply state that  $\text{bnd}_{\mathcal{R}}(A) = \text{res}_{\mathcal{R}}(A)$ .

## 6 Further Structures Derived from Relators

By using Definition 1, we may also naturally introduce the following

**Definition 3** If  $\mathcal{R}$  is a relator on  $X$ , then for any  $A \subseteq X$  we also define:

- (1)  $A \in \tau_{\mathcal{R}}$  if  $A \in \text{Int}_{\mathcal{R}}(A)$ ;                      (2)  $A \in \varepsilon_{\mathcal{R}}$  if  $A^c \notin \text{Cl}_{\mathcal{R}}(A)$ ;
- (3)  $A \in \mathcal{I}_{\mathcal{R}}$  if  $A \subseteq \text{int}_{\mathcal{R}}(A)$ ;                      (4)  $A \in \mathcal{F}_{\mathcal{R}}$  if  $\text{cl}_{\mathcal{R}}(A) \subseteq A$ ;
- (5)  $A \in \mathcal{N}_{\mathcal{R}}$  if  $\text{cl}_{\mathcal{R}}(A) \notin \mathcal{E}_{\mathcal{R}}$ ;                      (6)  $A \in \mathcal{M}_{\mathcal{R}}$  if  $\text{int}_{\mathcal{R}}(A) \in \mathcal{D}_{\mathcal{R}}$ .

*Remark 7* The members of the families,  $\tau_{\mathcal{R}}$  and  $\mathcal{I}_{\mathcal{R}}$  and  $\mathcal{N}_{\mathcal{R}}$  are called the *proximally open, topologically open and rare (or nowhere dense) subsets* of the relator space  $X(\mathcal{R})$ , respectively.

The family  $\tau_{\mathcal{R}}$  was first introduced by the second author in [75, 76]. While, the practical notation  $\varepsilon_{\mathcal{R}}$  was suggested by János Kurdics who first noticed that “connectedness” is a particular case of “well-chainedness”. (See [46, 47, 65, 69].)

By using the corresponding results of Section 5, we can easily establish the following theorems.

**Theorem 9** *If  $\mathcal{R}$  is a relator on  $X$ , then for any  $A \subseteq X$ , we have (1)  $A \in \varepsilon_{\mathcal{R}} \iff A^c \in \tau_{\mathcal{R}}$ ;                      (2)  $A \in \tau_{\mathcal{R}} \iff A^c \in \varepsilon_{\mathcal{R}}$ .*

**Theorem 10** *For any relator  $\mathcal{R}$  on  $X$ , we have (1)  $\varepsilon_{\mathcal{R}} = \tau_{\mathcal{R}^{-1}}$ ;                      (2)  $\tau_{\mathcal{R}} = \varepsilon_{\mathcal{R}^{-1}}$ .*

**Theorem 11** *If  $\mathcal{R}$  is a relator on  $X$ , then for any  $A \subseteq X$ , we have (1)  $A \in \mathcal{F}_{\mathcal{R}} \iff A^c \in \mathcal{I}_{\mathcal{R}}$ ;                      (2)  $A \in \mathcal{I}_{\mathcal{R}} \iff A^c \in \mathcal{F}_{\mathcal{R}}$ .*

**Corollary 2** *If  $\mathcal{R}$  is a relator on  $X$  and  $A \subseteq X$  and  $V \in \mathcal{I}_{\mathcal{R}}$  such that  $A \cap V = \emptyset$ , then  $\text{cl}_{\mathcal{R}}(A) \cap V = \emptyset$  also hold.*

**Proof** By Theorem 11, we have  $V^c \in \mathcal{F}_{\mathcal{R}}$ . Thus, by Definition 3, we also have  $V^{c-} \subseteq V^c$ . Hence, by using the increasingness of the operation  $-$ , we can see that  $A \cap V = \emptyset \implies A \subseteq V^c \implies A^- \subseteq V^{c-} \implies A^- \subseteq V^c \implies A^- \cap V = \emptyset$ .

*Remark 8* Note that if  $\mathcal{R}$  is a reflexive relator on  $X$ , then  $A \subseteq A^-$  for any  $A \subseteq X$ . Therefore,  $A^- \cap V = \emptyset$  trivially implies  $A \cap V = \emptyset$  for any  $A, V \subseteq X$ .

**Theorem 12** For any relator  $\mathcal{R}$  on  $X$ , we have (1)  $\tau_{\mathcal{R}} \subseteq \mathcal{T}_{\mathcal{R}}$ ; (2)  $\varepsilon_{\mathcal{R}} \subseteq \mathcal{F}_{\mathcal{R}}$ .

*Remark 9* In particular, for any relation  $R$  on  $X$ , we have

(1)  $\tau_R = \mathcal{T}_R$ ; (2)  $\varepsilon_R = \mathcal{F}_R$ .

**Theorem 13** For any relator  $\mathcal{R}$  on  $X$ , we have (1)  $\mathcal{T}_{\mathcal{R}} \setminus \{\emptyset\} \subseteq \mathcal{E}_{\mathcal{R}}$ ; (2)

$\mathcal{D}_{\mathcal{R}} \cap \mathcal{F}_{\mathcal{R}} \subseteq \{X\}$ .

*Remark 10* Hence, by using global complementations, we can easily infer that  $\mathcal{F}_{\mathcal{R}} \subseteq (\mathcal{D}_{\mathcal{R}})^c \cup \{X\}$  and  $\mathcal{D}_{\mathcal{R}} \subseteq (\mathcal{F}_{\mathcal{R}})^c \cup \{X\}$ .

**Theorem 14** If  $\mathcal{R}$  is a relator on  $X$ , then for any  $A \subseteq X$  we have (1)  $A \in \mathcal{E}_{\mathcal{R}}$  if  $V \subseteq A$  for some  $V \in \mathcal{T}_{\mathcal{R}} \setminus \{\emptyset\}$ ;

(2)  $A \in \mathcal{D}_{\mathcal{R}}$  only if  $A \setminus W \neq \emptyset$  for all  $W \in \mathcal{F}_{\mathcal{R}} \setminus \{X\}$ .

*Remark 11* The fat sets are frequently more convenient tools than the topologically open ones. For instance, if  $\leq$  is a relation on  $X$ , then  $\mathcal{T}_{\leq}$  and  $\mathcal{E}_{\leq}$  are the families of all ascending and residual subsets of the goset  $X(\leq)$ , respectively.

Moreover, if in particular  $X = \mathbb{R}$  and  $R(x) = \{x - 1\} \cup [x, +\infty[$  for all  $x \in X$ , then  $R$  is a reflexive relation on  $X$  such that  $\mathcal{T}_R = \{\emptyset, X\}$ , but  $\mathcal{E}_R$  is quite a large family. Namely, the supersets of each  $R(x)$  are also contained in  $\mathcal{E}_R$ .

However, the importance of fat and dense lies mainly in the following

**Definition 4** If  $\mathcal{R}$  is a relator on  $X$  to  $Y$ , and  $\varphi$  and  $\psi$  are functions of a relator space  $\Gamma(\mathcal{U})$  to  $X$  and  $Y$ , respectively, then by using the notation

$$(\varphi, \psi)(\gamma) = (\varphi(\gamma), \psi(\gamma))$$

for all  $\gamma \in \Gamma$ , we may also naturally define (1)  $\varphi \in \text{Lim}_{\mathcal{R}}(\psi)$  if  $(\varphi, \psi)^{-1}[R] \in \mathcal{E}_{\mathcal{U}}$  for all  $R \in \mathcal{R}$ ,

(2)  $\varphi \in \text{Adh}_{\mathcal{R}}(\psi)$  if  $(\varphi, \psi)^{-1}[R] \in \mathcal{D}_{\mathcal{U}}$  for all  $R \in \mathcal{R}$ .

Moreover, for any  $x \in X$ , we may also naturally define:

(3)  $x \in \lim_{\mathcal{R}}(\psi)$  if  $x_{\Gamma} \in \text{Lim}_{\mathcal{R}}(\psi)$ , (4)  $x \in \text{adh}_{\mathcal{R}}(\psi)$  if  $x_{\Gamma} \in \text{Adh}_{\mathcal{R}}(\psi)$ , where  $x_{\Gamma}$  is a function of  $\Gamma$  to  $X$  such that  $x_{\Gamma}(\gamma) = x$  for all  $\gamma \in \Gamma$ .

*Remark 12* Fortunately, the small limit and adherece relations are equivalent to the small closure and interior ones.

However, the big limit and adherence relations, suggested by Efremović and Švarc [25], are usually stronger tools than the big closure and interior ones.

In this respect, it seems convenient to only mention here the following



**Theorem 15** *If  $\mathcal{R}$  is a relator on  $X$  to  $Y$ , then for any  $A \subseteq X$  and  $B \subseteq Y$  the following assertions are equivalent:*

- (1)  $A \in \text{Cl}_{\mathcal{R}}(B)$ ;
- (2) *there exist functions  $\varphi$  and  $\psi$  of the poset  $\mathcal{R}(\supseteq)$  to  $A$  and  $B$ , respectively, such that  $\varphi \in \text{Lim}_{\mathcal{R}}(\psi)$ ;*
- (3) *there exist functions  $\varphi$  and  $\psi$  of a relator space  $\Gamma(\mathcal{U})$  to  $A$  and  $B$ , respectively, such that  $\varphi \in \text{Lim}_{\mathcal{R}}(\psi)$ .*

**Proof** For instance, if (1) holds, then for each  $R \in \mathcal{R}$ , we have  $R[A] \cap B \neq \emptyset$ . Therefore, there exist  $\varphi(R) \in A$  and  $\psi(R) \in B$  such that  $\psi(R) \in R(\varphi(R))$ . Hence, we can already infer that  $(\varphi, \psi)(R) = (\varphi(R), \psi(R)) \in R$ , and thus also  $R \in (\varphi, \psi)^{-1}[R]$ .

Therefore, if  $R \in \mathcal{R}$ , then for any  $S \in \mathcal{R}$ , with  $R \supseteq S$ , we have

$$S \in (\varphi, \psi)^{-1}[S] \subseteq (\varphi, \psi)^{-1}[R].$$

This shows that  $(\varphi, \psi)^{-1}[R]$  is a fat subset of  $\mathcal{R}(\supseteq)$ , and thus  $\varphi \in \text{Lim}_{\mathcal{R}}(\psi)$ .

*Remark 13* Finally, we note that if  $\mathcal{R}$  is a relator on  $X$  to  $Y$ , then according to [82], for any  $A \subseteq X$  and  $B \subseteq Y$ , we may also naturally write  $A \in \text{Lb}_{\mathcal{R}}(B)$  and  $B \in \text{Ub}_{\mathcal{R}}(A)$  if there exists  $R \in \mathcal{R}$  such that  $A \times B \subseteq R$ .

However, the algebraic structures  $\text{Lb}_{\mathcal{R}}$  and  $\text{Ub}_{\mathcal{R}}$ , and the structures derivable from them, are not independent of the former topological ones. Namely, it can be easily shown that  $\text{Lb}_{\mathcal{R}} = \text{Int}_{\mathcal{R}^c} \circ \mathcal{C}_Y$  and  $\text{Int}_{\mathcal{R}} = \text{Lb}_{\mathcal{R}^c} \circ \mathcal{C}_Y$ .

## 7 Closure Operations for Relators

Similar operations for relators have formerly been studied by Kenyon [42], Nakano–Nakano [58], Száz [77, 79] and Pataki [64].

**Definition 5** For any relator  $\mathcal{R}$  on  $X$  to  $Y$ , the relators

$$\begin{aligned} \mathcal{R}^* &= \{ S \subseteq X \times Y : \exists R \in \mathcal{R} : R \subseteq S \}; \\ \mathcal{R}^\# &= \{ S \subseteq X \times Y : \forall A \subseteq X : \exists R \in \mathcal{R} : R[A] \subseteq S[A] \}; \\ \mathcal{R}^\wedge &= \{ S \subseteq X \times Y : \forall x \in X : \exists R \in \mathcal{R} : R(x) \subseteq S(x) \}; \end{aligned}$$

and

$$\mathcal{R}^\Delta = \{ S \subseteq X \times Y : \forall x \in X : \exists u \in X : \exists R \in \mathcal{R} : R(u) \subseteq S(x) \}$$

are called the *uniform, proximal, topological and paratopological closures (or refinements)* of the relator  $\mathcal{R}$ , respectively.

*Remark 14* Thus, we evidently have  $\mathcal{R} \subseteq \mathcal{R}^* \subseteq \mathcal{R}^\# \subseteq \mathcal{R}^\wedge \subseteq \mathcal{R}^\Delta$ . Moreover, if  $\mathcal{R}$  is a relator on  $X$ , then it can be easily shown that  $\mathcal{R}^\infty \subseteq \mathcal{R}^{*\infty} \subseteq \mathcal{R}^{\infty*} \subseteq \mathcal{R}^*$ .

*Remark 15* However, it is now more important to note that, because of Definition 1, we also have

$$\begin{aligned} \mathcal{R}^\# &= \{ S \subseteq X \times Y : \forall A \subseteq X : A \in \text{Int}_{\mathcal{R}}(S[A]) \}, \\ \mathcal{R}^\wedge &= \{ S \subseteq X \times Y : \forall x \in X : x \in \text{int}_{\mathcal{R}}(S(x)) \}, \\ \mathcal{R}^\Delta &= \{ S \subseteq X \times Y : \forall x \in X : S(x) \in \mathcal{E}_{\mathcal{R}} \}. \end{aligned}$$

Moreover, by using Pataki connections [64, 69, 95], the following equivalences and their consequences can be proved in a unified way.

**Theorem 16**  $\#, \wedge$  and  $\Delta$  are closure operations for relators on  $X$  to  $Y$  such that, for any two relators  $\mathcal{R}$  and  $\mathcal{S}$  on  $X$  to  $Y$ , we have

- (1)  $\mathcal{S} \subseteq \mathcal{R}^\# \iff \mathcal{S}^\# \subseteq \mathcal{R}^\# \iff \text{Int}_{\mathcal{S}} \subseteq \text{Int}_{\mathcal{R}} \iff \text{Cl}_{\mathcal{R}} \subseteq \text{Cl}_{\mathcal{S}}$ ,
- (2)  $\mathcal{S} \subseteq \mathcal{R}^\wedge \iff \mathcal{S}^\wedge \subseteq \mathcal{R}^\wedge \iff \text{int}_{\mathcal{S}} \subseteq \text{int}_{\mathcal{R}} \iff \text{cl}_{\mathcal{R}} \subseteq \text{cl}_{\mathcal{S}}$ ,
- (3)  $\mathcal{S} \subseteq \mathcal{R}^\Delta \iff \mathcal{S}^\Delta \subseteq \mathcal{R}^\Delta \iff \mathcal{E}_{\mathcal{S}} \subseteq \mathcal{E}_{\mathcal{R}} \iff \mathcal{D}_{\mathcal{R}} \subseteq \mathcal{D}_{\mathcal{S}}$ .

**Corollary 3** For any relator  $\mathcal{R}$  on  $X$  to  $Y$ ,

- (1)  $\mathcal{S} = \mathcal{R}^\#$  is the largest relator on  $X$  to  $Y$  such that  $\text{Int}_{\mathcal{S}} \subseteq \text{Int}_{\mathcal{R}}$  ( $\text{Int}_{\mathcal{S}} = \text{Int}_{\mathcal{R}}$ ), or equivalently  $\text{Cl}_{\mathcal{R}} \subseteq \text{Cl}_{\mathcal{S}}$  ( $\text{Cl}_{\mathcal{S}} = \text{Cl}_{\mathcal{R}}$ );
- (2)  $\mathcal{S} = \mathcal{R}^\wedge$  is the largest relator on  $X$  to  $Y$  such that  $\text{int}_{\mathcal{S}} \subseteq \text{int}_{\mathcal{R}}$  ( $\text{int}_{\mathcal{S}} = \text{int}_{\mathcal{R}}$ ), or equivalently  $\text{cl}_{\mathcal{R}} \subseteq \text{cl}_{\mathcal{S}}$  ( $\text{cl}_{\mathcal{S}} = \text{cl}_{\mathcal{R}}$ );
- (3)  $\mathcal{S} = \mathcal{R}^\Delta$  is the largest relator on  $X$  to  $Y$  such that  $\mathcal{E}_{\mathcal{S}} \subseteq \mathcal{E}_{\mathcal{R}}$  ( $\mathcal{E}_{\mathcal{S}} = \mathcal{E}_{\mathcal{R}}$ ), or equivalently  $\mathcal{D}_{\mathcal{R}} \subseteq \mathcal{D}_{\mathcal{S}}$  ( $\mathcal{D}_{\mathcal{S}} = \mathcal{D}_{\mathcal{R}}$ ).

*Remark 16* To prove some similar statements for the operation  $*$ , the structures  $\text{Lim}_{\mathcal{R}}$  and  $\text{Adh}_{\mathcal{R}}$  have to be used [73].

Moreover, for instance, to investigate the structures  $\text{Lb}_{\mathcal{R}}$  and  $\text{Ub}_{\mathcal{R}}$  the compound operation  $\# \circledast = c \# c$  is needed [89].

Concerning the above basic closure operations, we can also prove the following two theorems.

**Theorem 17** For any relator  $\mathcal{R}$  on  $X$  to  $Y$ , we have

- (1)  $\mathcal{R}^\# = \mathcal{R}^{\diamond\#} = \mathcal{R}^{\#\diamond}$  with  $\diamond = *$  and  $\#$ ;
- (2)  $\mathcal{R}^\wedge = \mathcal{R}^{\diamond\wedge} = \mathcal{R}^{\wedge\diamond}$  with  $\diamond = *, \#$  and  $\wedge$ ;
- (3)  $\mathcal{R}^\Delta = \mathcal{R}^{\diamond\Delta} = \mathcal{R}^{\Delta\diamond}$  with  $\diamond = *, \#, \wedge$  and  $\Delta$ .

**Proof** To prove (1), note that, by Remark 14 and the closure properties, we have  $\mathcal{R}^\# \subseteq \mathcal{R}^{\#\#} \subseteq \mathcal{R}^{\#\#} = \mathcal{R}^\#$  and  $\mathcal{R}^\# \subseteq \mathcal{R}^{*\#} \subseteq \mathcal{R}^{\#\#} = \mathcal{R}^\#$ .

**Theorem 18** For any relator  $\mathcal{R}$  on  $X$  to  $Y$ , we have (1)  $\mathcal{R}^{*-1} = \mathcal{R}^{-1*}$ ; (2)  $\mathcal{R}^{\#-1} = \mathcal{R}^{-1\#}$ .

**Proof** To prove (2), note that by Theorems 4 and 16 we have

$$\text{Cl}_{\mathcal{R}^{\#-1}} = \text{Cl}_{\mathcal{R}^{\#}}^{-1} = \text{Cl}_{\mathcal{R}}^{-1} = \text{Cl}_{\mathcal{R}^{-1}}.$$

and thus in particular  $\text{Cl}_{\mathcal{R}^{-1}} \subseteq \text{Cl}_{\mathcal{R}^{\#-1}}$ . Hence, by using Theorem 16, we can infer that  $\mathcal{R}^{\#-1} \subseteq \mathcal{R}^{-1\#}$ .

Now, by writing  $\mathcal{R}^{-1}$  in place of  $\mathcal{R}$ , we can see that  $\mathcal{R}^{-1\#-1} \subseteq \mathcal{R}^{\#}$ , and thus  $\mathcal{R}^{-1\#} \subseteq \mathcal{R}^{\#-1}$ . Therefore, (2) is also true.

*Remark 17* For instance, the elementwise operations  $c$  and  $\infty$  are also inversion compatible. Moreover, the operation  $\partial$  is also inversion compatible.

However, unfortunately, the operations  $\wedge$  and  $\Delta$  are not inversion compatible. Therefore, in addition to Definition 5, we must also have the following

**Definition 6** For any relator  $\mathcal{R}$  on  $X$  to  $Y$ , we define

$$\mathcal{R}^{\vee} = \mathcal{R}^{\wedge-1} \quad \text{and} \quad \mathcal{R}^{\nabla} = \mathcal{R}^{\Delta-1}.$$

*Remark 18* The latter operations have very curious properties. For instance, if  $\mathcal{R}$  is nonvoid, then  $\mathcal{R}^{\vee\wedge} = \{\rho_{\mathcal{R}}\}^{\wedge}$  [54].

Thus, in particular, the relator  $\mathcal{R}^{\vee}$  is topologically simple in the sense that it is topologically equivalent to a singleton relator.

## 8 Some Further Theorems on the Operations $\wedge$ and $\Delta$

A preliminary form of the following theorem was already proved in [73].

**Theorem 19** If  $\mathcal{R}$  is nonvoid relator on  $X$  to  $Y$ , then for any  $B \subseteq Y$  we have: (1)

$$\text{Int}_{\mathcal{R}^{\wedge}}(B) = \mathcal{P}(\text{int}_{\mathcal{R}}(B)); \quad (2) \text{Cl}_{\mathcal{R}^{\wedge}}(B) = \mathcal{P}(\text{cl}_{\mathcal{R}}(B)^c)^c.$$

**Proof** If  $A \in \text{Int}_{\mathcal{R}^{\wedge}}(B)$ , then by Theorems 3 and 16 we have

$$A \subseteq \text{int}_{\mathcal{R}^{\wedge}}(B) = \text{int}_{\mathcal{R}}(B),$$

and thus  $A \in \mathcal{P}(\text{int}_{\mathcal{R}}(B))$ . Therefore,  $\text{Int}_{\mathcal{R}^{\wedge}}(B) \subseteq \mathcal{P}(\text{int}_{\mathcal{R}}(B))$ .

While, if  $A \in \mathcal{P}(\text{int}_{\mathcal{R}}(B))$ , then  $A \subseteq \text{int}_{\mathcal{R}}(B)$ . Therefore, for each  $x \in A$ , there exists  $R_x \in \mathcal{R}$  such that  $R_x(x) \subseteq B$ . Now, by defining

$$S(x) = R_x(x) \quad \text{for all } x \in A \quad \text{and} \quad S(x) = Y \quad \text{for all } x \in A^c,$$

we can easily see that  $S \in \mathcal{R}^\wedge$  such that  $S[A] \subseteq B$ . Therefore, we also have  $A \in \text{Int}_{\mathcal{R}^\wedge}(B)$ . Consequently,  $\mathcal{P}(\text{int}_{\mathcal{R}}(B)) \subseteq \text{Int}_{\mathcal{R}^\wedge}(B)$ , and thus (1) also holds.

Now, by using Theorems 1 and 2, we can also easily see that

$$\text{Cl}_{\mathcal{R}^\wedge}(B) = \text{Int}_{\mathcal{R}^\wedge}(B^c)^c = \mathcal{P}(\text{int}_{\mathcal{R}}(B^c))^c = \mathcal{P}(\text{cl}_{\mathcal{R}}(B)^c)^c.$$

*Remark 19* Thus, for any  $A \subseteq X$ , we have

$$A \in \text{Cl}_{\mathcal{R}^\wedge}(B) \iff A \cap \text{cl}_{\mathcal{R}}(B) \neq \emptyset.$$

From Theorem 19, by using Definition 3, we can immediately derived

**Corollary 4** *If  $\mathcal{R}$  is a nonvoid relator on  $X$ , then (1)  $\tau_{\mathcal{R}^\wedge} = \mathcal{T}_{\mathcal{R}}$ ; (2)  $\varepsilon_{\mathcal{R}^\wedge} = \mathcal{F}_{\mathcal{R}}$ .*

*Remark 20* Hence, since  $\tau_{\mathcal{R}} = \bigcup_{R \in \mathcal{R}} \tau_R = \bigcup_{R \in \mathcal{R}} \mathcal{T}_R$ , we can infer that

$$\mathcal{T}_{\mathcal{R}} = \bigcup_{R \in \mathcal{R}^\wedge} \mathcal{T}_R.$$

Unfortunately, in contrast to the structures  $\text{Int}$ ,  $\text{int}$ ,  $\mathcal{E}$  and  $\tau$ , the increasing structure  $\mathcal{T}$  is already not union-preserving.

*Example 1* If  $\text{card}(X) > 2$ ,  $x_1 \in X$  and  $x_2 \in X \setminus \{x_1\}$ , and

$$R_i = \{x_i\}^2 \cup (X \setminus \{x_i\})^2$$

for all  $i = 1, 2$ , then  $\mathcal{R} = \{R_1, R_2\}$  is an equivalence relator on  $X$  such that

$$\{x_1, x_2\} \in \mathcal{T}_{\mathcal{R}} \setminus (\mathcal{T}_{R_1} \cup \mathcal{T}_{R_2}), \quad \text{and thus} \quad \mathcal{T}_{\mathcal{R}} \not\subseteq \mathcal{T}_{R_1} \cup \mathcal{T}_{R_2}.$$

From Corollary 4, by using Theorem 17, we can immediately derive

**Corollary 5** *If  $\mathcal{R}$  is a nonvoid relator on  $X$ , then (1)  $\tau_{\mathcal{R}^\Delta} = \mathcal{T}_{\mathcal{R}^\Delta}$ ; (2)  $\varepsilon_{\mathcal{R}^\Delta} = \mathcal{F}_{\mathcal{R}^\Delta}$ .*

Concerning the operation  $\Delta$ , we can also prove the following

**Theorem 20** *If  $\mathcal{R}$  is a nonvoid relator on  $X$  to  $Y$ , then for any  $B \subseteq Y$  we have:*

- (1)  $\text{Int}_{\mathcal{R}^\Delta}(B) = \{\emptyset\}$  if  $B \notin \mathcal{E}_{\mathcal{R}}$  and  $\text{Int}_{\mathcal{R}^\Delta}(B) = \mathcal{P}(X)$  if  $B \in \mathcal{E}_{\mathcal{R}}$ ;
- (2)  $\text{Cl}_{\mathcal{R}^\Delta}(B) = \emptyset$  if  $B \notin \mathcal{D}_{\mathcal{R}}$  and  $\text{Cl}_{\mathcal{R}^\Delta}(B) = \mathcal{P}(X) \setminus \{\emptyset\}$  if  $B \in \mathcal{D}_{\mathcal{R}}$ .

**Proof** If  $A \in \text{Int}_{\mathcal{R}^\Delta}(B)$ , then there exists  $S \in \mathcal{R}^\Delta$  such that  $S[A] \subseteq B$ . Therefore, if  $A \neq \emptyset$ , then there exists  $x \in X$  such that  $S(x) \subseteq B$ . Hence, since  $S(x) \in \mathcal{E}_{\mathcal{R}}$ , it follows that  $B \in \mathcal{E}_{\mathcal{R}}$ . Therefore, the first part of (1) is true.

To prove the second part of (1), it is enough to note only that if  $B \in \mathcal{E}_{\mathcal{R}}$ , then  $R = X \times B \in \mathcal{R}^\Delta$  such that  $R[A] \subseteq B$ , and thus  $A \in \text{Int}_{\mathcal{R}^\Delta}(B)$  for all  $A \subseteq X$ .

Assertion (2) can again be derived from (1) by using Theorem 1.

From this theorem, by Definition 1, it is clear that in particular we also have

**Corollary 6** *If  $\mathcal{R}$  is nonvoid relator on  $X$  to  $Y$ , then for any  $B \subseteq Y$ :*

- (1)  $\text{cl}_{\mathcal{R}^\Delta}(B) = \emptyset$  if  $B \notin \mathcal{D}_{\mathcal{R}}$  and  $\text{cl}_{\mathcal{R}^\Delta}(B) = X$  if  $B \in \mathcal{D}_{\mathcal{R}}$ ;
- (2)  $\text{int}_{\mathcal{R}^\Delta}(B) = \emptyset$  if  $B \notin \mathcal{E}_{\mathcal{R}}$  and  $\text{int}_{\mathcal{R}^\Delta}(B) = X$  if  $B \in \mathcal{E}_{\mathcal{R}}$ .

Hence, by using Definitions 1 and 3, we can immediately derive

**Corollary 7** *If  $\mathcal{R}$  is a relator on  $X$ , then (1)  $\mathcal{I}_{\mathcal{R}^\Delta} = \mathcal{E}_{\mathcal{R}} \cup \{\emptyset\}$ ; (2)  $\mathcal{F}_{\mathcal{R}^\Delta} = (\mathcal{P}(X) \setminus \mathcal{D}_{\mathcal{R}}) \cup \{X\}$ .*

*Remark 21* Note that if in particular  $\mathcal{R} = \emptyset$ , then  $\mathcal{E}_{\mathcal{R}} = \emptyset$ . Moreover,  $\mathcal{R}^\Delta = \emptyset$  if  $X \neq \emptyset$ , and  $\mathcal{R}^\Delta = \{\emptyset\}$  if  $X = \emptyset$ . Therefore,  $\mathcal{I}_{\mathcal{R}^\Delta} = \{\emptyset\}$ , and thus (1) is still true.

Now, since  $\emptyset \notin \mathcal{E}_{\mathcal{R}}$  if  $\mathcal{R}$  is non-partial, we can also state

**Corollary 8** *If  $\mathcal{R}$  is a non-partial relator on  $X$ , then (1)  $\mathcal{E}_{\mathcal{R}} = \mathcal{I}_{\mathcal{R}^\Delta} \setminus \{\emptyset\}$ , (2)  $\mathcal{D}_{\mathcal{R}} = (\mathcal{P}(X) \setminus \mathcal{F}_{\mathcal{R}^\Delta}) \cup \{X\}$ .*

## 9 Projection Operations for Relators

By using the basic properties of the operation  $\infty$ , in addition to a particular case of Theorem 16, we can also prove the following

**Theorem 21**  *$\infty$  is a closure operation for relations on  $X$  such that, for any two relations  $R$  and  $S$  on  $X$ , we have*

$$S \subseteq R^\infty \iff S^\infty \subseteq R^\infty \iff \tau_R \subseteq \tau_S \iff \varepsilon_R \subseteq \varepsilon_S;$$

**Proof** To prove that  $\tau_R \subseteq \tau_S \iff S \subseteq R^\infty$ , note that if  $x \in X$ , then because of the inclusion  $R \subseteq R^\infty$  and the transitivity of  $R^\infty$  we have

$$R[R^\infty(x)] \subseteq R^\infty[R^\infty(x)] = (R^\infty \circ R^\infty)(x) \subseteq R^\infty(x).$$

Thus, by the definition of  $\tau_{\mathcal{R}}$ , we have  $R^\infty(x) \in \tau_{\mathcal{R}}$ . Now, if  $\tau_R \subseteq \tau_S$  holds, then we can see that  $R^\infty(x) \in \tau_S$ , and thus  $S[R^\infty(x)] \subseteq R^\infty(x)$ . Hence, by using the reflexivity of  $R^\infty$ , we can already infer that  $S(x) \subseteq R^\infty(x)$ . Therefore,  $S \subseteq R^\infty$  also holds.

While, if  $A \in \tau_{\mathcal{R}}$ , then by the definition of  $\tau_{\mathcal{R}}$  we have  $R[A] \subseteq A$ . Hence, by induction, we can see that  $R^n[A] \subseteq A$  for all  $n \in \mathbb{N}$ . Now, since  $R^0[A] = \Delta_X[A] = A$  also holds, we can already state that

$$R^\infty[A] = \left( \bigcup_{n=0}^\infty R^n \right)[A] = \bigcup_{n=0}^\infty R^n[A] \subseteq \bigcup_{n=0}^\infty A = A.$$

Therefore, if  $S \subseteq R^\infty$  holds, then we have  $S[A] \subseteq R^\infty[A] \subseteq A$ , and thus  $A \in \tau_S$  also holds.

Now, analogously to Corollary 3, we can also state

**Corollary 9** *For any relation  $R$  on  $X$ ,  $S = R^\infty$  is the largest relation on  $X$  such that  $\tau_R \subseteq \tau_S$  ( $\tau_R = \tau_S$ ), or equivalently  $\varepsilon_R \subseteq \varepsilon_S$  ( $\varepsilon_R = \varepsilon_S$ ).*

*Remark 22* Preliminary forms of the above theorem and its corollary were first proved by Mala [53].

Moreover, he also proved that  $R^\infty(x) = \bigcap \{ A \in \tau_R : x \in A \}$  for all  $x \in X$ , and thus  $R^\infty = \bigcap \{ R_A : A \in \tau_R \}$ .

By using Theorem 21, as an analogue of Theorem 16, we can also prove

**Theorem 22**  *$\# \partial$  is a closure operation for relators on  $X$  such that, for any two relators  $\mathcal{R}$  and  $\mathcal{S}$  on  $X$ , we have*

$$\mathcal{S} \subseteq \mathcal{R}^{\# \partial} \iff \mathcal{S}^{\# \partial} \subseteq \mathcal{R}^{\# \partial} \iff \tau_{\mathcal{S}} \subseteq \tau_{\mathcal{R}} \iff \varepsilon_{\mathcal{S}} \subseteq \varepsilon_{\mathcal{R}};$$

Thus, analogously to Corollary 3, we can also state

**Corollary 10** *For any relator  $\mathcal{R}$  on  $X$ ,  $\mathcal{S} = \mathcal{R}^{\# \partial}$  is the largest relator on  $X$  such that  $\tau_{\mathcal{S}} \subseteq \tau_{\mathcal{R}}$  ( $\tau_{\mathcal{S}} = \tau_{\mathcal{R}}$ ), or equivalently  $\varepsilon_{\mathcal{S}} \subseteq \varepsilon_{\mathcal{R}}$  ( $\varepsilon_{\mathcal{S}} = \varepsilon_{\mathcal{R}}$ ).*

By using the Galois property of the operations  $\infty$  and  $\partial$ , Theorem 22 can be reformulated in the following more convenient form.

**Theorem 23**  *$\# \infty$  is a projection operation for relators on  $X$  such that, for any two relators  $\mathcal{R}$  and  $\mathcal{S}$  on  $X$ , we have*

$$\mathcal{S}^\infty \subseteq \mathcal{R}^\# \iff \mathcal{S}^{\# \infty} \subseteq \mathcal{R}^{\# \infty} \iff \tau_{\mathcal{S}} \subseteq \tau_{\mathcal{R}} \iff \varepsilon_{\mathcal{S}} \subseteq \varepsilon_{\mathcal{R}}.$$

*Remark 23* Moreover, it can be shown that the inclusions  $\mathcal{S}^\infty \subseteq \mathcal{R}^\#$ ,  $\mathcal{S}^{\# \infty} \subseteq \mathcal{R}^\#$  and  $\mathcal{S}^\infty \subseteq \mathcal{R}^{\# \infty}$  are also equivalent.

Now, analogously to our former corollaries, we can also state

**Corollary 11** *For any relator  $\mathcal{R}$  on  $X$ ,  $\mathcal{S} = \mathcal{R}^{\# \infty}$  is the largest preorder relator on  $X$  such that  $\tau_{\mathcal{S}} \subseteq \tau_{\mathcal{R}}$  ( $\tau_{\mathcal{S}} = \tau_{\mathcal{R}}$ ), or equivalently  $\varepsilon_{\mathcal{S}} \subseteq \varepsilon_{\mathcal{R}}$  ( $\varepsilon_{\mathcal{S}} = \varepsilon_{\mathcal{R}}$ ).*

*Remark 24* The advantage of the projection operation  $\# \infty$  over the closure operation  $\# \partial$  lies mainly in the fact that, in contrast to  $\# \partial$ , it is *stable* in the sense  $\{X^2\}^{\# \infty} = \{X^2\}$ .

Since the structure  $\mathcal{T}$  is not union-preserving, by using some parts of the theory of Pataki connections [64, 69, 95], we can only prove the following

**Theorem 24**  $\wedge \partial$  is a preclosure operation for relators such that, for any two relators  $\mathcal{R}$  and  $\mathcal{S}$  on  $X$ , we have

$$\mathcal{I}_{\mathcal{S}} \subseteq \mathcal{I}_{\mathcal{R}} \iff \mathcal{F}_{\mathcal{S}} \subseteq \mathcal{F}_{\mathcal{R}} \implies \mathcal{S}^{\wedge} \subseteq \mathcal{R}^{\wedge \partial} \implies \mathcal{S}^{\wedge \partial} \subseteq \mathcal{R}^{\wedge \partial}.$$

*Remark 25* If  $\text{card}(X) > 2$ , then by using the equivalence relation  $\mathcal{R} = \{X^2\}$  Mala [53, Example 5.3] proved that there does not exist a largest relator  $\mathcal{S}$  on  $X$  such that  $\mathcal{I}_{\mathcal{R}} = \mathcal{I}_{\mathcal{S}}$ .

Moreover, Pataki [64, Example 7.2] proved that  $\mathcal{I}_{\mathcal{R}^{\wedge \partial}} \not\subseteq \mathcal{I}_{\mathcal{R}}$  and  $\wedge \partial$  is not idempotent. (Actually, it can be proved that  $\mathcal{R}^{\wedge \partial \wedge} \not\subseteq \mathcal{R}^{\wedge \partial}$  also holds [84, Example 10.11].)

Fortunately, as an analogue of Theorem 23, we can also prove

**Theorem 25**  $\wedge \infty$  is a projection operation for relators on  $X$  such that, for any two nonvoid relators  $\mathcal{R}$  and  $\mathcal{S}$  on  $X$ , we have

$$\mathcal{S}^{\wedge \infty} \subseteq \mathcal{R}^{\wedge} \iff \mathcal{S}^{\wedge \infty} \subseteq \mathcal{R}^{\wedge \infty} \iff \mathcal{I}_{\mathcal{S}} \subseteq \mathcal{I}_{\mathcal{R}} \iff \mathcal{F}_{\mathcal{S}} \subseteq \mathcal{F}_{\mathcal{R}}.$$

Thus, in particular, we can also state

**Corollary 12** For any nonvoid relator  $\mathcal{R}$  on  $X$ ,  $\mathcal{S} = \mathcal{R}^{\wedge \infty}$  is the largest preorder relator on  $X$  such that  $\mathcal{I}_{\mathcal{S}} \subseteq \mathcal{I}_{\mathcal{R}}$  ( $\mathcal{I}_{\mathcal{S}} = \mathcal{I}_{\mathcal{R}}$ ), or equivalently  $\mathcal{F}_{\mathcal{S}} \subseteq \mathcal{F}_{\mathcal{R}}$  ( $\mathcal{F}_{\mathcal{S}} = \mathcal{F}_{\mathcal{R}}$ ).

*Remark 26* In the light of the several disadvantages of the structure  $\mathcal{I}$ , it is rather curious that most of the works in general topology and abstract analysis have been based on open sets suggested by Tietze [98] and Alexandroff [2], and standardized by Bourbaki [6] and Kelley [41]. (See Thron [97, p. 18].)

Moreover, it also a very striking fact that, despite the results of Davis [20], Pervin [66], Hunsaker and Lindgren [35] and the second author [76, 85], generalized proximities and closures, minimal structures, generalized topologies and stacks (ascending systems) are still intensively investigated by a great number of mathematicians without using generalized uniformities.

## 10 Reflexive, Non-partial and Non-degenerated Relators

**Definition 7** A relator  $\mathcal{R}$  on  $X$  is called *reflexive* if each member  $R$  of  $\mathcal{R}$  is a reflexive relation on  $X$ .

*Remark 27* Thus, the following assertions are equivalent:

- (1)  $\mathcal{R}$  is reflexive;
- (2)  $x \in R(x)$  for all  $x \in X$  and  $R \in \mathcal{R}$ ;
- (3)  $A \subseteq R[A]$  for all  $A \subseteq X$  and  $R \in \mathcal{R}$ .

The importance of reflexive relators is also apparent from the following two obvious theorems.

**Theorem 26** For a relator  $\mathcal{R}$  on  $X$ , the following assertions are equivalent: (1)

- $\rho_{\mathcal{R}}$  is reflexive;
- (2)  $\mathcal{R}$  is reflexive;
- (3)  $A \subseteq \text{cl}_{\mathcal{R}}(A)$  ( $\text{int}_{\mathcal{R}}(A) \subseteq A$ ) for all  $A \subseteq X$ .

**Proof** To prove the equivalence of (1) and (2), recall that  $\rho_{\mathcal{R}} = (\bigcap \mathcal{R})^{-1}$ .

**Remark 28** Thus, the relator  $\mathcal{R}$  is reflexive if and only if  $A^\circ \subseteq A$  ( $A \subseteq A^-$ ) for all  $A \subseteq X$ .

Therefore, if  $\mathcal{R}$  is a reflexive relator on  $X$ , then for any  $A \subseteq X$  we have  $A \in \mathcal{T}_{\mathcal{R}}$  ( $A \in \mathcal{F}_{\mathcal{R}}$ ) if and only if  $A^\circ = A$  ( $A^- = A$ ).

**Theorem 27** For a relator  $\mathcal{R}$  on  $X$ , the following assertions are equivalent:

- (1)  $\mathcal{R}$  is reflexive;
- (2)  $A \in \text{Int}_{\mathcal{R}}(B)$  implies  $A \subseteq B$  for all  $A, B \subseteq X$ ;
- (3)  $A \cap B \neq \emptyset$  implies  $A \in \text{Cl}_{\mathcal{R}}(B)$  for all  $A, B \subseteq X$ .

**Remark 29** In addition to the above two theorems, it is also worth mentioning that if  $\mathcal{R}$  is a reflexive relator on  $X$ , then

- (1)  $\text{Int}_{\mathcal{R}}$  is a transitive relation on  $\mathcal{P}(X)$ ;
- (2)  $B \in \text{Cl}_{\mathcal{R}}(A)$  implies  $\mathcal{P}(X) = \text{Cl}_{\mathcal{R}}(A)^c \cup \text{Cl}_{\mathcal{R}}^{-1}(B)$ ;
- (3)  $\text{int}_{\mathcal{R}}(\text{bor}_{\mathcal{R}}(A)) = \emptyset$  and  $\text{int}_{\mathcal{R}}(\text{res}_{\mathcal{R}}(A)) = \emptyset$  for all  $A \subseteq X$ .

Thus, for instance, for any  $A \subseteq X$  we have  $\text{res}_{\mathcal{R}}(A) \in \mathcal{T}_{\mathcal{R}}$  if and only if  $A \in \mathcal{F}_{\mathcal{R}}$ .

Analogously to Definition 7, we may also naturally have the following

**Definition 8** A relator  $\mathcal{R}$  on  $X$  to  $Y$  is called *non-partial* if each member  $R$  of  $\mathcal{R}$  is a non-partial relation on  $X$  to  $Y$ .

**Remark 30** Thus, the following assertions are equivalent:

- (1)  $\mathcal{R}$  is non-partial;
- (2)  $R^{-1}[Y] = X$  for all  $R \in \mathcal{R}$ ;
- (3)  $R(x) \neq \emptyset$  for all  $x \in X$  and  $R \in \mathcal{R}$ .

The importance of non-partial relators is apparent from the following

**Theorem 28** For a relator  $\mathcal{R}$  on  $X$  to  $Y$ , the following assertions are equivalent:

- (1)  $\mathcal{R}$  is non-partial;
- (2)  $\emptyset \notin \mathcal{E}_{\mathcal{R}}$ ; (3)  $\mathcal{D}_{\mathcal{R}} \neq \emptyset$ ; (4)  $Y \in \mathcal{D}_{\mathcal{R}}$ ; (5)  $\mathcal{E}_{\mathcal{R}} \neq \mathcal{P}(Y)$ .

Sometimes, we also need the following localized form of Definition 8.

**Definition 9** A relator  $\mathcal{R}$  on  $X$  is called *locally non-partial* if for each  $x \in X$  there exists  $R \in \mathcal{R}$  such that for any  $y \in R(x)$  and  $S \in \mathcal{R}$  we have  $S(y) \neq \emptyset$ .



*Remark 31* Thus, if either  $X = \emptyset$  or  $\mathcal{R}$  is nonvoid and non-partial, then  $\mathcal{R}$  is locally non-partial.

Moreover, by using the corresponding definitions, we can also easily prove

**Theorem 29** For a relator  $\mathcal{R}$  on  $X$ , the following assertions are equivalent: (1)  $\mathcal{R}$  is locally non-partial; (2)  $X = \text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(X))$ .

*Proof* To prove the implication (1)  $\implies$  (2), note that if (1) holds, then for each  $x \in X$  there exists  $R \in \mathcal{R}$  such that for any  $y \in R(x)$  and for any  $S \in \mathcal{R}$  we have  $S(y) \cap X = S(y) \neq \emptyset$ , and thus  $y \in \text{cl}_{\mathcal{R}}(X)$ .

Therefore, for each  $x \in X$  there exists  $R \in \mathcal{R}$  such that  $R(x) \subseteq \text{cl}_{\mathcal{R}}(X)$ , and thus  $x \in \text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(X))$ . Hence, we can already see that  $X \subseteq \text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(X))$ , and thus (2) also holds. (Therefore, by a former notation,  $X \in \mathcal{I}_{\mathcal{R}}^r$ .)

In addition to Definition 8, it is also worth introducing the following

**Definition 10** A relator  $\mathcal{R}$  on  $X$  to  $Y$  is called *non-degerated* if both  $X \neq \emptyset$  and  $\mathcal{R} \neq \emptyset$ .

Thus, analogously to Theorem 28, we can also easily establish the following

**Theorem 30** For a relator  $\mathcal{R}$  on  $X$  to  $Y$ , the following assertions are equivalent: (1)  $\mathcal{R}$  is non-degenerated;

(2)  $\emptyset \notin \mathcal{D}_{\mathcal{R}}$ ; (3)  $\mathcal{E}_{\mathcal{R}} \neq \emptyset$ ; (4)  $Y \in \mathcal{E}_{\mathcal{R}}$ ; (5)  $\mathcal{D}_{\mathcal{R}} \neq \mathcal{P}(Y)$ .

*Remark 32* In addition to Theorems 28 and 30, it is also worth mentioning that if a relator  $\mathcal{R}$  on  $X$  to  $Y$  is paratopologically simple in the sense that  $\mathcal{E}_{\mathcal{R}} = \mathcal{E}_S$  for some relation  $S$  on  $X$  to  $Y$ , then the stack  $\mathcal{E}_{\mathcal{R}}$  has a base  $\mathcal{B}$  with  $\text{card}(\mathcal{B}) \leq \text{card}(X)$ . (See [63, Theorem 5.9] of Pataki.)

The existence of a non-paratopologically simple (actually finite equivalence) relator, proved first by Pataki [63, Example 5.11], shows that in our definitions of the relations  $\text{Lim}_{\mathcal{R}}$  and  $\text{Adh}_{\mathcal{R}}$  we cannot restrict ourselves to functions of gosets (generalized ordered sets) without some loss of generality.

## 11 Topological and Quasi-Topological Relators

The following improvement of [74, Definition 2.1] was first considered in [75].

**Definition 11** A relator  $\mathcal{R}$  on  $X$  is called:

- (1) *quasi-topological* if  $x \in \text{int}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(R(x)))$  for all  $x \in X$  and  $R \in \mathcal{R}$ ;
- (2) *topological* if for any  $x \in X$  and  $R \in \mathcal{R}$  there exists  $V \in \mathcal{I}_{\mathcal{R}}$  such that  $x \in V \subseteq R(x)$ .

The appropriateness of these definitions is already quite obvious from the following four theorems.

**Theorem 31** For a relator  $\mathcal{R}$  on  $X$ , the following assertions are equivalent:

- (1)  $\mathcal{R}$  is quasi-topological;
- (2)  $\text{int}_{\mathcal{R}}(R(x)) \in \mathcal{T}_{\mathcal{R}}$  for all  $x \in X$  and  $R \in \mathcal{R}$ .
- (3)  $\text{cl}_{\mathcal{R}}(A) \in \mathcal{F}_{\mathcal{R}}$  ( $\text{int}_{\mathcal{R}}(A) \in \mathcal{T}_{\mathcal{R}}$ ) for all  $A \subseteq X$ .

*Remark 33* Hence, by Definition 3, we can see that the relator  $\mathcal{R}$  is quasi-topological if and only if  $A^\circ \subseteq A^{\circ\circ}$  ( $A^{--} \subseteq A^-$ ) for all  $A \subseteq X$ .

**Theorem 32** For a relator  $\mathcal{R}$  on  $X$ , the following assertions are equivalent:

- (1)  $\mathcal{R}$  is topological;
- (2)  $\mathcal{R}$  is reflexive and quasi-topological.

*Remark 34* By Theorem 31, the relator  $\mathcal{R}$  may be called weakly (strongly) quasi-topological if  $\rho_{\mathcal{R}}(x) \in \mathcal{F}_{\mathcal{R}}$  ( $R(x) \in \mathcal{T}_{\mathcal{R}}$ ) for all  $x \in X$  and  $R \in \mathcal{R}$ .

Moreover, by Theorem 32, the relator  $\mathcal{R}$  may be called weakly (strongly) topological if it is reflexive and weakly (strongly) quasi-topological.

The following theorem shows that in a topological relator space  $X(\mathcal{R})$ , the relation  $\text{int}_{\mathcal{R}}$  and the family  $\mathcal{T}_{\mathcal{R}}$  are also equivalent tools.

**Theorem 33** For a relator  $\mathcal{R}$  on  $X$ , the following assertions are equivalent:

- (1)  $\mathcal{R}$  is topological;
- (2)  $\text{int}_{\mathcal{R}}(A) = \bigcup \mathcal{T}_{\mathcal{R}} \cap \mathcal{P}(A)$  for all  $A \subseteq X$ ;
- (3)  $\text{cl}_{\mathcal{R}}(A) = \bigcap \mathcal{F}_{\mathcal{R}} \cap \mathcal{P}^{-1}(A)$  for all  $A \subseteq X$ .

Now, as an immediate consequence of Theorems 14 and 33, we can also state

**Corollary 13** If  $\mathcal{R}$  is topological relator on  $X$ , then for any  $A \subset X$ , we have

- (1)  $A \in \mathcal{E}_{\mathcal{R}}$  if and only if there exists  $V \in \mathcal{T}_{\mathcal{R}} \setminus \{\emptyset\}$  such that  $V \subseteq A$ ;
- (2)  $A \in \mathcal{D}_{\mathcal{R}}$  if and only if for all  $W \in \mathcal{F}_{\mathcal{R}} \setminus \{X\}$  we have  $A \setminus W \neq \emptyset$ .

However, it is now more important to note that we can also prove the following

**Theorem 34** For a relator  $\mathcal{R}$  on  $X$ , the following assertions are equivalent:

- (1)  $\mathcal{R}$  is topological;
- (2)  $\mathcal{R}$  is topologically equivalent to  $\mathcal{R}^{\wedge\infty}$ ;
- (3)  $\mathcal{R}$  is topologically equivalent to a preorder relator on  $X$ .

**Proof** To prove the implication (1)  $\implies$  (3), note that if (1) holds, then by Definition 11, for any  $x \in X$  and  $R \in \mathcal{R}$ , there exists  $V \in \mathcal{T}_{\mathcal{R}}$  such that  $x \in V \subseteq R(x)$ . Thus, by using the Pervin preorder relator

$$\mathcal{S} = \mathcal{R}_{\mathcal{T}_{\mathcal{R}}} = \{R_V : V \in \mathcal{T}_{\mathcal{R}}\}, \quad \text{where} \quad R_V = V^2 \cup V^c \times X,$$

we can show that  $\text{int}_{\mathcal{R}}(A) = \text{int}_{\mathcal{S}}(A)$  for all  $A \subseteq X$ , and thus (3) also holds.

For this, we have to note that

$$R_V(x) = V \quad \text{if } x \in V \quad \text{and} \quad R_V(x) = X \quad \text{if } x \in V^c.$$

In addition to Theorem 31, it is also worth proving the following

**Theorem 35** *For a relator  $\mathcal{R}$  on  $X$ , the following assertions are equivalent: (1)  $\mathcal{R}$  is quasi-topological; (2)  $\mathcal{R} \subseteq (\mathcal{R} \wedge \circ \mathcal{R})^\wedge$ ; (3)  $\mathcal{R}^\wedge \subseteq (\mathcal{R} \wedge \circ \mathcal{R}^\wedge)^\wedge$ .*

*Remark 35* By [74], a relator  $\mathcal{R}$  on  $X$  may be naturally called *topologically transitive* if, for each  $x \in X$  and  $R \in \mathcal{R}$  there exist  $S, T \in \mathcal{R}$  such that  $T[S(x)] \subseteq R(x)$ .

This property can be easily reformulated in the more concise form that  $\mathcal{R} \subseteq (\mathcal{R} \circ \mathcal{R})^\wedge$ . Thus, the equivalence (1) and (3) can be expressed by saying that  $\mathcal{R}$  is quasi-topological if and only if  $\mathcal{R}^\wedge$  is topologically transitive.

## 12 Proximal and Quasi-Proximal Relators

Analogously to Definition 11, we may also naturally have the following

**Definition 12** A relator  $\mathcal{R}$  on  $X$  is called

- (1) *quasi-proximal* if  $A \in \text{Int}_{\mathcal{R}} [\tau_{\mathcal{R}} \cap \text{Int}_{\mathcal{R}} (R[A])]$  for all  $A \subseteq X$  and  $R \in \mathcal{R}$ ;
- (2) *proximal* if for any  $A \subseteq X$  and  $R \in \mathcal{R}$  there exists  $V \in \tau_{\mathcal{R}}$  such that  $A \subseteq V \subseteq R[A]$ .

*Remark 36* Thus, the relator  $\mathcal{R}$  is quasi-proximal if and only if, for any  $A \subseteq X$  and  $R \in \mathcal{R}$ , there exists  $V \in \tau_{\mathcal{R}}$  such that  $A \in \text{Int}_{\mathcal{R}}(V)$  and  $V \in \text{Int}_{\mathcal{R}}(R[A])$ .

Now, by using the corresponding definitions, we can also easily prove the following analogues of Theorems 32 and 33.

**Theorem 36** *For a relator  $\mathcal{R}$  on  $X$ , the following assertions are equivalent: (1)  $\mathcal{R}$  is proximal; (2)  $\mathcal{R}$  is reflexive and quasi-proximal.*

**Proof** To prove the implication (1)  $\implies$  (2), note that if (1) holds, then for any  $A \subseteq X$  and  $R \in \mathcal{R}$ , there exists  $V \in \tau_{\mathcal{R}}$  such that  $A \subseteq V \subseteq R[A]$ . Hence, by taking  $A = \{x\}$  for  $x \in X$ , we can see that  $\mathcal{R}$  is reflexive.

Moreover, since  $V \in \tau_{\mathcal{R}}$ , we can also note that  $V \in \text{Int}_{\mathcal{R}}(V)$ . Hence, by using that  $A \subseteq V$  and  $V \subseteq R[A]$ , we can already infer that  $A \in \text{Int}_{\mathcal{R}}(V)$  and  $V \in \text{Int}_{\mathcal{R}}(R[A])$ . Therefore, by Remark 36,  $\mathcal{R}$  is quasi-proximal.

*Remark 37* Note that if  $\mathcal{R}$  is only a *weakly proximal* relator on  $X$  in the sense that, for any  $x \in X$  and  $R \in \mathcal{R}$ , there exists  $V \in \tau_{\mathcal{R}}$  such that  $x \in V \subseteq R(x)$ , then because of  $\tau_{\mathcal{R}} \subseteq \mathcal{I}_{\mathcal{R}}$  we can already state that  $\mathcal{R}$  is topological.

**Theorem 37** For a relator  $\mathcal{R}$  on  $X$ , the following assertions are equivalent:

- (1)  $\mathcal{R}$  is proximal;
- (2)  $\text{Int}_{\mathcal{R}}(A) = \mathcal{P}[\tau_{\mathcal{R}} \cap \mathcal{P}(A)]$  for all  $A \subseteq X$ ;
- (3)  $\text{Cl}_{\mathcal{R}}(A) = \bigcap \{ \mathcal{P}(W^c)^c : W \in \tau_{\mathcal{R}} \cap \mathcal{P}^{-1}(A) \}$  for all  $A \subseteq X$ .

**Proof** Note that if  $A \subseteq X$  and  $B \in \mathcal{P}[\tau_{\mathcal{R}} \cap \mathcal{P}(A)]$ , then there exists  $V \in \tau_{\mathcal{R}}$  such that  $B \in \mathcal{P}(V)$  and  $V \in \mathcal{P}(A)$ , and thus  $B \subseteq V \subseteq A$ . Hence, by using that  $V \in \text{Int}_{\mathcal{R}}(V)$  we can already infer that  $B \in \text{Int}_{\mathcal{R}}(A)$ . Thus, the inclusion  $\mathcal{P}[\tau_{\mathcal{R}} \cap \mathcal{P}(A)] \subseteq \text{Int}_{\mathcal{R}}(A)$  is always true.

Therefore, to obtain (1), it is enough to assume only the converse inclusion. For this, note that if  $A \subseteq X$  and  $R \in \mathcal{R}$ , then because of  $R[A] \subseteq R[A]$ , we always have  $A \in \text{Int}_{\mathcal{R}}(R[A])$ . Therefore, if  $\text{Int}_{\mathcal{R}}(R[A]) \subseteq \mathcal{P}[\tau_{\mathcal{R}} \cap \mathcal{P}(R[A])]$ , then we also have  $A \in \mathcal{P}[\tau_{\mathcal{R}} \cap \mathcal{P}(R[A])]$ . Thus, there exists  $V \in \tau_{\mathcal{R}}$  such that  $A \in \mathcal{P}(V)$  and  $V \in \mathcal{P}(R[A])$ , and thus  $A \subseteq V \subseteq R[A]$ .

**Remark 38** Note that  $\mathcal{P}(A) = \text{Int}_{\Delta_X}(A)$  for all  $A \subseteq X$ . Therefore, instead of (2) we may write that  $\text{Int}_{\mathcal{R}}(A) = \text{Int}_{\Delta_X}[\tau_{\mathcal{R}} \cap \text{Int}_{\Delta_X}(A)]$  for all  $A \subseteq X$ .

However, it is now more important to note that we also have the following

**Theorem 38** For a relator  $\mathcal{R}$  on  $X$ , the following assertions are equivalent:

- (1)  $\mathcal{R}$  is proximal;
- (2)  $\mathcal{R}$  is proximally equivalent to  $\mathcal{R}^\infty$  or  $\mathcal{R}^{\#\infty}$ ;
- (3)  $\mathcal{R}$  is proximally equivalent to a preorder relator on  $X$ .

In principle, each theorem on topological and quasi-topological relators can be immediately derived from a corresponding theorem on proximal and quasi-proximal relators by using the following two theorems.

**Theorem 39** For a relator  $\mathcal{R}$  on  $X$ , the following assertions are equivalent: (1)  $\mathcal{R}$  is quasi-topological; (2)  $\mathcal{R}^\wedge$  is quasi-proximal.

**Proof** To prove the implication (1)  $\implies$  (2), assume that (1) holds, and moreover  $A \subseteq X$  and  $S \in \mathcal{R}^\wedge$ . Define  $V = \text{int}_{\mathcal{R}}(S[A])$ . Then, if  $\mathcal{R} \neq \emptyset$ , by Theorem 31 and Corollary 4, we have  $V \in \mathcal{I}_{\mathcal{R}} = \tau_{\mathcal{R}^\wedge}$ . Moreover, since  $V \subseteq \text{int}_{\mathcal{R}}(S[A])$ , by Theorem 19 we also have  $V \in \text{Int}_{\mathcal{R}^\wedge}(S[A])$ . Therefore,  $V \in \tau_{\mathcal{R}^\wedge} \cap \text{Int}_{\mathcal{R}^\wedge}(S[A])$ .

Furthermore, since  $S[A] \subseteq S[A]$ , we can also note that  $A \in \text{Int}_{\mathcal{R}^\wedge}(S[A])$ . Hence, by Theorem 19, we can infer that  $A \subseteq \text{int}_{\mathcal{R}}(S[A]) = V$ . Moreover, since  $V \in \tau_{\mathcal{R}^\wedge}$ , we can also note that  $V \in \text{Int}_{\mathcal{R}^\wedge}(V)$ . Hence, since  $A \subseteq V$ , we can infer that  $A \in \text{Int}_{\mathcal{R}^\wedge}(V)$ . Therefore, since  $V \in \tau_{\mathcal{R}^\wedge} \cap \text{Int}_{\mathcal{R}^\wedge}(S[A])$ , we also have  $A \in \text{Int}_{\mathcal{R}^\wedge}[\tau_{\mathcal{R}^\wedge} \cap \text{Int}_{\mathcal{R}^\wedge}(S[A])]$ .

This shows that (1) implies (2) whenever  $\mathcal{R} \neq \emptyset$ . However, if  $\mathcal{R} = \emptyset$ , then it can be easily seen that  $\mathcal{R}$  is topological and  $\mathcal{R}^\wedge$  is proximal.

*Remark 39* If assertion (2) holds, then  $\mathcal{R}^\wedge$  is *semi-proximal* in the sense that  $A \in \text{Int}_{\mathcal{R}^\wedge} [\text{Int}_{\mathcal{R}^\wedge} (S[A])]$  for all  $A \subseteq X$  and  $S \in \mathcal{R}^\wedge$ .

Moreover, if in particular  $\{x\} \in \text{Int}_{\mathcal{R}^\wedge} [\text{Int}_{\mathcal{R}^\wedge} (R(x))]$  for all  $x \in X$  and  $R \in \mathcal{R}$ , then we can already prove that assertion (1) also holds.

From Theorem 39, by using Theorems 32 and 36, we can immediately derive

**Theorem 40** *For a relator  $\mathcal{R}$  on  $X$ , the following assertions are equivalent: (1)  $\mathcal{R}$  is topological, (2)  $\mathcal{R}^\wedge$  is proximal.*

*Remark 40* By the corresponding definitions, it is clear that the relator  $\mathcal{R}^\wedge$  is reflexive if and only if  $\mathcal{R}$  is reflexive.

However, if  $\mathcal{R} \not\subseteq \{X^2\}$ , then there exists  $R \in \mathcal{R}$  such that  $R \neq X^2$ . Therefore, there exist  $x, y \in X$  such that  $x \notin R(y)$ . Thus,  $S = \{x\} \times R(y) \cup \{x\}^c \times X$  is a non-reflexive relation on  $X$  such that  $S \in \mathcal{R}^\Delta$ . Therefore,  $\mathcal{R}^\Delta$  cannot be reflexive.

Note that if in particular either  $\mathcal{R} = \emptyset$  or  $\mathcal{R} = \{X^2\}$ , then  $\mathcal{R}^\Delta$  is also reflexive.

### 13 A Few Basic Facts on Filtred Relators

Intersection properties of relators were also first investigated in [74, 75].

**Definition 13** A relator  $\mathcal{R}$  on  $X$  to  $Y$  is called

- (1) *properly filtered* if for any  $R, S \in \mathcal{R}$  we have  $R \cap S \in \mathcal{R}$ ;
- (2) *uniformly filtered* if for any  $R, S \in \mathcal{R}$  there exists  $T \in \mathcal{R}$  such that  $T \subseteq R \cap S$ ;
- (3) *proximally filtered* if for any  $A \subseteq X$  and  $R, S \in \mathcal{R}$  there exists  $T \in \mathcal{R}$  such that  $T[A] \subseteq R[A] \cap S[A]$ ;
- (4) *topologically filtered* if for any  $x \in X$  and  $R, S \in \mathcal{R}$  there exists  $T \in \mathcal{R}$  such that  $T(x) \subseteq R(x) \cap S(x)$ .

*Remark 41* By using the binary operation  $\wedge$  and the basic closure operations on relators, the above properties can be reformulated in some more concise forms.

For instance, we can see that  $\mathcal{R}$  is topologically filtered if and only if any one of the properties  $\mathcal{R} \wedge \mathcal{R} \subseteq \mathcal{R}^\wedge$ ,  $(\mathcal{R} \wedge \mathcal{R})^\wedge = \mathcal{R}^\wedge$  and  $\mathcal{R}^\wedge \wedge \mathcal{R}^\wedge = \mathcal{R}^\wedge$  holds.

However, in general, we only have  $(R \cap S)[A] \subseteq R[A] \cap S[A]$ . Therefore, the corresponding proximal filteredness properties are, unfortunately, not equivalent.

Despite this, we can easily prove the following theorem which shows the appropriateness of the above proximal filteredness property.

**Theorem 41** *For a relator  $\mathcal{R}$  on  $X$  to  $Y$ , the following assertions are equivalent:*

- (1)  $\mathcal{R}$  is proximally filtered;
- (2)  $\text{Cl}_{\mathcal{R}}(A \cup B) = \text{Cl}_{\mathcal{R}}(A) \cup \text{Cl}_{\mathcal{R}}(B)$  for all  $A, B \subseteq Y$ ;

(3)  $\text{Int}_{\mathcal{R}}(A \cap B) = \text{Int}_{\mathcal{R}}(A) \cap \text{Int}_{\mathcal{R}}(B)$  for all  $A, B \subseteq Y$ .

**Proof** To prove the implication (3)  $\implies$  (1), note that if  $A \subseteq X$  and  $R, S \in \mathcal{R}$ , then by the definition of  $\text{Int}_{\mathcal{R}}$  we trivially have  $A \in \text{Int}_{\mathcal{R}}(R[A])$  and  $A \in \text{Int}_{\mathcal{R}}(S[A])$ . Therefore, if (3) holds, then we also have  $A \in \text{Int}_{\mathcal{R}}(R[A] \cap S[A])$ . Thus, by the definition of  $\text{Int}_{\mathcal{R}}$ , there exists  $T \in \mathcal{R}$  such that  $T[A] \subseteq R[A] \cap S[A]$ .

Now, as an immediate consequence of this theorem, we can also state

**Corollary 14** *If  $\mathcal{R}$  is a proximally filtered relator on  $X$ , then the families  $\varepsilon_{\mathcal{R}}$  and  $\tau_{\mathcal{R}}$  are closed under binary unions and intersections, respectively.*

From Theorem 41, we can also easily derive the following

**Theorem 42** *For a relator  $\mathcal{R}$  on  $X$  to  $Y$ , the following assertions are equivalent:*

- (1)  $\mathcal{R}$  is topologically filtered;
- (2)  $\text{cl}_{\mathcal{R}}(A \cup B) = \text{cl}_{\mathcal{R}}(A) \cup \text{cl}_{\mathcal{R}}(B)$  for all  $A, B \subseteq Y$ ;
- (3)  $\text{int}_{\mathcal{R}}(A \cap B) = \text{int}_{\mathcal{R}}(A) \cap \text{int}_{\mathcal{R}}(B)$  for all  $A, B \subseteq Y$ .

Thus, in particular, we can also state the following

**Corollary 15** *If  $\mathcal{R}$  is a topologically filtered relator on  $X$ , then the families  $\mathcal{F}_{\mathcal{R}}$  and  $\mathcal{I}_{\mathcal{R}}$  are closed under binary unions and intersections, respectively.*

The following example shows that, for a non-topological relator  $\mathcal{R}$ , the converse of the above corollary need not be true.

*Example 2* If  $X = \{1, 2, 3\}$  and  $R_i$  is relation on  $X$ , for each  $i = 1, 2$ , such that

$$R_i(1) = \{1, i + 1\} \quad \text{and} \quad R_i(2) = R_i(3) = \{2, 3\},$$

then  $\mathcal{R} = \{R_1, R_2\}$  is a reflexive relator on  $X$  such that  $\mathcal{I}_{\mathcal{R}}$  is closed under arbitrary intersections, but  $\mathcal{R}$  is still not topologically filtered.

By the corresponding definitions, it is clear that  $\mathcal{I}_{\mathcal{R}} = \{\emptyset, \{2, 3\}, X\}$ . Moreover, we can note that  $R_i(1) \not\subseteq R_1(1) \cap R_2(1)$  for each  $i = 1, 2$ . Thus, by Definition 13, the relator  $\mathcal{R}$  is not topologically filtered.

In addition to Theorem 42, we can also prove the following generalization of [48, Lemma 7] of Levine.

**Theorem 43** *If  $\mathcal{R}$  is a topologically filtered relator on  $X$ ,  $A, B \subseteq X$  and there exists  $V \in \mathcal{I}_{\mathcal{R}} \cap \mathcal{F}_{\mathcal{R}}$  such that  $A \subseteq V$  and  $B \subseteq V^c$ , then*

$$\text{int}_{\mathcal{R}}(A \cup B) = \text{int}_{\mathcal{R}}(A) \cup \text{int}_{\mathcal{R}}(B).$$

**Proof** Because of the increasingness of  $\circ$ , we evidently have  $A^\circ \cup B^\circ \subseteq (A \cup B)^\circ$ . Therefore, we need actually prove the converse inclusion.

For this, note that if  $x \in (A \cup B)^\circ$ , then by the definition of the operation  $\circ$ , there exists  $R \in \mathcal{R}$  such that  $R(x) \subseteq A \cup B$ .

Now, if  $x \in V$ , then by the definition of  $\mathcal{T}_{\mathcal{R}}$ , we can see that there exists  $S \in \mathcal{R}$  such that  $S(x) \subseteq V$ . Moreover, since  $\mathcal{R}$  is topologically filtered, there exists  $T \in \mathcal{R}$  such that  $T(x) \subseteq R(x) \cap S(x)$ . Hence, we can already infer that

$$T(x) \subseteq R(x) \cap S(x) \subseteq (A \cup B) \cap V = (A \cap V) \cup (B \cap V) = A \cup \emptyset = A.$$

Therefore, if  $x \in V$ , then  $x \in A^\circ$ .

A quite similar argument shows that if  $x \in V^c$ , then  $x \in B^\circ$ . Therefore, in both cases, we have  $x \in A^\circ \cup B^\circ$ . Thus, the inclusion  $(A \cup B)^\circ \subseteq A^\circ \cup B^\circ$  is also true.

*Remark 42* More difficult conditions for the dual equality  $(A \cap B)^- = A^- \cap B^-$  to hold were given by Gottschalk [34] and Jung and Nam [39, 40].

Concerning the latter problem, we shall only mention here the following

**Theorem 44** *If  $\mathcal{R}$  is a nonvoid, reflexive relator on  $X$  such that*

$$\text{cl}_{\mathcal{R}}(A \cap B) = \text{cl}_{\mathcal{R}}(A) \cap \text{cl}_{\mathcal{R}}(B)$$

*for all  $A, B \subseteq X$ , then  $\mathcal{R}^{\wedge\infty} = \mathcal{P}(X^2)^\infty$ .*

**Proof** For this, it is enough to prove only that  $\mathcal{T}_{\mathcal{R}} = \mathcal{P}(X)$ . Namely, in this case we have  $\mathcal{T}_{\mathcal{R}} = \mathcal{T}_{\Delta_X}$ . Hence, by using Theorem 25 and the corresponding definitions, we can already infer that  $\mathcal{R}^{\wedge\infty} = \{\Delta_X\}^{\wedge\infty} = \{\Delta_X\}^{*\infty} = \mathcal{P}(X^2)^\infty$ .

To prove the equality  $\mathcal{T}_{\mathcal{R}} = \mathcal{P}(X)$ , note that if this not true, then there exists  $A \subseteq X$  such that  $A \notin \mathcal{T}_{\mathcal{R}}$ , and thus  $B = A^c \notin \mathcal{F}_{\mathcal{R}}$ . Therefore,  $B^- \not\subseteq B$ , and thus there exists  $x \in B^-$  such that  $x \notin B$ . Hence, by using the assumptions of the theorem, we can arrive at the contradiction that  $x \in \{x\}^- \cap B^- = (\{x\} \cap B)^- = \emptyset^- = \emptyset$ .

## 14 A Few Basic Facts on Quasi-Filtered Relators

Since  $R \subseteq R^\infty$  for every relation  $R$  on  $X$ , in addition to Definition 13, we may also naturally introduce the following

**Definition 14** A relator  $\mathcal{R}$  on  $X$  is called

- (1) *quasi-uniformly filtered* if for any  $R, S \in \mathcal{R}$  there exists  $T \in \mathcal{R}$  such that  $T \subseteq R^\infty \cap S^\infty$ ;
- (2) *quasi-proximally filtered* if for any  $A \subseteq X$  and  $R, S \in \mathcal{R}$  there exists  $T \in \mathcal{R}$  such that  $T[A] \subseteq R^\infty[A] \cap S^\infty[A]$ ;

(3) *quasi-topologically filtered* if for any  $x \in X$  and  $R, S \in \mathcal{R}^\wedge$  there exists  $T \in \mathcal{R}$  such that  $T(x) \subseteq R^\infty(x) \cap S^\infty(x)$ .

*Remark 43* Analogously to Remark 41, the above quasi-filteredness properties can also be reformulated in some more concise forms.

For instance, we can see that  $\mathcal{R}$  is quasi-topologically filtered if and only if  $\mathcal{R}^{\wedge\infty} \wedge \mathcal{R}^{\wedge\infty} \subseteq \mathcal{R}^\wedge$ ,  $(\mathcal{R}^{\wedge\infty} \wedge \mathcal{R}^{\wedge\infty})^{\wedge\infty} = \mathcal{R}^{\wedge\infty}$  or  $\mathcal{R}^{\wedge\infty} \wedge \mathcal{R}^{\wedge\infty} = \mathcal{R}^{\wedge\infty}$ .

However, it is now more important to note that, by using some former results, we can also prove the following two theorems which show the appropriateness of the above quasi-proximal and quasi-topological filteredness properties.

**Theorem 45** *For any relator  $\mathcal{R}$  on  $X$ , the following assertions are equivalent:*

- (1)  $\mathcal{R}$  is a quasi-proximally filtered;
- (2)  $\varepsilon_{\mathcal{R}}$  is closed under binary unions;
- (3)  $\tau_{\mathcal{R}}$  is closed under binary intersections.

**Theorem 46** *For any relator  $\mathcal{R}$  on  $X$ , the following assertions are equivalent:*

- (1)  $\mathcal{R}$  is a quasi-topologically filtered;
- (2)  $\mathcal{F}_{\mathcal{R}}$  is closed under binary unions;
- (3)  $\mathcal{I}_{\mathcal{R}}$  is closed under binary intersections.

*Remark 44* In this respect it is also worth mentioning that if  $\mathcal{R}$  is a relator on  $X$  to  $Y$ , then the family  $\mathcal{E}_{\mathcal{R}}$  is closed under binary intersections if and only if  $\mathcal{R}$  is *quasi-directed* in the sense that for any  $x, y \in X$  and  $R, S \in \mathcal{R}$  we have  $R(x) \cap S(y) \in \mathcal{E}_{\mathcal{R}}$ .

From the above two theorems, by using Corollaries 14 and 15, we can derive

**Corollary 16** *If  $\mathcal{R}$  is a proximally (topologically) filtered relator on  $X$ , then  $\mathcal{R}$  is also quasi-proximally (quasi-topologically) filtered.*

Now, by using Theorem 45, we can also easily prove the following

**Theorem 47** *If  $\mathcal{R}$  is a quasi-proximally filtered, proximal relator on  $X$ , then  $\mathcal{R}$  is proximally filtered.*

**Proof** Suppose that  $A \subseteq X$  and  $R, S \in \mathcal{R}$ . Then, by Definition 12, there exist  $U, V \in \tau_{\mathcal{R}}$  such that  $A \subseteq U \subseteq R[A]$  and  $A \subseteq V \subseteq S[A]$ .

Moreover, by Theorem 45, we can state that  $U \cap V \in \tau_{\mathcal{R}}$ . Therefore, by the definition of  $\tau_{\mathcal{R}}$ , there exists  $T \in \mathcal{R}$  such that  $T[U \cap V] \subseteq U \cap V$ . Hence, we can already see that  $T[A] \subseteq T[U \cap V] \subseteq U \cap V \subseteq R[A] \cap S[A]$ .

Moreover, by using Theorem 46, we can quite similarly prove the following

**Theorem 48** *If  $\mathcal{R}$  is a quasi-topologically filtered, topological relator on  $X$ , then  $\mathcal{R}$  is topologically filtered.*



*Remark 45* Our former Example 2 shows that even a quasi-proximally filtered, reflexive relator need not be topologically filtered.

Namely, if  $X$  and  $\mathcal{R}$  are as in Example 2, then by the corresponding definitions it is clear that  $\tau_{\mathcal{R}} = \{\emptyset, \{2, 3\}, X\}$ , and thus by Theorem 45 the relator  $\mathcal{R}$  is quasi-proximally filtered.

## 15 Some Further Theorems on Topologically Filtered Relators

In our former papers [67, 68], by using the arguments of Kuratowski [45, pp. 39, 45], we have also proved the following two basic theorems.

**Theorem 49** *If  $\mathcal{R}$  is a topologically filtered relator on  $X$  to  $Y$ , then for any  $A, B \subseteq Y$  we have*

$$\text{cl}_{\mathcal{R}}(A) \setminus \text{cl}_{\mathcal{R}}(B) = \text{cl}_{\mathcal{R}}(A \setminus B) \setminus \text{cl}_{\mathcal{R}}(B).$$

*Proof* By using Theorem 42, we can see that

$$A^- \cup B^- = (A \cup B)^- = ((A \setminus B) \cup B)^- = (A \setminus B)^- \cup B^-.$$

Hence, because of the identity  $(U \cup V) \setminus V = U \setminus V$ , the required equality follows.

Thus, in particular we can also state the following

**Corollary 17** *If  $\mathcal{R}$  is a topologically filtered relator on  $X$  to  $Y$ , then for any  $A, B \subseteq Y$  we have  $\text{cl}_{\mathcal{R}}(A) \setminus \text{cl}_{\mathcal{R}}(B) \subseteq \text{cl}_{\mathcal{R}}(A \setminus B)$ .*

This corollary already allows us to easily prove the following

**Theorem 50** *If  $\mathcal{R}$  is a topologically filtered relator on  $X$ , then for any  $A \subseteq X$  and  $U \in \mathcal{F}_{\mathcal{R}}$  we have*

$$\text{cl}_{\mathcal{R}}(A) \cap U = \text{cl}_{\mathcal{R}}(A \cap U) \cap U.$$

*Proof* By Definition 3 and Theorem 2, we have  $U \subseteq U^\circ = U^{c-c}$ . Hence, by using Corollary 17, we can infer that

$$A^- \cap U \subseteq A^- \cap U^{c-c} = A^- \setminus U^{c-c} \subseteq (A \setminus U^c)^- = (A \cap U)^-.$$

Therefore,  $A^- \cap U = A^- \cap (U \cap U) = (A^- \cap U) \cap U \subseteq (A \cap U)^- \cap U$ .

Moreover, by using the increasingness of  $-$ , we can see that  $(A \cap U)^- \subseteq A^-$ , and thus  $(A \cap U)^- \cap U \subseteq A^- \cap U$  is always true. Therefore, we actually have  $A^- \cap U = (A \cap U)^- \cap U$ .

This theorem can also be easily derived from its subsequent corollary which can also be easily proved directly, without using Corollary 17. (See [68].)

**Corollary 18** *If  $\mathcal{R}$  is a topologically filtered relator on  $X$ , then for any  $A \subseteq X$  and  $U \in \mathcal{T}_{\mathcal{R}}$  we have  $\text{cl}_{\mathcal{R}}(A) \cap U \subseteq \text{cl}_{\mathcal{R}}(A \cap U)$ .*

*Remark 46* The importance of the closure space counterpart of Corollary 18 was also recognized Császár [13–18] and Sivagami [70] who assumed it as an axiom for an increasing set-to-set function  $\gamma$ .

Now, as a dual form of Theorem 50, it is also worth proving the following

**Theorem 51** *If  $\mathcal{R}$  is a topologically filtered relator on  $X$ , then for any  $A \subseteq X$  and  $V \in \mathcal{F}_{\mathcal{R}}$  we have*

$$\text{int}_{\mathcal{R}}(A) \cup V = \text{int}_{\mathcal{R}}(A \cup V) \cup V.$$

**Proof** By using Theorems 2 and 50, we can see that

$$\begin{aligned} A^{\circ} \cup V &= A^{c-c} \cup V^{cc} = (A^{c-} \cap V^c)^c = ((A^c \cap V^c)^- \cap V^c)^c \\ &= ((A \cup V)^{c-} \cap V^c)^c = (A \cup V)^{c-c} \cup V = (A \cup V)^{\circ} \cup V. \end{aligned}$$

Thus, in particular we can also state the following

**Corollary 19** *If  $\mathcal{R}$  is a topologically filtered relator on  $X$ , then for any  $A \subseteq X$  and  $V \in \mathcal{F}_{\mathcal{R}}$  we have  $\text{int}_{\mathcal{R}}(A \cup V) \subseteq \text{int}_{\mathcal{R}}(A) \cup V$ .*

The importance of this corollary is apparent from the following

**Theorem 52** *If  $\mathcal{R}$  is a topologically filtered, quasi-topological relator on  $X$ , then for any  $A, B \in \mathcal{N}_{\mathcal{R}}$  we have  $A \cup B \in \mathcal{N}_{\mathcal{R}}$ .*

**Proof** By using Theorem 42 and Corollary 19, we can see that

$$(A \cup B)^{-\circ} = (A^- \cup B^-)^{\circ} \subseteq A^{-\circ} \cup B^- = \emptyset \cup B^- = B^-.$$

Hence, by using Theorem 31 and the definition of  $\mathcal{T}_{\mathcal{R}}$ , we can infer that

$$(A \cup B)^{-\circ} \subseteq (A \cup B)^{-\circ\circ} \subseteq B^{-\circ} = \emptyset.$$

Therefore,  $(A \cup B)^{-\circ} = \emptyset$ , and thus  $A \cup B \in \mathcal{N}_{\mathcal{R}}$ .

Now, by using this theorem, we can also easily establish the following

**Corollary 20** *If  $\mathcal{R}$  is a nonvoid, non-partial, topologically filtered, quasi-topological relator on  $X$ , then  $\mathcal{N}_{\mathcal{R}}$  is an ideal on  $X$ .*

**Proof** By the definition of  $\mathcal{N}_{\mathcal{R}}$  and the increasingness of  $-$  and  $\circ$ , it is clear that  $\mathcal{N}_{\mathcal{R}}$  is always descending. Moreover, since  $\mathcal{R}$  is nonvoid and non-partial, we can also see that  $\emptyset^{-\circ} = \emptyset^{\circ} = \emptyset$ . Therefore,  $\emptyset \in \mathcal{N}_{\mathcal{R}}$ , and thus  $\mathcal{N}_{\mathcal{R}} \neq \emptyset$ . Furthermore, from Theorem 52, we know that  $\mathcal{N}_{\mathcal{R}}$  is closed under pairwise unions.

*Remark 47* Note that if  $\mathcal{R}$  is a locally non-partial relator on  $X$ , then by Theorem 29 we have  $X^{-\circ} = X$ . Therefore, if  $X \neq \emptyset$ , then we can also state that  $X \notin \mathcal{N}_{\mathcal{R}}$ , and thus  $\mathcal{N}_{\mathcal{R}} \neq \mathcal{P}(X)$ .

While, if  $\mathcal{R}$  is a quasi-topological relator on  $X$  and  $A \in \mathcal{N}_{\mathcal{R}}$ , then by using Theorem 31 and the increasingness of  $\circ$  we can also see that  $A^{-\circ} \subseteq A^{-\circ} = \emptyset$ . Therefore,  $A^{-\circ} = \emptyset$ , and thus  $A^{-} \in \mathcal{N}_{\mathcal{R}}$ .

## 16 Some More Particular Theorems on Topologically Filtered Relators

By using Corollary 18, we can also easily prove the following

**Theorem 53** *If  $\mathcal{R}$  is a topologically filtered, topological relator on  $X$ , then for any  $A \subseteq X$  and  $U \in \mathcal{I}_{\mathcal{R}}$  we have*

$$\text{cl}_{\mathcal{R}}(A \cap U) = \text{cl}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A) \cap U).$$

**Proof** By Corollary 18 we have  $A^{-} \cap B \subseteq (A \cap B)^{-}$ . Hence, by using Theorems 32 and 31, we can infer that

$$(A^{-} \cap B)^{-} \subseteq (A \cap B)^{-\circ} \subseteq (A \cap B)^{-}.$$

On the other hand, by Theorem 26, we have  $A \subseteq A^{-}$ , and thus also  $A \cap B \subseteq A^{-} \cap B$ . Hence, we can infer that  $(A \cap B)^{-} \subseteq (A^{-} \cap B)^{-}$ . Therefore, the corresponding equality is also true.

From this theorem, we can immediate derive

**Corollary 21** *If  $\mathcal{R}$  is a topologically filtered, topological relator on  $X$ , then for any  $A \in \mathcal{D}_{\mathcal{R}}$  and  $U \in \mathcal{I}_{\mathcal{R}}$  we have  $\text{cl}_{\mathcal{R}}(U) = \text{cl}_{\mathcal{R}}(A \cap U)$ .*

**Proof** By Definition 1 and Theorem 53, we evidently have

$$U^{-} = (X \cap U)^{-} = (A^{-} \cap U)^{-} = (A \cap U)^{-}.$$

Now, by modifying an argument of Levine [50], we can also prove

**Theorem 54** *If  $\mathcal{R}$  is a nonvoid, topological relator on  $X$  and  $A \subseteq X$  such that  $\text{cl}_{\mathcal{R}}(U) = \text{cl}_{\mathcal{R}}(A \cap U)$  for all  $U \in \mathcal{I}_{\mathcal{R}}$ , then  $A \in \mathcal{D}_{\mathcal{R}}$ .*

**Proof** Assume on the contrary that  $A \notin \mathcal{D}_{\mathcal{R}}$ . Then, by Definition 1, there exists  $x \in X$  such that  $x \notin A^-$ . Thus, by Definition 1, there exists  $R \in \mathcal{R}$  such that  $A \cap R(x) = \emptyset$ . Moreover, by Definition 11, there exists  $U \in \mathcal{T}_{\mathcal{R}}$  such that  $x \in U \subseteq R(x)$ . Thus, in particular we also have  $A \cap U = \emptyset$ .

Hence, by using the assumptions of the theorem, we can already infer that  $U^- = (A \cap U)^- = \emptyset^- = \emptyset$ . Note that the latter equality already requires that  $\mathcal{R} \neq \emptyset$ .

On the other hand, from the inclusion  $x \in U$ , by using Theorems 32 and 26 and the increasingness of  $-$ , we can infer that  $x \in \{x\}^- \subseteq U^-$ , and thus  $U^- \neq \emptyset$ . This contradiction proves that  $A \in \mathcal{D}_{\mathcal{R}}$ .

*Remark 48* If  $\mathcal{R}$  is a nonvoid, reflexive relator on  $X$  and  $A \subseteq X$  such that  $\text{cl}_{\mathcal{R}}(R(x)) = \text{cl}_{\mathcal{R}}(A \cap R(x))$  for all  $x \in X$  and  $R \in \mathcal{R}$ , then we can even more easily prove that  $A \in \mathcal{D}_{\mathcal{R}}$ .

In addition to Theorem 54, we can also prove the following

**Theorem 55** *If  $\mathcal{R}$  is a topologically filtered, topological relator on  $X$ , then for any  $U \in \mathcal{T}_{\mathcal{R}}$  we have*

$$\text{res}_{\mathcal{R}}(U) \in \mathcal{F}_{\mathcal{R}} \setminus \mathcal{E}_{\mathcal{R}}.$$

**Proof** By Theorem 11, we have  $U^c \in \mathcal{F}_{\mathcal{R}}$ . Moreover, by Theorems 32 and 31, we have  $U^- \in \mathcal{F}_{\mathcal{R}}$ . Hence, by using the notation  $U^\dagger = \text{res}_{\mathcal{R}}(U)$ , we can infer that

$$U^\dagger = U^- \setminus U = U^- \cap U^c \in \mathcal{F}_{\mathcal{R}}.$$

Moreover, by using Theorems 42, 2, 32, and 26, we can also see that

$$U^{\dagger\circ} = (U^- \setminus U)^\circ = U^{-\circ} \cap U^{c\circ} = U^{-\circ} \cap U^{-c} \subseteq U^- \cap U^{-c} = \emptyset,$$

and thus  $U^{\dagger\circ} = \emptyset$ . Therefore,  $U^\dagger \notin \mathcal{E}_{\mathcal{R}}$ , and thus  $U^\dagger \in \mathcal{F}_{\mathcal{R}} \setminus \mathcal{E}_{\mathcal{R}}$ .

Now, as an immediate consequence of this theorem, we can also state

**Corollary 22** *If  $\mathcal{R}$  is a topologically filtered, topological relator on  $X$ , then  $\text{res}_{\mathcal{R}}(U) \in \mathcal{N}_{\mathcal{R}}$  for all  $U \in \mathcal{T}_{\mathcal{R}}$ .*

*Remark 49* Note that if  $\mathcal{R}$  is a topological relator on  $X$  and  $U \in \mathcal{T}_{\mathcal{R}}$ , then by Definition 3 and Theorems 32 and 26 we have  $U = U^\circ$ . Therefore, under the notation  $U^\ddagger = \text{bnd}_{\mathcal{R}}(U)$ , we have  $U^\dagger = U^- \setminus U = U^- \setminus U^\circ = U^\ddagger$ .

Moreover, it is also worth noticing that in Theorem 55 and Corollary 22, it is enough to assume only that  $\mathcal{R}$  is a quasi-topologically filtered, topological relator on  $X$ . Namely, in this case,  $\mathcal{R}$  is already topologically filtered by Theorem 48.

### 17 Some Generalized Topologically Open Sets

**Notation 1** In the sequel, we shall always assume that  $X$  is a set and  $\mathcal{R}$  is a relator on  $X$ .

Moreover, to shorten the subsequent proofs, we shall again use the notations

$$A^- = \text{cl}_{\mathcal{R}}(A), \quad A^\circ = \text{int}_{\mathcal{R}}(A) \quad \text{and} \quad A^\dagger = \text{res}_{\mathcal{R}}(A).$$

In our first joint paper [67], motivated by the corresponding definitions on generalized open subsets of topological spaces mentioned in the Motivations, we have introduced the following

**Definition 15** For a subset  $A$  of the relator space  $X(\mathcal{R})$ , we write

- (1)  $A \in \mathcal{T}_{\mathcal{R}}^s$  if  $A \subseteq \text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A))$ ;
- (2)  $A \in \mathcal{T}_{\mathcal{R}}^p$  if  $A \subseteq \text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))$ ;
- (3)  $A \in \mathcal{T}_{\mathcal{R}}^\alpha$  if  $A \subseteq \text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A)))$ ;
- (4)  $A \in \mathcal{T}_{\mathcal{R}}^\beta$  if  $A \subseteq \text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)))$ ;
- (5)  $A \in \mathcal{T}_{\mathcal{R}}^a$  if  $A \subseteq \text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A)) \cap \text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))$ ;
- (6)  $A \in \mathcal{T}_{\mathcal{R}}^b$  if  $A \subseteq \text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A)) \cup \text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))$ .

And, the members of the above families are called the *topologically semi-open, preopen,  $\alpha$ -open,  $\beta$ -open,  $a$ -open and  $b$ -open subsets* of the relator space  $X(\mathcal{R})$ , respectively.

*Remark 50* Note that, for instance,  $\circ -$  is always an increasing operation on  $\mathcal{P}(X)$ . Moreover, if  $\mathcal{R}$  is nonvoid and non-partial, then  $\emptyset^{\circ-} = \emptyset$  and  $X^{\circ-} = X$ .

Therefore, by [76, Theorem 9.4], there exists a nonvoid and non-partial relator  $\mathcal{S}$  on  $X$  such that  $A^{\circ-} = \text{int}_{\mathcal{S}}(A)$  for all  $A \subseteq X$ , and thus  $\mathcal{T}_{\mathcal{R}}^s = \mathcal{T}_{\mathcal{S}}$ .

However, this fact can be used to establish only that if  $\mathcal{R}$  is a nonvoid and non-partial relator on  $X$ , then  $\mathcal{T}_{\mathcal{R}}^s$  is a generalized topology on  $X$  [85].

In [67], by using the enormous literature on generalized open sets in topological spaces and their generalizations, we have, for instance, proved the following theorems.

**Theorem 56** We have (1)  $\mathcal{T}_{\mathcal{R}}^a = \mathcal{T}_{\mathcal{R}}^s \cap \mathcal{T}_{\mathcal{R}}^p$ ; (2)  $\mathcal{T}_{\mathcal{R}}^s \cup \mathcal{T}_{\mathcal{R}}^p \subseteq \mathcal{T}_{\mathcal{R}}^b$ .

**Theorem 57** If  $\mathcal{R}$  is a reflexive relator on  $X$ , then (1)  $\mathcal{T}_{\mathcal{R}}^\alpha \subseteq \mathcal{T}_{\mathcal{R}}^a$ ; (2)

$$\mathcal{T}_{\mathcal{R}}^b \subseteq \mathcal{T}_{\mathcal{R}}^\beta.$$

(3)  $\mathcal{T}_{\mathcal{R}} \subseteq \mathcal{T}_{\mathcal{R}}^\kappa$  for all  $\kappa = s, p, \alpha, \beta, a$  and  $b$ .

*Remark 51* Note that, by the former inclusions, it is enough to prove the inclusion  $\mathcal{T}_{\mathcal{R}} \subseteq \mathcal{T}_{\mathcal{R}}^\kappa$  only for  $\kappa = \alpha$ .

**Theorem 58** If  $\mathcal{R}$  is a topological relator on  $X$ , then  $\mathcal{T}_{\mathcal{R}}^\alpha = \mathcal{T}_{\mathcal{R}}^a$ .

**Theorem 59** If  $\mathcal{R}$  is a topological relator on  $X$ , then for any  $A \subseteq X$ , the following assertions are equivalent:

- (1)  $A \in \mathcal{T}_{\mathcal{R}}^s$ ; (2)  $\text{cl}_{\mathcal{R}}(A) \subseteq \text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A))$ ; (3)  $\text{cl}_{\mathcal{R}}(A) = \text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A))$ .

**Theorem 60** If  $\mathcal{R}$  is a topological relator on  $X$ , then for any  $A \subseteq X$  the following assertions are equivalent:

- (1)  $A \in \mathcal{T}_{\mathcal{R}}^s$ ;  
 (2) there exists  $V \in \mathcal{T}_{\mathcal{R}}$  such that  $V \subseteq A$  and  $\text{cl}_{\mathcal{R}}(A) = \text{cl}_{\mathcal{R}}(V)$ ;  
 (3) there exist  $V \in \mathcal{T}_{\mathcal{R}}$  and  $B \subseteq X$  such that  $A = V \cup B$  and  $B \subseteq \text{res}_{\mathcal{R}}(V)$ .

**Corollary 23** If  $\mathcal{R}$  is a topologically filtered, topological relator on  $X$  and  $A \in \mathcal{T}_{\mathcal{R}}^s$ , then there exist  $V \in \mathcal{T}_{\mathcal{R}}$  and  $B \in \mathcal{N}_{\mathcal{R}}$  such that  $A = V \cup B$  and  $V \cap B = \emptyset$ .

**Theorem 61** If  $\mathcal{R}$  is a topologically filtered, topological relator on  $X$ , then for any  $A \subseteq X$  the following assertions are equivalent:

- (1)  $A \in \mathcal{T}_{\mathcal{R}}^p$ ;  
 (2) there exist  $V \in \mathcal{T}_{\mathcal{R}}$  and  $B \in \mathcal{D}_{\mathcal{R}}$  such that  $A = V \cap B$ ;  
 (3) there exists  $V \in \mathcal{T}_{\mathcal{R}}$  such that  $A \subseteq V$  and  $\text{cl}_{\mathcal{R}}(A) = \text{cl}_{\mathcal{R}}(V)$ .

**Theorem 62** If  $\mathcal{R}$  is a topologically filtered, topological relator on  $X$ , then for any  $A \subseteq X$  the following assertions are equivalent:

- (1)  $A \in \mathcal{T}_{\mathcal{R}}^{\alpha}$ ;  
 (2) there exist  $V \in \mathcal{T}_{\mathcal{R}}$  and  $B \in \mathcal{N}_{\mathcal{R}}$  such that  $A = V \setminus B$ ;  
 (3) there exist  $V \in \mathcal{T}_{\mathcal{R}}$  and  $B \subseteq \text{res}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A))$  such that  $A = V \setminus B$ .

**Theorem 63** If  $\mathcal{R}$  is a topological relator on  $X$ , then for any  $A \subseteq X$  the following assertions are equivalent: (1)  $A \in \mathcal{T}_{\mathcal{R}}^{\beta}$ ; (2)  $\text{cl}_{\mathcal{R}}(A) \in \mathcal{T}_{\mathcal{R}}^s$ ;

- (3) there exists  $V \in \mathcal{T}_{\mathcal{R}}$  such that  $\text{cl}_{\mathcal{R}}(A) = \text{cl}_{\mathcal{R}}(V)$ .

*Remark 52* Hence, by using Theorems 59 and 60, we can derive some further characterizations of the family  $\mathcal{T}_{\mathcal{R}}^{\beta}$ .

On the other hand, from Theorem 56, by using Theorems 59, 60, and 61, we can derive several characterizations of the family  $\mathcal{T}_{\mathcal{R}}^{\alpha}$ .

**Theorem 64** If  $\mathcal{R}$  is a topologically filtered, topological relator on  $X$ , then for any  $A \subseteq X$  the following assertions are equivalent:

- (1)  $A \in \mathcal{T}_{\mathcal{R}}^b$ ;  
 (2) there exist  $B \in \mathcal{T}_{\mathcal{R}}^s$  and  $C \in \mathcal{T}_{\mathcal{R}}^p$  such that  $A = B \cup C$ .

## 18 Some Further Theorems on Generalized Topologically Open Sets

In our second joint paper [68], by using the enormous literature on generalized open sets in topological spaces and their generalizations, we have, for instance, proved the following theorems.

**Theorem 65** *The families  $\mathcal{T}_{\mathcal{R}}^{\kappa}$ , with  $\kappa = s, p, \alpha, \beta, a$  and  $b$ , are closed under arbitrary unions.*

*Remark 53* Thus, in particular we have  $\emptyset \in \mathcal{T}_{\mathcal{R}}^{\kappa}$  for all  $\kappa = s, p, \alpha, \beta, a$  and  $b$ .

**Theorem 66** *The following assertions are equivalent: (1)  $X \in \mathcal{T}_{\mathcal{R}}^s$ ; (2)*

*$X \in \mathcal{T}_{\mathcal{R}}^{\beta}$ ; (3)  $\mathcal{R}$  is non-partial.*

**Theorem 67** *The following assertions are equivalent: (1)  $X \in \mathcal{T}_{\mathcal{R}}^p$ ; (2)*

*$X \in \mathcal{T}_{\mathcal{R}}^{\alpha}$ ; (3)  $\mathcal{R}$  is locally non-partial.*

**Corollary 24** *The following assertions are equivalent: (1)  $X \in \mathcal{T}_{\mathcal{R}}^a$ ; (2)  $\mathcal{R}$*

*is non-partial and locally non-partial.*

*Remark 54* If  $\mathcal{R}$  is either non-partial or locally non-partial, then by Theorems 56, 66, and 67, we can also state that  $X \in \mathcal{T}_{\mathcal{R}}^b$ .

**Theorem 68** *If  $\mathcal{R}$  is a nonvoid, non-partial relator on  $X$ , then  $\mathcal{T}_{\mathcal{R}}^{\kappa}$  is a generalized topology on  $X$  for all  $\kappa = s, p, \alpha, \beta, a$  and  $b$ .*

**Theorem 69** *If  $\mathcal{R}$  is a topologically filtered relator on  $X$  and  $U \in \mathcal{T}_{\mathcal{R}}$ , then  $U \cap A \in \mathcal{T}_{\mathcal{R}}^{\kappa}$  for all  $A \in \mathcal{T}_{\mathcal{R}}^{\kappa}$  with  $\kappa = s, p, \alpha, \beta, a$  and  $b$ .*

**Theorem 70** *If  $\mathcal{R}$  is a topologically filtered, topological relator on  $X$  and  $A \in \mathcal{T}_{\mathcal{R}}^{\alpha}$ , then  $A \cap B \in \mathcal{T}_{\mathcal{R}}^{\kappa}$  for all  $B \in \mathcal{T}_{\mathcal{R}}^{\kappa}$  with  $\kappa = s, p, \alpha, \beta, a$  and  $b$ .*

**Corollary 25** *If  $\mathcal{R}$  is a nonvoid, topologically filtered, topological relator on  $X$ , then  $\mathcal{T}_{\mathcal{R}}^{\alpha}$  is a topology on  $X$ .*

In [68], by using an argument of Njåstad [62, p. 962], we have also proved

**Theorem 71** *If  $\mathcal{R}$  is a nonvoid, topologically filtered, topological relator on  $X$  and  $A \subseteq X$  such that  $A \cap B \in \mathcal{T}_{\mathcal{R}}^s$  for all  $B \in \mathcal{T}_{\mathcal{R}}^s$ , then  $A \in \mathcal{T}_{\mathcal{R}}^{\alpha}$ .*

Moreover, to introduce the corresponding generalized topologically closed sets, we have also used the following

**Definition 16** For any  $\kappa = s, p, \alpha, \beta, a$  and  $b$ , we define

$$\mathcal{F}_{\mathcal{R}}^{\kappa} = \{ A \subseteq X : A^c \in \mathcal{T}_{\mathcal{R}}^{\kappa} \}.$$

Thus, by using Theorem 2, we could easily prove the following

**Theorem 72** For any  $A \subseteq X$ , we have

- (1)  $A \in \mathcal{F}_{\mathcal{R}}^s$  if and only if  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)) \subseteq A$ ;
- (2)  $A \in \mathcal{F}_{\mathcal{R}}^p$  if and only if  $\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A)) \subseteq A$ ;
- (3)  $A \in \mathcal{F}_{\mathcal{R}}^\alpha$  if and only if  $\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))) \subseteq A$ ;
- (4)  $A \in \mathcal{F}_{\mathcal{R}}^\beta$  if and only if  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A))) \subseteq A$ ;
- (5)  $A \in \mathcal{F}_{\mathcal{R}}^a$  if and only if  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)) \cup \text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A)) \subseteq A$ ;
- (6)  $A \in \mathcal{F}_{\mathcal{R}}^b$  if and only if  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)) \cap \text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A)) \subseteq A$ .

In [67] and [68], following Kuratowski [44], we have also introduced

**Definition 17** A subset  $A$  of the relator space  $X(\mathcal{R})$  is called *topologically regular open* if

$$A = \text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)).$$

And, the family of all such subsets of  $X(\mathcal{R})$  is denoted by  $\mathcal{T}_{\mathcal{R}}^r$ .

Thus, by using several papers on regular open sets in topological spaces, we have, for instance, proved the following theorems.

**Theorem 73** If  $\mathcal{R}$  is a quasi-topological relator on  $X$ , then  $\mathcal{T}_{\mathcal{R}}^r \subseteq \mathcal{T}_{\mathcal{R}}$ .

**Theorem 74** If  $\mathcal{R}$  is a reflexive relator on  $X$ , then  $\mathcal{T}_{\mathcal{R}} \cap \mathcal{F}_{\mathcal{R}} \subseteq \mathcal{T}_{\mathcal{R}}^r$ .

**Theorem 75** We always have  $\mathcal{T}_{\mathcal{R}}^r = \mathcal{T}_{\mathcal{R}}^p \cap \mathcal{F}_{\mathcal{R}}^s$ .

**Theorem 76** If  $\mathcal{R}$  is a topological relator on  $X$ , then (1)  $\mathcal{T}_{\mathcal{R}}^r = \mathcal{T}_{\mathcal{R}} \cap \mathcal{F}_{\mathcal{R}}^s$ ;

(2)  $\mathcal{T}_{\mathcal{R}}^r = \mathcal{T}_{\mathcal{R}}^\alpha \cap \mathcal{F}_{\mathcal{R}}^\beta$ .

**Theorem 77** If  $\mathcal{R}$  is a topological relator on  $X$ , then (1)  $\text{cl}_{\mathcal{R}}(A) \in \mathcal{T}_{\mathcal{R}}^r$  for all  $A \in \mathcal{T}_{\mathcal{R}}^s$ ;

(2)  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)) \in \mathcal{T}_{\mathcal{R}}^r$  for all  $A \subseteq X$ .

**Theorem 78** If  $\mathcal{R}$  is a topologically filtered, topological relator on  $X$ , then  $\mathcal{T}_{\mathcal{R}}^r$  is closed under pairwise intersections.

The following example shows that the counterparts of Theorems 65 and 69 fail to hold for the family  $\mathcal{T}_{\mathcal{R}}^r$ .

*Example 3* If  $X = \mathbb{R}$  and

$$R_n = \{(x, y) \in X^2 : d(x, y) < n^{-1}\}$$



for all  $n \in \mathbb{N}$ , then  $\mathcal{R} = \{R_n : n \in \mathbb{N}\}$  is a properly filtered, strongly topological, tolerance relator on  $X$  such that, for the sets

$$A = ]0, 1[, \quad B = ]1, 2[, \quad C = ]0, 2[ \quad \text{and} \quad U = \{1\}^c,$$

we have  $A \cup B = C \cap U \notin \mathcal{T}_{\mathcal{R}}^r$  despite that  $A, B, C \in \mathcal{T}_{\mathcal{R}}^r$  and  $U \in \mathcal{T}_{\mathcal{R}}$ .

### 19 A Further Family of Generalized Topologically Open Sets

The following definitions and some of the forthcoming theorems have been mainly suggested to us by Levine [48], Chattopadhyay and Bandyopadhyay [8] and Császár [15, 16].

In [48], a subset  $A$  of a topological space is said to have *property Q* if  $A^{-\circ} = A^{\circ-}$ . While, in [8] and [15], the set  $A$  is called a  $\delta$ -set and a *quasi-open set*, respectively, if  $A^{-\circ} \subseteq A^{\circ-}$ .

**Definition 18** If  $\Phi$  and  $\Psi$  are relation on  $\mathcal{P}(X)$  to  $X$ , then for any  $A \subseteq X$  we shall write

$$A \in \mathcal{A}(\Phi, \Psi) \quad \text{if} \quad \Psi(\Phi(A)) \subseteq \Phi(\Psi(A)).$$

*Remark 55* Thus, for any  $A \subseteq X$ , we have

$$A \in \mathcal{A}(\Phi, \Psi) \cap \mathcal{A}(\Psi, \Phi) \quad \text{if and only if} \quad \Psi(\Phi(A)) = \Phi(\Psi(A)).$$

That is, the associated set-valued functions commute at the set  $A$ .

In the sequel, we shall also restrict ourselves to the particular case when  $\Phi = \text{cl}_{\mathcal{R}}$  and  $\Psi = \text{int}_{\mathcal{R}}$ . And, to shorten the subsequent statements and proofs, we shall use the following

**Definition 19** In particular, we define

$$\mathcal{A}_{\mathcal{R}} = \mathcal{A}_{\mathcal{R}}(-, \circ) = \mathcal{A}(\text{cl}_{\mathcal{R}}, \text{int}_{\mathcal{R}}).$$

*Remark 56* Thus, for any  $A \subseteq X$ , we have

$$A \in \mathcal{A}_{\mathcal{R}} \quad \text{if and only if} \quad A^{-\circ} \subseteq A^{\circ-}.$$

By using the latter property, we can easily prove the following theorems which could, to some extent, be generalized to the more general case in Definition 18.

**Theorem 79** We have  $\mathcal{N}_{\mathcal{R}} \subseteq \mathcal{A}_{\mathcal{R}}$ .

**Proof** If  $A \in \mathcal{N}_{\mathcal{R}}$ , then by Definition 3, we have  $A^{-\circ} = \emptyset$ . Thus,  $A^{-\circ} \subseteq A^{\circ-}$  trivially holds. Hence, by Remark 56, we can see that  $A \in \mathcal{A}_{\mathcal{R}}$ . Therefore, the required inclusion is true.

*Remark 57* If  $\mathcal{R}$  is a nonvoid, reflexive relator on  $X$ , then we can also state that  $\mathcal{N}_{\mathcal{R}} \subseteq \mathcal{A}_{\mathcal{R}}(-, \circ) \cap \mathcal{A}_{\mathcal{R}}(\circ, -)$ .

Namely, if  $\mathcal{R}$  is reflexive, then by Theorem 26 we have  $A \subseteq A^{-}$ , and thus also  $A^{\circ} \subseteq A^{-\circ}$ . Therefore, if  $A \in \mathcal{N}_{\mathcal{R}}$ , i.e.,  $A^{-\circ} = \emptyset$ , then we also have  $A^{\circ} = \emptyset$ . Moreover, if  $\mathcal{R}$  is nonvoid, then we also have  $\emptyset^{-} = \emptyset$ . Therefore, in this particular case, we actually have  $A^{\circ-} = \emptyset^{-} = \emptyset = A^{-\circ}$ , and thus also  $A \in \mathcal{A}_{\mathcal{R}}(\circ, -)$ .

**Theorem 80** *If  $\mathcal{R}$  is a nonvoid relator on  $X$ , then*

- (1)  $\mathcal{A}_{\mathcal{R}} \setminus \mathcal{N}_{\mathcal{R}} \subseteq \mathcal{E}_{\mathcal{R}}$ ;      (2)  $A \in \mathcal{A}_{\mathcal{R}} \cap \mathcal{D}_{\mathcal{R}}$  implies  $\text{int}_{\mathcal{R}}(A) \in \mathcal{D}_{\mathcal{R}}$ .

**Proof** If  $A \in \mathcal{A}_{\mathcal{R}} \setminus \mathcal{N}_{\mathcal{R}}$ , then  $A \in \mathcal{A}_{\mathcal{R}}$  and  $A \notin \mathcal{N}_{\mathcal{R}}$ . Therefore, by Remark 56 and Definition 3, we have  $A^{-\circ} \subseteq A^{\circ-}$  and  $A^{-\circ} \neq \emptyset$ , and thus  $A^{\circ-} \neq \emptyset$ . Hence, since  $\mathcal{R} \neq \emptyset$ , and thus  $\emptyset^{-} = \emptyset$ , we can infer that  $A^{\circ} \neq \emptyset$ . Therefore,  $A \in \mathcal{E}_{\mathcal{R}}$ , and thus assertion (1) is true.

While, if  $A \in \mathcal{A}_{\mathcal{R}} \cap \mathcal{D}_{\mathcal{R}}$ , then we have  $A \in \mathcal{A}_{\mathcal{R}}$  and  $A \in \mathcal{D}_{\mathcal{R}}$ . Therefore, by Remark 56 and Definition 1, we have  $A^{-\circ} \subseteq A^{\circ-}$  and  $A^{-} = X$ . Hence, since  $\mathcal{R} \neq \emptyset$ , and thus  $X^{\circ} = X$ , we can infer that  $X = X^{\circ} = A^{-\circ} \subseteq A^{\circ-}$ , and thus  $A^{\circ-} = X$ . Therefore,  $A^{\circ} \in \mathcal{D}_{\mathcal{R}}$ , and thus assertion (2) is also true.

*Remark 58* Assertions (1) and (2) can be reformulated in the forms that:

- (1)  $\mathcal{A}_{\mathcal{R}} \setminus \mathcal{E}_{\mathcal{R}} \subseteq \mathcal{N}_{\mathcal{R}}$ ;      (2)  $(\mathcal{A}_{\mathcal{R}} \cap \mathcal{D}_{\mathcal{R}})^{\circ} \subseteq \mathcal{D}_{\mathcal{R}}$ .

**Theorem 81** *If  $\mathcal{R}$  is a reflexive relator on  $X$ , then  $\mathcal{T}_{\mathcal{R}} \subseteq \mathcal{A}_{\mathcal{R}}$ .*

**Proof** If  $A \in \mathcal{T}_{\mathcal{R}}$ , then by Definition 3 and Theorem 26 we have  $A \subseteq A^{\circ}$  and  $A^{\circ} \subseteq A$ , and thus also  $A = A^{\circ}$ .

Now, by Theorem 26, we can also see that  $A^{-\circ} \subseteq A^{-} = A^{\circ-}$ , and thus by Remark 56 we also have  $A \in \mathcal{A}_{\mathcal{R}}$ .

**Theorem 82** *If  $\mathcal{R}$  is a topological relator on  $X$ , then  $\mathcal{T}_{\mathcal{R}}^s \subseteq \mathcal{A}_{\mathcal{R}}$ .*

**Proof** If  $A \in \mathcal{T}_{\mathcal{R}}^s$ , then by Definition 15, we have  $A \subseteq A^{\circ-}$ . Hence, by using the increasingness of the operation  $- \circ$  and Theorems 32, 31, and 26, we can see that

$$A^{-\circ} \subseteq A^{\circ--\circ} \subseteq A^{\circ-\circ} \subseteq A^{\circ-}.$$

Therefore, by Remark 56, we also have  $A \in \mathcal{A}_{\mathcal{R}}$ .

**Theorem 83** *We have*

- (1)  $\mathcal{A}_{\mathcal{R}} \cap \mathcal{T}_{\mathcal{R}}^p \subseteq \mathcal{T}_{\mathcal{R}}^s$ ;  
 (2)  $\mathcal{A}_{\mathcal{R}} \cap \mathcal{T}_{\mathcal{R}}^p = \mathcal{A}_{\mathcal{R}} \cap \mathcal{T}_{\mathcal{R}}^a$ ;      (3)  $\mathcal{A}_{\mathcal{R}} \cap \mathcal{T}_{\mathcal{R}}^s = \mathcal{A}_{\mathcal{R}} \cap \mathcal{T}_{\mathcal{R}}^b$ .

**Proof** If  $A \in \mathcal{A}_{\mathcal{R}} \cap \mathcal{T}_{\mathcal{R}}^p$ , then  $A \in \mathcal{T}_{\mathcal{R}}^p$  and  $A \in \mathcal{A}_{\mathcal{R}}$ . Hence, by using Definition 15 and Remark 56, we can see that  $A \subseteq A^{-\circ} \subseteq A^{\circ-}$ . Thus, by Remark 56, we also have  $A \in \mathcal{T}_{\mathcal{R}}^s$ . Therefore, inclusion (1) is true.

By Remark 56, it is clear that, for any  $A \subseteq X$ , the following assertions are equivalent:

$$(a) A \in \mathcal{A}_{\mathcal{R}}; \quad (b) A^{-\circ} = A^{-\circ} \cap A^{\circ-}; \quad (c) A^{\circ-} = A^{-\circ} \cup A^{\circ-}.$$

Hence, by Definition 15, it is clear that, for any  $A \in \mathcal{A}_{\mathcal{R}}$ , we have

$$A \in \mathcal{T}_{\mathcal{R}}^p \iff A \subseteq A^{-\circ} \iff A \subseteq A^{-\circ} \cap A^{\circ-} \iff A \in \mathcal{T}_{\mathcal{R}}^a,$$

and quite similarly

$$A \in \mathcal{T}_{\mathcal{R}}^s \iff A \subseteq A^{\circ-} \iff A \subseteq A^{-\circ} \cup A^{\circ-} \iff A \in \mathcal{T}_{\mathcal{R}}^b.$$

Therefore, equalities (2) and (3) are also true.

*Remark 59* By using Theorem 56 and assertion (3), we can also see that

$$\mathcal{A}_{\mathcal{R}} \cap \mathcal{T}_{\mathcal{R}}^p \subseteq \mathcal{A}_{\mathcal{R}} \cap (\mathcal{T}_{\mathcal{R}}^s \cup \mathcal{T}_{\mathcal{R}}^p) \subseteq \mathcal{A}_{\mathcal{R}} \cap \mathcal{T}_{\mathcal{R}}^b = \mathcal{A}_{\mathcal{R}} \cap \mathcal{T}_{\mathcal{R}}^s \subseteq \mathcal{T}_{\mathcal{R}}^s.$$

Therefore, assertion (3) is somewhat stronger than assertion (1).

## 20 Some Further Theorems on the Family $\mathcal{A}_{\mathcal{R}}$

The following theorem shows that, in contrast to  $\mathcal{T}_{\mathcal{R}}$  and  $\mathcal{T}_{\mathcal{R}}^k$ , the family  $\mathcal{A}_{\mathcal{R}}$  is closed under elementwise complementation.

**Theorem 84** For any  $A \in \mathcal{A}_{\mathcal{R}}$ , we have  $A^c \in \mathcal{A}_{\mathcal{R}}$ .

**Proof** By using Remark 56 and Theorem 2, we can see that

$$\begin{aligned} A \in \mathcal{A}_{\mathcal{R}} &\implies A^{-\circ} \subseteq A^{\circ-} \implies A^{-c-c} \subseteq A^{\circ c \circ c} \\ &\implies A^{\circ c \circ} \subseteq A^{-c-} \implies A^{c-\circ} \subseteq A^{c \circ-} \implies A^c \in \mathcal{A}_{\mathcal{R}}. \end{aligned}$$

*Remark 60* By using the practical, but ambiguous notation

$$\mathcal{A}^c = [\mathcal{A}]^c = \mathcal{C}_X[\mathcal{A}] = \{ \mathcal{C}_X(A) : A \in \mathcal{A} \}$$

the above theorem can be reformulated in the instructive form that  $\mathcal{A}_{\mathcal{R}}^c = \mathcal{A}_{\mathcal{R}}$ .

From our former Theorems 79 and 81–83, by using Theorem 84, we can immediately derive the following four theorems.

**Theorem 85** *We have*

$$\mathcal{N}_{\mathcal{R}} \cup \mathcal{N}_{\mathcal{R}}^c \subseteq \mathcal{A}_{\mathcal{R}}.$$

**Theorem 86** *If  $\mathcal{R}$  is a reflexive relator on  $X$ , then*

$$\mathcal{T}_{\mathcal{R}} \cup \mathcal{F}_{\mathcal{R}} \subseteq \mathcal{A}_{\mathcal{R}}.$$

*Remark 61* If  $\mathcal{R}$  is a reflexive relator on  $X$ , then we can also state that  $\mathcal{T}_{\mathcal{R}} \cap \mathcal{F}_{\mathcal{R}} \subseteq \mathcal{A}_{\mathcal{R}}(-, \circ) \cap \mathcal{A}_{\mathcal{R}}(\circ, -)$ .

Namely, if  $A \in \mathcal{T}_{\mathcal{R}} \cap \mathcal{F}_{\mathcal{R}}$ , then by using Definition 3 and Theorem 26 we can see that

$$A^{\circ-} = A^{-} = A = A^{\circ} = A^{-\circ}.$$

Therefore,  $A \in \mathcal{A}_{\mathcal{R}}(\circ, -)$ , and thus  $\mathcal{T}_{\mathcal{R}} \cap \mathcal{F}_{\mathcal{R}} \subseteq \mathcal{A}_{\mathcal{R}}(\circ, -)$  also holds.

**Theorem 87** *If  $\mathcal{R}$  is a topological relator on  $X$ , then*

$$\mathcal{T}_{\mathcal{R}}^s \cup \mathcal{F}_{\mathcal{R}}^s \subseteq \mathcal{A}_{\mathcal{R}}.$$

**Theorem 88** *We have*

- (1)  $\mathcal{A}_{\mathcal{R}} \cap \mathcal{F}_{\mathcal{R}}^p \subseteq \mathcal{F}_{\mathcal{R}}^s$ ;
- (2)  $\mathcal{A}_{\mathcal{R}} \cap \mathcal{F}_{\mathcal{R}}^p = \mathcal{A}_{\mathcal{R}} \cap \mathcal{F}_{\mathcal{R}}^a$ ;      (3)  $\mathcal{A}_{\mathcal{R}} \cap \mathcal{F}_{\mathcal{R}}^s = \mathcal{A}_{\mathcal{R}} \cap \mathcal{F}_{\mathcal{R}}^b$ .

Now, analogously to Theorem 84, we can also prove the following

**Theorem 89** *If  $\mathcal{R}$  is a topological relator on  $X$ , then for any  $A \in \mathcal{A}_{\mathcal{R}}$  we have*

- (1)  $\text{cl}_{\mathcal{R}}(A) \in \mathcal{A}_{\mathcal{R}}$ ;
- (2)  $\text{int}_{\mathcal{R}}(A) \in \mathcal{A}_{\mathcal{R}}$ .

**Proof** If  $A \in \mathcal{A}_{\mathcal{R}}$ , then by Remark 56 we have

$$A^{-\circ} \subseteq A^{\circ-}.$$

Moreover, by Theorems 32, 31 and 26, we also have

$$A^{--} \subseteq A^{-} \quad \text{and} \quad A \subseteq A^{-}.$$

Hence, by using the increasingness of the operations  $\circ$  and  $\circ -$ , we can infer that

$$A^{--\circ} \subseteq A^{-\circ} \subseteq A^{\circ-} \subseteq A^{-\circ-}.$$

Therefore, by the definition of  $\mathcal{A}_{\mathcal{R}}$ , we also have  $A^{-} \in \mathcal{A}_{\mathcal{R}}$ .

Assertion (2) can now be easily derived from (1) by using Theorems 2 and 84. Namely, by Theorem 2, we have  $A^\circ = A^{c-c}$ .

*Remark 62* Assertions (1) and (2) can be reformulated in the form that

$$\mathcal{A}_{\mathcal{R}}^- \cup \mathcal{A}_{\mathcal{R}}^\circ \subseteq \mathcal{A}_{\mathcal{R}}.$$

The following theorem shows that, analogously to the families  $\mathcal{T}_{\mathcal{R}}^\alpha$  and  $\mathcal{T}_{\mathcal{R}}^r$ , the family  $\mathcal{A}_{\mathcal{R}}$  may also be frequently closed under pairwise intersections.

**Theorem 90** *If  $\mathcal{R}$  is a topologically filtered, topological relator on  $X$ , then for any  $A, B \in \mathcal{A}_{\mathcal{R}}$  we have  $A \cap B \in \mathcal{A}_{\mathcal{R}}$ .*

*Proof* By Remark 56, we have

$$A^{-\circ} \subseteq A^{\circ-} \quad \text{and} \quad B^{-\circ} \subseteq B^{\circ-}.$$

Hence, by using the increasingness of the operations—and  $\circ$ , Theorems 32 and 31, Corollary 18 and Theorem 42, we can see that

$$\begin{aligned} (A \cap B)^{-\circ} &\subseteq A^{-\circ} \cap B^{-\circ} \subseteq A^{\circ-} \cap B^{\circ-} \subseteq (A^\circ \cap B^{\circ-})^- \\ &\subseteq (A^\circ \cap B^{\circ-})^- \subseteq (A^\circ \cap B^\circ)^{-\circ} \subseteq (A^\circ \cap B^\circ)^- = (A \cap B)^{\circ-}. \end{aligned}$$

Thus, by Remark 56, we also have  $A \cap B \in \mathcal{A}_{\mathcal{R}}$ .

Now, as a useful consequence of Theorems 84 and 90, we can also state

**Corollary 26** *If  $\mathcal{R}$  is a topologically filtered, topological relator on  $X$ , then for any  $A, B \in \mathcal{A}_{\mathcal{R}}$  we have (1)  $A \setminus B \in \mathcal{A}_{\mathcal{R}}$ ; (2)  $A \cup B \in \mathcal{A}_{\mathcal{R}}$ ; (3)  $A \Delta B \in \mathcal{A}_{\mathcal{R}}$ .*

*Proof* To prove these, recall that

$$A \setminus B = A \cap B^c, \quad A \cup B = (A^c \cap B^c)^c, \quad A \Delta B = (A \setminus B) \cup (B \setminus A).$$

*Remark 63* Thus, if  $\mathcal{R}$  is a topologically filtered, topological relator on  $X$ , then  $\mathcal{A}_{\mathcal{R}}$  forms an algebra of subsets of  $X$ .

The countable infinite set  $\mathbb{Q}$  of all rational numbers, used also by Levine [48] and Császár [15], shows that  $\mathcal{A}_{\mathcal{R}}$  is not in general a  $\sigma$ -algebra.

*Example 4* If  $X$  and  $\mathcal{R}$  are as in Example 3 and  $A_r = \{r\}$  for all  $r \in \mathbb{Q}$ , then  $A_r \in \mathcal{A}_{\mathcal{R}}$  for all  $r \in \mathbb{Q}$  such that  $\bigcup_{r \in \mathbb{Q}} A_r = \mathbb{Q} \notin \mathcal{A}_{\mathcal{R}}$ .

To check this, note that  $A_r \in \mathcal{F}_{\mathcal{R}}$ , and thus by Theorem 86 we also have  $A^r \in \mathcal{A}_{\mathcal{R}}$  for all  $r \in \mathbb{Q}$ . However,  $\mathbb{Q}^{-\circ} = \mathbb{R}^\circ = \mathbb{R}$ , but  $\mathbb{Q}^{\circ-} = \emptyset^- = \emptyset$ . Therefore, by Remark 56, we have  $\mathbb{Q} \notin \mathcal{A}_{\mathcal{R}}$ .

## 21 Characterizations of the Family $\mathcal{A}_{\mathcal{R}}$

By using the corresponding definitions, one can easily establish the following

**Theorem 91** *For any  $A \subseteq X$ , the following assertions are equivalent:*

- (1)  $A \in \mathcal{A}_{\mathcal{R}}$ ;
- (2) if  $x \in X$  and there exists  $R \in \mathcal{R}$  such that for any  $y \in R(x)$  and  $S \in \mathcal{R}$  we have  $S(y) \cap A \neq \emptyset$ , then for every  $U \in \mathcal{R}$  there exist  $z \in U(x)$  and  $V \in \mathcal{R}$  such that  $V(z) \subseteq A$ .

However, this intrinsic characterization of the family  $\mathcal{A}_{\mathcal{R}}$  cannot, certainly, be used to prove our former and subsequent theorems on  $\mathcal{A}_{\mathcal{R}}$ .

Therefore, analogously to the first part of [15, Theorem 2.3] of Császár, we shall also prove a more particular, but much deeper characterization theorem.

**Theorem 92** *If  $\mathcal{R}$  is a topologically filtered, topological relator on  $X$ , then for any  $A \subseteq X$  the following assertions are equivalent:*

- (1)  $A \in \mathcal{A}_{\mathcal{R}}$ ;
- (2) there exist  $V \in \mathcal{T}_{\mathcal{R}}$  and  $B \in \mathcal{N}_{\mathcal{R}}$  such that  $A = V \cup B$ ;
- (3) there exist  $V \in \mathcal{T}_{\mathcal{R}}^s$  and  $B \in \mathcal{N}_{\mathcal{R}}$  such that  $A = V \cup B$ .

**Proof** If (1) holds, then by Remark 56 we have  $A^{-\circ} \subseteq A^{\circ-}$ . Moreover, by Theorems 32 and 26, we can also note that  $A^{\circ} \subseteq A$ , and thus

$$A = A^{\circ} \cup (A \setminus A^{\circ}).$$

Therefore, by taking

$$V = A^{\circ} \quad \text{and} \quad B = A \setminus A^{\circ}$$

we can state that  $A = V \cup B$  with  $V \cap B = \emptyset$ . Moreover, by Theorem 32 and 31, we can also note that  $V = A^{\circ} \in \mathcal{T}_{\mathcal{R}}$ , and thus  $V \subseteq V^{\circ}$ .

On the other hand, by using Corollary 2, Theorem 42 and the inclusions

$$V^{-} = A^{\circ-} \supseteq A^{-\circ} \supseteq B^{-\circ},$$

we can see that

$$\begin{aligned} B = A \setminus V &\implies V \cap B = \emptyset \implies V \cap B^{-} = \emptyset \\ &\implies (V \cap B^{-})^{\circ} = \emptyset \implies V^{\circ} \cap B^{-\circ} = \emptyset \implies V \cap B^{-\circ} = \emptyset \\ &\implies V^{-} \cap B^{-\circ} = \emptyset \implies B^{-\circ} = \emptyset \implies B \in \mathcal{N}_{\mathcal{R}}. \end{aligned}$$

Therefore, (2) also holds.

Now, since the implication (2)  $\implies$  (3) is immediate from the inclusion  $\mathcal{I}_{\mathcal{R}} \subseteq \mathcal{I}_{\mathcal{R}}^s$ , we need only show that (3) also implies (1).

For this, note that if (3) holds, then by Theorems 82 and 79 we also have  $V \in \mathcal{A}_{\mathcal{R}}$  and  $B \in \mathcal{A}_{\mathcal{R}}$ . Hence, by Corollary 26, we can already infer that  $A = V \cup B \in \mathcal{A}_{\mathcal{R}}$ .

*Remark 64* From the inclusions  $\mathcal{I}_{\mathcal{R}} \subseteq \mathcal{I}_{\mathcal{R}}^\alpha \subseteq \mathcal{I}_{\mathcal{R}}^a \subseteq \mathcal{I}_{\mathcal{R}}^s$ , valid for any reflexive relator  $\mathcal{R}$  on  $X$ , it clear that in the above theorem we may write  $\mathcal{I}_{\mathcal{R}}^\alpha$  and  $\mathcal{I}_{\mathcal{R}}^a$  instead of  $\mathcal{I}_{\mathcal{R}}^s$ .

Moreover, if  $\mathcal{R}$  is a reflexive relator on  $X$ , then from assertion (2) of Theorem 92, by using the inclusion  $\mathcal{I}_{\mathcal{R}} \subseteq \mathcal{I}_{\mathcal{R}}^p$ , we can also infer that there exist  $V \in \mathcal{I}_{\mathcal{R}}^p$  and  $B \in \mathcal{N}_{\mathcal{R}}$  such that  $A = V \cup B$ .

However, this statement cannot certainly be used to obtain assertion (1) of Theorem 92. Namely, if  $\mathcal{R}$  is as in Theorem 92 and  $V \in \mathcal{I}_{\mathcal{R}}^p$ , then by Theorem 61 we can only state that there exist  $W \in \mathcal{I}_{\mathcal{R}}$  and  $C \in \mathcal{D}_{\mathcal{R}}$  such that  $V = W \cap C$ . And,  $C$  need not belong to  $\mathcal{A}_{\mathcal{R}}$ .

Namely, if  $\mathcal{R}$  is a nonvoid relator on  $X$  and in addition to  $C \in \mathcal{D}_{\mathcal{R}}$  we also have  $C \in \mathcal{A}_{\mathcal{R}}$ , then  $X = X^\circ = C^{\circ-} \subseteq C^{\circ-}$ . Therefore,  $X = C^{\circ-}$ , and thus  $C^\circ \in \mathcal{D}_{\mathcal{R}}$  which usually does not hold.

Now, analogously to [9, Lemma 1.3] of Chattopadhyay and Roy, we can also prove the following

**Theorem 93** *If  $\mathcal{R}$  is a topologically filtered, topological relator on  $X$ , then the following assertions are equivalent:*

- (1)  $\mathcal{A}_{\mathcal{R}} = \mathcal{I}_{\mathcal{R}}$ ;                      (2)  $\mathcal{A}_{\mathcal{R}} \subseteq \mathcal{I}_{\mathcal{R}}$ ;
- (3)  $\mathcal{I}_{\mathcal{R}} = \mathcal{F}_{\mathcal{R}}$ ;                      (4)  $\mathcal{I}_{\mathcal{R}} \subseteq \mathcal{F}_{\mathcal{R}}$ ;                      (5)  $\mathcal{F}_{\mathcal{R}} \subseteq \mathcal{I}_{\mathcal{R}}$ .

**Proof** By Theorem 81, it is clear that assertions (1) and (2) are equivalent even if  $\mathcal{R}$  is assumed only to be reflexive. Moreover, by using Theorem 11, we can easily see that assertions (3), (4) and (5) are always equivalent. Therefore, it is sufficient to prove only the equivalence of (2) and (4).

For this, note that if (4) does not hold, then there exists  $V \in \mathcal{I}_{\mathcal{R}}$  such that  $V \notin \mathcal{F}_{\mathcal{R}}$ . Hence, by using Theorem 11, we can infer that  $V^c \in \mathcal{F}_{\mathcal{R}}$  and  $V^c \notin \mathcal{I}_{\mathcal{R}}$ . Now, by Theorem 86, we can also state that  $V^c \in \mathcal{A}_{\mathcal{R}}$ . Thus, assertion (2) does not also hold. Therefore, (2) implies (4) even if  $\mathcal{R}$  is assumed to be only reflexive.

On the other hand, if  $A \in \mathcal{A}_{\mathcal{R}}$ , then by Theorem 92 there exist  $V \in \mathcal{I}_{\mathcal{R}}$  and  $B \in \mathcal{N}_{\mathcal{R}}$  such that

$$A = V \cup B.$$

Moreover, if (4) holds, then (5) also holds. Therefore, by using Theorems 32, 26 and 31, we can also see that

$$B \subseteq B^- = B^{-\circ} = \emptyset,$$

and thus  $B = \emptyset$ . Hence, we can already infer that  $A = V \cup B = V \cup \emptyset = V \in \mathcal{T}_{\mathcal{R}}$ . Therefore, (2) also holds.

## 22 Some Further Characterizations of the Family $\mathcal{A}_{\mathcal{R}}$

Analogously to [15, Theorem 2.3] of Császár and [29, Theorem 2.2] of Ganster, Reilly and Vamanamurthy, we can also prove the following two theorems.

**Theorem 94** *If  $\mathcal{R}$  is a nonvoid, topologically filtered, topological relator on  $X$ , then for any  $A \subseteq X$  the following assertions are equivalent:*

- (1)  $A \in \mathcal{A}_{\mathcal{R}}$ ;
- (2) there exists  $V \in \mathcal{T}_{\mathcal{R}}$  such that  $A \Delta V \in \mathcal{N}_{\mathcal{R}}$ ;
- (3) there exists  $V \in \mathcal{T}_{\mathcal{R}}$  such that  $A \setminus V \in \mathcal{N}_{\mathcal{R}}$  and  $V \setminus A \in \mathcal{N}_{\mathcal{R}}$ .

**Proof** If (1) holds, then by Theorem 92 there exist  $V \in \mathcal{T}_{\mathcal{R}}$  and  $B \in \mathcal{N}_{\mathcal{R}}$  such that  $A = V \cup B$ . Hence, we can infer that

$$A \setminus V \subseteq B \quad \text{and} \quad V \setminus A = \emptyset.$$

Now, since  $B \in \mathcal{N}_{\mathcal{R}}$  and  $\mathcal{N}_{\mathcal{R}}$  is descending, we can see that  $A \setminus V \in \mathcal{N}_{\mathcal{R}}$ . Moreover, since  $\mathcal{R}$  is nonvoid and reflexive, we can also see that  $\emptyset^{-\circ} = \emptyset^{\circ} = \emptyset$ , and thus  $V \setminus A = \emptyset \in \mathcal{N}_{\mathcal{R}}$ . Therefore, (3) also holds.

Conversely, if (3) holds, then by defining

$$W = V \setminus (V \setminus A)^{-}$$

and using Theorems 32, 32, and 11 and Corollary 15 we can see that

$$W = V \cap (V \cap A^c)^{-c} \in \mathcal{T}_{\mathcal{R}}.$$

Moreover, by using Theorems 2, 32, and 26, we can also see that

$$W = V \cap (V \cap A^c)^{c\circ} = V \cap (V^c \cup A)^{\circ} \subseteq V \cap (V^c \cup A) = V \cap A \subseteq A,$$

and thus  $A = W \cup (A \setminus W)$ . Furthermore, we can also note that

$$\begin{aligned} A \setminus W &= A \cap W^c = A \cap (V \cap (V \cap A^c)^{-c})^c = A \cap (V^c \cup (V \cap A^c)^{-}) = \\ &= (A \cap V^c) \cup (A \cap (V \cap A^c)^{-}) = (A \setminus V) \cup (A \cap (V \setminus A)^{-}) \subseteq (A \setminus V) \cup (V \setminus A)^{-}. \end{aligned}$$

Hence, by using the assumptions  $A \setminus V \in \mathcal{N}_{\mathcal{R}}$  and  $V \setminus A \in \mathcal{N}_{\mathcal{R}}$  and three basic properties of  $\mathcal{N}_{\mathcal{R}}$  established at the end of Section 15, we can already infer that  $A \setminus W \in \mathcal{N}_{\mathcal{R}}$ . Therefore, by Theorem 92, assertion (1) also holds.



Now, to complete the prove it remains only to note that, because of the equality  $A \Delta V = (A \setminus V) \cup (V \setminus A)$  and two basic properties of  $\mathcal{N}_{\mathcal{R}}$ , assertions (2) and (3) are also equivalent.

*Remark 65* Note that if,  $V \in \mathcal{T}_{\mathcal{R}}^{\kappa}$ , with  $\kappa = s, \alpha$  or  $a$ , then by Theorems 31, 11, and 69, we have  $W = V \cap (V \cap A^c)^{-c} \in \mathcal{T}_{\mathcal{R}}^{\kappa}$ . Therefore, by Theorem 92 and Remark 64, in Theorem 94 we can write  $\mathcal{T}_{\mathcal{R}}^{\kappa}$  instead of  $\mathcal{I}_{\mathcal{R}}$ .

However, it is now more important to note that, by a general definition of Kuratowski [45, p. 11], Theorem 94 can be briefly reformulated by stating that  $A \in \mathcal{A}_{\mathcal{R}}$  if and only if  $A$  is congruent to a member  $V$  of  $\mathcal{I}_{\mathcal{R}}$  modulo  $\mathcal{N}_{\mathcal{R}}$ .

Thus, under the assumptions of Theorem 94, by [45, p. 12], we can also state that  $A \in \mathcal{A}_{\mathcal{R}}$  if and only if there exist  $V \in \mathcal{I}_{\mathcal{R}}$  and  $B, C \in \mathcal{N}_{\mathcal{R}}$  such that  $A = (V \setminus B) \cup C$ . However, by Theorem 62, we have  $V \setminus B \in \mathcal{I}_{\mathcal{R}}^{\alpha}$ .

**Theorem 95** *If  $\mathcal{R}$  is a nonvoid, topologically filtered, topological relator on  $X$ , then for any  $A \subseteq X$  the following assertions are equivalent:*

- (1)  $A \in \mathcal{A}_{\mathcal{R}}$ ;
- (2) there exist  $V \in \mathcal{I}_{\mathcal{R}}^s$  and  $W \in \mathcal{F}_{\mathcal{R}}^s$  such that  $A = V \cap W$ ;
- (3) there exist  $V \in \mathcal{I}_{\mathcal{R}}^{\alpha}$  and  $W \in \mathcal{F}_{\mathcal{R}}^{\alpha}$  such that  $A = V \cap W$ .

**Proof** Because of the inclusion  $\mathcal{I}_{\mathcal{R}}^{\alpha} \subseteq \mathcal{I}_{\mathcal{R}}^s$  and its consequence  $\mathcal{F}_{\mathcal{R}}^{\alpha} \subseteq \mathcal{F}_{\mathcal{R}}^s$ , it is clear that (3) implies (2).

Moreover, if (2) holds, then by Definition 15 and Theorem 72, we have

$$V \subseteq V^{\circ-} \quad \text{and} \quad W^{\circ-} \subseteq W.$$

Hence, by using the inclusions  $A \subseteq V$  and  $A \subseteq W$ , and the corresponding properties of the operations—and  $\circ$ , we can infer that

$$A^{\circ-} \subseteq V^{\circ-} \subseteq V^{\circ--\circ} = V^{\circ-\circ} \subseteq V^{\circ-}$$

and

$$A^{\circ-} \subseteq W^{\circ-} \subseteq W, \quad \text{and thus} \quad A^{\circ-} = A^{\circ-\circ} \subseteq W^{\circ}.$$

Now, by using Corollary 18 and Theorem 42, we can already see that

$$A^{\circ-} \subseteq V^{\circ-} \cap W^{\circ} \subseteq (V^{\circ} \cap W^{\circ})^{-} = (V \cap W)^{\circ-} = A^{\circ-}.$$

Thus, by Remark 56, assertion (1) also holds.

Finally, if (1) holds, then by Theorem 84 we also have  $A^c \in \mathcal{A}_{\mathcal{R}}$ . Thus, by Theorem 92, there exist  $B \in \mathcal{N}_{\mathcal{R}}$  and  $V \in \mathcal{I}_{\mathcal{R}}$  such that  $A^c = B \cap V$ , and thus

$$A = B^c \cap V^c.$$

Moreover, by using Theorems 62 and 57, we can see that  $B^c \in \mathcal{T}_{\mathcal{R}}^\alpha$  and  $V \in \mathcal{F}_{\mathcal{R}}^\alpha$ , and hence also  $V^c \in \mathcal{F}_{\mathcal{R}}^\alpha$ . Therefore, assertion (3) also holds.

*Remark 66* If  $\mathcal{R}$  is nonvoid, topologically filtered topological relator on  $X$ , then necessary and sufficient conditions for a subset  $A$  of  $X$ , in order that  $A^{-\circ} = A^{\circ-}$ , i.e.,  $A \in \mathcal{A}(-, \circ) \cap \mathcal{A}(\circ, -)$  could hold, can be derived from the results of Levine [48], Choda and Matoba [10] and Chapman [7].

### 23 Lower and Upper Nearness Relations for Sets

If  $A, B \subseteq X$ , then  $A$  is said to be *near* to  $B$ , with respect to the relator  $\mathcal{R}$ , if  $A \in \text{Cl}_{\mathcal{R}}(B)$ . That is,  $R[A] \cap B \neq \emptyset$  for all  $R \in \mathcal{R}$ .

Note that if  $A \cap \text{cl}_{\mathcal{R}}(B) \neq \emptyset$ , then by Theorem 3 we also have  $A \in \text{Cl}_{\mathcal{R}}(B)$ . Therefore,  $A$  is near to  $B$  with respect to  $\mathcal{R}$ .

Moreover, if in particular  $\mathcal{R}$  is nonvoid and topologically fine, then by Theorem 19 the converse of the above implication is also true.

Therefore, for some particular purposes, the relation  $\text{cl}_{\mathcal{R}}$  can also be naturally used to define a reasonable nearness relation for sets.

However, our main motivation for the subsequent definition has mainly come from the famous observations of Levine [49] and Corson and Michael [11] that for a subset  $A$  of a topological space  $X(\mathcal{T})$  we can state that;

- (1)  $A$  is preopen if and only if there exists  $V \in \mathcal{T}$  such that  $A \subseteq V \subseteq A^-$ ;
- (2)  $A$  is semi-open if and only if there exists  $V \in \mathcal{T}$  such that  $V \subseteq A \subseteq V^-$ .

Now, by an observation of Andrijević [3], we can also state that  $A$  is  $\beta$ -open if and only if there exists a preopen subset  $V$  of  $X(\mathcal{T})$  such that  $A \subseteq V \subseteq A^-$ . Therefore, for the approximations of sets, we may also use generalized open sets.

The above observations strongly suggest that, analogously to the definition of the big lower and upper bound relations  $\text{Lb}_{\mathcal{R}}$  and  $\text{Ub}_{\mathcal{R}}$  mentioned in Remark 13, we may also naturally introduce the following

**Definition 20** If  $A, B \subseteq X$ , such that

$$A \subseteq B \subseteq \text{cl}_{\mathcal{R}}(A),$$

then we shall write

$$A \in \text{Ln}_{\mathcal{R}}(B) \qquad \text{and} \qquad B \in \text{Un}_{\mathcal{R}}(A).$$

Moreover, in this case, we shall say that  $A$  is near to  $B$  from below and  $B$  is near to  $A$  from above with respect to the relator  $\mathcal{R}$ .

*Remark 67* Thus, the relations  $\text{Ln}_{\mathcal{R}}$  and  $\text{Un}_{\mathcal{R}}$  are also not independent of each other. Namely, by the above definition, evidently have  $\text{Un}_{\mathcal{R}} = \text{Ln}_{\mathcal{R}}^{-1}$ .

Therefore, in the sequel, we shall only establish some basic properties of the relation  $\text{Ln}_{\mathcal{R}}$ .

**Theorem 96** *For any  $A, B \subseteq X$ , the following assertions are equivalent:*

- (1)  $A \in \text{Ln}_{\mathcal{R}}(B)$ ;
- (2)  $A \subseteq B$  and for each  $x \in B$  and  $R \in \mathcal{R}$  we have  $R(x) \cap A \neq \emptyset$ .

**Proof** By the corresponding definitions, we have

$$A \in \text{Ln}_{\mathcal{R}}(B) \iff A \subseteq B \subseteq A^-$$

and

$$B \subseteq A^- \iff \forall x \in B : x \in A^- \iff \forall x \in B : \forall R \in \mathcal{R} : R(x) \cap A \neq \emptyset.$$

Therefore, assertions (1) and (2) are also equivalent.

*Remark 68* In principle, this intrinsic characterization of the relation  $\text{Ln}_{\mathcal{R}}$  can be used to establish every possible properties of  $\text{Ln}_{\mathcal{R}}$ .

However, in most of the forthcoming proofs, it will be more convenient to use Definition 20 and the basic properties of the relation  $\text{cl}_{\mathcal{R}}$ .

**Theorem 97** *For any  $A, B \subseteq X$ , the following assertions are equivalent:*

- (1)  $A^c \in \text{Ln}_{\mathcal{R}}(B^c)$ ;
- (2)  $\text{int}_{\mathcal{R}}(A) \subseteq B \subseteq A$ .

**Proof** If (1) holds, then by Definition 20 we have  $A^c \subseteq B^c \subseteq A^{c-}$ . Hence, by using Theorem 2, we can infer that

$$A^\circ = A^{c-c} \subseteq B \subseteq A.$$

Thus, (2) also holds. The converse implication (2)  $\implies$  (1) can be proved quite similarly, by reversing the above argument.

*Remark 69* Analogously to Remark 2, assertion (1) can also be expressed in the more instructive form that  $A \in (\text{Ln}_{\mathcal{R}} \circ \mathcal{C}_X)^c(B)$ .

**Theorem 98** *If  $\mathcal{R}$  is a topological relator on  $X$ , then for any  $A, B \subseteq X$ , the following assertions are equivalent:*

- (1)  $A \in \text{Ln}_{\mathcal{R}}(B)$ ;
- (2)  $A \subseteq B$  and  $\text{cl}_{\mathcal{R}}(A) = \text{cl}_{\mathcal{R}}(B)$ .

**Proof** If (2) holds, then by Theorems 32 and 26, we can see that

$$A \subseteq B \subseteq B^- = A^-.$$

Therefore, by Definition 20, assertion (1) also holds even if  $\mathcal{R}$  is assumed to be only reflexive.

Conversely if (1) holds, then by Definition 20 we have  $A \subseteq B \subseteq A^-$ . Hence, by using the increasingness of  $-$  and Theorem 31, we can infer that

$$A^- \subseteq B^- \subseteq A^{--} \subseteq A^-.$$

Therefore,  $A^- = B^-$ , and thus (2) also holds even if  $\mathcal{R}$  is assumed to be only quasi-topological.

## 24 Some Set-Theoretic Properties of the Relation $\text{Ln}_{\mathcal{R}}$

To prove the following six theorems, it will be convenient to use Theorem 96.

**Theorem 99** For any  $A \subseteq X$ , the following assertions are equivalent: (1)  $A \in \text{Ln}_{\mathcal{R}}(\emptyset)$ ; (2)  $A = \emptyset$ .

**Proof** By Theorem 96, assertion (1) holds if and only if

$$(a) A \subseteq \emptyset; \quad (b) \forall x \in \emptyset : \forall R \in \mathcal{R} : R(x) \cap A \neq \emptyset.$$

Hence, since (a) is equivalent to (2), and (b) automatically holds, it is clear that (1) and (2) are equivalent.

**Theorem 100** For any  $B \subseteq X$ , the following assertions are equivalent:

$$(1) \emptyset \in \text{Ln}_{\mathcal{R}}(B); \quad (2) \text{ either } B = \emptyset \text{ or } \mathcal{R} = \emptyset.$$

**Proof** By Theorem 96, assertion (1) holds if and only if

$$(a) \emptyset \subseteq B; \quad (b) \forall x \in B : \forall R \in \mathcal{R} : R(x) \cap \emptyset \neq \emptyset.$$

Hence, since (a) automatically holds, and (b) can only hold if and only if (2) holds, it is clear that (1) and (2) are equivalent.

**Theorem 101** For any  $a \in X$  and  $B \subseteq X$ , the following assertions are equivalent:

$$(1) \{a\} \in \text{Ln}_{\mathcal{R}}(B); \\ (2) a \in B \text{ and for any } x \in B \text{ and } R \in \mathcal{R} \text{ we have } a \in R(x).$$

**Proof** By Theorem 96, assertion (1) holds if and only if

$$(a) \{a\} \subseteq B; \quad (b) \forall x \in B : \forall R \in \mathcal{R} : R(x) \cap \{a\} \neq \emptyset.$$

Hence, since (a) is equivalent to  $a \in B$ , and  $R(x) \cap \{a\} \neq \emptyset$  is equivalent to  $a \in R(x)$ , it is clear that (1) and (2) are equivalent.

**Theorem 102** For any  $A \subseteq X$  and  $b \in X$ , the following assertions are equivalent:

$$(1) A \in \text{Ln}_{\mathcal{R}}(\{b\}); \\ (2) \text{ either } A = \emptyset \text{ and } \mathcal{R} = \emptyset, \text{ or } A = \{b\} \text{ and for any } R \in \mathcal{R} \text{ we have } b \in R(b).$$

**Proof** By Theorem 96, assertion (1) holds if and only if

(a)  $A \subseteq \{b\}$ ;                      (b)  $\forall x \in \{b\} : \forall R \in \mathcal{R} : R(x) \cap A \neq \emptyset$ .

Hence, since (a) can hold if and only if either  $A = \emptyset$  or  $A = \{b\}$ , and moreover  $x \in \{b\}$  is equivalent to  $x = b$ , and  $R(b) \cap A \neq \emptyset$  can hold if and only if  $A = \{b\}$  and  $b \in R(b)$ , it is clear that (1) and (2) are equivalent.

**Theorem 103** *For any  $A \subseteq X$ , the following assertions are equivalent:*

- (1)  $A \in \text{Ln}_{\mathcal{R}}(X)$ ;  
 (2) for any  $x \in X$  and  $R \in \mathcal{R}$  we have  $R(x) \cap A \neq \emptyset$ .

**Proof** By Theorem 96, assertion (1) holds if and only if

(a)  $A \subseteq X$ ;                      (b)  $\forall x \in X : \forall R \in \mathcal{R} : R(x) \cap A \neq \emptyset$ .

Hence, since (a) automatically holds, and (b) coincides with (2), it is clear that (1) and (2) are equivalent.

**Theorem 104** *For any  $B \subseteq X$ , the following assertions are equivalent:*

- (1)  $X \in \text{Ln}_{\mathcal{R}}(B)$ ;                      (2)  $B = X$  and  $\mathcal{R}$  is non-partial.

**Proof** By Theorem 96, assertion (1) holds if and only if

(a)  $X \subseteq B$ ;                      (b)  $\forall x \in B : \forall R \in \mathcal{R} : R(x) \cap X \neq \emptyset$ .

Hence, since (a) is equivalent to  $B = X$ , and  $R(x) \cap X = R(x)$ , it is clear that now (b) means only that for any  $x \in X$  and  $R \in \mathcal{R}$  we have  $R(x) \neq \emptyset$ , i.e.,  $\mathcal{R}$  is non-partial. Therefore, (1) and (2) are equivalent.

To prove the following two theorems, we shall again use Definition 20 and some basic properties of the relation  $\text{cl}_{\mathcal{R}}$ .

**Theorem 105** *If  $(A_i)_{i \in I}$  and  $(B_i)_{i \in I}$  are families of subsets of  $X$  such that  $A_i \in \text{Ln}_{\mathcal{R}}(B_i)$  for all  $i \in I$ , then*

$$\bigcup_{i \in I} A_i \in \text{Ln}_{\mathcal{R}} \left( \bigcup_{i \in I} B_i \right).$$

**Proof** By Definition 20, we have  $A_i \subseteq B_i \subseteq A_i^-$  for all  $i \in I$ . Hence, by using the increasingness of the corresponding operations, we can infer that

$$\bigcup_{i \in I} A_i \subseteq \bigcup_{i \in I} B_i \subseteq \bigcup_{i \in I} A_i^- \subseteq \left( \bigcup_{i \in I} A_i \right)^-.$$

Therefore, by Definition 20, the required assertion is also true.

**Theorem 106** *If  $\mathcal{R}$  is a topologically filtered relator on  $X$  and  $A, B \subseteq X$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$ , then for any  $V \in \mathcal{T}_{\mathcal{R}}$  we have*

$$A \cap V \in \text{Ln}_{\mathcal{R}}(B \cap V).$$

**Proof** By Definition 20, we have  $A \subseteq B \subseteq A^-$ . Hence, by using Corollary 18, we can infer that

$$A \cap V \subseteq B \cap V \subseteq A^- \cap V \subseteq (A \cap V)^-.$$

Thus, by Definition 20, the required assertion is also true.

The next simple example shows that the condition  $V \in \mathcal{I}_{\mathcal{R}}$  of the above theorem cannot be omitted or replaced with  $V \in \mathcal{F}_{\mathcal{R}}$ .

*Example 5* If  $X$  and  $\mathcal{R}$  are as in Example 3, then for the sets

$$A = [0, 1[, \quad B = [0, 1] \quad \text{and} \quad V = \{1\},$$

we have  $A \in \text{Ln}_{\mathcal{R}}(B)$  and  $A \cap V \notin \text{Ln}_{\mathcal{R}}(B \cap V)$ .

## 25 Some Topological Properties of the Relation $\text{Ln}_{\mathcal{R}}$

**Theorem 107** *If  $A, B \subseteq X$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$  and  $\Phi$  is an increasing relation on  $\mathcal{P}(X)$  to  $X$ , then*

- (1)  $\Phi(A) \in \text{Ln}_{\mathcal{R}}(\Phi(B))$  if  $A \in \mathcal{A}(\text{cl}_{\mathcal{R}}, \Phi)$ ;
- (2)  $\Phi(B) \in \text{Ln}_{\mathcal{R}}(\Phi(\text{cl}_{\mathcal{R}}(A)))$  if  $B \in \mathcal{A}(\text{cl}_{\mathcal{R}}, \Phi)$ .

**Proof** By Definition 20 and the increasingness of  $-$ , we have

$$A \subseteq B \subseteq A^- \subseteq B^-.$$

Hence, by using the increasingness of  $\Phi$ , we can infer that

$$\Phi(A) \subseteq \Phi(B) \subseteq \Phi(A^-) \subseteq \Phi(B^-).$$

Moreover, if  $A \in \mathcal{A}(-, \Phi)$ , then by Definition 18 we have

$$\Phi(A^-) \subseteq \Phi(A)^-.$$

Hence, we can that see that

$$\Phi(A) \subseteq \Phi(B) \subseteq \Phi(A)^-.$$

Thus, by Definition 20, we have  $\Phi(A) \in \text{Ln}_{\mathcal{R}}(\Phi(B))$ .

While, if  $B \in \mathcal{A}(-, \Phi)$ , then by Definition 18 we have

$$\Phi(B^-) \subseteq \Phi(B)^-.$$

Hence, we can see that

$$\Phi(B) \subseteq \Phi(A^-) \subseteq \Phi(B)^-$$

Thus, by Definition 20, we have  $\Phi(B) \in \text{Ln}_{\mathcal{R}}(\Phi(A^-))$ .

*Remark 70* Note that if in particular  $\Phi(A) = A$  for all  $A \subseteq X$ , or  $\Phi(A) = \text{cl}_{\mathcal{R}}(A)$  for all  $A \subseteq X$ , then  $\Phi$  is increasing, and by Definition 18 we have  $A(\text{cl}_{\mathcal{R}}, \Phi) = \mathcal{P}(X)$ .

Therefore, as an immediate consequence of the above theorem, we can state

**Corollary 27** *If  $A, B \subseteq X$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$  then*

- (1)  $B \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))$ ;
- (2)  $\text{cl}_{\mathcal{R}}(A) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B))$ ;      (3)  $\text{cl}_{\mathcal{R}}(B) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)))$ .

*Remark 71* Note that if in particular  $\Phi(A) = \text{int}_{\mathcal{R}}(A)$  for all  $A \subseteq X$ , then  $\Phi$  is increasing, and by Definition 19 we have  $\mathcal{A}(\text{cl}_{\mathcal{R}}, \Phi) = \mathcal{A}_{\mathcal{R}}$ .

Therefore, as an immediate consequence of Theorem 107, we can also state

**Corollary 28** *If  $A, B \subseteq X$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$ , then*

- (1)  $\text{int}_{\mathcal{R}}(A) \in \text{Ln}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(B))$  if  $A \in \mathcal{A}_{\mathcal{R}}$ ;
- (2)  $\text{int}_{\mathcal{R}}(B) \in \text{Ln}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)))$  if  $B \in \mathcal{A}_{\mathcal{R}}$ .

From the above two corollaries, we can easily derive the following two corollaries.

**Corollary 29** *If  $A \in \mathcal{A}_{\mathcal{R}}$  and  $B \subseteq X$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$ , then*

- (1)  $\text{int}_{\mathcal{R}}(B) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A)))$ ;
- (2)  $\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(B)))$ ;
- (3)  $\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(B)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A))))$ .

**Proof** By Corollary 28, we have

$$A^\circ \in \text{Ln}_{\mathcal{R}}(B^\circ).$$

Hence, by using Corollary 27, we can infer that

$$B^\circ \in \text{Ln}_{\mathcal{R}}(A^{\circ-}), \quad A^{\circ-} \in \text{Ln}_{\mathcal{R}}(B^{\circ-}), \quad B^{\circ-} \in \text{Ln}_{\mathcal{R}}(A^{\circ--}).$$

Therefore, assertions (1), (2) and (3) are true.

**Corollary 30** *If  $A \subseteq X$  and  $B \in \mathcal{A}_{\mathcal{R}}$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$ , then*

- (1)  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(B)))$ ;
- (2)  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))))$ ;
- (3)  $\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(B))))$ .

**Proof** By Corollary 28, we have

$$B^\circ \in \text{Ln}_{\mathcal{R}}(A^{-\circ}).$$

Hence, by using Corollary 27, we can infer that

$$A^{-\circ} \in \text{Ln}_{\mathcal{R}}(B^{\circ-}), \quad B^{\circ-} \in \text{Ln}_{\mathcal{R}}(A^{-\circ-}), \quad A^{-\circ-} \in \text{Ln}_{\mathcal{R}}(B^{\circ--}).$$

Therefore, assertions (1), (2) and (3) are true.

*Remark 72* Note that if  $\mathcal{R}$  is a topological relator on  $X$ , then by Theorems 32, 26, and 31, we have  $A^{--} = A^-$  for all  $A \subseteq X$ .

Therefore, in this particular case, assertions (3) in Corollaries 27, 29, and 30 can be simplified by writing  $\text{cl}_{\mathcal{R}}$  instead of  $\text{cl}_{\mathcal{R}} \text{cl}_{\mathcal{R}}$ .

## 26 Some Further Topological Properties of the Relation $\text{Ln}_{\mathcal{R}}$

From the above corollaries, by using the results of Section 20, we can easily derive several further topological properties of the relation  $\text{Ln}_{\mathcal{R}}$ .

For instance, by Corollaries 28 and 29 and Theorem 85, we can at once state the following two theorems.

**Theorem 108** *If  $A \in \mathcal{N}_{\mathcal{R}} \cup \mathcal{N}_{\mathcal{R}}^c$  and  $B \subseteq X$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$ , then*

- (1)  $\text{int}_{\mathcal{R}}(A) \in \text{Ln}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(B))$ ;
- (2)  $\text{int}_{\mathcal{R}}(B) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A)))$ ;
- (3)  $\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(B)))$ ;
- (4)  $\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(B)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A))))$ .

**Theorem 109** *If  $A \subseteq X$  and  $B \in \mathcal{N}_{\mathcal{R}} \cup \mathcal{N}_{\mathcal{R}}^c$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$ , then*

- (1)  $\text{int}_{\mathcal{R}}(B) \in \text{Ln}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)))$ ;
- (2)  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(B)))$ ;
- (3)  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))))$ ;
- (4)  $\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(B))))$ .

*Remark 73* Note that if in particular  $\mathcal{R}$  is reflexive (resp. topological), then in the above two theorems, by Theorem 86 (resp. 87), we may write  $\mathcal{I}_{\mathcal{R}} \cup \mathcal{F}_{\mathcal{R}}$  (resp.  $\mathcal{I}_{\mathcal{R}}^s \cup \mathcal{F}_{\mathcal{R}}^s$ ) instead of  $\mathcal{N}_{\mathcal{R}} \cup \mathcal{N}_{\mathcal{R}}^c$ .

Moreover, if  $\mathcal{R}$  is topological, then assertions (3) in the above two theorems can also be simplified by writing  $\text{cl}_{\mathcal{R}}$  instead of  $\text{cl}_{\mathcal{R}} \text{cl}_{\mathcal{R}}$ .

More importantly, concerning topological relators, we can also prove the following two theorems.

**Theorem 110** *If  $\mathcal{R}$  is a topological relator on  $X$ ,  $A \in \mathcal{A}$  and  $B \subseteq X$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$ , then*

- (1)  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)) \in \text{Ln}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B)))$ ;
- (2)  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))))$ ;



- (3)  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B))))$ ;
- (4)  $\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B))))$ ;
- (5)  $\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B))) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))))$ .

**Proof** By Corollary 27 and Theorem 89, we have

$$A^- \in \text{Ln}_{\mathcal{R}}(B^-) \quad \text{and} \quad A^- \in \mathcal{A}_{\mathcal{R}}.$$

Hence, by using Corollary 28 and 29, we can infer that

$$A^{-\circ} \in \text{Ln}_{\mathcal{R}}(B^{-\circ})$$

and

$$B^{-\circ} \in \text{Ln}_{\mathcal{R}}(A^{-\circ-}); \quad A^{-\circ-} \in \text{Ln}_{\mathcal{R}}(B^{-\circ-}); \quad B^{-\circ-} \in \text{Ln}_{\mathcal{R}}(A^{-\circ-}).$$

Thus, assertions (1), (2), (4) and (5) are true.

Moreover, by Corollary 27, we now also have  $B \in \text{Ln}_{\mathcal{R}}(A^-)$ . Hence, by using Corollary 30, we can infer that  $A^{-\circ} \in \text{Ln}_{\mathcal{R}}(B^{-\circ-})$ . Thus, assertion (3) is also true.

**Theorem 111** *If  $\mathcal{R}$  is a topological relator on  $X$ ,  $A \subseteq X$  and  $B \in \mathcal{A}_{\mathcal{R}}$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$ , then*

- (1)  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B)) \in \text{Ln}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)))$ ;
- (2)  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))))$ ;
- (3)  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B))))$ ;
- (4)  $\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B))))$ ;
- (5)  $\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B))) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))))$ .

**Proof** By Corollary 27 and Theorem 89, we have

$$A^- \in \text{Ln}_{\mathcal{R}}(B^-) \quad \text{and} \quad B^- \in \mathcal{A}_{\mathcal{R}}.$$

Hence, by using Corollary 28 and 30, we can infer that

$$B^{-\circ} \in \text{Ln}_{\mathcal{R}}(A^{-\circ})$$

and

$$A^{-\circ} \in \text{Ln}_{\mathcal{R}}(B^{-\circ-}); \quad B^{-\circ} \in \text{Ln}_{\mathcal{R}}(A^{-\circ-}); \quad A^{-\circ-} \in \text{Ln}_{\mathcal{R}}(B^{-\circ-}).$$

Thus, assertions (1), (3), (2) and (4) are true.

Moreover, by Corollary 27, we now also have  $B^- \in \text{Ln}_{\mathcal{R}}(A^-)$ . Hence, by using Corollary 29, we can infer that  $B^{-\circ-} \in \text{Ln}_{\mathcal{R}}(A^{-\circ-})$ . Thus, assertion (5) is also true.

## 27 Some Relation-Theoretic Properties of the Relation $\text{Ln}_{\mathcal{R}}$

By using Definition 20, we can also easily establish the following two theorems.

**Theorem 112** *If  $\mathcal{R}$  is a reflexive relator on  $X$ , then  $\text{Ln}_{\mathcal{R}}$  is a reflexive relation on  $\mathcal{P}(X)$ .*

**Proof** Namely, if  $A \subseteq X$ , then we trivially have  $A \subseteq A$ . Moreover, by Theorem 26, we also have  $A \subseteq A^-$ . Thus, by Definition 20, we have  $A \in \text{Ln}_{\mathcal{R}}(A)$ .

*Remark 74* Actually, by the corresponding definitions and Theorem 26, the converse of the above theorem is also true.

For, if the relation  $\text{Ln}_{\mathcal{R}}$  is reflexive, then for any  $A \subseteq X$  we have  $A \in \text{Ln}_{\mathcal{R}}(A)$ , and thus  $A \subseteq A^-$ . Therefore, by Theorem 26, the relator  $\mathcal{R}$  is reflexive.

**Theorem 113** *If  $\mathcal{R}$  is a quasi-topological relator on  $X$ , then  $\text{Ln}_{\mathcal{R}}$  is a transitive relation on  $\mathcal{P}(X)$ .*

**Proof** If  $A, B, C \subseteq X$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$  and  $B \in \text{Ln}_{\mathcal{R}}(C)$ , then by Definition 20 we have

$$A \subseteq B \subseteq A^- \quad \text{and} \quad B \subseteq C \subseteq B^-.$$

Moreover, by using Theorem 31 and the increasingness of  $-$ , we can see that

$$B^- \subseteq A^{--} \subseteq A^-.$$

Hence, we can already see that  $A \subseteq C \subseteq A^-$ , and thus  $A \in \text{Ln}_{\mathcal{R}}(C)$  also holds.

Thus, since the relation  $\text{Ln}_{\mathcal{R}}$  is always antisymmetric, we can also state

**Corollary 31** *If  $\mathcal{R}$  is a topological relator on  $X$ , then  $\text{Ln}_{\mathcal{R}}$  is a partial order relation on  $\mathcal{P}(X)$ .*

*Remark 75* Note that this corollary can also be immediately derived from Theorem 98.

Now, from Theorem 107 and its corollaries, by using Theorem 113, we can easily derive the following five theorems.

**Theorem 114** *If  $\mathcal{R}$  is a quasi-topological relator on  $X$  and  $A, B \in \mathcal{A}(\text{cl}_{\mathcal{R}}, \Phi)$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$ , then for any increasing relation on  $\mathcal{P}(X)$  to  $X$ , we have*

$$\Phi(A) \in \text{Ln}_{\mathcal{R}}(\Phi(\text{cl}_{\mathcal{R}}(A))).$$

**Proof** By Theorem 107, we have both

$$\Phi(A) \in \text{Ln}_{\mathcal{R}}(\Phi(B)) \quad \text{and} \quad \Phi(B) \in \text{Ln}_{\mathcal{R}}(\Phi(\text{cl}_{\mathcal{R}}(A))).$$

Moreover, by Theorem 113, the relation  $\text{Ln}_{\mathcal{R}}$  is transitive. Hence, it is clear that the required assertion is true.

**Theorem 115** *If  $\mathcal{R}$  is a quasi-topological relator on  $X$  and  $A, B \subseteq X$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$  then*

- (1)  $B \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B))$ ;      (2)  $A \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B))$ ;  
 (3)  $\text{cl}_{\mathcal{R}}(A) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)))$ .

**Proof** Assertion (3) is the  $\Phi = \text{cl}_{\mathcal{R}}$  particular case of Theorem 114. Assertion (1) can be derived from Corollary 27 by using Theorem 113. While, assertion (2) is an immediate consequence of the assumption  $A \in \text{Ln}_{\mathcal{R}}(B)$ , assertion (1) and Theorem 113.

The following three theorems can be derived quite similarly from Corollaries 28, 29, and 30 by using Theorem 113.

**Theorem 116** *If  $\mathcal{R}$  is a quasi-topological relator on  $X$  and  $A \in \mathcal{A}_{\mathcal{R}}$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$  for some  $B \in \mathcal{A}_{\mathcal{R}}$ , then*

$$\text{int}_{\mathcal{R}}(A) \in \text{Ln}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))).$$

**Theorem 117** *If  $\mathcal{R}$  is a quasi-topological relator on  $X$  and  $A \in \mathcal{A}_{\mathcal{R}}$  and  $B \subseteq X$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$ , then*

- (1)  $\text{int}_{\mathcal{R}}(B) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(B)))$ ;  
 (2)  $\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(A))))$ .

**Theorem 118** *If  $\mathcal{R}$  is a quasi-topological relator on  $X$  and  $A \subseteq X$  and  $B \in \mathcal{A}_{\mathcal{R}}$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$ , then*

- (1)  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))))$ ;  
 (2)  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(B))))$ .

Moreover, by using Theorem 113, from the results of Section 26, we can also derive some similar theorems.

For instance, from Theorems 110 and 111, by using Theorem 113, we can immediately derive the following two theorems.

**Theorem 119** *If  $\mathcal{R}$  is a topological relator on  $X$  and  $A \in \mathcal{A}_{\mathcal{R}}$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$  for some  $B \subseteq X$ , then*

$$\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)))).$$

**Theorem 120** *If  $\mathcal{R}$  is a topological relator on  $X$ ,  $A \subseteq X$  and  $B \in \mathcal{A}_{\mathcal{R}}$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$ , then*

- (1)  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))))$ ;  
 (2)  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B))))$ .

From the latter two theorems, by using Theorem 113, we can immediately derive

**Theorem 121** *If  $\mathcal{R}$  is a topological relator on  $X$  and  $A, B \subseteq X$  such that such that  $A \in \text{Ln}_{\mathcal{R}}(B)$  and either  $A \in \mathcal{A}_{\mathcal{R}}$  or  $B \in \mathcal{A}_{\mathcal{R}}$ , then*

- (1)  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(A))))$ ;
- (2)  $\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B)) \in \text{Ln}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(\text{int}_{\mathcal{R}}(\text{cl}_{\mathcal{R}}(B))))$ .

## 28 Lower and Upper Nearness Closures of Families of Sets

By using our former nearness relations and a very particular case of the induced topological closure  $\text{cl}_{\mathcal{R}}$ , we may naturally introduce

**Definition 21** For any  $\mathcal{A} \subseteq \mathcal{P}(X)$ , the families

$$\mathcal{A}^{\ell} = \mathcal{A}^{\ell_{\mathcal{R}}} = \text{cl}_{\text{Ln}_{\mathcal{R}}}(\mathcal{A}) \quad \text{and} \quad \mathcal{A}^u = \mathcal{A}^{u_{\mathcal{R}}} = \text{cl}_{\text{Un}_{\mathcal{R}}}(\mathcal{A})$$

will be called the *lower and upper nearness closures* of the family  $\mathcal{A}$  with respect to the relator  $\mathcal{R}$ , respectively.

Thus, by Theorem 5, Remark 67 and Definition 1, we can at once state the following two theorems.

**Theorem 122** *For any  $\mathcal{A} \subseteq \mathcal{P}(X)$ , we have*

- (1)  $\mathcal{A}^{\ell} = \text{Ln}_{\mathcal{R}}^{-1}[\mathcal{A}]$ ;
- (2)  $\mathcal{A}^u = \text{Ln}_{\mathcal{R}}[\mathcal{A}]$ .

**Theorem 123** *For any  $A \subseteq X$  and  $\mathcal{A} \subseteq \mathcal{P}(X)$ , we have*

- (1)  $A \in \mathcal{A}^{\ell}$  if and only if  $\mathcal{A} \cap \text{Ln}_{\mathcal{R}}(A) \neq \emptyset$ ;
- (2)  $A \in \mathcal{A}^u$  if and only if  $\mathcal{A} \cap \text{Ln}_{\mathcal{R}}^{-1}(A) \neq \emptyset$ .

Hence, by the corresponding definitions, it is clear that we also have the following

**Theorem 124** *For any  $A \subseteq X$  and  $\mathcal{A} \subseteq \mathcal{P}(X)$ , we have*

- (1)  $A \in \mathcal{A}^{\ell}$  if and only if there exists  $V \in \mathcal{A}$  such that  $V \in \text{Ln}_{\mathcal{R}}(A)$ ;
- (2)  $A \in \mathcal{A}^u$  if and only if there exists  $V \in \mathcal{A}$  such that  $A \in \text{Ln}_{\mathcal{R}}(V)$ .

From this theorem, by using Definition 20 and Theorem 96, we can immediately derive the following two theorems.

**Theorem 125** *For any  $A \subseteq X$  and  $\mathcal{A} \subseteq \mathcal{P}(X)$ , we have*

- (1)  $A \in \mathcal{A}^{\ell}$  if and only if there exists  $V \in \mathcal{A}$  such that  $V \subseteq A \subseteq \text{cl}_{\mathcal{R}}(V)$ ;
- (2)  $A \in \mathcal{A}^u$  if and only if there exists  $V \in \mathcal{A}$  such that  $A \subseteq V \subseteq \text{cl}_{\mathcal{R}}(A)$ .

**Theorem 126** For any  $A \subseteq X$  and  $\mathcal{A} \subseteq \mathcal{P}(X)$ , we have

- (1)  $A \in \mathcal{A}^\ell$  if and only if there exists  $V \in \mathcal{A}$  such that  $V \subseteq A$ , and for any  $x \in A$  and  $R \in \mathcal{R}$  we have  $R(x) \cap V \neq \emptyset$ ;
- (2)  $A \in \mathcal{A}^u$  if and only if there exists  $V \in \mathcal{A}$  such that  $A \subseteq V$ , and for any  $x \in V$  and  $R \in \mathcal{R}$  we have  $R(x) \cap A \neq \emptyset$ .

**Proof** To prove the “only if part” of (2), note that if  $A \in \mathcal{A}^u$ , then by Theorem 124 there exists  $V \in \mathcal{A}$  such that  $A \in \text{Ln}_{\mathcal{R}}(V)$ . Hence, by Theorem 96, we can see that  $A \subseteq V$  and for any  $x \in V$  and  $R \in \mathcal{R}$  we have  $R(x) \cap A \neq \emptyset$ .

*Remark 76* The most important particular cases of Definition 21 are when  $\mathcal{A} = \tau_{\mathcal{S}}$  or  $\mathcal{I}_{\mathcal{S}}$  for some nonvoid relator  $\mathcal{S}$  on  $X$ . That is,  $\mathcal{A}$  is a minimal structure or a generalized topology on  $X$ . (See [85].)

In our former papers [67, 68], the members of the families  $\mathcal{F}_{\mathcal{R}}^\ell$ ,  $\mathcal{F}_{\mathcal{R}}^u$ ,  $\mathcal{F}_{\mathcal{R}}^{su}$  and  $\mathcal{T}_{\mathcal{R}}^{p\ell}$  have been called the *topologically quasi-open, pseudo-open,  $\gamma$ -open and  $\delta$ -open subsets* of the relator space  $X(\mathcal{R})$ , respectively.

## 29 Some Set-Theoretic Properties of the Families $\mathcal{A}^\ell$ and $\mathcal{A}^u$

By using Theorems 124 and 99–104, we can easily prove the following three theorems.

**Theorem 127** For any  $A \subseteq X$  and  $\mathcal{A} \subseteq \mathcal{P}(X)$  we have

- (1)  $\emptyset \in \mathcal{A}^\ell$  if and only if  $\emptyset \in \mathcal{A}$ ;
- (2)  $\emptyset \in \mathcal{A}^u$  if and only if either  $\emptyset \in \mathcal{A}$ , or  $\mathcal{A} \neq \emptyset$  and  $\mathcal{R} = \emptyset$ .

**Proof** By Theorems 124 and 99, we have

$$\emptyset \in \mathcal{A}^\ell \iff \exists V \in \mathcal{A} : V \in \text{Ln}_{\mathcal{R}}(\emptyset) \iff \exists V \in \mathcal{A} : V = \emptyset \iff \emptyset \in \mathcal{A}.$$

While, by Theorems 124 and 100, we have

$$\begin{aligned} \emptyset \in \mathcal{A}^u &\iff \exists V \in \mathcal{A} : \emptyset \in \text{Ln}_{\mathcal{R}}(V) \iff \\ &\iff \exists V \in \mathcal{A} : V = \emptyset \text{ or } \mathcal{R} = \emptyset \iff \emptyset \in \mathcal{A} \text{ or } (\mathcal{A} \neq \emptyset, \mathcal{R} = \emptyset). \end{aligned}$$

**Theorem 128** For any  $a \in X$  and  $\mathcal{A} \subseteq \mathcal{P}(X)$  we have

- (1)  $\{a\} \in \mathcal{A}^\ell$  if  $\{a\} \in \mathcal{A}$  and for any  $R \in \mathcal{R}$  we have  $a \in R(a)$ ;
- (2)  $\{a\} \in \mathcal{A}^u$  if and only if there exists  $V \in \mathcal{A}$  such that  $a \in V$  and for any  $x \in V$  and  $R \in \mathcal{R}$  we have  $a \in R(x)$ .

**Proof** By Theorems 124 and 101, we have

$$\{a\} \in \mathcal{A}^u \iff \exists V \in \mathcal{A} : \{a\} \in \text{Ln}_{\mathcal{R}}(V) \iff \\ \exists V \in \mathcal{A} : a \in V, \forall x \in V : \forall R \in \mathcal{R} : a \in R(x).$$

While, by Theorems 124 and 102, we have

$$\{a\} \in \mathcal{A}^l \iff \exists V \in \mathcal{A} : V \in \text{Ln}_{\mathcal{R}}(\{a\}) \iff \\ \exists V \in \mathcal{A} : (V = \emptyset, \mathcal{R} = \emptyset) \text{ or } (V = \{a\}, \forall R \in \mathcal{R} : a \in R(a)).$$

**Theorem 129** For any  $A \subseteq X$  and  $\mathcal{A} \subseteq \mathcal{P}(X)$  we have

- (1)  $X \in \mathcal{A}^l$  if  $X \in \mathcal{A}$  and  $\mathcal{R}$  non-partial;
- (2)  $X \in \mathcal{A}^u$  if and only if  $X \in \mathcal{A}$  and  $\mathcal{R}$  non-partial.

**Proof** By Theorems 124 and 104, we have

$$X \in \mathcal{A}^u \iff \exists V \in \mathcal{A} : X \in \text{Ln}_{\mathcal{R}}(V) \iff \exists V \in \mathcal{A} : V = \emptyset, \mathcal{R} \text{ is non-partial.}$$

While, by Theorems 124 and 103, we have

$$X \in \mathcal{A}^l \iff \exists V \in \mathcal{A} : V \in \text{Ln}_{\mathcal{R}}(X) \iff \\ \exists V \in \mathcal{A} : \forall x \in X : \forall R \in \mathcal{R} : R(x) \cap V \neq \emptyset.$$

From Theorems 127 and 129, we can obtain the following

**Corollary 32** If  $\mathcal{R}$  is a non-partial relator on  $X$  and  $\mathcal{A}$  is a minimal structure on  $X$ , then  $\mathcal{A}^l$  and  $\mathcal{A}^u$  are also minimal structures on  $X$ .

In addition to this corollary, by using Theorems 124 and 105, we can also prove

**Theorem 130** If  $\mathcal{A} \subseteq \mathcal{P}(X)$  such that  $\mathcal{A}$  is closed under arbitrary unions, then  $\mathcal{A}^l$  and  $\mathcal{A}^u$  are also closed under arbitrary unions.

**Proof** If  $A_i \in \mathcal{A}^l$  for all  $i \in I$ , then by Theorem 124 for each  $i \in I$  there exists  $V_i \in \mathcal{A}$  such that  $V_i \in \text{Ln}_{\mathcal{R}}(A_i)$ . Hence, by using the assumed property of  $\mathcal{A}$  and Theorem 105, we can infer that

$$\bigcup_{i \in I} A_i \in \mathcal{A} \quad \text{and} \quad \bigcup_{i \in I} A_i \in \text{Ln}_{\mathcal{R}}\left(\bigcup_{i \in I} B_i\right).$$

Therefore, by Theorem 124, we also have  $\bigcup_{i \in I} A_i \in \mathcal{A}^l$ .

This proves the first statement of the theorem. The second statement can be proved quite similarly.

Now, as an immediate consequence of Corollary 32 and Theorem 130, we can also state

**Corollary 33** *If  $\mathcal{R}$  is a non-partial relator on  $X$  and  $\mathcal{A}$  is a generalized topology on  $X$ , then  $\mathcal{A}^\ell$  and  $\mathcal{A}^u$  are also generalized topologies on  $X$ .*

In this respect, it is also worth proving the following

**Theorem 131** *If  $\mathcal{A}$  is a stack on  $X$ , then*

$$(1) \mathcal{A}^\ell \subseteq \mathcal{A}; \quad (2) \mathcal{A}^{u-} \subseteq \mathcal{A}.$$

**Proof** For instance, if  $A \in \mathcal{A}^u$ , then by Theorem 125, there exists  $V \in \mathcal{A}$  such that  $A \subseteq V \subseteq A^-$ . Thus, since  $\mathcal{A}$  is ascending,  $A^- \in \mathcal{A}$  also holds. Therefore, assertion (2) is true.

Moreover, by using Theorems 124 and 97, we can also prove

**Theorem 132** *For any  $A \subseteq X$  and  $\mathcal{A} \subseteq \mathcal{P}(X)$ , we have*

- (1)  $A \in \mathcal{A}^{\ell c}$  if and only if there exists  $W \in \mathcal{A}^c$  such that  $\text{int}_{\mathcal{R}}(W) \subseteq A \subseteq W$ ;
- (2)  $A \in \mathcal{A}^{uc}$  if and only if there exists  $W \in \mathcal{A}^c$  such that  $\text{int}_{\mathcal{R}}(A) \subseteq W \subseteq A$ .

**Proof** If  $A \in \mathcal{A}^{\ell c}$ , then  $A^c \in \mathcal{A}^\ell$ . Thus, by Theorem 124, there exists  $V \in \mathcal{A}$  such that  $V \in \text{Ln}_{\mathcal{R}}(A^c)$ . Hence, by taking  $W = V^c$ , we can see that  $W \in \mathcal{A}^c$  such that  $W^c \in \text{Ln}_{\mathcal{R}}(A^c)$ . Therefore, by Theorem 97, we can also state that  $W^\circ \subseteq A \subseteq W$ .

Thus, we have proved the “only if part” of (1). The if part of (1) and assertion (2) can be proved quite similarly.

### 30 Intersection Properties of the Families $\mathcal{A}^\ell$ and $\mathcal{A}^u$

**Theorem 133** *If  $\mathcal{R}$  is a topologically filtered relator on  $X$ , and moreover  $U \in \mathcal{T}_{\mathcal{R}}$  and  $\mathcal{A} \subseteq \mathcal{P}(X)$  such that  $U \cap V \in \mathcal{A}$  for all  $V \in \mathcal{A}$ , then*

- (1)  $U \cap A \in \mathcal{A}^\ell$  for all  $A \in \mathcal{A}^\ell$ ;
- (2)  $U \cap A \in \mathcal{A}^u$  for all  $A \in \mathcal{A}^u$ .

**Proof** If  $A \in \mathcal{A}^\ell$ , then by Theorem 124 there exists  $V \in \mathcal{A}$  such that  $V \in \text{Ln}_{\mathcal{R}}(A)$ . Hence, by assumption, we can infer that  $U \cap V \in \mathcal{A}$ . Moreover, by using Theorem 106, we can see that

$$U \cap V \in \text{Ln}_{\mathcal{R}}(U \cap A).$$

Thus, by Theorem 124, we also have  $U \cap A \in \mathcal{A}^\ell$ . Therefore, assertion (1) is true.

While, if  $A \in \mathcal{A}^u$ , then by Theorem 124 there exists  $V \in \mathcal{A}$  such that  $A \in \text{Ln}_{\mathcal{R}}(V)$ . Hence, by assumption, we can infer that  $U \cap V \in \mathcal{A}$ . Moreover, by using Theorem 106, we can see that

$$U \cap A \in \text{Ln}_{\mathcal{R}}(U \cap V).$$

Thus, by Theorem 124, we also have  $U \cap A \in \mathcal{A}^u$ . Therefore, assertion (2) is also true.

Repeated applications of this theorem give the following

**Corollary 34** *If  $\mathcal{R}$  is a topologically filtered relator on  $X$ , and moreover  $U \in \mathcal{T}_{\mathcal{R}}$  and  $\mathcal{A} \subseteq \mathcal{P}(X)$  such that  $U \cap V \in \mathcal{A}$  for all  $V \in \mathcal{A}$ , then  $U \cap A \in \mathcal{A}^\kappa$  for all  $A \in \mathcal{A}^\kappa$  with  $\kappa = \ell\ell, uu, \ell u$  and  $u\ell$ .*

**Proof** By Theorem 133, for instance we have  $U \cap A \in \mathcal{A}^\ell$  for all  $A \in \mathcal{A}^\ell$ . Hence, by applying Theorem 133, to the family  $\mathcal{A}^\ell$  instead of  $\mathcal{A}$ , for instance we can infer that  $U \cap A \in \mathcal{A}^{\ell\ell}$  for all  $A \in \mathcal{A}^{\ell\ell}$ .

From Theorem 133 and Corollary 34, by Corollary 15, it is clear that in particular we also have

**Corollary 35** *If  $\mathcal{R}$  is a topologically filtered relator on  $X$  and  $U \in \mathcal{T}_{\mathcal{R}}$ , then  $U \cap A \in \mathcal{T}_{\mathcal{R}}^\kappa$  for all  $A \in \mathcal{T}_{\mathcal{R}}^\kappa$  with  $\kappa = \ell, u, \ell\ell, uu, \ell u$  and  $u\ell$ .*

**Proof** By Corollary 15, we have  $U \cap V \in \mathcal{T}_{\mathcal{R}}$  for all  $V \in \mathcal{T}_{\mathcal{R}}$ . Thus, by Theorem 133, for instance we have  $U \cap A \in \mathcal{T}_{\mathcal{R}}^\ell$  for all  $A \in \mathcal{T}_{\mathcal{R}}^\ell$ .

Moreover, from Theorem 133, by using Theorem 69, we can derive

**Corollary 36** *If  $\mathcal{R}$  is a topologically filtered relator on  $X$  and  $U \in \mathcal{T}_{\mathcal{R}}$ , then  $U \cap A \in \mathcal{T}_{\mathcal{R}}^{\kappa\nu}$  for all  $A \in \mathcal{T}_{\mathcal{R}}^{\kappa\mu}$  with  $\kappa = s, p, \alpha, \beta, a, b$  and  $\mu = \ell, u$ .*

**Proof** By Theorem 69, for instance we have  $U \cap A \in \mathcal{T}_{\mathcal{R}}^s$  for all  $A \in \mathcal{T}_{\mathcal{R}}^s$ . Thus, by Theorem 133, for instance we have  $U \cap A \in \mathcal{T}_{\mathcal{R}}^{s\ell}$  for all  $A \in \mathcal{T}_{\mathcal{R}}^{s\ell}$ .

**Remark 77** If  $\mathcal{R}$  is a nonvoid, topologically filtered, topological relator on  $X$ , then by Theorems 70, 71, and 144, we can see that

$$\mathcal{T}_{\mathcal{R}}^\alpha = \{ A \subseteq X : \forall B \in \mathcal{T}_{\mathcal{R}}^\ell : A \cap B \in \mathcal{T}_{\mathcal{R}}^\ell \}.$$

This fact can also be used to prove that now  $\mathcal{T}_{\mathcal{R}}^\alpha$  is an ordinary topology on  $X$ .

Other classes of generalized topologically open sets do not, in general, form topologies. For instance, the following example shows that the family  $\mathcal{T}_{\mathcal{R}}^p$  need not be a topology.

**Example 6** If  $X$  and  $\mathcal{R}$  are as in Example 3 and

$$A = \mathbb{Q} \quad \text{and} \quad B = \{1\} \cup \mathbb{Q}^c,$$

then  $A, B \in \mathcal{T}_{\mathcal{R}}^p$ , but  $A \cap B \notin \mathcal{T}_{\mathcal{R}}^p$ .

To check this, recall that families of all rational and irrational numbers are dense in  $\mathbb{R}$ . Therefore,

$$A^{-\circ} = \mathbb{R}^\circ = \mathbb{R} = X$$



and quite similarly  $B^{-\circ} = X$ . Thus, in particular  $A, B \in \mathcal{T}_{\mathcal{R}}^P$ . However,

$$(A \cap B)^{-\circ} = \{1\}^{-\circ} = \{1\}^{\circ} = \emptyset,$$

and thus  $A \cap B \notin \mathcal{T}_{\mathcal{R}}^P$ .

*Remark 78* In this respect, it is also worth mentioning that if  $\mathcal{R}$  is a topologically filtered, topological relator on  $X$ , then by Theorem 78 the family  $\mathcal{T}_{\mathcal{R}}^r$  is closed under pairwise intersections.

However, the family  $\mathcal{T}_{\mathcal{R}}^r$  does not also form a topology. Namely, in contrast to the various families of generalized topological open sets, it is not, in general, closed even under pairwise unions.

### 31 Some Algebraic and Topological Properties of the Operations $\ell$ and $u$

By using Theorems 122, 112, and 113, we can easily establish the following basic properties of the operations  $\ell$  and  $u$  introduced in Definition 21.

**Theorem 134** *The operations  $\ell$  and  $u$  are union-preserving.*

Thus, in particular, we can also state

**Corollary 37** *The operations  $\ell$  and  $u$  are increasing.*

**Theorem 135** *If  $\mathcal{R}$  is a reflexive relator on  $X$ , then the operations  $\ell$  and  $u$  are expansive.*

Hence, by using Theorem 131, we can derive

**Corollary 38** *If  $\mathcal{R}$  is a reflexive relator on  $X$  and  $\mathcal{A}$  is a stack on  $X$ , then  $\mathcal{A}^{\ell} = \mathcal{A}$ .*

**Theorem 136** *If  $\mathcal{R}$  is a quasi-topological relator on  $X$ , then the operations  $\ell$  and  $u$  are upper quasi-idempotent.*

*Proof* By Theorem 113,  $\text{Ln}_{\mathcal{R}}$  is a transitive relation on  $\mathcal{P}(X)$ . Hence, it is clear that  $\text{Ln}_{\mathcal{R}}^{-1}$  is also a transitive relation on  $\mathcal{P}(X)$ . Therefore,  $\text{Ln}_{\mathcal{R}}^{-1} \circ \text{Ln}_{\mathcal{R}}^{-1} \subseteq \text{Ln}_{\mathcal{R}}^{-1}$ . Hence, by using Theorem 122, we can already see that

$$\mathcal{A}^{\ell\ell} = \text{Ln}_{\mathcal{R}}^{-1} [ \text{Ln}_{\mathcal{R}}^{-1} [ \mathcal{A} ] ] = ( \text{Ln}_{\mathcal{R}}^{-1} \circ \text{Ln}_{\mathcal{R}}^{-1} ) [ \mathcal{A} ] \subseteq \text{Ln}_{\mathcal{R}}^{-1} [ \mathcal{A} ] = \mathcal{A}^{\ell}$$

for all  $\mathcal{A} \subseteq \mathcal{P}(X)$ . Thus, the operation  $\ell$  is quasi-idempotent. The corresponding assertion for  $u$  can be proved even more easily.

Now, as an immediate consequence of the above three theorems, we can state

**Corollary 39** *If  $\mathcal{R}$  is a topological relator on  $X$ , then  $\ell$  and  $u$  are union-preserving closure operations on  $\mathcal{P}(X)$ .*

By using Theorem 124, Corollaries 27, and 28 and Theorem 115, we can prove the following three theorems.

**Theorem 137** *For any  $\mathcal{A} \subseteq \mathcal{P}(X)$ , we have*

- (1)  $\mathcal{A}^\ell \subseteq \mathcal{A}^{-u}$ ;                      (2)  $\mathcal{A}^{u-} \subseteq \mathcal{A}^\ell$ ;
- (3)  $\mathcal{A}^{\ell-} \subseteq \mathcal{A}^{-\ell}$ ;                      (4)  $\mathcal{A}^{u-} \subseteq \mathcal{A}^{-u}$ .

**Proof** If  $A \in \mathcal{A}^\ell$ , then by Theorem 124 there exists  $V \in \mathcal{A}$  such that  $V \in \text{Ln}_{\mathcal{R}}(A)$ . Hence, by using Corollary 27, we can infer that  $A \in \text{Ln}_{\mathcal{R}}(V^-)$ . Now, since  $V^- \in \mathcal{A}^-$ , by Theorem 124 we can see that  $A \in \mathcal{A}^{-u}$ . Therefore,  $\mathcal{A}^\ell \subseteq \mathcal{A}^{-u}$ .

While, if  $A \in \mathcal{A}^{u-}$ , then by Theorem 124 there exists  $V \in \mathcal{A}$  such that  $A \in \text{Ln}_{\mathcal{R}}(V)$ . Hence, by using Corollary 27, we can infer that  $A^- \in \text{Ln}_{\mathcal{R}}(V^-)$ . Now, since  $V^- \in \mathcal{A}^-$ , by Theorem 124 we can see that  $A^- \in \mathcal{A}^{-u}$ . Therefore,  $\mathcal{A}^{u-} \subseteq \mathcal{A}^{-u}$  also holds.

Thus, we have proved assertions (1) and (4). The proof of assertions (2) and (3) are quite similar.

**Theorem 138** *We have*

- (1)  $\mathcal{A}_{\mathcal{R}}^{\ell\circ} \subseteq \mathcal{A}_{\mathcal{R}}^{\circ\ell}$ ;                      (2)  $\mathcal{A}_{\mathcal{R}}^{u-\circ} \subseteq \mathcal{A}_{\mathcal{R}}^{\circ\ell}$ .

**Proof** If  $A \in \mathcal{A}_{\mathcal{R}}^{\ell\circ}$ , then by Theorem 124 there exists  $V \in \mathcal{A}_{\mathcal{R}}$  such that  $V \in \text{Ln}_{\mathcal{R}}(A)$ . Hence, by using Corollary 28 we can infer that  $V^\circ \in \text{Ln}_{\mathcal{R}}(A^\circ)$ . Now, since  $V^\circ \in \mathcal{A}_{\mathcal{R}}^\circ$ , by Theorem 124 we can see that  $A^\circ \in \mathcal{A}_{\mathcal{R}}^{\circ\ell}$ . Therefore,  $\mathcal{A}_{\mathcal{R}}^{\ell\circ} \subseteq \mathcal{A}_{\mathcal{R}}^{\circ\ell}$ .

While, if  $A \in \mathcal{A}_{\mathcal{R}}^{u-\circ}$ , then by Theorem 124 there exists  $V \in \mathcal{A}_{\mathcal{R}}$  such that  $A \in \text{Ln}_{\mathcal{R}}(V)$ . Hence, by using Corollary 28, we can infer that  $V^\circ \in \text{Ln}_{\mathcal{R}}(A^{-\circ})$ . Now, since  $V^\circ \in \mathcal{A}_{\mathcal{R}}^\circ$ , by Theorem 124 we can see that  $A^{-\circ} \in \mathcal{A}_{\mathcal{R}}^{\circ\ell}$ . Therefore,  $\mathcal{A}_{\mathcal{R}}^{u-\circ} \subseteq \mathcal{A}_{\mathcal{R}}^{\circ\ell}$  also holds.

**Theorem 139** *If  $\mathcal{R}$  is a quasi-topological relator on  $X$ , then for any  $\mathcal{A} \subseteq \mathcal{P}(X)$  we have  $\mathcal{A}^{\ell-} \subseteq \mathcal{A}^\ell$ .*

**Proof** If  $A \in \mathcal{A}^{\ell-}$ , then by Theorem 124 there exists  $V \in \mathcal{A}$  such that  $V \in \text{Ln}_{\mathcal{R}}(A)$ . Hence, by using Theorem 115, we can infer that  $V \in \text{Ln}_{\mathcal{R}}(A^-)$ . Now, by Theorem 124, we can see that  $A^- \in \mathcal{A}^\ell$ . Therefore,  $\mathcal{A}^{\ell-} \subseteq \mathcal{A}^\ell$ .

**Theorem 140** *If  $\mathcal{R}$  is a quasi-topological relator on  $X$ , then for any  $\mathcal{A} \subseteq \mathcal{P}(X)$  we have*

- (1)  $\mathcal{A}^{\ell u} \subseteq \mathcal{A}^{-u}$ ;                      (2)  $\mathcal{A}^{\ell u-} \subseteq \mathcal{A}^\ell$ ;
- (3)  $\mathcal{A}^{u\ell} \subseteq \mathcal{A}^{-u}$ ;                      (4)  $\mathcal{A}^{u\ell-} \subseteq \mathcal{A}^\ell$ .

**Proof** If  $A \in \mathcal{A}^{\ell u}$ , then by Theorem 124 there exists  $B \in \mathcal{A}^{\ell}$  such that  $A \in \text{Ln}_{\mathcal{R}}(B)$ . Moreover, also by Theorem 124, there exists  $V \in \mathcal{A}$  such that  $V \in \text{Ln}_{\mathcal{R}}(B)$ . Hence, by using Corollary 27, we can infer that  $B \in \text{Ln}_{\mathcal{R}}(V^-)$ . Now, by using Theorem 113, we can see that  $A \in \text{Ln}_{\mathcal{R}}(V^-)$ . Thus, since  $V^- \in \mathcal{A}^-$ , by using Theorem 124, we can see that  $A^- \in \mathcal{A}^{-u}$ . Therefore,  $\mathcal{A}^{\ell u} \subseteq \mathcal{A}^{-u}$ .

While, if  $A \in \mathcal{A}^{u\ell}$ , then by Theorem 124 there exists  $B \in \mathcal{A}^u$  such that  $B \in \text{Ln}_{\mathcal{R}}(A)$ . Moreover, also by Theorem 124, there exists  $V \in \mathcal{A}$  such that  $B \in \text{Ln}_{\mathcal{R}}(V)$ . Hence, by using Corollary 27, we can infer that

$$V \in \text{Ln}_{\mathcal{R}}(B^-) \quad \text{and} \quad B^- \in \text{Ln}_{\mathcal{R}}(A^-).$$

Now, by using Theorem 113, we can see that  $V \in \text{Ln}_{\mathcal{R}}(A^-)$ . Thus, by using Theorem 124, we can see that  $A^- \in \mathcal{A}^{-\ell}$ . Therefore,  $\mathcal{A}^{u\ell} \subseteq \mathcal{A}^{-\ell}$ .

Thus, we have proved assertions (1) and (4). The proofs of assertions (2) and (3) are quite similar.

### 32 Nearness Closures of the Families $\mathcal{T}_{\mathcal{R}}$ , $\mathcal{T}_{\mathcal{R}}^s$ and $\mathcal{T}_{\mathcal{R}}^p$

The importance of Definition 21 lies mainly in the following theorems.

**Theorem 141** *We have*

$$(1) \mathcal{T}_{\mathcal{R}}^{\ell} \subseteq \mathcal{T}_{\mathcal{R}}^s; \quad (2) \mathcal{T}_{\mathcal{R}}^u \subseteq \mathcal{T}_{\mathcal{R}}^p.$$

**Proof** If  $A \in \mathcal{T}_{\mathcal{R}}^{\ell}$ , then by Theorem 125 there exists  $V \in \mathcal{T}_{\mathcal{R}}$  such that

$$V \subseteq A \subseteq V^-.$$

Hence, by using the definition of  $\mathcal{T}_{\mathcal{R}}$  and the increasingness of  $\circ$ , we can infer that  $V \subseteq V^{\circ} \subseteq A^{\circ}$ . Now, by using the increasingness of  $-$ , we can also see that

$$A \subseteq V^- \subseteq A^{\circ-}.$$

Therefore, by Definition 15, we also have  $A \in \mathcal{T}_{\mathcal{R}}^s$ .

While, if  $A \in \mathcal{T}_{\mathcal{R}}^u$ , then by Theorem 125, there exists  $V \in \mathcal{T}_{\mathcal{R}}$  such that

$$A \subseteq V \subseteq A^-.$$

Hence, by using the definition of  $\mathcal{T}_{\mathcal{R}}$  and the increasingness of  $\circ$ , we can infer that

$$A \subseteq V \subseteq V^{\circ} \subseteq A^{\circ-}.$$

Therefore, by Definition 15, we also have  $A \in \mathcal{T}_{\mathcal{R}}^p$ .

**Theorem 142** *If  $\mathcal{R}$  is a reflexive relator on  $X$ , then*

$$(1) \mathcal{T}_{\mathcal{R}}^{\ell} \subseteq \mathcal{T}_{\mathcal{R}}^{s\ell}; \quad (2) \mathcal{T}_{\mathcal{R}}^u \subseteq \mathcal{T}_{\mathcal{R}}^{pu}.$$

**Proof** By Theorems 135 and 141 and Corollary 37, we have

$$\mathcal{T}_{\mathcal{R}}^{\ell} \subseteq \mathcal{T}_{\mathcal{R}}^{\ell\ell} \subseteq \mathcal{T}_{\mathcal{R}}^{s\ell} \quad \text{and} \quad \mathcal{T}_{\mathcal{R}}^u \subseteq \mathcal{T}_{\mathcal{R}}^{uu} \subseteq \mathcal{T}_{\mathcal{R}}^{pu}.$$

**Theorem 143** *If  $\mathcal{R}$  is a quasi-topological relator on  $X$ , then*

$$(1) \mathcal{T}_{\mathcal{R}}^{s\ell} \subseteq \mathcal{T}_{\mathcal{R}}^s; \quad (2) \mathcal{T}_{\mathcal{R}}^{pu} \subseteq \mathcal{T}_{\mathcal{R}}^p.$$

**Proof** If  $A \in \mathcal{T}_{\mathcal{R}}^{s\ell}$ , then by Theorem 125 there exists  $V \in \mathcal{T}_{\mathcal{R}}^s$  such that

$$V \subseteq A \subseteq V^{-}.$$

Moreover, by Definition 15, we have  $V \subseteq V^{\circ-}$ . Hence, by using the increasingness of–and Theorem 31, we can infer that

$$V^{-} \subseteq V^{\circ--} \subseteq V^{\circ-}.$$

On the other hand, because of  $V \subseteq A$  and the increasingness of  $\circ-$ , we also have  $V^{\circ-} \subseteq A^{\circ-}$ . Therefore,  $V^{-} \subseteq A^{\circ-}$  also holds. Hence, because of  $A \subseteq V^{-}$ , we can already see that  $A \subseteq A^{\circ-}$ . Thus, by Definition 15, we also have  $A \in \mathcal{T}_{\mathcal{R}}^s$ . Therefore, assertion (1) is true.

While, if  $A \in \mathcal{T}_{\mathcal{R}}^{pu}$ , then by Theorem 125 there exists  $V \in \mathcal{T}_{\mathcal{R}}^p$  such that

$$A \subseteq V \subseteq A^{-}.$$

Moreover, by Definition 15, we have  $A \subseteq A^{-\circ}$ . Hence, by using that  $B \subseteq A$ , we can see that  $B \subseteq A^{-\circ}$ . Moreover, from the inclusion  $A \subseteq B^{-}$ , by using the increasingness of–and Theorem 31, we can infer that

$$A^{-} \subseteq B^{--} \subseteq B^{-}.$$

Hence, by using the increasingness of  $\circ$ , we can infer that  $A^{-\circ} \subseteq B^{-\circ}$ . Therefore, because of  $B \subseteq A^{-\circ}$ , we also have  $B \subseteq B^{-\circ}$ . Hence, by Definition 15, we can see that  $B \in \mathcal{T}_{\mathcal{R}}^p$  also holds. Therefore, assertion (2) is also true.

**Theorem 144** *If  $\mathcal{R}$  is a topological relator on  $X$ , then*

$$(1) \mathcal{T}_{\mathcal{R}}^{\ell} = \mathcal{T}_{\mathcal{R}}^s; \quad (2) \mathcal{T}_{\mathcal{R}}^u = \mathcal{T}_{\mathcal{R}}^p.$$

**Proof** If  $A \in \mathcal{T}_{\mathcal{R}}^s$ , then by Definition 15 we have  $A \subseteq A^{\circ-}$ . Hence, by taking  $V = A^{\circ}$ , we get  $A \subseteq V^{-}$ . Moreover, by using Theorems 32, 26, and 31, we can see that

$$V = A^\circ \subseteq A \quad \text{and} \quad V = A^\circ \in \mathcal{I}_{\mathcal{R}}.$$

Thus, by Theorem 125, we also have  $A \in \mathcal{I}_{\mathcal{R}}^\ell$ . This proves that  $\mathcal{I}_{\mathcal{R}}^s \subseteq \mathcal{I}_{\mathcal{R}}^\ell$ .

While, if  $A \in \mathcal{I}_{\mathcal{R}}^p$ , then by Definition 15 we have  $A \subseteq A^{-\circ}$ . Hence, by taking  $V = A^{-\circ}$ , we get  $A \subseteq V$ . Moreover, by using Theorems 32, 26, and 31, we can see that

$$V = A^{-\circ} \subseteq A^- \quad \text{and} \quad V = A^{-\circ} \in \mathcal{I}_{\mathcal{R}}.$$

Thus, by Theorem 125, we also have  $A \in \mathcal{I}_{\mathcal{R}}^u$ . This proves that  $\mathcal{I}_{\mathcal{R}}^p \subseteq \mathcal{I}_{\mathcal{R}}^u$ .

Now, by Theorem 141, we can see that assertions (1) and (2) are also true.

From Theorems 142 and 143, by using Theorem 144, we can immediately derive

**Corollary 40** *If  $\mathcal{R}$  is a topological relator on  $X$ , then*

$$(1) \mathcal{I}_{\mathcal{R}}^{s\ell} = \mathcal{I}_{\mathcal{R}}^s; \quad (2) \mathcal{I}_{\mathcal{R}}^{pu} = \mathcal{I}_{\mathcal{R}}^p.$$

*Remark 79* Note that this corollary can also be immediately derived from Theorem 144 by using Corollary 39.

### 33 Some Further Theorems on the Families $\mathcal{I}_{\mathcal{R}}^s$ and $\mathcal{I}_{\mathcal{R}}^p$

In addition to Corollary 40, we can also prove the following

**Theorem 145** *If  $\mathcal{R}$  is a topological relator, then*

$$(1) \mathcal{I}_{\mathcal{R}} = \mathcal{I}_{\mathcal{R}}^{s^\circ}; \quad (2) \mathcal{I}_{\mathcal{R}} = \mathcal{I}_{\mathcal{R}}^{p^\circ}.$$

**Proof** By Theorem 32,  $\mathcal{R}$  is reflexive and quasi-topological. Thus, if in particular  $V \in \mathcal{I}_{\mathcal{R}}$ , then by Theorem 57 we also have  $V \in \mathcal{I}_{\mathcal{R}}^s$ .

Moreover, by Definition 3 and Theorem 26, we also have  $V \subseteq V^\circ$  and  $V^\circ \subseteq V$ , and thus also  $V = V^\circ$ . Therefore,  $V \in \mathcal{I}_{\mathcal{R}}^{s^\circ}$ , and thus  $\mathcal{I}_{\mathcal{R}} \subseteq \mathcal{I}_{\mathcal{R}}^{s^\circ}$ .

On the other hand, by Theorem 31 we have  $A^\circ \in \mathcal{I}_{\mathcal{R}}$  for all  $A \subseteq X$ . Thus,  $\mathcal{P}(X)^\circ \subseteq \mathcal{I}_{\mathcal{R}}$ , and thus in particular  $\mathcal{I}_{\mathcal{R}}^{s^\circ} \subseteq \mathcal{I}_{\mathcal{R}}$  also holds.

Therefore, assertion (1) is true. The proof of assertion (2) is quite similar.

However, it is now more important to note that, by using Corollary 40, we can prove the following theorems.

**Theorem 146** *If  $\mathcal{R}$  is a topological relator on  $X$ , then*

$$(1) \mathcal{I}_{\mathcal{R}}^{s^-} \subseteq \mathcal{I}_{\mathcal{R}}^s; \quad (2) A^- \in \mathcal{I}_{\mathcal{R}}^p \text{ implies } A \in \mathcal{I}_{\mathcal{R}}^p \text{ for all } A \subseteq X.$$

**Proof** If  $A \subseteq X$ , then by Theorems 32 and 31, we have  $A \subseteq A^- \subseteq A^-$ . Hence, if  $A \in \mathcal{I}_{\mathcal{R}}^s$ , then by using Theorem 125 we can see that  $A^- \in \mathcal{I}_{\mathcal{R}}^{s\ell}$ . Thus, by Corollary 40, we also have  $A^- \in \mathcal{I}_{\mathcal{R}}^s$ . Therefore, assertion (1) is true.

On the other hand, if  $A \subseteq X$  such that  $A^- \in \mathcal{T}_{\mathcal{R}}^p$ , then from the inclusions  $A \subseteq A^- \subseteq A^-$  by using Theorem 125 we can see that  $A \in \mathcal{T}_{\mathcal{R}}^{pu}$ . Thus, by Corollary 40, we also have  $A \in \mathcal{T}_{\mathcal{R}}^p$ . Therefore, assertion (2) is also true.

From this theorem, by using Theorem 57, we can immediately derive

**Corollary 41** *If  $\mathcal{R}$  is a topological relator on  $X$ , then*

- (1)  $\mathcal{T}_{\mathcal{R}}^- \subseteq \mathcal{T}_{\mathcal{R}}^s$ ;      (2)  $A^- \in \mathcal{T}_{\mathcal{R}}$  implies  $A \in \mathcal{T}_{\mathcal{R}}^p$  for all  $A \subseteq X$ .

**Proof** Namely, by Theorem 32 and 57, we have  $\mathcal{T}_{\mathcal{R}} \subseteq \mathcal{T}_{\mathcal{R}}^s$  and  $\mathcal{T}_{\mathcal{R}} \subseteq \mathcal{T}_{\mathcal{R}}^p$ . Hence, by Theorem 146, it is clear that in particular assertions are also true.

*Remark 80* If  $\mathcal{R}$  is a topological relator on  $X$ , then in particular we also have  $\mathcal{P}(X)^{\circ-} \subseteq \mathcal{T}_{\mathcal{R}}^s$ .

Namely, by Theorems 32 and 31, we have  $\mathcal{P}(X)^{\circ} \subseteq \mathcal{T}_{\mathcal{R}}$ . Hence, by using Corollary 41, we can infer that  $\mathcal{P}(X)^{\circ-} \subseteq \mathcal{T}_{\mathcal{R}}^- \subseteq \mathcal{T}_{\mathcal{R}}^s$ .

**Theorem 147** *If  $\mathcal{R}$  is a topological relator on  $X$ , then  $\mathcal{A} = \mathcal{T}_{\mathcal{R}}^s$  is the smallest subset of  $\mathcal{P}(X)$  such that*

- (1)  $\mathcal{T}_{\mathcal{R}} \subseteq \mathcal{A}$ ;      (2)  $\mathcal{A}^l \subseteq \mathcal{A}$ .

**Proof** If  $\mathcal{A} = \mathcal{T}_{\mathcal{R}}^s$ , then from Theorems 57 and 141, we can see that (1) and (2) hold.

While, if  $\mathcal{A} \subseteq \mathcal{P}(X)$  such that (1) and (2) hold, then by using Theorem 144 and Corollary 37 we can see that

$$\mathcal{T}_{\mathcal{R}}^s = \mathcal{T}_{\mathcal{R}}^l \subseteq \mathcal{A}^l \subseteq \mathcal{A}.$$

Therefore, the stated minimality property of  $\mathcal{T}_{\mathcal{R}}^s$  is also true.

*Remark 81* Assertion (1) of Theorem 144 showed that, for a topological relator  $\mathcal{R}$  on  $X$ , the family  $\mathcal{T}_{\mathcal{R}}^s$  is the  $l$ -closure of  $\mathcal{T}_{\mathcal{R}}$  in  $\mathcal{P}(X)$ .

While, assertion (1) of Theorem 147 shows that, for a topological relator  $\mathcal{R}$  on  $X$ , the family  $\mathcal{T}_{\mathcal{R}}^s$  is the smallest  $l$ -closed subset of  $\mathcal{P}(X)$  containing  $\mathcal{T}_{\mathcal{R}}$ .

**Theorem 148** *If  $\mathcal{R}$  is a topological relator on  $X$ , then  $\mathcal{A} = \mathcal{T}_{\mathcal{R}}^p$  is the smallest subset of  $\mathcal{P}(X)$  such that*

- (1)  $\mathcal{T}_{\mathcal{R}} \subseteq \mathcal{A}$ ;      (2)  $\mathcal{A}^u \subseteq \mathcal{A}$ .

**Proof** If  $\mathcal{A} = \mathcal{T}_{\mathcal{R}}^p$ , then from Theorems 57 and 141, we can see that (1) and (2) hold.

While, if  $\mathcal{A} \subseteq \mathcal{P}(X)$  such that (1) and (2) hold, then by using Theorem 144 and Corollary 37 we can see that

$$\mathcal{T}_{\mathcal{R}}^p = \mathcal{T}_{\mathcal{R}}^u \subseteq \mathcal{A}^u \subseteq \mathcal{A}.$$

Therefore, the stated minimality property of  $\mathcal{T}_{\mathcal{R}}^p$  is also true.

*Remark 82* In the above theorems, by Theorem 144, we may also write  $\mathcal{T}_{\mathcal{R}}^{\ell}$  and  $\mathcal{T}_{\mathcal{R}}^u$  instead of  $\mathcal{T}_{\mathcal{R}}^s$  and  $\mathcal{T}_{\mathcal{R}}^p$ , respectively.

### 34 Some Basic Properties of the Families $\mathcal{T}_{\mathcal{R}}^{su}$ and $\mathcal{T}_{\mathcal{R}}^{p\ell}$

The following theorems will show that the families  $\mathcal{T}_{\mathcal{R}}^{su}$  and  $\mathcal{T}_{\mathcal{R}}^{p\ell}$  are usually more important than the families  $\mathcal{T}_{\mathcal{R}}^{sl}$  and  $\mathcal{T}_{\mathcal{R}}^{pu}$  considered in Section 32.

**Theorem 149** *We have*

$$\mathcal{T}_{\mathcal{R}}^{su} \cup \mathcal{T}_{\mathcal{R}}^{p\ell} \subseteq \mathcal{T}_{\mathcal{R}}^{\beta}.$$

*Proof* If  $A \in \mathcal{T}_{\mathcal{R}}^{su}$ , then by Theorem 125 there exists  $V \in \mathcal{T}_{\mathcal{R}}^s$  such that

$$A \subseteq V \subseteq A^{-}.$$

Hence, by using Definition 15 and the increasingness of  $\circ-$ , we can infer that

$$A \subseteq V \subseteq V^{\circ-} \subseteq A^{-\circ-}.$$

Thus, by Definition 15, we also have  $A \in \mathcal{T}_{\mathcal{R}}^{\beta}$ . Therefore,  $\mathcal{T}_{\mathcal{R}}^{su} \subseteq \mathcal{T}_{\mathcal{R}}^{\beta}$ .

While, if  $A \in \mathcal{T}_{\mathcal{R}}^{p\ell}$ , then by Theorem 125 there exists  $V \in \mathcal{T}_{\mathcal{R}}^p$  such that

$$V \subseteq A \subseteq V^{-}.$$

Hence, by using Definition 15 and the increasingness of— and  $- \circ -$ , we can infer that

$$A \subseteq V^{-} \subseteq V^{-\circ-} \subseteq A^{-\circ-}.$$

Thus, by Definition 15, we also have  $A \in \mathcal{T}_{\mathcal{R}}^{\beta}$ . Therefore,  $\mathcal{T}_{\mathcal{R}}^{p\ell} \subseteq \mathcal{T}_{\mathcal{R}}^{\beta}$ .

From this theorem, by using Theorem 141 and Corollary 37, we can derive

**Corollary 42** *We have  $\mathcal{T}_{\mathcal{R}}^{\ell u} \cup \mathcal{T}_{\mathcal{R}}^{u\ell} \subseteq \mathcal{T}_{\mathcal{R}}^{\beta}$ .*

**Theorem 150** *If  $\mathcal{R}$  is a reflexive relator on  $X$ , then*

$$(1) \mathcal{T}_{\mathcal{R}}^s \cup \mathcal{T}_{\mathcal{R}}^u \subseteq \mathcal{T}_{\mathcal{R}}^{su}; \quad (2) \mathcal{T}_{\mathcal{R}}^p \cup \mathcal{T}_{\mathcal{R}}^{\ell} \subseteq \mathcal{T}_{\mathcal{R}}^{p\ell}.$$

*Proof* By Theorem 135, we have  $\mathcal{T}_{\mathcal{R}}^s \subseteq \mathcal{T}_{\mathcal{R}}^{su}$ . Moreover, by Theorem 57, we have  $\mathcal{T}_{\mathcal{R}}^u \subseteq \mathcal{T}_{\mathcal{R}}^{su}$ . Hence, by using Corollary 37, we can infer that  $\mathcal{T}_{\mathcal{R}}^u \subseteq \mathcal{T}_{\mathcal{R}}^{su}$ .

Therefore, assertion (1) is true. Assertion (2) can be proved quite similarly.

From this theorem, by using Theorem 141, we can derive

**Corollary 43** *If  $\mathcal{R}$  is a reflexive relator on  $X$ , then  $\mathcal{T}_{\mathcal{R}}^{\ell} \cup \mathcal{T}_{\mathcal{R}}^u \subseteq \mathcal{T}_{\mathcal{R}}^{su} \cap \mathcal{T}_{\mathcal{R}}^{pl}$ .*

However, this is of no particular importance since we can now also prove a stronger statement.

**Theorem 151** *If  $\mathcal{R}$  is a reflexive relator on  $X$ , then*

$$\mathcal{T}_{\mathcal{R}}^{\ell} \cup \mathcal{T}_{\mathcal{R}}^u \subseteq \mathcal{T}_{\mathcal{R}}^{\ell u} \cap \mathcal{T}_{\mathcal{R}}^{u\ell}.$$

**Proof** By Theorem 135, we have  $\mathcal{T}_{\mathcal{R}}^{\ell} \subseteq \mathcal{T}_{\mathcal{R}}^{\ell u}$  and  $\mathcal{T}_{\mathcal{R}} \subseteq \mathcal{T}_{\mathcal{R}}^{\ell}$ . Hence, by using Corollary 37, we can infer that  $\mathcal{T}_{\mathcal{R}}^u \subseteq \mathcal{T}_{\mathcal{R}}^{\ell u}$ .

Therefore,  $\mathcal{T}_{\mathcal{R}}^{\ell} \cup \mathcal{T}_{\mathcal{R}}^u \subseteq \mathcal{T}_{\mathcal{R}}^{\ell u}$ . The inclusion  $\mathcal{T}_{\mathcal{R}}^{\ell} \cup \mathcal{T}_{\mathcal{R}}^u \subseteq \mathcal{T}_{\mathcal{R}}^{u\ell}$  can be proved quite similarly.

From this theorem, by using Theorem 135, we can derive

**Corollary 44** *If  $\mathcal{R}$  is a reflexive relator on  $X$ , then  $\mathcal{T}_{\mathcal{R}} \subseteq \mathcal{T}_{\mathcal{R}}^{\ell u} \cap \mathcal{T}_{\mathcal{R}}^{u\ell}$ .*

Moreover, by using Theorem 135, we can also prove

**Theorem 152** *If  $\mathcal{R}$  is a reflexive relator on  $X$ , then*

$$\mathcal{T}_{\mathcal{R}}^a \subseteq \mathcal{T}_{\mathcal{R}}^{su} \cap \mathcal{T}_{\mathcal{R}}^{pl}.$$

**Proof** By Theorems 56 and 135, we have

$$\mathcal{T}_{\mathcal{R}}^a = \mathcal{T}_{\mathcal{R}}^s \cap \mathcal{T}_{\mathcal{R}}^p \subseteq \mathcal{T}_{\mathcal{R}}^{su} \cap \mathcal{T}_{\mathcal{R}}^{pl}.$$

However, it is now more important to note that, in addition to Theorem 149, we can also prove the following two theorems.

**Theorem 153** *If  $\mathcal{R}$  is a topological relator on  $X$ , then  $\mathcal{T}_{\mathcal{R}}^{su} = \mathcal{T}_{\mathcal{R}}^{\beta}$ .*

**Proof** By Theorem 149, we always have  $\mathcal{T}_{\mathcal{R}}^{su} \subseteq \mathcal{T}_{\mathcal{R}}^{\beta}$ . Therefore, we need only prove that now  $\mathcal{T}_{\mathcal{R}}^{\beta} \subseteq \mathcal{T}_{\mathcal{R}}^{su}$  also holds.

For this, note that if  $A \in \mathcal{T}_{\mathcal{R}}^{\beta}$ , then by Theorem 63 we have also  $A^- \in \mathcal{T}_{\mathcal{R}}^s$ . Hence, by defining  $V = A^-$ , we can note that  $V \in \mathcal{T}_{\mathcal{R}}^s$  such that  $V \subseteq A^-$ . Moreover, by Theorems 32 and 26, we can see that  $A \subseteq A^- = V$  is also true. Therefore, by Theorem 125, we also have  $A \in \mathcal{T}_{\mathcal{R}}^u$ .

**Theorem 154** *If  $\mathcal{R}$  is a topologically filtered, topological relator on  $X$ , then  $\mathcal{T}_{\mathcal{R}}^{pl} = \mathcal{T}_{\mathcal{R}}^{\beta}$ .*

**Proof** By Theorem 149, we always have  $\mathcal{T}_{\mathcal{R}}^{pl} \subseteq \mathcal{T}_{\mathcal{R}}^{\beta}$ . Therefore, we need only prove that now  $\mathcal{T}_{\mathcal{R}}^{\beta} \subseteq \mathcal{T}_{\mathcal{R}}^{pl}$  also holds.



For this, note that if  $A \in \mathcal{T}_{\mathcal{R}}^{\beta}$ , then by Theorem 63 we have  $A^{-} \in \mathcal{T}_{\mathcal{R}}^s$ . Hence, by using Theorems 32, 26, 31, and 59, we can infer that

$$A^{-} = A^{-\circ} = A^{-\circ-}.$$

Now, by defining  $V = A \cap A^{-\circ}$ , we can note that  $V \subseteq A$ . Moreover, by using Corollary 18 and Theorems 32, 31, and 26, we can see that

$$V^{-} = (A \cap A^{-\circ})^{-} \supseteq A^{-} \cap A^{-\circ} = A^{-\circ}.$$

Hence, by using the increasingness of  $\circ$  and Theorems 32, 31, and 26, we can infer that

$$V^{-\circ} \supseteq A^{-\circ\circ} = A^{-\circ} \supseteq A \cap A^{-\circ} = V.$$

Thus, by Definition 15, we also have  $V \in \mathcal{T}_{\mathcal{R}}^p$ .

Moreover, quite similarly, we can now also note that

$$A \subseteq A^{-} = A^{-\circ-} \subseteq V^{-\circ} = V^{-}.$$

Therefore, by Theorem 125, we also have  $A \in \mathcal{T}_{\mathcal{R}}^{p\ell}$ .

From Theorems 153 and 154, by using Corollary 39, we can immediately derive the following two corollaries.

**Corollary 45** *If  $\mathcal{R}$  is a topological relator on  $X$ , then  $\mathcal{T}_{\mathcal{R}}^{\beta u} = \mathcal{T}_{\mathcal{R}}^{\beta}$ .*

**Corollary 46** *If  $\mathcal{R}$  is a topologically filtered, topological relator on  $X$ , then  $\mathcal{T}_{\mathcal{R}}^{\beta\ell} = \mathcal{T}_{\mathcal{R}}^{\beta}$ .*

### 35 Some Further Theorems on the Closure Operation $\Delta$

By Theorems 17 and 16, we evidently have the following

**Theorem 155** *For any  $\diamond = *, \#$  and  $\wedge$ , we have*

- (1)  $\mathcal{A}_{\mathcal{R}} = \mathcal{A}_{\mathcal{R}\diamond}$ ,
- (2)  $\text{Ln}_{\mathcal{R}} = \text{Ln}_{\mathcal{R}\diamond}$  and  $\text{Un}_{\mathcal{R}} = \text{Un}_{\mathcal{R}\diamond}$ ;
- (3)  $\ell_{\mathcal{R}} = \ell_{\mathcal{R}\diamond}$  and  $u_{\mathcal{R}} = u_{\mathcal{R}\diamond}$ .

**Proof** By Theorem 17, we have  $\mathcal{R}^{\wedge} = \mathcal{R}^{\diamond\wedge}$ . Hence, by using Theorem 16, we can infer that

$$\text{cl}_{\mathcal{R}} = \text{cl}_{\mathcal{R}\diamond} \qquad \text{and} \qquad \text{int}_{\mathcal{R}} = \text{int}_{\mathcal{R}\diamond}.$$

Thus, by Definitions 19, 20, and 21, the required assertions are also true.

Moreover, analogously to the results of [67, Sections 33–35], we can also easily prove the following theorems.

**Theorem 156** *If  $\mathcal{R}$  is nonvoid, non-partial relator on  $X$ , then*

$$\mathcal{E}_{\mathcal{R}} \cup \mathcal{E}_{\mathcal{R}}^c \subseteq \mathcal{A}_{\mathcal{R}\Delta}.$$

**Proof** By Theorem 28, we have

$$\emptyset \notin \mathcal{E}_{\mathcal{R}} \quad \text{and} \quad X \in \mathcal{D}_{\mathcal{R}}.$$

Therefore, if  $A \in \mathcal{E}_{\mathcal{R}}$ , then by Corollary 6 we have

$$\text{cl}_{\mathcal{R}\Delta}(\text{int}_{\mathcal{R}\Delta}(A)) = \text{cl}_{\mathcal{R}\Delta}(X) = X.$$

Thus,  $\text{int}_{\mathcal{R}\Delta}(\text{cl}_{\mathcal{R}\Delta}(A)) \subseteq \text{cl}_{\mathcal{R}\Delta}(\text{int}_{\mathcal{R}\Delta}(A))$ , i.e.,  $A \in \mathcal{A}_{\mathcal{R}\Delta}$  trivially holds.

While, if  $A \in \mathcal{E}_{\mathcal{R}}^c$ , i.e.,  $A^c \in \mathcal{E}_{\mathcal{R}}$ , then by Theorem 7 we have  $A \notin \mathcal{D}_{\mathcal{R}}$ . Therefore, by Corollary 6, we have

$$\text{int}_{\mathcal{R}\Delta}(\text{cl}_{\mathcal{R}\Delta}(A)) = \text{int}_{\mathcal{R}\Delta}(\emptyset) = \emptyset.$$

Thus,  $\text{int}_{\mathcal{R}\Delta}(\text{cl}_{\mathcal{R}\Delta}(A)) \subseteq \text{cl}_{\mathcal{R}\Delta}(\text{int}_{\mathcal{R}\Delta}(A))$ , i.e.,  $A \in \mathcal{A}_{\mathcal{R}\Delta}$  trivially holds.

**Theorem 157** *If  $\mathcal{R}$  is non-degenerated relator on  $X$ , then*

$$\mathcal{D}_{\mathcal{R}} \cap \mathcal{A}_{\mathcal{R}\Delta} \subseteq \mathcal{E}_{\mathcal{R}}.$$

**Proof** By Definition 10 and Theorem 30, we have

$$X \neq \emptyset, \quad \mathcal{R} \neq \emptyset \quad \text{and} \quad \emptyset \notin \mathcal{D}_{\mathcal{R}}, \quad X \in \mathcal{E}_{\mathcal{R}}.$$

Thus, if  $A \in \mathcal{D}_{\mathcal{R}}$ , then by Corollary 6 we have

$$\text{int}_{\mathcal{R}\Delta}(\text{cl}_{\mathcal{R}\Delta}(A)) = \text{int}_{\mathcal{R}\Delta}(X) = X.$$

Now, if  $A \in \mathcal{A}_{\mathcal{R}\Delta}$ , i.e.,  $\text{int}_{\mathcal{R}\Delta}(\text{cl}_{\mathcal{R}\Delta}(A)) \subseteq \text{cl}_{\mathcal{R}\Delta}(\text{int}_{\mathcal{R}\Delta}(A))$ , then we also have

$$\text{cl}_{\mathcal{R}\Delta}(\text{int}_{\mathcal{R}\Delta}(A)) = X \neq \emptyset.$$

Hence, by using Corollary 6, we can infer that  $\text{int}_{\mathcal{R}\Delta}(A) \in \mathcal{D}_{\mathcal{R}}$ . Thus, since  $\emptyset \notin \mathcal{D}_{\mathcal{R}}$ , we can also state that  $\text{int}_{\mathcal{R}\Delta}(A) \neq \emptyset$ . Hence, by using Corollary 6, we can already infer that  $A \in \mathcal{E}_{\mathcal{R}}$ .

**Theorem 158** *If  $\mathcal{R}$  is a nonvoid relator on  $X$ , then for any  $A \subseteq X$  we have*

- (1)  $\text{Un}_{\mathcal{R}\Delta}(A) = \mathcal{P}^{-1}(A)$  if  $A \in \mathcal{D}_{\mathcal{R}}$ ;
- (2)  $\text{Un}_{\mathcal{R}\Delta}(A) = \emptyset$  if  $A \notin \mathcal{D}_{\mathcal{R}} \setminus \{\emptyset\}$ ;
- (3)  $\mathcal{U}_{\mathcal{R}\Delta}(\emptyset) = \{\emptyset\}$  if  $\emptyset \notin \mathcal{D}_{\mathcal{R}}$ .

**Proof** If  $A \in \mathcal{D}_{\mathcal{R}}$ , then by Definition 20 and Corollary 6, for any  $B \subseteq X$  we have

$$\begin{aligned}
 B \in \text{Un}_{\mathcal{R}\Delta}(A) &\iff A \in \text{Ln}_{\mathcal{R}\Delta}(B) \iff A \subseteq B \subseteq \text{cl}_{\mathcal{R}\Delta}(A) \\
 &\iff A \subseteq B \subseteq X \iff A \in \mathcal{P}(B) \iff B \in \mathcal{P}^{-1}(A).
 \end{aligned}$$

Therefore, assertion (1) is true.

While, if  $A \notin \mathcal{D}_{\mathcal{R}}$ , then by Definition 20 and Corollary 6, for any  $B \subseteq X$  we have

$$\begin{aligned}
 B \in \text{Un}_{\mathcal{R}\Delta}(A) &\iff A \in \text{Ln}_{\mathcal{R}\Delta}(B) \\
 &\iff A \subseteq B \subseteq \text{cl}_{\mathcal{R}\Delta}(A) \iff A \subseteq B \subseteq \emptyset \iff A = B = \emptyset.
 \end{aligned}$$

Hence, it is clear that assertions (2) and (3) are also true.

From this theorem, by using Theorem 30, we can immediately derive

**Corollary 47** *If  $\mathcal{R}$  is a non-degenerated relator on  $X$ , then for any  $A \subseteq X$  we have*

- (1)  $\text{Un}_{\mathcal{R}\Delta}(A) = \emptyset$  if  $A \notin \mathcal{D}_{\mathcal{R}}$ ;
- (2)  $\mathcal{U}_{\mathcal{R}\Delta}(A) = \mathcal{P}^{-1}(A)$  if  $A \in \mathcal{D}_{\mathcal{R}}$ .

Hence, by using Remark 67 and Theorem 122, we can easily derive

**Theorem 159** *If  $\mathcal{R}$  is a non-degenerated relator on  $X$ , then for any  $\mathcal{A} \subseteq \mathcal{P}(X)$  we have*

$$\mathcal{A}^{\ell_{\mathcal{R}\Delta}} = \bigcup_{A \in \mathcal{A} \cap \mathcal{D}_{\mathcal{R}}} \mathcal{P}^{-1}(A).$$

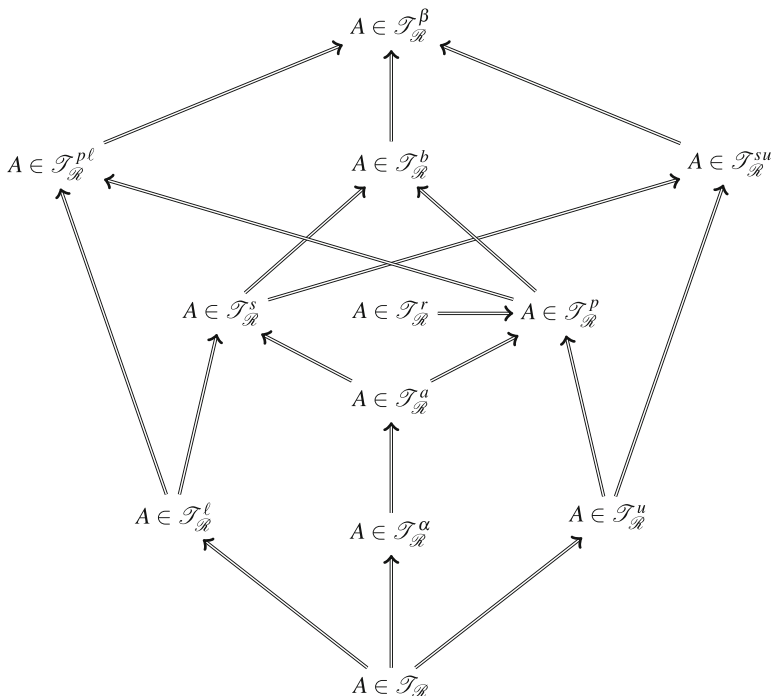
**Proof** By Remark 67, Theorem 122 and Corollary 47, we can see that

$$\mathcal{A}^{\ell_{\mathcal{R}\Delta}} = \text{Ln}_{\mathcal{R}\Delta}^{-1}[\mathcal{A}] = \text{Un}_{\mathcal{R}\Delta}[\mathcal{A}] = \bigcup_{A \in \mathcal{A}} \text{Un}_{\mathcal{R}\Delta}(A) = \bigcup_{A \in \mathcal{A} \cap \mathcal{D}_{\mathcal{R}}} \mathcal{P}^{-1}(A).$$

### 36 An Illustrating Diagram and Two Related Examples

The following diagram and the subsequent two examples have been constructed by Muwafaq Mahdi Salih, a PhD student of the second author from the University of Duhok, Kurdistan Region, Iraq.

**Diagram 1** For a reflexive relation  $\mathcal{R}$  on  $X$ , the following implications hold:



*Remark 83* By our former theorems, nine from the above implications do not require the relation  $\mathcal{R}$  to be reflexive.

Moreover, the following two examples established in [67] show that seventeen implications in Diagram 1 are not reversible.

*Example 7* If  $X = \{1, 2, 3\}$  and  $R$  is a relation on  $X$  such that

$$R(1) = \{1, 2\} \quad \text{and} \quad R(2) = R(3) = X,$$

then  $\mathcal{R} = \{R\}$  is a reflexive relation on  $X$  such that:

- (1)  $\mathcal{T}_{\mathcal{R}} = \mathcal{T}_{\mathcal{R}}^r = \mathcal{T}_{\mathcal{R}}^l = \{\emptyset, X\}$ ;
- (2)  $\mathcal{T}_{\mathcal{R}}^s = \mathcal{T}_{\mathcal{R}}^a = \mathcal{T}_{\mathcal{R}}^{\alpha} = \{\emptyset, \{1, 2\}, X\}$ ;
- (3)  $\mathcal{T}_{\mathcal{R}}^u = \mathcal{T}_{\mathcal{R}}^p = \mathcal{T}_{\mathcal{R}}^b = \mathcal{T}_{\mathcal{R}}^{\beta} = \mathcal{T}_{\mathcal{R}}^{su} = \mathcal{T}_{\mathcal{R}}^{pl} = \mathcal{P}(X) \setminus \{\{3\}\}$ .

*Example 8* If  $X = \{1, 2, 3, 4\}$  and  $R_1$  and  $R_2$  are relations on  $X$  such that

$$\begin{aligned} R_1(1) = R_1(2) = \{1, 2, 3\}, & \quad R_1(3) = R_1(4) = \{1, 3, 4\}; \\ R_2(1) = \{1, 2, 3\}, & \quad R_2(2) = \{1, 2\}, \quad R_2(3) = R_2(4) = \{3, 4\}; \end{aligned}$$

then  $\mathcal{R} = \{R_1, R_2\}$  is a reflexive relator on  $X$  such that:

- (1)  $\mathcal{T}_{\mathcal{R}} = \mathcal{T}_{\mathcal{R}}^{\alpha} = \{\emptyset, \{3, 4\}, X\}$ ;      (2)  $\mathcal{T}_{\mathcal{R}}^r = \mathcal{T}_{\mathcal{R}} \cup \{\{2\}\}$ ;  
 (3)  $\mathcal{T}_{\mathcal{R}}^{\ell} = \mathcal{T}_{\mathcal{R}}^a = \mathcal{T}_{\mathcal{R}} \cup \{\{1, 3, 4\}\}$ ;      (4)  $\mathcal{T}_{\mathcal{R}}^s = \mathcal{T}_{\mathcal{R}}^{\ell} \cup \{\{1, 2\}\}$ ;  
 (5)  $\mathcal{T}_{\mathcal{R}}^p = \mathcal{P}(X) \setminus \{\{1\}, \{1, 2\}\}$ ;      (6)  $\mathcal{T}_{\mathcal{R}}^u = \mathcal{T}_{\mathcal{R}}^p \setminus \{\{2\}\}$ ;  
 (7)  $\mathcal{T}_{\mathcal{R}}^b = \mathcal{T}_{\mathcal{R}}^{p\ell} = \mathcal{P}(X) \setminus \{\{1\}\}$ ;      (8)  $\mathcal{T}_{\mathcal{R}}^{\beta} = \mathcal{T}_{\mathcal{R}}^{su} = \mathcal{P}(X)$ .

*Remark 84* Unfortunately, the above two examples cannot be used to show that the implication  $A \in \mathcal{T}_{\mathcal{R}}^{su} \implies A \in \mathcal{T}_{\mathcal{R}}^{\beta}$  is also not reversible.

Note that, by Theorem 149, this implication does not also require the relator  $\mathcal{R}$  to be reflexive. Moreover, if  $\mathcal{R}$  is topological, then by Theorem 153 the reverse implication is also true.

*Note 1* Some of the results of this paper can be generalized according to the ideas of Gargouri and Rezgui [31] and the second author [96].

## References

1. M.E. Abd El-Monsef, S.N. El-Deeb, R.A. Mahmoud,  $\beta$ -open sets and  $\beta$ -continuous mappings. *Bull. Fac. Sci. Assiut Univ.* **12**, 77–90 (1983)
2. P. Alexandroff, Zur Begründung der  $n$ -dimensionalen mengentheoretischen Topologie. *Math. Ann.* **94**, 296–308 (1925)
3. D. Andrijević, Semi-preopen sets. *Mat. Vesnik* **38**, 24–32 (1986)
4. D. Andrijević, On  $b$ -open sets. *Mat. Vesnik* **48**, 59–64 (1996)
5. N. Biswas, On some mappings in topological spaces. *Bull. Cal. Math. Soc.* **61**, 127–135 (1969)
6. N. Bourbaki, *General Topology, Chapters 1–4* (Springer, Berlin, 1989)
7. T.A. Chapman, A further note on closure and interior operators. *Am. Math. Monthly* **69**, 524–529 (1962)
8. Ch. Chattopadhyay, Ch. Bandyopadhyay, On structure of  $\delta$ -sets. *Bull. Calcutta Math. Soc.* **83**, 281–290 (1991)
9. Ch. Chattopadhyay, U.K. Roy,  $\delta$ -sets, irresolvable and resolvable spaces. *Math. Slovaca* **42**, 371–378 (1992)
10. H. Choda, K. Matoba, On a theorem of Levine. *Proc. Jpn. Acad.* **37**, 462–463 (1961)
11. H.H. Corson, E. Michael, Metrizable unions. *Ill. J. Math.* **8**, 351–360 (1964)
12. Á. Császár, *Foundations of General Topology* (Pergamon Press, London, 1963)
13. Á. Császár, Generalized open sets. *Acta Math. Hungar.* **75**, 65–87 (1997)
14. Á. Császár, On the  $\gamma$ -interior and  $\gamma$ -closure of a set. *Acta Math. Hungar.* **80**, 89–93 (1998)
15. Á. Császár,  $\gamma$ -quasi-open sets. *Stud. Sci. Math. Hungar.* **38**, 171–176 (2001)
16. Á. Császár, Remarks on  $\gamma$ -quasi-open sets. *Stud. Sci. Math. Hungar.* **39**, 137–141 (2002)
17. Á. Császár, Further remarks on the formula for the  $\gamma$ -interior. *Acta Math. Hungar.* **113**, 325–332 (2006)
18. Á. Császár, Remarks on quasi-topologies. *Acta Math. Hungar.* **119**, 197–200 (2008)
19. B.A. Davey, H.A. Priestley, *Introduction to Lattices and Order* (Cambridge University Press, Cambridge, 2002)
20. A.S. Davis, Indexed systems of neighborhoods for general topological spaces. *Am. Math. Monthly* **68**, 886–893 (1961)
21. K. Dłaska, N. Ergun, M. Ganster, On the topology generated by semi-regular sets. *Indian J. Pure Appl. Math.* **25**, 1163–1170 (1994)

22. J. Dontchev, Survey on preopen sets. Meetings on Topological Spaces, Theory and Applications, Yatsushiro College of Technology, Kumamoto, Japan (1998), 18 pp.
23. Z. Duszynski, T. Noiri, Semi-open, semi-closed sets and semi-continuity of functions. *Math. Pannon.* **23**, 195–200 (2012)
24. V.A. Efremovič, The geometry of proximity. *Mat. Sb.* **31**, 189–200 (1952) (Russian)
25. V.A. Efremovič, A.S. Švarc, A new definition of uniform spaces. Metrization of proximity spaces. *Dokl. Acad. Nauk. SSSR* **89**, 393–396 (1953) (Russian)
26. N. Elez, O. Papaz, The new operators in topological spaces. *Math. Moravica* **17**, 63–68 (2013)
27. P. Fletcher, W.F. Lindgren, *Quasi-Uniform Spaces* (Marcel Dekker, New York, 1982)
28. M. Ganster, Preopen sets and resolvable spaces. *Kyungpook J.* **27**, 135–143 (1987)
29. M. Ganster, I.L. Reilly, M.K. Vamanamurthy, Remarks on locally closed sets. *Math. Pannon.* **3**, 107–113 (1992)
30. B. Ganter, R. Wille, *Formal Concept Analysis* (Springer, Berlin, 1999)
31. R. Gargouri, A. Rezgui, A unification of weakening of open and closed subsets in a topological spaces. *Bull. Malays. Math. Sci. Soc.* **40**, 1219–1230 (2017)
32. S. Givant, P. Halmos, *Introduction to Boolean Algebras* (Springer, Berlin, 2009)
33. T. Glavosits, Generated preorders and equivalences. *Acta Acad. Paed. Agrienses, Sect. Math.* **29**, 95–103 (2002)
34. W.H. Gottschalk, Intersection and closure. *Proc. Am. Math. Soc.* **4**, 470–473 (1953)
35. W. Hunsaker, W. Lindgren, Construction of quasi-uniformities. *Math. Ann.* **188**, 39–42 (1970)
36. D.H. Hyers, On the stability of the linear functional equation. *Proc. Nat. Acad. Sci. U.S.A* **27**, 222–224 (1941)
37. Y. Isomichi, New concept in the theory of topological spaces—Supercondensed set, subcondensed set, and condensed set. *Pac. J. Math.* **38**, 657–668 (1971)
38. Y.B. Jun, S.W. Jeong, H.j. Lee, J.W. Lee, Applications of pre-open sets. *Appl. Gen. Top.* **9**, 213–228 (2008)
39. S.-M. Jung, Interiors and closure of sets and applications. *Int. J. Pure Math.* **3**, 41–45 (2016)
40. S.-M. Jung, D. Nam, Some properties of interior and closure in general topology. *Mathematics* **7**, 624 (2019)
41. J.L. Kelley, *General Topology* (Van Nostrand Reinhold Company, New York, 1955)
42. H. Kenyon, Two theorems on relations. *Trans. Am. Math. Soc.* **107**, 1–9 (1963)
43. V.L. Kljushin, Al bayati J.H. Hussein, On simply-open sets. *Vestnik UDC* **3**, 34–38 (2011). (Russian)
44. K. Kuratowski, Sur l'opération  $\bar{A}$  de l'analysis situs. *Fund. Math.* **3**, 182–199 (1922) (An English translation: On the operation  $\bar{A}$  in analysis situs, prepared by M. Bowron in 2010, is available on the Internet)
45. K. Kuratowski, *Topology I* (Academic Press, New York, 1966)
46. J. Kurdics, A note on connection properties. *Acta Math. Acad. Paedagog. Nyházi.* **12**, 57–59 (1990)
47. J. Kurdics, Á. Száz, Well-chainedness characterizations of connected relators. *Math. Pannon.* **4**, 37–45 (1993)
48. N. Levine, On the commutivity of the closure and interior operators in topological spaces. *Am. Math. Monthly* **68**, 474–477 (1961)
49. N. Levine, Semi-open sets and semi-continuity in topological spaces. *Am. Math. Monthly* **70**, 36–41 (1963)
50. N. Levine, Some remarks on the closure operator in topological spaces. *Am. Math. Monthly* **70**, 553 (1963)
51. N. Levine, On uniformities generated by equivalence relations. *Rend. Circ. Mat. Palermo* **18**, 62–70 (1969)
52. N. Levine, On Pervin's quasi uniformity. *Math. J. Okayama Univ.* **14**, 97–102 (1970)
53. J. Mala, Relators generating the same generalized topology. *Acta Math. Hungar.* **60**, 291–297 (1992)
54. J. Mala, Á. Száz, Properly topologically conjugated relators. *Pure Math. Appl. Ser. B* **3**, 119–136 (1992)

55. J. Mala, Á. Szász, Modifications of relators. *Acta Math. Hungar.* **77**, 69–81 (1997)
56. A.S. Mashhour, M.E. Abd El-Monsef, S.N. El-Deeb, On precontinuous and weak precontinuous mappings. *Proc. Math. Phys. Soc. Egypt* **53**, 47–53 (1982)
57. S.A. Naimpally, B.D. Warrack, *Proximity Spaces* (Cambridge University Press, Cambridge, 1970)
58. H. Nakano, K. Nakano, Connector theory. *Pac. J. Math.* **56**, 195–213 (1975)
59. A.A. Nasef, R. Mareay, More on simply open sets and its applications. *South Asian J. Math.* **5**, 100–108. (2015)
60. A.A. Nasef, R. Mareay, Ideals and some applications of simply open sets. *J. Adv. Math.* **13**, 7264–7271 (2017)
61. A. Neubrunnová, On transfinite sequences of certain types of functions. *Acta Fac. Rer. Natur. Univ. Commun. Math.* **30**, 121–126 (1975)
62. O. Njåstad, On some classes of nearly open sets. *Pac. J. Math.* **15**, 195–213 (1965)
63. G. Pataki, Supplementary notes to the theory of simple relators. *Radovi Mat.* **9**, 101–118 (1999)
64. G. Pataki, On the extensions, refinements and modifications of relators. *Math. Balk.* **15**, 155–186 (2001)
65. G. Pataki, Á. Szász, A unified treatment of well-chainedness and connectedness properties. *Acta Math. Acad. Paedagog. Nyházi. (N.S.)* **19**, 101–165 (2003)
66. W.J. Pervin, Quasi-uniformization of topological spaces. *Math. Ann.* **147**, 316–317 (1962)
67. Th.M. Rassias, M. Salih, Á. Szász, Characterizations of generalized topologically open sets in relator spaces, in *Recent Trends on Pure and Applied Mathematics, Special Issue of the Montes Taurus*, ed. by G.V. Milovanovic, Thm. M. Rassias, Y. Simsek. *J. Pure Appl. Math., Dedicated to Professor Hari Mohan Srivastava on the occasion of his 80th Birthday, Montes Taurus J. Pure Appl. Math.* **3**, 39–94 (2021)
68. Th.M. Rassias, M. Salih, Á. Szász, Set-theoretic properties of generalized topologically open sets in relator spaces, in *Mathematical Analysis in Interdisciplinary Research*, ed. by I.N. Parasidis, E. Providas, Th.M. Rassias, to appear
69. M. Salih, Á. Szász, Generalizations of some ordinary and extreme connectedness properties of topological spaces to relator spaces. *Electron. Res. Arch.* **28**, 471–548 (2020)
70. P. Sivagami, Remarks on  $\gamma$ -interior. *Acta Math. Hungar.* **119**, 81–94 (2008)
71. Yu.M. Smirnov, On proximity spaces. *Math. Sb.* **31**, 543–574 (1952) (Russian)
72. M.H. Stone, Application of the theory of Boolean rings to general topology. *Trans. Am. Math. Soc.* **41**, 374–481 (1937)
73. Á. Szász, Basic tools and mild continuities in relator spaces. *Acta Math. Hungar.* **50**, 177–201 (1987)
74. Á. Szász, Directed, topological and transitive relators. *Publ. Math. Debrecen* **35**, 179–196 (1988)
75. Á. Szász, Relators, Nets and Integrals. Unfinished doctoral thesis, Debrecen (1991), 126 pp.
76. Á. Szász, Structures derivable from relators. *Singularité* **3**, 14–30 (1992)
77. Á. Szász, Refinements of relators. *Tech. Rep., Inst. Math., Univ. Debrecen*, vol. 76 (1993), 19 pp.
78. Á. Szász, Cauchy nets and completeness in relator spaces. *Colloq. Math. Soc. János Bolyai* **55**, 479–489 (1993)
79. Á. Szász, Neighbourhood relators. *Bolyai Soc. Math. Stud.* **4**, 449–465 (1995)
80. Á. Szász, Uniformly, proximally and topologically compact relators. *Math. Pannon.* **8**, 103–116 (1997)
81. Á. Szász, Somewhat continuity in a unified framework for continuities of relations. *Tatra Mt. Math. Publ.* **24**, 41–56 (2002)
82. Á. Szász, Upper and lower bounds in relator spaces. *Serdica Math. J.* **29**, 239–270 (2003)
83. Á. Szász, Rare and meager sets in relator spaces. *Tatra Mt. Math. Publ.* **28**, 75–95 (2004)
84. Á. Szász, Galois-type connections on power sets and their applications to relators. *Tech. Rep., Inst. Math., Univ. Debrecen 2005/2* (2005), 38 pp.
85. Á. Szász, Minimal structures, generalized topologies, and ascending systems should not be studied without generalized uniformities. *Filomat* **21**, 87–97 (2007)

86. Á. Száz, Galois type connections and closure operations on preordered sets. *Acta Math. Univ. Comenian. (N.S.)* **78**, 1–21 (2009)
87. Á. Száz, Inclusions for compositions and box products of relations. *J. Int. Math. Virt. Inst.* **3**, 97–125 (2013)
88. Á. Száz, A particular Galois connection between relations and set functions. *Acta Univ. Sapientiae, Math.* **6**, 73–91 (2014)
89. Á. Száz, Generalizations of Galois and Pataki connections to relator spaces. *J. Int. Math. Virtual Inst.* **4**, 43–75 (2014)
90. Á. Száz, Basic tools, increasing functions, and closure operations in generalized ordered sets, in *Contributions in Mathematics and Engineering*, ed. by P.M. Pardalos, Th.M. Rassias. In Honor of Constantine Caratheodory (Springer, Berlin, 2016), pp. 551–616
91. Á. Száz, Four general continuity properties, for pairs of functions, relations and relators, whose particular cases could be investigated by hundreds of mathematicians. *Tech. Rep., Inst. Math., Univ. Debrecen, 2017/1* (2017), 17 pp.
92. Á. Száz, The closure-interior Galois connection and its applications to relational equations and inclusions. *J. Int. Math. Virt. Inst.* **8**, 181–224 (2018)
93. Á. Száz, Corelations are more powerful tools than relations, in *Applications of Nonlinear Analysis*, ed. by Th.M. Rassias. Optimization and Its Applications, vol. 134 (Springer, Berlin, 2018), pp. 711–779
94. Á. Száz, Relationships between inclusions for relations and inequalities for corelations. *Math. Pannon.* **26**, 15–31 (2018)
95. Á. Száz, Galois and Pataki connections on generalized ordered sets. *Earthline J. Math. Sci.* **2**, 283–323 (2019)
96. Á. Száz, Birelator spaces are natural generalizations of not only bitopological spaces, but also ideal topological spaces, in *Mathematical Analysis and Applications, Springer Optimization and Its Applications*, ed. by Th.M. Rassias, P.M. Pardalos, vol. 154 (Springer, Switzerland, 2019), pp. 543–586
97. W.J. Thron, *Topological Structures* (Holt, Rinehart and Winston, New York, 1966)
98. H. Tietze, Beiträge zur allgemeinen Topologie I. Axiome für verschiedene Fassungen des Umgebungsbegriffs. *Math. Ann.* **88**, 290–312 (1923)
99. A. Weil, Sur les espaces á structure uniforme et sur la topologie générale. *Actual. Sci. Ind.*, vol. 551 (Herman and Cie, Paris 1937)



# Graphical Mean Curvature Flow



Andreas Savas-Halilaj

**Abstract** In this survey article, we discuss recent developments on the mean curvature flow of graphical submanifolds, generated by smooth maps between Riemannian manifolds. We will see interesting applications of this technique, in the understanding of the homotopy type of maps between manifolds.<sup>12</sup>

## 1 Introduction

Let  $f: M \rightarrow N$  be a smooth map between two manifolds  $M$  and  $N$ . It is a fundamental problem to find *canonical representatives* in the homotopy class of  $f$ . By a canonical representative is usually meant a map in the homotopy class of the given map  $f$  which is a critical point of a suitable functional. In the mid-1960s, Eells and Sampson [34] introduced the *harmonic maps* as critical points of the energy density, to attack the aforementioned problem.

One possible approach to construct harmonic maps is via the *harmonic map heat flow*. If  $M$  is compact and  $N$  is negatively curved, in [34] Eells and Sampson were able to prove long-time existence and convergence of the flow, showing that under these assumptions one finds harmonic representatives in a given homotopy class. In general, one can neither expect long-time existence nor convergence of this flow. For example, the situation is very complicated in the case of maps between spheres.

---

<sup>1</sup>These notes are based partly on a series of lectures delivered by the author at the Chern Institute of Mathematics held in Tianjin-China in November 2019.

<sup>2</sup>The author would like to acknowledge support by the General Secretariat for Research and Technology (GSRT) and the Hellenic Foundation for Research and Innovation (HFRI) Grant No:133.

---

A. Savas-Halilaj (✉)

Department of Mathematics, University of Ioannina, Ioannina, Greece

e-mail: [ansavas@uoi.gr](mailto:ansavas@uoi.gr)

There is another important functional that we may consider in the space of smooth maps. Given a map  $f : M \rightarrow N$  between Riemannian manifolds, let us denote its *graph* in the product space  $M \times N$  by

$$\Gamma(f) = \{(x, f(x)) \in M \times N : x \in M\}.$$

Following the terminology introduced by Schoen [79], a map whose graph is minimal submanifold is called *minimal map*. Therefore, minimal maps are critical points of the volume functional.

In this survey, among others, we will discuss deformation of graphical submanifolds via the *mean curvature flow*. Before stating the problems that we would like to deal with, let us provide some basic facts and definitions. Let  $M$  be a smooth  $m$ -dimensional manifold,  $T > 0$  a positive number and  $F : M \times [0, T) \rightarrow P$  a smooth time-dependent family of immersions of  $M$  into a Riemannian manifold  $P$ . We say that  $F$  evolves in time under the *mean curvature flow* if it satisfies the evolution equation

$$dF(\partial_t)(x, t) = H(x, t)$$

for any  $(x, t) \in M \times [0, T)$ , where  $H(x, t)$  stands for the mean curvature vector at the point  $x$  of the immersion  $F(\cdot, t) : M \rightarrow P$ . It is a well-known fact that if  $M$  is compact and  $F_0 : M \rightarrow P$  is an immersion, then the initial value problem for the mean curvature flow admits a unique smooth solution on a maximal time interval  $[0, T_{\max})$ , where  $0 < T_{\max} \leq \infty$ . Suppose now that  $P$  is the product manifold  $M \times N$  and  $F_0$  is the graph of a map  $f : M \rightarrow N$ . Notice that long as the submanifolds deformed under mean curvature flow remain graphical, one obtains a smooth family of maps which belong to the homotopy class of the map  $f$ . In the case of long-time existence and convergence of the flow, we obtain a smooth homotopy from  $f$  to a minimal map.

The first result regarding evolutions by mean curvature of graphical submanifolds is due to Ecker and Huisken [33]. They proved long-time existence of entire graphical hypersurfaces in  $\mathbb{R}^{n+1}$ . Moreover, Ecker and Huisken proved convergence to a flat subspace, if the growth rate at infinity of the initial graphical submanifold is linear. On the other hand, in higher codimensions, the complexity of the normal bundle makes the situation more complicated. Results analogous to that of Ecker and Huisken are not available any more without further assumptions. However, the ideas developed in the paper of Ecker and Huisken opened a new era for the study of the mean curvature flow of submanifolds in Riemannian manifolds of arbitrary codimension; see for example [12, 13, 16–18, 60–62, 64, 75, 77, 78, 85, 87–91, 94–96, 98–100].

This new deformation of maps between Riemannian manifolds via the mean curvature flow has been used in order to have a better understanding of the relation between the  $k$ -dilation  $Dil_k$  and the homotopy type of maps. In order to be precise, let us recall at first the following definition:

**Definition 1** Let  $f : M \rightarrow N$  be a map between two Riemannian manifolds. We say that  $Dil_k(f) \leq \alpha$  if  $f$  maps each  $k$ -dimensional submanifold  $\Sigma \subset M$  to an image with  $k$ -dimensional volume at most  $\alpha \cdot \mathcal{H}^k(\Sigma)$ , where  $\mathcal{H}^k(\Sigma)$  stands for the  $k$ -dimensional Hausdorff measure of  $\Sigma$ . In particular, we say that  $f$  is *area decreasing* if  $Dil_2(f) \leq 1$ , *strictly area decreasing* if  $Dil_2(f) = 1$ , and *area preserving* if  $Dil_2(f) = 1$ .

Roughly speaking, the  $k$ -dilation measures how much the map  $f : M \rightarrow N$  contracts  $k$ -dimensional volumes. Gromov in [38] realized that there is a close relationship between the 1-dilation of a map and its homotopy type. For instance, he proved that if  $f$  is a map from  $\mathbb{S}^m$  to  $\mathbb{S}^m$ , then its degree is at most  $Dil_1^m(f)$  and this bound is sharp up to a constant factor. Motivated by this result, in [40, 41] Gromov proposed the following:

**Problem 1** Let  $f : \mathbb{S}^m \rightarrow \mathbb{S}^n$  be a smooth map between euclidean spheres. Is there a number  $\varepsilon(k, m, n)$  such that if  $Dil_k(f) < \varepsilon$  would imply that  $f$  is null-homotopic?

Tsui and Wang in [91] proved using the mean curvature flow that smooth strictly area decreasing maps  $f : \mathbb{S}^m \rightarrow \mathbb{S}^n$  can be smoothly deformed to a constant map. Guth [42] proved this result cannot be extended in the case of maps with  $k$ -dilation strictly less than 1, if  $k \geq 3$ . The result of Tsui and Wang was generalized by Lee and Lee in [60]. In the matter of fact, they proved that any strictly area decreasing map between compact Riemannian manifolds  $M$  and  $N$  whose sectional curvatures are bounded by  $sec_M \geq \sigma_1$  and  $\sigma_2 \geq sec_N$ , where  $\sigma_1, \sigma_2$  are two real constants such that  $\sigma_1 \geq \sigma_2 > 0$  or  $\sigma_1 > 0 \geq \sigma_2$ , is homotopic by mean curvature flow to a constant map. We would like to point out here that the curvature assumptions can be relaxed to

$$sec_M > -\sigma \quad \text{and} \quad Ric_M \geq (m - 1)\sigma \geq (m - 1) sec_N,$$

where  $\sigma$  is a positive constant number, as it was shown in [75] by Savas-Halilaj and Smoczyk.

In the case of a smooth area decreasing map  $f : M \rightarrow N$  between two compact Riemann surfaces  $M$  and  $N$  of the same constant sectional curvature  $\sigma$ , we have a complete picture of the behaviour of the mean curvature flow. *It turns out that, under the mean curvature flow, such a map either instantly becomes strictly area decreasing or it was and remains an area preserving map. Moreover, the mean curvature flow preserves the graphical property, exists for all time, and converges to a minimal surface  $\Sigma_\infty$  of the product  $M \times N$ . Additionally:*

- (I) *If the evolved graphs are generated by strictly area decreasing maps then:*
  - (a) *If  $\sigma > 0$ , then  $\Sigma_\infty$  is the graph of a constant map.*
  - (b) *If  $\sigma = 0$ , then  $\Sigma_\infty$  is the graph of an affine minimal map.*
- (II) *If the evolved graphs are generated by area preserving maps then:*

- (a) If  $\sigma > 0$ , then  $\Sigma_\infty$  is the graph of an isometry.
- (b) If  $\sigma = 0$ , then  $\Sigma_\infty$  is the graph of an affine minimal diffeomorphism.

The first steps in the proof of the above result were made in the seminal works of Smoczyk [85] and Wang [94, 95, 100], where the area preserving case was investigated. The strictly area decreasing case was first treated by Tsui and Wang [91], in the positive case, and completed recently by Savas-Halilaj and Smoczyk in [78]. The primary goal of this survey is to present a unified proof of this result, based in [78].

From the results of Wang [94, 95, 100], we get another proof of Smale's Theorem [84] which says that any diffeomorphism of  $\mathbb{S}^2$  can be smoothly deformed into an isometry. Let us mention here that, according to a deep theorem of Hatcher [49], any diffeomorphism of  $\mathbb{S}^3$  can be deformed into an isometry of  $\mathbb{S}^3$ . Such a result is not expected for spheres of dimension greater or equal than 4; see for example [28]. However, the following problem is challenging:

**Problem 2** Let  $f : \mathbb{S}^m \rightarrow \mathbb{S}^m$ ,  $m \geq 4$ , be a smooth diffeomorphism. Under which conditions  $f$  can be smoothly deformed into an isometry of the sphere?

Another interesting problem is the investigation of the symplectomorphism group of the complex projective space  $\mathbb{C}\mathbb{P}^m$ . Gromov [39] proved that the bi-holomorphic group of  $\mathbb{C}\mathbb{P}^2$  is a deformation retract of its symplectomorphism group. A natural problem is to determine whether a similar result holds for  $\mathbb{C}\mathbb{P}^m$  with  $m \geq 3$ . In the matter of fact, the following problem is still open:

**Problem 3** Let  $f : \mathbb{C}\mathbb{P}^m \rightarrow \mathbb{C}\mathbb{P}^m$ ,  $m \geq 3$ , be a smooth symplectomorphism. Is it true that the mean curvature flow deforms the graph  $\Gamma(f)$  of  $f$  smoothly to the graph  $\Gamma(g)$  of bi-holomorphic map  $g : \mathbb{C}\mathbb{P}^m \rightarrow \mathbb{C}\mathbb{P}^m$ ? Is it true that any minimal symplectomorphism  $f : \mathbb{C}\mathbb{P}^m \rightarrow \mathbb{C}\mathbb{P}^m$  is a bi-holomorphic isometry?

Medoš and Wang in [64] made some contribution by giving an affirmative answer to the above problem under the additional assumption that the singular values of the differential of the symplectomorphism are close to 1.

The paper is structured as follows. In Section 2 we set up the notation and recall basic facts from submanifold geometry. In Section 3 we discuss minimal submanifolds in euclidean spaces. We introduce the generalized Gauss map and prove the Ruh-Vilms Theorem. Section 4 describes the class of graphical submanifolds and review some Bernstein-type theorems. Section 5 is devoted to the maximum principle for scalar and systems of PDEs. In Section 6, we introduce the mean curvature flow, prove short-time existence, and derive various basic evolution equations. Section 7 describes how to built smooth singularity models for the mean curvature flow. Section 8 combines results from the previous sections to prove our main result.

## 2 Riemannian Submanifolds

In this section we set up the notation and recall some basic facts from submanifold geometry. We closely follow the exposition in [5, 29, 55, 59, 92, 102].

### 2.1 Notation and Conventions

Let  $M$  be a  $m$ -dimensional manifold and  $(E, \pi, M)$  a vector bundle of rank  $k$  over  $M$ . We often denote the bundles only by its total space  $E$ . The fiber of  $E$  at a point  $x \in M$  is denoted by  $E_x$ , the tangent space of  $M$  at a point  $x \in M$  will be denoted by  $T_x M$  and the space of sections of  $E$  is denoted by  $\Gamma(E)$ . For the tangent bundle of  $M$ , we use the symbol  $TM$ . Sections of the tangent bundle are called *vector fields* and usually  $\Gamma(TM)$  is denoted by  $\mathfrak{X}(M)$ . A smooth map  $T : E \rightarrow V$  between two vector bundles  $E$  and  $V$  over  $M$  which maps the fiber  $E_x$  linearly to  $V_x$ , for any  $x \in M$  is called *bundle morphism*. If additionally,  $T$  is bijective we call it *bundle isomorphism*.

**Definition 2** A (linear) *connection* on a vector bundle  $E$  is a map  $\nabla^E : \mathfrak{X}(M) \times \Gamma(E) \rightarrow \Gamma(E)$ , written  $\nabla^E(v, \phi) = \nabla_v^E \phi$ , satisfying the properties:

(a) For any  $v_1, v_2 \in \mathfrak{X}(M)$  and  $\phi \in \Gamma(E)$ , it holds

$$\nabla_{v_1+v_2}^E \phi = \nabla_{v_1}^E \phi + \nabla_{v_2}^E \phi.$$

(b) For any  $v \in \mathfrak{X}(M)$ ,  $f \in C^\infty(M)$  and  $\phi \in \Gamma(E)$ , it holds

$$\nabla_{fv}^E \phi = f \nabla_v^E \phi.$$

(c) For any  $v \in \mathfrak{X}(M)$ ,  $f \in C^\infty(M)$  and  $\phi_1, \phi_2 \in \Gamma(E)$ , it holds

$$\nabla_v^E (\phi_1 + \phi_2) = \nabla_v^E \phi_1 + \nabla_v^E \phi_2.$$

(d) For any  $v \in \mathfrak{X}(M)$ ,  $\phi \in \Gamma(E)$  and  $f \in C^\infty(M)$ , it holds

$$\nabla_v^E (f\phi) = (vf)\phi + f \nabla_v^E \phi.$$

For any  $\phi \in \Gamma(E)$  and  $x_0 \in M$ , the value  $\nabla^E \phi_x|_{x_0}$  of the quantity  $\nabla_v^E \phi$  at  $x_0 \in M$  depends only on the value of  $v$  at  $x_0$  and on the restriction of  $\phi$  along a curve passing through  $x_0$  with speed  $v$ . If  $\phi_1, \phi_2 \in \Gamma(E)$  coincide on a neighbourhood of  $x_0 \in M$ , then

$$\nabla_{v_1}^E \phi_1|_{x_0} = \nabla_{v_2}^E \phi_2|_{x_0},$$

for any pair of vector fields  $v_1, v_2 \in \mathfrak{X}(M)$  with  $v_1|_{x_0} = v_2|_{x_0}$ .

**Definition 3** A section  $\phi \in \Gamma(E)$  is said to be *parallel* with respect to  $\nabla^E$  if, for any vector field  $v$  on  $M$ , it holds  $\nabla_v^E \phi = 0$ .

We can define higher derivatives of sections of a vector bundle over a manifold  $M$  whose tangent bundle  $TM$  is equipped with a connection.

**Definition 4** Suppose that  $M$  is a smooth manifold and  $E$  a vector bundle over  $M$ . Let  $\nabla^M$  be a connection of  $TM$  and  $\nabla^E$  a connection of  $E$ . For any pair  $v_1, v_2 \in \mathfrak{X}(M)$ , the map  $\nabla_{v_1, v_2}^2 : \Gamma(E) \rightarrow \Gamma(E)$ , given by

$$\nabla_{v_1, v_2}^2 \phi = \nabla_{v_1}^E \nabla_{v_2}^E \phi - \nabla_{\nabla_{v_1}^M v_2}^E \phi,$$

is called the *second covariant derivative* of  $\phi$ , with respect to the directions  $v_1$  and  $v_2$ . By coupling the connections  $\nabla^M$  and  $\nabla^E$ , one may define, the  $k$ -th derivative  $\nabla^k$  of a section  $\phi$  in  $\Gamma(E)$ .

To each connection there is associated an important operator, which measures the non commutativity of the covariant derivatiation.

**Definition 5** The operator  $R^E : \mathfrak{X}(M) \times \mathfrak{X}(M) \times \Gamma(E) \rightarrow \Gamma(E)$ , defined by

$$R^E(v_1, v_2, \phi) = \nabla_{v_1, v_2}^2 \phi - \nabla_{v_2, v_1}^2 \phi,$$

for any  $v_1, v_2 \in \mathfrak{X}(M)$  and  $\phi \in \Gamma(E)$ , is called the *curvature operator* of  $\nabla^E$ .

Now let us turn our attention to vector bundles equipped with a Riemannian metric structure.

**Definition 6** A *Riemannian metric* on a vector bundle  $E$  of rank  $k$  over the manifold  $M$  is a smooth map  $g_E : \Gamma(E) \times \Gamma(E) \rightarrow C^\infty(M)$ , such that its restriction to the fibers is a positive definite inner product.

**Definition 7** A connection  $\nabla^E$  is called *compatible with the Riemannian metric*  $g_E$  or *metric compatible* if it satisfies

$$v g_E(\phi_1, \phi_2) = g_E(\nabla_v^E \phi_1, \phi_2) + g_E(\phi_1, \nabla_v^E \phi_2),$$

for any  $v \in \mathfrak{X}(M)$  and  $\phi_1, \phi_2 \in \Gamma(E)$ . A vector bundle  $E$  endowed with these structures is called *Riemannian vector bundle endowed with a compatible linear connection*.

We say that a set of sections  $\{\phi_1, \dots, \phi_k\}$  forms an *orthonormal frame*, with respect to  $g_E$  if and only if  $g_E(\phi_i, \phi_j) = \delta_{ij}$ , for any  $i, j \in \{1, \dots, k\}$ . In particular, around any point  $x_0$  of  $M$  it is possible to find a local orthonormal frame  $\{\phi_1, \dots, \phi_k\}$  such that

$$\nabla_v \phi_i|_{x_0} = 0$$

for any tangent vector  $v$ . Such frames are called *normal* or *geodesic frames*.

Let us restrict ourselves at the tangent bundle  $TM$  of  $M$ . Given a Riemannian metric  $g$  on  $M$ , there is a unique connection  $\nabla$ , referred as the *Levi-Civita connection*, which is compatible with the Riemannian metric. More precisely,  $\nabla$  is given by the *Koszul formula*

$$2g(\nabla_{v_1} v_2, v_3) = v_1(g(v_2, v_3)) + v_2(g(v_1, v_3)) - v_3(g(v_1, v_2)) + g([v_1, v_2], v_3) - g([v_1, v_3], v_2) - g([v_2, v_3], v_1),$$

for all  $v_1, v_2, v_3 \in \mathfrak{X}(M)$ . The Levi-Civita also satisfy

$$\nabla_{v_1} v_2 - \nabla_{v_2} v_1 = [v_1, v_2],$$

for any  $v_1, v_2 \in \mathfrak{X}(M)$ .

Denote by  $R$  the curvature operator with respect to the connection  $\nabla$ . Combining  $R$  with  $g$  we obtain a  $(4,0)$ -tensor which, for simplicity, we again denote with the letter  $R$ . More precisely,

$$R(v_1, v_2, v_3, v_4) = -g(R(v_1, v_2, v_3), v_4),$$

for any  $v_1, v_2, v_3, v_4 \in \mathfrak{X}(M)$ . If  $v_1, v_2$  are linearly independent vectors, then

$$\text{sec}(v_1, v_2) = \frac{R(v_1, v_2, v_1, v_2)}{g(v_1, v_1)g(v_2, v_2) - g(v_1, v_2)^2},$$

is called the *sectional curvature* of the plane spanned by the vectors  $v_1$  and  $v_2$ . By contracting the operator  $R$  with  $g$  we obtain the *Ricci operator*  $\text{Ric}$  and the *scalar curvature*  $\text{scal}$ , i.e.,

$$\text{Ric}(v_1, v_2) = \sum_{i=1}^m R(v_1, e_i, v_2, e_i) \quad \text{and} \quad \text{scal} = \sum_{i=1}^m \text{Ric}(e_i, e_i),$$

where  $v \in \mathfrak{X}(M)$  and  $\{e_1, \dots, e_m\}$  is a local orthonormal frame on  $M$ .

*Remark 1* One can use the operations of Linear Algebra to produce new vector bundles from given ones. For example, if  $E$  and  $V$  are vector bundles over a manifold  $M$ , then  $E^*, E \times V, E \otimes V, E \oplus V, \text{Hom}(E; V), \Lambda^r(V)$  and  $\text{Sym}^r(V)$  gives rise to new bundles over  $M$ . If  $M$  is endowed with a Riemannian metric, then this metric and its Levi-Civita connections extends in a natural way to all the aforementioned bundles; for more details see [59] or [102].

## 2.2 The Pull-back Bundle

Let  $M$  and  $N$  be two manifolds,  $(E, \pi, N)$  is a vector bundle of rank  $k$  over  $N$  and suppose that  $f: M \rightarrow N$  is a smooth map. The map  $f$  induces a new vector bundle of rank  $k$  over  $M$ . Indeed, take as total space

$$f^*E = \{(x, \xi) : x \in M, \xi \in E_{f(x)}\}$$

and as projection the map  $\pi_f: f^*E \rightarrow M$  given by  $\pi_f(x, \xi) = x$ . The space  $\Gamma(f^*E)$  contains all sections of  $E$  with base point at  $f(M)$  and inherit naturally a vector space structure from  $E_{f(x)}$ , given by

$$(x, \xi) + (x, \eta) = (x, \xi + \eta) \quad \text{and} \quad \lambda(x, \xi) = (x, \lambda\xi).$$

The triple  $(f^*E, \pi_f, M)$  carries the structure of a vector bundle over  $M$ . This bundle is called the *pull-back* or the *induced by  $f$  vector bundle* on  $M$ .

Suppose that  $h$  is a Riemannian metric on  $E$  and  $\nabla^E$  is a metric compatible connection. The map  $f$  induces a connection  $\nabla^{f^*E}$  on the pull-back bundle which is defined as follows: Let  $\{\vartheta_1, \dots, \vartheta_k\}$  be a local orthonormal frame field of  $E$  in a neighbourhood of the point  $f(x) \in N$ . Then, any section  $\phi \in \Gamma(f^*E)$  can be written in the form

$$\phi|_x = \left(x, \sum_{\alpha=1}^k \phi^\alpha(x) \vartheta_\alpha|_{f(x)}\right) \cong \sum_{\alpha=1}^k \phi^\alpha(x) \vartheta_\alpha|_{f(x)},$$

where  $\phi^\alpha$ ,  $\alpha \in \{1, \dots, k\}$ , are the *components* of  $\phi$  with respect to the given orthonormal frame field. Define now the induced connection by

$$\nabla_v^{f^*E} \phi|_x = \sum_{\alpha=1}^k (v\phi^\alpha) \vartheta_\alpha|_{f(x)} + \sum_{\alpha=1}^k \phi^\alpha \nabla_{df(v)}^E \vartheta_\alpha|_{f(x)},$$

for  $x \in M$  and  $v \in T_x M$ . One can easily show that the curvature operator  $R^{f^*E}$  of  $\nabla^{f^*E}$  is given by

$$R^{f^*E}(v_1, v_2, \phi|_x) = R^E(df(v_1), df(v_2), \phi|_x),$$

for any  $x \in M$ ,  $v_1, v_2 \in T_x M$  and  $\phi \in E|_{f(x)}$ .

Let us discuss the case where  $f: (M, g, \nabla^g) \rightarrow (N, h, \nabla^h)$  is a map between Riemannian manifolds. The restriction of  $h$  on  $f^*TN$ , induces a Riemannian metric on  $f^*TN$ , which is compatible with the pull-back connection, that is

$$vh(\phi_1, \phi_2) = h(\nabla_v^{f^*TN} \phi_1, \phi_2) + h(\phi_1, \nabla_v^{f^*TN} \phi_2).$$

Moreover, for  $v_1, v_2 \in \mathfrak{X}(M)$ , it holds



$$\nabla_{v_1}^{f^*TN} df(v_2) - \nabla_{v_2}^{f^*TN} df(v_1) = df([v_1, v_2]).$$

**Definition 8** The *Hessian* of a map  $f : (M, g, \nabla^g) \rightarrow (N, h, \nabla^h)$  is defined to be the symmetric tensor  $B : \mathfrak{X}(M) \times \mathfrak{X}(M) \rightarrow \Gamma(f^*E)$  given by

$$B(v_1, v_2) = \nabla_{v_1}^{f^*TN} df(v_2) - df(\nabla_{v_1}^g v_2),$$

for any  $v_1, v_2 \in \mathfrak{X}(M)$ . The trace of  $B$  with respect to  $g$  is denoted by  $\Delta_{g,h} f$  and is called the *Laplacian* of  $f$ . If the Laplacian of  $f$  is zero, then  $f$  is called a *harmonic map*.

### 2.3 The Second Fundamental Form

Consider Riemannian manifolds  $(M, g, \nabla^g)$  and  $(N, h, \nabla^h)$  of dimension  $m$  and  $n$ , respectively, with  $m \leq n$ . A map  $f : M \rightarrow N$  is called an *isometric immersion* if and only if  $f^*h = g$ . For simplicity, we often denote both metrics  $g$  and  $h$  by  $\langle \cdot, \cdot \rangle$ . At every  $x \in M$ , we have the orthogonal decomposition

$$T_{f(x)}N = df_x(T_xM) \oplus N_{f(x)}M,$$

where  $N_{f(x)}M$  is the orthogonal complement of  $df_x(T_xM)$  with respect to  $h$ . The union  $NM$  of all normal spaces form a vector bundle of rank  $n - m$  over  $M$  which is called the *normal bundle*. According to the above decomposition, any section  $v \in \Gamma(f^*TN)$  can be decomposed in a unique way in the form

$$v = v^\top + v^\perp,$$

where  $v^\top$  is the *tangent* and  $v^\perp$  is the *normal* part of  $v$  along the submanifold. A well known fact in submanifold theory is that

$$(\nabla_{v_1}^{f^*TN} df(v_2))^\top = df(\nabla_{v_1}^g v_2), \tag{1}$$

for any  $v_1, v_2 \in \mathfrak{X}(M)$ . In submanifold theory, the Hessian of  $f$  is denoted by the letter  $A$ , i.e., we have

$$A(v_1, v_2) = \nabla_{v_1}^{f^*TN} df(v_2) - df(\nabla_{v_1}^g v_2).$$

The tensor  $A$  is called the *second fundamental form* of  $f$ . If  $\xi$  is a normal vector, then the tensor  $A^\xi$  given by

$$A^\xi(v_1, v_2) = \langle A(v_1, v_2), \xi \rangle,$$

for any tangent vectors  $v_1, v_2$ , is called *shape operator with respect to  $\xi$* . The Weingarten operator  $A_\xi$  associated with  $\xi$  is defined by

$$\langle A_\xi v_1, v_2 \rangle = A^\xi(v_1, v_2) = \langle A(v_1, v_2), \xi \rangle.$$

The Laplacian of  $f$  or, equivalently, the trace of  $A$  with respect to  $g$  is called the (*unnormalized*) *mean curvature* and is denoted by the letter  $H$ , that is

$$H = \text{trace}_g A.$$

**Definition 9** A submanifold with zero mean curvature is called *minimal*.

The restriction of  $h$  on  $NM$  gives rise to a Riemannian metric on the normal bundle. Moreover, the restriction of  $\nabla^h$  on  $NM$  induces a connection  $\nabla^\perp$  on  $NM$  which is compatible with the metric; i.e., just define

$$\nabla_v^\perp \xi = (\nabla_v^N \xi)^\perp.$$

The curvature tensor of the normal bundle is denoted by  $R^\perp$  and is given by

$$R^\perp(v_1, v_2, \xi) = \nabla_{v_1}^\perp \nabla_{v_2}^\perp \xi - \nabla_{v_2}^\perp \nabla_{v_1}^\perp \xi - \nabla_{[v_1, v_2]}^\perp \xi.$$

As usual, we can form from  $R^\perp$  a  $C^\infty(M)$ -valued tensor which we denote again by  $R^\perp$ , that is

$$R^\perp(v_1, v_2, \xi, \eta) = -\langle R^\perp(v_1, v_2, \xi), \eta \rangle.$$

The Riemann curvature operator  $R$  of  $M$ , the curvature operator  $\tilde{R}$  of  $N$  and the normal curvature  $R^\perp$  are related to the second fundamental form  $A$  through the *Gauss-Codazzi-Ricci equations*:

(a) **Gauss equation:**

$$R(v_1, v_2, v_3, v_4) = \tilde{R}(df(v_1), df(v_2), df(v_3), df(v_4)) + \langle A(v_1, v_3), A(v_2, v_4) \rangle - \langle A(v_2, v_3), A(v_1, v_4) \rangle.$$

(b) **Codazzi equation:**

$$(\nabla_{v_1}^\perp A)(v_2, v_3) - (\nabla_{v_2}^\perp A)(v_1, v_3) = (\tilde{R}(df(v_1)df(v_2), df(v_3)))^\perp.$$

(c) **Ricci equation:**

$$R^\perp(v_1, v_2, \xi, \eta) = \tilde{R}(df(v_1), df(v_2), \xi, \eta)$$

$$+ \sum_k (A^\xi(v_1, e_k)A^\eta(v_2, e_k) - A^\eta(v_1, e_k)A^\xi(v_2, e_k)),$$

where  $v_1, v_2, v_3, v_4 \in \mathfrak{X}(M)$ ,  $\xi, \eta \in NM$  and  $\{e_1, \dots, e_m\}$  is a local orthonormal frame field with respect to  $g$ .

### 2.4 Local Representations

Let  $f : (M, g) \rightarrow (N, h)$  be a smooth map between Riemannian manifolds. For computational reasons, we need expressions for components of various tensorial quantities. We can express coordinates with respect to local charts or with respect to orthonormal frames.

Let discuss at first the notation with respect to a local coordinate system. Choose a chart  $(U, \varphi)$  around a point  $x \in M$  and a chart  $(V, \psi)$  around  $f(x) \in N$ . Assume that  $\varphi : U \rightarrow \mathbb{R}^m$  is represented as  $\varphi = (x_1, \dots, x_m)$  and suppose that  $\psi : V \rightarrow \mathbb{R}^n$  is represented as  $\psi = (y_1, \dots, y_n)$ . We use Latin indices to describe quantities on  $M$  and Greek indices for quantities on  $N$ . From the charts  $\varphi$  and  $\psi$ , we obtain for  $f$  the local expression expression

$$\psi \circ f \circ \varphi^{-1} = (f^1, \dots, f^n),$$

where

$$f^\alpha = y^\alpha \circ f \circ \varphi^{-1}.$$

Denote now the basic vector fields associated with the charts  $(U, \varphi)$  and  $(V, \psi)$  by  $\{\partial_{x_i}, \dots, \partial_{x_m}\}$  and  $\{\partial_{y_1}, \dots, \partial_{y_n}\}$ , respectively. Moreover, denote their corresponding dual forms by  $\{dx_1, \dots, dx_m\}$  and  $\{dy_1, \dots, dy_n\}$ . With respect to these conventions, the Riemannian metrics  $g$  and  $h$  can be written in the form

$$g = \sum_{i,j} g_{ij} dx_i \otimes dx_j \quad \text{and} \quad h = \sum_{\alpha,\beta} h_{\alpha\beta} dy_\alpha \otimes dy_\beta.$$

The *Christoffel symbols*  $\Gamma_{ij}^k$  of the metric  $g$ , are defined by the formula

$$\nabla_{\partial_{x_i}}^g \partial_{x_j} = \sum_k \Gamma_{ij}^k \partial_{x_k}$$

and they can be expressed in terms of the metric as

$$\Gamma_{ij}^k = \frac{1}{2} \sum_l g^{kl} (-\partial_{x_l} g_{ij} + \partial_{x_i} g_{jl} + \partial_{x_j} g_{li}),$$

where  $g^{ij}$  are the components of the inverse of the matrix of the metric  $g$ , with respect to the basis  $\{\partial_{x_1}, \dots, \partial_{x_m}\}$ . Similarly, are defined the Christoffel symbols  $\Gamma_{\beta\gamma}^\alpha$  of  $h$ . The differential of the map  $f$  and the pull-back via  $f$  of the metric  $h$  are given by

$$df = \sum_{\alpha} f_{x_i}^{\alpha} \partial_{y_{\alpha}} \otimes dx_i \quad \text{and} \quad f^*h = \sum_{\alpha, \beta} h_{\alpha\beta} f_{x_i}^{\alpha} f_{x_j}^{\beta}.$$

By a straightforward computation, we see that the Hessian  $B$  of  $f$  can be represented in the form

$$B(\partial_{x_i}, \partial_{x_j}) = \sum_{\alpha} B_{ij}^{\alpha} \partial_{y_{\alpha}} = \sum_{\alpha} (f_{x_i x_j}^{\alpha} - \sum_k \Gamma_{ij}^k f_{x_k}^{\alpha} + \sum_{\beta, \gamma} \Gamma_{\beta\gamma}^{\alpha} f_{x_i}^{\beta} f_{x_j}^{\gamma}) \partial_{y_{\alpha}}.$$

Suppose now that  $f : M \rightarrow N$  is an isometric immersion. Then, the second fundamental form  $A$  and the mean curvature  $H$  are represented, respectively, as

$$A(\partial_{x_i}, \partial_{x_j}) = \sum_{\alpha} A_{ij}^{\alpha} \partial_{y_{\alpha}} = \sum_{\alpha} (f_{x_i x_j}^{\alpha} - \sum_k \Gamma_{ij}^k f_{x_k}^{\alpha} + \sum_{\beta, \gamma} \Gamma_{\beta\gamma}^{\alpha} f_{x_i}^{\beta} f_{x_j}^{\gamma}) \partial_{y_{\alpha}}$$

and

$$\begin{aligned} H &= \sum_{\alpha} H^{\alpha} \partial_{y_{\alpha}} = \sum_{i, j, \alpha} g^{ij} A_{ij}^{\alpha} \partial_{y_{\alpha}} \\ &= \sum_{i, j} g^{ij} (f_{x_i x_j}^{\alpha} - \sum_k \Gamma_{ij}^k f_{x_k}^{\alpha} + \sum_{\beta, \gamma} \Gamma_{\beta\gamma}^{\alpha} f_{x_i}^{\beta} f_{x_j}^{\gamma}) \partial_{y_{\alpha}}. \end{aligned} \quad (2)$$

Let us discuss now expressions of tensors in local orthonormal frames. Let  $\{e_1, \dots, e_m\}$  be a local orthonormal frame of the tangent bundle and  $\{\xi_{m+1}, \dots, \xi_n\}$  a local orthonormal frame of the normal bundle. Here we use Latin indices for components on the tangent bundle and Greek indices for components on the normal bundle. For example, we write:

$$\begin{aligned} A_{ij}^{\alpha} &= \langle A(e_i, e_j), \xi_{\alpha} \rangle = \langle A_{ij}, \xi_{\alpha} \rangle, \\ \tilde{R}_{ijkl} &= \tilde{R}(df(e_i), df(e_j), df(e_k), df(e_l)), \\ \tilde{R}_{ij\alpha\beta} &= \tilde{R}(df(e_i), df(e_j), \xi_{\alpha}, \xi_{\beta}). \end{aligned}$$

Now the Gauss-Codazzi-Ricci equations can be written as:

(a) **Gauss equation:**

$$R_{ijkl} = \tilde{R}_{ijkl} + \sum_{\alpha} (A_{ik}^{\alpha} A_{jl}^{\alpha} - A_{jk}^{\alpha} A_{il}^{\alpha}). \quad (3)$$

(b) **Codazzi equation:**

$$(\nabla_{e_i}^\perp A)_{jk}^\alpha - (\nabla_{e_j}^\perp A)_{ik}^\alpha = - \sum_{\alpha} \tilde{R}_{ijk\alpha}. \quad (4)$$

(c) **Ricci equation:**

$$R_{ij\alpha\beta}^\perp = \tilde{R}_{ij\alpha\beta} + \sum_k (A_{ik}^\alpha A_{jk}^\beta - A_{ik}^\beta A_{jk}^\alpha). \quad (5)$$

### 3 Minimal Submanifolds

The theory of minimal submanifolds is one of the most active subjects of differential geometry. There is a vast of literature, but here we will present rather basic facts concerning higher codimensional minimal submanifolds. For more details we refer to [21, 22, 70].

#### 3.1 The Gauss Map of a Minimal Submanifold

One of the most important objects in the submanifold geometry is the Gauss map. For codimension one oriented submanifolds in the euclidean space, the *Gauss map* associates to every point of the hypersurface its oriented unit normal vector. This concept can be generalized to higher codimensional oriented submanifolds. Let  $f : M \rightarrow \mathbb{R}^n$  be an isometric immersion of a  $m$ -dimensional oriented Riemannian manifold  $M$  into the euclidean space. The image  $df(T_x M)$ , can be taken after a suitable parallel displacement in  $\mathbb{R}^n$ , into a point  $\mathcal{G}(x)$  of the *oriented Grassmann space*  $\mathbb{G}_+(m, n)$  of  $m$ -dimensional oriented subspaces of  $\mathbb{R}^n$ . The map  $\mathcal{G} : M \rightarrow \mathbb{G}_+(m, n)$  defined in this way, is called the *generalized Gauss map*.

There is a natural way to visualize the Grassmann space  $\mathbb{G}_+(m, n)$ . Let us denote by  $\Lambda^m(\mathbb{R}^n)$  the dual space of all alternative multilinear forms of degree  $m$ . Elements of  $\Lambda^m(\mathbb{R}^n)$  are called  *$m$ -vectors*. Hence, given vectors  $v_1, \dots, v_m$  on  $\mathbb{R}^n$ , the *exterior product*  $v_1 \wedge \dots \wedge v_m$  is the linear map which on an alternating form  $\Omega$  of degree  $m$  takes the value

$$(v_1 \wedge \dots \wedge v_m)(\Omega) = \Omega(v_1, \dots, v_m).$$

The exterior product is linear in each variable separately. Interchanging two elements the sign of the product changes and if two variables are the same the exterior product vanishes. An  $m$ -vector  $\xi$  is called *simple* or *decomposable* if it can be written as a single wedge product of vectors, that is

$$\xi = v_1 \wedge \dots \wedge v_m.$$

Note that are  $m$ -vectors which are not simple. Using standard techniques from Linear Algebra one can verify that the exterior product  $v_1 \wedge \dots \wedge v_m$  is zero if and only if the vectors are linearly dependent. Moreover, if  $\{e_1, \dots, e_n\}$  consists a basis for  $\mathbb{R}^n$ , then the  $m$ -vectors

$$\{e_{i_1} \wedge \dots \wedge e_{i_m} : 1 \leq i_1 < \dots < i_m \leq n\}$$

consists a basis of  $\Lambda^m(\mathbb{R}^n)$ . Therefore, the dimension of the vector space of  $m$ -vectors is

$$\dim \Lambda^m(\mathbb{R}^n) = \binom{n}{m} = \frac{n!}{m!(n-m)!}.$$

Each simple vector represents a unique  $m$ -dimensional subspace of  $\mathbb{R}^n$ . Moreover, if  $\xi$  and  $\eta$  are simple vectors representing the same subspace, then there exists a non-zero real number such that  $\xi = \lambda\eta$ . Therefore, there is an obvious equivalence relation on the space of simple vectors such that the space of equivalence classes is to an one to one correspondence with the space of  $m$ -dimensional subspaces of  $\mathbb{R}^n$ . Additionally, we can consider the following relation on the set of non-zero simple  $m$ -vectors:  $\xi$  and  $\eta$  are called equivalent if and only if  $\xi = \lambda\eta$  for some positive number  $\lambda$ . Denote by  $[\xi]$  the class containing all simple  $m$ -vectors that are equivalent to  $\xi$ . The equivalence classes now obtained are called *oriented  $m$ -dimensional subspaces* of  $\mathbb{R}^n$ .

We can equip  $\Lambda^m(\mathbb{R}^n)$  with a natural inner product, which for simplicity we denote again by  $\langle \cdot, \cdot \rangle$ . Indeed, define

$$\langle v_1 \wedge \dots \wedge v_m, w_1 \wedge \dots \wedge w_m \rangle = \det((v_i, w_j))_{1 \leq i, j \leq m}$$

on simple  $m$ -vectors and then extend linearly. Moreover, if  $\{e_1, \dots, e_n\}$  is an orthonormal basis of  $\mathbb{R}^n$  then, the  $m$ -vectors

$$\{e_{i_1} \wedge \dots \wedge e_{i_m} : 1 \leq i_1 < \dots < i_m \leq n\}$$

consists an orthonormal basis for the exterior power  $\Lambda^m(\mathbb{R}^n)$ . Moreover, it turns out that for vectors  $v_1, \dots, v_m$  in  $\mathbb{R}^n$ , the norm

$$|v_1 \wedge \dots \wedge v_m|$$

gives the  $m$ -volume of the parallelepiped spanned by these vectors.

We can equip  $\mathbb{G}_+(m, n)$  with a natural differentiable structure. For every  $m$ -dimensional subspace  $V_0$  of  $\mathbb{G}_+(m, n)$ , consider the open neighbourhood  $U(V_0)$  of oriented  $m$ -dimensional subspaces whose orthogonal projection into  $V_0$  is one-to-one. Let  $\{e_1, \dots, e_m\}$  be an orthonormal base spanning  $V_0$  and  $\{\eta_{m+1}, \dots, \eta_n\}$  an orthonormal base spanning its orthogonal complement  $V_0^\perp$  in  $\mathbb{R}^n$ . Then, we may parametrize  $U(V_0)$  via  $\xi : \mathbb{R}^{m(n-m)} \rightarrow U(V_0)$  given by

$$\begin{aligned} (x_{1m+1}, \dots, x_{i\alpha}, \dots, x_{mn}) &\rightarrow \xi(x_{1m+1}, \dots, x_{i\alpha}, \dots, x_{mn}) \\ &= (e_1 + \sum_{\alpha} x_{1\alpha} \eta_{\alpha}) \wedge \cdots \wedge (e_m + \sum_{\alpha} x_{m\alpha} \eta_{\alpha}). \end{aligned} \quad (6)$$

Two charts  $U(V_i), U(V_j)$  with distinct  $V_i, V_j$  are analytically compatible.

**Definition 10** The map  $\Psi : \mathbb{G}_+(m, n) \rightarrow \mathbb{S}^{\binom{n}{m}-1}$  given by

$$\Psi([v_1 \wedge \cdots \wedge v_m]) = \frac{v_1 \wedge \cdots \wedge v_m}{|v_1 \wedge \cdots \wedge v_m|}$$

is called the *Plücker embedding*. We regard the Grassmann space  $\mathbb{G}_+(m, n)$  as a Riemannian manifold with the induced by  $\Psi$  Riemannian metric.

**Theorem 1** *The Plücker embedding is minimal.*

**Proof** Fix a  $m$ -dimensional linear space  $V_0 \in \mathbb{G}_+(m, n)$  and consider the parametrization  $\xi : \mathbb{R}^{m(n-m)} \rightarrow U(V_0) \subset \mathbb{G}_+(m, n)$  described in (6). Now

$$\Psi = W\xi = \frac{(e_1 + \sum_{\alpha} x_{1\alpha} \eta_{\alpha}) \wedge \cdots \wedge (e_m + \sum_{\alpha} x_{m\alpha} \eta_{\alpha})}{|(e_1 + \sum_{\alpha} x_{1\alpha} \eta_{\alpha}) \wedge \cdots \wedge (e_m + \sum_{\alpha} x_{m\alpha} \eta_{\alpha})|},$$

where the index  $\alpha$  run from  $m+1$  to  $n$  and

$$W = \frac{1}{|\xi|} = \frac{1}{|(e_1 + \sum_{\alpha} x_{1\alpha} \eta_{\alpha}) \wedge \cdots \wedge (e_m + \sum_{\alpha} x_{m\alpha} \eta_{\alpha})|}.$$

Note that

$$\begin{aligned} \xi_{x_{i\alpha}} &= (e_1 + \sum_{\alpha} x_{1\alpha} \eta_{\alpha}) \wedge \cdots \wedge (e_{i-1} + \sum_{\alpha} x_{i-1\alpha} \eta_{\alpha}) \wedge \eta_{\alpha} \\ &\quad \wedge (e_{i+1} + \sum_{\alpha} x_{i+1\alpha} \eta_{\alpha}) \wedge \cdots \wedge (e_m + \sum_{\alpha} x_{m\alpha} \eta_{\alpha}) \end{aligned}$$

and

$$\begin{aligned} \xi_{x_{i\alpha} x_{j\beta}} &= (e_1 + \sum_{\alpha} x_{1\alpha} \eta_{\alpha}) \wedge \cdots \wedge (e_{i-1} + \sum_{\alpha} x_{i-1\alpha} \eta_{\alpha}) \wedge \eta_{\alpha} \\ &\quad \wedge (e_{i+1} + \sum_{\alpha} x_{i+1\alpha} \eta_{\alpha}) \wedge \cdots \wedge (e_{j-1} + \sum_{\alpha} x_{j-1\alpha} \eta_{\alpha}) \wedge \eta_{\beta} \\ &\quad \wedge (e_{j+1} + \sum_{\alpha} x_{j+1\alpha} \eta_{\alpha}) \wedge \cdots \wedge (e_m + \sum_{\alpha} x_{m\alpha} \eta_{\alpha}) (1 - \delta_{ij}), \end{aligned}$$

where  $i, j \in \{1, \dots, m\}$  and  $\alpha, \beta \in \{m+1, \dots, n\}$ . In particular,

$$\xi_{x_{i\alpha}}(0) = e_1 \wedge \cdots \wedge e_{i-1} \wedge \eta_{\alpha} \wedge e_{i+1} \wedge \cdots \wedge e_m \quad (7)$$

and

$$\xi_{x_{i\alpha}x_{j\beta}}(0) = e_1 \wedge \cdots \wedge e_{i-1} \wedge \eta_\alpha \wedge e_{i+1} \wedge \cdots \wedge e_{j-1} \wedge \eta_\beta \wedge e_{j+1} \wedge \cdots \wedge e_m (1 - \delta_{ij}). \tag{8}$$

Additionally,

$$W_{x_{i\alpha}} = -W^3(\xi, \xi_{x_{i\alpha}})$$

and

$$W_{x_{i\alpha}x_{j\beta}} = -3W^5(\xi, \xi_{x_{i\alpha}})(\xi, \xi_{x_{j\beta}}) - W^3(\xi_{x_{i\alpha}}, \xi_{x_{j\beta}}) - W^3(\xi, \xi_{x_{i\alpha}x_{j\beta}}).$$

Moreover,

$$W(0) = 1, \quad W_{x_{i\alpha}}(0) = 0 \quad \text{and} \quad W_{x_{i\alpha}x_{j\beta}}(0) = -\delta_{ij}\delta_{\alpha\beta}. \tag{9}$$

From (7) and (9) we see that

$$\Psi_{x_{i\alpha}}(0) = e_1 \wedge \cdots \wedge e_{i-1} \wedge \eta_\alpha \wedge e_{i+1} \wedge \cdots \wedge e_m.$$

Hence, the vectors

$$\{\partial_{x_{1m+1}}|0, \dots, \partial_{x_{i\alpha}}|0, \dots, \partial_{x_{mn}}|0\}$$

form an orthonormal basis of  $T_{V_0}\mathbb{G}_+(m, n)$  with respect to the induced by  $\Psi$  Riemannian metric. Moreover, from (7), (8), and (9) we deduce that

$$\Psi_{x_{i\alpha}x_{j\beta}}(0) = -\delta_{ij}\delta_{\alpha\beta}\Psi(0) + \xi_{x_{i\alpha}x_{j\beta}}(0).$$

According to (8), the second fundamental form  $A$  of the Plücker embedding at the point  $V_0$  is equal to

$$A(\partial_{x_{i\alpha}}, \partial_{x_{j\beta}}) = e_1 \wedge \cdots \wedge e_{i-1} \wedge \eta_\alpha \wedge e_{i+1} \wedge \cdots \wedge e_{j-1} \wedge \eta_\beta \wedge e_{j+1} \wedge \cdots \wedge e_m (1 - \delta_{ij}) \tag{10}$$

and, in particular,

$$A(\partial_{x_{i\alpha}}, \partial_{x_{i\alpha}}) = 0 \tag{11}$$

for any  $i \in \{1, \dots, m\}$  and  $\alpha \in \{m + 1, \dots, n\}$ . Thus, the mean curvature  $H$  of the embedding  $\Psi$  at  $V_0$  is given by

$$H(V_0) = \sum_{i,\alpha} A(\partial_{x_{i\alpha}}, \partial_{x_{i\alpha}}) = 0.$$

Consequently,  $\Psi$  gives rise to a minimal submanifold of the sphere. □



In 1970, Ruh and Vilms [74] obtained an important link between minimality of a submanifold and the harmonicity of its generalized Gauss map. More precisely, the following result holds:

**Theorem 2** *Let  $f : M \rightarrow \mathbb{R}^n$  be a minimal isometric immersion. Then, the generalized Gauss map  $\mathcal{G}$  of  $f$  is a harmonic map.*

**Proof** Consider the map

$$F = \Psi \circ \mathcal{G} : M \rightarrow \mathbb{S}^{\binom{n}{m}-1} \subset A^m(\mathbb{R}^n)$$

where  $\Psi$  is the Plücker embedding. From the composition formula, we have

$$B_F(v_1, v_2) = d\Psi(B_{\mathcal{G}}(v_1, v_2)) + A_{\Psi}(d\mathcal{G}(v_1), d\mathcal{G}(v_2)) \tag{12}$$

for any  $v_1, v_2 \in \mathfrak{X}(M)$ , where  $B_F$  and  $B_{\mathcal{G}}$  are the Hessians of  $F$  and  $\mathcal{G}$  and  $A_{\Psi}$  the second fundamental form of  $\Psi$ , respectively. Fix a local orthonormal frame field  $\{e_1, \dots, e_m\}$  defined on an open neighbourhood  $U$  of  $M$  and a local orthonormal frame  $\{\eta_{m+1}, \dots, \eta_n\}$  in the normal bundle of the immersion. Note that since  $f$  is isometric immersion, for any  $x \in U$ , we have

$$F(x) = df_x(e_1) \wedge df_x(e_2) \wedge \dots \wedge df_x(e_m).$$

Fix now a point  $x_0 \in U$  and for simplicity suppose that the frame  $\{e_1, \dots, e_m\}$  is normal at  $x_0$ . By straightforward computations we see that at  $x_0$  we have

$$dF(e_j) = A(e_j, e_1) \wedge \dots \wedge df(e_m) + \dots + df(e_1) \wedge \dots \wedge A(e_j, e_m),$$

where  $A$  is the second fundamental form of  $f$ . Hence, in view of (7), we obtain that the differential of  $\mathcal{G}$  at  $x_0$  is equal to

$$d\mathcal{G}(e_j) = \sum_{\alpha} A_{1j}^{\alpha} \eta_{\alpha} \wedge e_2 \wedge \dots \wedge e_m + \dots + e_1 \wedge e_2 \wedge \dots \wedge \sum_{\alpha} A_{mj}^{\alpha} \eta_{\alpha}.$$

Recall that, from the Codazzi equations (4), we have at  $x_0$  that

$$(\nabla_{e_j}^{f^*T\mathbb{R}^n} A)_{kl} = (\nabla_{e_k}^{\perp} A)_{jl} - \sum_{i,\alpha} A_{kl}^{\alpha} A_{ij}^{\alpha} df(e_i),$$

for any  $j, k, l \in \{1, \dots, m\}$ . Differentiating  $dF$  and estimating at  $x_0$  we get

$$\begin{aligned} \nabla_{e_j}^{F^*T A^m(\mathbb{R}^n)} dF(e_j) &= -|A|^2 F + (\nabla_{e_1}^{\perp} A)_{jj} \wedge df(e_2) \wedge \dots \wedge df(e_m) \\ &\quad + df(e_1) \wedge (\nabla_{e_2}^{\perp} A)_{jj} \wedge \dots \wedge df(e_m) + \dots + df(e_1) \wedge \dots \wedge (\nabla_{e_m}^{\perp} A)_{jj} \\ &\quad + A_{1j} \wedge A_{2j} \wedge \dots \wedge df(e_m) + \dots + df(e_1) \wedge \dots \wedge A_{jm-1} \wedge A_{jm}. \end{aligned}$$



The idea to obtain such a parametrization is the following: Let  $f : M \rightarrow \mathbb{R}^n$  be a minimal immersion of a 2-manifold. Choose a local isothermal system of coordinates  $(U \subset \mathbb{C}, z = x + iy)$ , where  $U$  is simply connected; see [54]. Then, the induced by  $f$  metric  $g$  on  $M$  has the form  $g = E|dz|^2$ , where  $E$  is a smooth positive function. Moreover, in these coordinates, the Laplace–Beltrami operator  $\Delta$  with respect to  $g$  is expressed by

$$\Delta = E^{-1}(\partial_x \partial_x + \partial_y \partial_y).$$

With respect to such parameters, minimality is equivalent to harmonicity. Consider now the map  $\varphi = (\varphi_1, \dots, \varphi_n) : U \rightarrow \mathbb{C}^n$ ,  $\varphi = f_x - if_y$ . One can readily check that  $\varphi$  is holomorphic and its components satisfy

$$\varphi_1^2 + \dots + \varphi_n^2 = 0 \quad \text{and} \quad |\varphi_1|^2 + \dots + |\varphi_n|^2 = 2E > 0.$$

By fixing a point  $z_0 \in U$  it is clear that, up to a parallel transport,

$$f(z) = \operatorname{Re} \int_{z_0}^z \varphi(\zeta) d\zeta, \quad z \in U.$$

The map  $\varphi$  has also a very important geometric interpretation. At first we observe that the variety

$$Q_{n-2} = \{[z_1, \dots, z_n] \in \mathbb{C}\mathbb{P}^{n-1} : z_1^2 + \dots + z_n^2 = 0\}$$

is diffeomorphic with  $\mathbb{G}_+(2, n)$ . To see this, consider a 2-plane  $\Pi \subset \mathbb{R}^n$  that is spanned by  $u \wedge v$ , where the vectors  $u$  and  $v$  satisfy  $|u| - |v| = 0$  and  $\langle u, v \rangle = 0$ . Then, vector  $w = u + iv$  belongs to  $Q_{n-2}$ . Hence, to each oriented 2-plane we associate a point in  $Q_{n-2}$ . In fact, this correspondence actually is a diffeomorphism. Consequently, the map  $\bar{\varphi} : U \rightarrow Q_{n-2}$ ,  $\bar{\varphi} = f_x + if_y$ , is exactly the generalized Gauss map of the minimal surface.

Let  $M$  be a manifold of dimension  $2m$  endowed with a Riemannian metric  $g$  and a metric connection  $\nabla$ . An *almost complex structure* on  $M$  is by definition a bundle isomorphism  $J : TM \rightarrow TM$  satisfying  $J \circ J = -I$ . The pair  $(M, J)$  is called an *almost complex manifold*. If  $J$  is an isometry with respect to  $g$  and parallel with respect to  $\nabla$ , then the triple  $(M, g, J)$  is called *Kähler manifold*. In this case, the 2-form  $\Omega$  given by

$$\Omega(v_1, v_2) = g(Jv_1, v_2),$$

where  $v_1, v_2 \in \mathfrak{X}(M)$ , is closed and is called the *Kähler form*. A smooth map  $f : (M, J_M) \rightarrow (N, J_N)$  between Kähler manifolds is called *holomorphic* if  $df \circ J_M = J_N \circ df$  and *anti-holomorphic* if  $df \circ J_M = -J_N \circ df$ . If the map  $f$  is a holomorphic or anti-holomorphic isometric immersion, then  $f(M)$  will be called an *immersed complex submanifold* of  $N$ . Such immersions are automatically minimal.

With the terminology we just introduced and the discussion above, we can now state the following result which was originally proved by Chern [20].

**Theorem 4** *An oriented surface of the euclidean space is minimal if and only if its generalized Gauss map is anti-holomorphic.*

We will present now another interesting category of submanifolds, the so called Lagrangian submanifolds.

**Definition 11** Let  $M^m$  be a Riemannian manifold,  $(N^{2m}, g_N, \Omega)$  be a Kähler manifold and  $f: M^m \rightarrow N^{2m}$  an isometric immersion. The immersion  $f$  will be called *Lagrangian* if and only if  $f^*\Omega = 0$ .

Let us conclude this section with the following parametrization of minimal Lagrangian surfaces in  $\mathbb{R}^4$ ; see Chen and Morvan [19] and Aiyama [1, 2].

**Theorem 5** *Suppose that  $f, g: U \rightarrow \mathbb{C}$  are two holomorphic maps defined in a simply connected domain  $U$  of the complex plane satisfying  $|f_z|^2 + |g_z|^2 > 0$ . Then the map*

$$F = \frac{e^{i\beta/2}}{\sqrt{2}}(f - i\bar{g}, g + i\bar{f}),$$

where  $\beta$  is a real number, is a minimal conformal Lagrangian immersion in  $\mathbb{C}^2$ . The generalized Gauss map  $\mathcal{G}$  takes values in  $\mathbb{S}^2 \times \{(e^{i\beta}, 0)\} \simeq \mathbb{C} \cup \{\infty\}$  and is given by the formula

$$\mathcal{G} = f_z/g_z.$$

Conversely, every minimal Lagrangian immersion  $f: M \rightarrow \mathbb{C}^2$  can be, at least locally, parametrized as above.

## 4 Scalar and Vectorial Maximum Principles

The maximum principle is one of the most useful tools employed in the study of PDEs. All maximum principles rely on the following elementary result of calculus: *Suppose that  $\Omega$  is an open, bounded domain of  $\mathbb{R}^m$  and let  $u: \bar{\Omega} \rightarrow \mathbb{R}$  be a continuous function which is  $C^2$ -smooth in  $\Omega$ . If  $u$  attains its maximum at interior point  $x_0$ , then*

$$\nabla u(x_0) = 0 \quad \text{and} \quad \nabla^2 u(x_0) \leq 0.$$

As an immediate consequence of this fact is that any continuous and  $C^2$ -smooth up to the boundary strictly convex function must attain its maximum at the boundary of  $\Omega$ . In the matter of fact, one can show a little bit more: *Any continuous and  $C^2$ -*

smooth up to the boundary weakly convex function either attain its maximum at the boundary of  $\Omega$  or otherwise is constant. The above principle holds for a large class of solutions of partial differential inequalities.

### 4.1 Hopf's Maximum Principles

Suppose that  $\Omega$  is a bounded, open and connected domain of  $\mathbb{R}^m$ . We wish to study operators  $\mathcal{L}: C^2(\Omega) \rightarrow C^0(\Omega)$  of the form

$$\mathcal{L} = \sum_{i=1}^m a_{ij} \partial_{x_i} \partial_{x_j} + \sum_{i=1}^m b_i \partial_{x_i}, \quad (13)$$

where here  $a_{ij} = a_{ij}, b_j: \Omega \rightarrow \mathbb{R}, i, j \in \{1, \dots, m\}$ , are uniformly bounded functions and  $\partial_{x_i}, i \in \{1, \dots, m\}$ , the partial derivatives with respect the cartesian coordinates of  $\mathbb{R}^m$ . The symmetric matrix  $\mathcal{A}$  with coefficients the functions  $a^{ij}$  is called the *representative matrix* of  $\mathcal{L}$ .

**Definition 12** The operator  $\mathcal{L}$  is called *elliptic* if the matrix  $\mathcal{A}$  is positive at each point of  $\Omega$ . Moreover,  $\mathcal{L}$  is called *uniformly elliptic* if the smallest eigenvalue of its matrix  $\mathcal{A}$  is a function which is bounded away from zero.

**Theorem 6 (Hopf's Strong Maximum Principle)** *Let  $\Omega \subset \mathbb{R}^m$  be an open, connected and bounded domain. Suppose that  $u \in C^2(\Omega) \cap C^0(\overline{\Omega})$  is a solution of the differential inequality*

$$\mathcal{L}u + hu \geq 0,$$

where  $\mathcal{L}$  is a uniformly elliptic differential operator with uniformly bounded coefficients and  $h$  a continuous function on  $\overline{\Omega}$ .

- (a) Suppose that  $h = 0$  and that  $u$  attains its maximum at an interior point of  $\Omega$ . Then,  $u$  is constant.
- (b) Suppose that  $h \leq 0$  and that  $u$  attains a non-negative maximum at an interior point of  $\Omega$ . Then,  $u$  is constant.

For the proof see [35] or [72].

### 4.2 Maximum Principles for Systems

We would like to have a form of the maximum principle that is applicable for sections in vector bundles. To generalize, first note that Hopf's maximum principle for functions can be re-formulated as follows:

Let  $\Omega$  be an open subset of  $\mathbb{R}^m$  and  $u : \Omega \subset \mathbb{R}^m \rightarrow [a, b]$  a  $C^2$ -smooth function satisfying the uniformly elliptic differential equation

$$\sum_{i,j=1}^m a_{ij}u_{x_i x_j} + \sum_{i=1}^m b_i u_{x_i} = 0.$$

If a point of  $\Omega$  is mapped into a boundary point of  $[a, b]$ , then any point of  $\Omega$  is mapped into the boundary.

From this point of view of the statement of Hopf's maximum principle, one can guess how the generalization of the maximum principle for vector valued maps should be. The interval is replaced by a convex set  $K$  and the statement reads:

Let  $\Omega \subset \mathbb{R}^m$  be open,  $K \subset \mathbb{R}^n$  closed convex and  $u : \Omega \rightarrow K$  a  $C^2$ -smooth vector valued map satisfying the uniformly elliptic differential system

$$\sum_{i,j=1}^m a_{ij}u_{x_i x_j} + \sum_{i=1}^m b_i u_{x_i} = 0.$$

If a point of  $\Omega$  is mapped into a boundary point of  $K$  then every point is mapped into the boundary.

#### 4.2.1 Convexity and Distance Functions

A crucial role in the proof of the vectorial maximum principle plays the geometry of the (signed) distance function from the boundary of a convex set. In this subsection, we review the basic definitions of the geometry of convex sets in euclidean space such as supporting half-spaces, tangent cones, normal vectors and distance functions.

**Definition 13** A subset  $K$  of  $\mathbb{R}^n$  is called *convex* if for any pair of points  $z, w \in K$ , the segment

$$\mathcal{E}_{z,w} = \{tz + (1-t)w \in \mathbb{R}^n : t \in [0, 1]\}$$

is contained in  $K$ . The set  $K$  is said to be *strictly convex*, if for any pair  $z, w \in K$  the segment  $\mathcal{E}_{z,w}$  belongs to the interior of  $K$ .

A convex set  $K \subset \mathbb{R}^n$  may have non-smooth boundary. It is a well-known fact in Convex Geometry that the boundary  $\partial K$  is a continuous hypersurface of  $\mathbb{R}^n$ . In fact, according to a result of Reidemeister [73], the boundary  $\partial K$  is Lipschitz continuous and so almost everywhere differentiable. In particular, there is no well-defined tangent or normal space of  $K$  in the classical sense. However, there is a way to generalize these notions for convex subsets of  $\mathbb{R}^n$ .

**Definition 14** Let  $K$  be a closed convex subset of the euclidean space  $\mathbb{R}^n$ . A *supporting half-space* of the set  $K$  is a closed half-space of  $\mathbb{R}^n$  which contains  $K$  and has points of  $K$  on its boundary. A *supporting hyperplane* of  $K$  is a hyperplane which is the boundary of a *supporting half-space* of  $K$ . The *tangent cone*  $C_{y_0}K$  of  $K$  at  $y_0 \in \partial K$  is defined as the intersection of all supporting half-spaces of  $K$  that contain  $y_0$ .

**Definition 15** Let  $K \subset \mathbb{R}^n$  be a closed convex subset and  $y_0 \in \partial K$ . Then:

- (a) A non-zero vector  $\xi$  is called *normal vector* of  $\partial K$  at  $y_0$ , if  $\xi$  is normal to a supporting hyperplane of  $K$  passing through  $y_0$ . This normal vector is called *inward pointing*, if it points into the half-space containing  $K$ .
- (b) A vector  $\eta$  is called *inward pointing* at  $y_0 \in \partial K$ , if

$$\langle \xi, \eta \rangle \geq 0$$

for any inward pointing normal vector  $\xi$  at  $y_0$ .

Let  $K \subset \mathbb{R}^n$  be a closed convex set and  $d : \mathbb{R}^n \rightarrow \mathbb{R}$  the function given by

$$d(z) = \begin{cases} +\text{dist}(z, \partial K), & \text{if } z \in K, \\ -\text{dist}(z, \partial K), & \text{if } z \notin K. \end{cases}$$

Note that for each  $x \in \mathbb{R}^n$  there is at least one point  $y \in \partial K$  such that

$$\text{dist}(z, \partial K) = |y - z|.$$

Moreover, the function  $d$  is Lipschitz continuous. For a better understanding of the properties of  $d$ , let us suppose that  $\partial K$  is  $C^2$ -smooth. Denote by  $\xi$  the inward pointing unit normal vector field along  $\partial K$  and by the  $A$  the corresponding Weingarten operator. Because  $K$  is convex, from Hadamard's Theorem,  $A$  is non-negative definite. In particular,  $K$  is strictly convex if and only if  $A$  is positive definite; see for example [29]. Fix a point  $y_0 \in \partial K$ . In an open neighbourhood  $U \subset \mathbb{R}^n$  of  $y_0$ , the part  $U \cap \partial K$  can be parametrized via an embedding  $f : \Omega = U \cap T_{y_0}K \rightarrow \mathbb{R}^n$ , which assigns to each point of  $\Omega$  the height of  $\partial K$  from its tangent plane at  $y_0$ . Recall from multi-variable calculus that the distance of any point  $z \in K^0$  to  $\partial K$  is realized as the intersection of a straight line passing through  $z$  and meeting  $\partial K$  orthogonally. Hence, the *level set*

$$K_t = \{z \in K : d(z) = t\},$$

of  $d$  is parametrized locally via the map  $f_t : \Omega \rightarrow \mathbb{R}^n$  given by

$$f_t = f + t\xi.$$

**Proposition 1** *There exists a positive real number  $\varepsilon$ , such that  $f_t$  is an immersion for all  $t \in (-\varepsilon, \varepsilon)$ . Moreover, the unit normal along  $f_t$  coincide with the unit normal  $\xi$  of  $f$ . Additionally, the Weingarten operator  $A_t$  of  $f_t$  is related to the Weingarten operator  $A$  of  $f$  by the formula*

$$A_t = (I - tA)^{-1} \circ A.$$

*In particular,  $K_t$  is strictly convex if and only if  $\partial K$  is strictly convex.*

**Proof** We have

$$df_t = df + t d\xi = df \circ (I - tA).$$

Hence,  $\xi$  is a unit normal vector field along  $f_t$ . Therefore,

$$-df \circ A = d\xi = -df_t \circ A_t = -df \circ (I - tA) \circ A_t$$

and so

$$A = (I - tA) \circ A_t.$$

From the above formula we deduce that there exists a positive constant  $\varepsilon$  such that  $K_t$  is convex for all  $t \in (-\varepsilon, \varepsilon)$ . In addition, if  $\partial K$  is strictly convex, the level sets close to the boundary are also strictly convex. □

**Proposition 2** *Let  $K$  be a closed and convex set in  $\mathbb{R}^n$ .*

- (a) *For any  $y_0 \in \partial K$  there exists a neighbourhood  $U \subset \mathbb{R}^n$  containing  $y_0$ , such that  $d$  is  $C^2$ -smooth function on  $U \cap K^0$ .*
- (b) *Let  $v, w$  tangent vectors of  $K_{d(z)}$  at  $z \in K^0$ . Then, the Hessian  $\nabla^2 d$  of  $d$  satisfies*

$$\nabla^2 d(v, w) = -\frac{A(v, w)}{1 - d(z)A(v, w)},$$

*where  $A$  is the shape operator of  $\partial K$  associated to the inward pointing unit normal, and*

$$\nabla^2 d(v, \xi) = 0.$$

**Proof** Parametrize, locally, the boundary  $\partial K$  as the image of an embedding  $f : \Omega \subset \mathbb{R}^{n-1} \rightarrow \mathbb{R}^n$ . Define the map  $F : \Omega \times \mathbb{R} \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ , given by

$$F(x, t) = f(x) + t\xi(x).$$

Then,  $dF_{(y_0,t)}(\partial_t) = \xi_{y_0}$  and  $dF_{(y_0,t)}(v) = df_{x_0}(v - tAv)$ , for any index  $i \in \{1, \dots, n - 1\}$  and  $v \in T_{y_0}\Omega$ . From the inverse mapping theorem, there exists



an open subset  $D \subset \Omega$  and a positive number  $\varepsilon$  such that the map  $F$  is a  $C^2$ -diffeomorphism for any  $(x, t) \in D \times (-\varepsilon, \varepsilon)$ . Hence, under a change of coordinates,  $d$  may be regarded as a  $C^2$ -smooth function defined on  $D \times (-\varepsilon, \varepsilon)$ . In the matter of fact, in these coordinates, we have

$$d(x, t) = \langle F(x, t) - f(x), \xi(x) \rangle = t.$$

Therefore,  $\nabla d_{(x,t)} = \xi_x$ . Because  $|\nabla d| = 1$ , we deduce that  $\nabla^2 d$  vanishes on the normal bundle of any level set  $K_t$ . Moreover,  $\nabla^2 d = -A_t$  on the tangent space of any level set  $K_t$ . Now the desired result follows from Proposition 1. This completes the proof.  $\square$

### 4.2.2 Weinberger’s Maximum Principle

Weinberger [101] established a strong maximum principle for  $C^2$ -smooth maps  $u : \Omega \subset \mathbb{R}^m \rightarrow K \subset \mathbb{R}^n$  with values in a closed convex set  $K \subset \mathbb{R}^n$ , whose boundary  $\partial K$  satisfies some regularity conditions that he called “slab conditions”. Inspired by the ideas of Weinberger, Wang [93] gave a geometric proof of the strong maximum principle, in the case where the boundary  $\partial K$  of  $K$  is of class  $C^2$ . The idea of Wang was to apply Hopf’s maximum principle to the function  $d \circ u : \Omega \rightarrow \mathbb{R}$ , whose value at  $x$  is equal to the distance of  $u(x)$  from the boundary  $\partial K$ . Later, Evans [36] removed all additional regularity requirements on the boundary of  $K$ .

**Theorem 7 (Weinberger-Evans)** *Let  $K$  be a closed, convex set of  $\mathbb{R}^n$  and  $u : \Omega \subset \mathbb{R}^m \rightarrow K \subset \mathbb{R}^n$ ,  $u = (u^1, \dots, u^n)$ , a solution of the uniformly elliptic system of partial differential equations*

$$\mathcal{L}u(x) + \Psi(x, u(x)) = 0, \quad x \in \Omega,$$

where  $\Omega$  is a connected open domain of  $\mathbb{R}^m$  and  $\Psi : \Omega \times \mathbb{R}^n \rightarrow \mathbb{R}$  a continuous map that is Lipschitz continuous in the second variable. Suppose further that  $\Psi$  is pointing into  $K$ .

- (a) *If there is a point  $x_0 \in \Omega$  such that  $u(x_0) \in \partial K$ , then  $u(x) \in \partial K$  for any point  $x \in \Omega$ .*
- (b) *Assume additionally that the boundary  $\partial K$  is strictly convex. If there is a point  $x_0 \in \Omega$  such that  $u(x_0) \in \partial K$ , then  $u$  is constant.*

**Proof** Let us give the proof in the case where the boundary of  $K$  is smooth of class  $C^2$ , following the ideas in [93]. Consider the function  $f = d \circ u : \Omega \rightarrow \mathbb{R}$ . We compute,

$$f_{x_i} = \sum_{\alpha=1}^n d_{u^\alpha} u_{x_i}^\alpha, \tag{14}$$

and

$$f_{x_i x_j} = \sum_{\alpha=1}^n d_{u^\alpha} u_{x_i x_j}^\alpha + \sum_{\alpha, \beta=1}^n d_{u^{\alpha\beta}} u_{x_i}^\alpha u_{x_j}^\beta.$$

Consider now the uniformly differential operator  $\tilde{\mathcal{L}}$  given by

$$\tilde{\mathcal{L}} = \mathcal{L} - \sum_{i=1}^m b_i \partial_{x_i}.$$

By a straightforward computation, we get

$$\tilde{\mathcal{L}}f = \sum_{\alpha=1}^n d_{u^\alpha} \sum_{i,j=1}^m a_{ij} u_{x_i x_j}^\alpha + \sum_{i,j=1}^m a_{ij} \sum_{\alpha, \beta=1}^n d_{u^{\alpha\beta}} u_{x_i}^\alpha u_{x_j}^\beta. \quad (15)$$

Denote the first sum in the right hand side of (15) by I and the second sum by II. Observe at first that

$$I(x) = -\langle \nabla d_{u(x)}, \Psi(x, u(x)) \rangle.$$

We restrict our selves in a sufficiently small neighbourhood  $U \subset \mathbb{R}^n$  around  $u(x_0)$  and in an neighbourhood  $V$  of  $x_0$  such that  $u(V) \subset U$ . For each  $x \in V$ , denote by  $\hat{u}(x)$  the unique point on  $\partial K$  with the property

$$f(x) = d(u(x)) = |u(x) - \hat{u}(x)|.$$

Recall that the integral curves  $\nabla d$  are straight lines perpendicular to each level set of  $d$ . Thus,  $\nabla d_{u(x)} = \nabla d_{\hat{u}(x)}$ . Since  $\Psi$  is inward pointing, we get that

$$\langle \nabla d_{u(x)}, \Psi(x, \hat{u}(x)) \rangle = \langle \nabla d_{\hat{u}(x)}, \Psi(x, \hat{u}(x)) \rangle \geq 0.$$

Therefore, exploiting the Lipschitz property of  $\Psi$ , we get that

$$\begin{aligned} I(x) &= -\langle \nabla d_{u(x)}, \Psi(x, u(x)) - \Psi(x, \hat{u}(x)) + \Psi(x, \hat{u}(x)) \rangle \\ &= -\langle \nabla d_{u(x)}, \Psi(x, u(x)) - \Psi(x, \hat{u}(x)) \rangle - \langle \nabla d_{u(x)}, \Psi(x, \hat{u}(x)) \rangle \\ &\leq |\nabla d_{u(x)}| \cdot |\Psi(x, u(x)) - \Psi(x, \hat{u}(x))| \\ &\leq h(x) |u(x) - \hat{u}(x)| \\ &= h(x) f(x), \end{aligned}$$

where  $h$  is a non-negative bounded function. Recall from Proposition 1 that  $U$  is foliated by level sets of  $d$ . Thus, we can decompose  $u_{x_i}$  in the form

$$u_{x_i} = u_{x_i}^\top + u_{x_i}^\perp$$

where  $(\cdot)^\top$  denotes the orthogonal projection into the tangent space and  $(\cdot)^\perp$  the orthogonal projection into the normal space of the foliation. Bearing in mind the conclusions of Proposition 2, we deduce that

$$\sum_{\alpha, \beta=1}^n d_{u^\alpha u^\beta} u_{x_i}^\alpha u_{x_j}^\beta = \nabla^2 d(u_{x_i}^\top, u_{x_j}^\top) = \frac{A(u_{x_i}^\top, u_{x_j}^\top)}{1 - f(x)A(u_{x_i}^\top, u_{x_j}^\top)}.$$

Since,  $A$  is non-negative definite and  $\mathcal{A} = (a_{ij})$  is positive definite, we deduce that

$$\Pi = \text{trace}(\mathcal{A} \cdot \nabla^2 d) \leq 0.$$

In addition, for any  $x$  such that  $u(x) \in \partial K$ , we have that

$$\Pi(x) = \begin{cases} \text{strictly negative,} & \text{if } \partial K \text{ is strictly convex in close to } u(x), \\ \text{zero} & \text{if } \partial K \text{ is flat in a neighbourhood of } u(x). \end{cases} \quad (16)$$

Putting everything together, we get

$$\tilde{\mathcal{L}}f(x) - h(x)f(x) \leq 0.$$

Since  $f \geq 0$  and there exists a point  $x_0$  such that  $f(x_0) = 0$ , from Hopf's strong maximum principle we deduce that  $f \equiv 0$ . This implies now that all the values of  $u$  lie in the boundary of  $K$ . Moreover, going back to the original equation for  $f$ , we see that  $\Pi \equiv 0$ . Consequently, if  $\partial K$  is strictly convex, from (16) it follows that  $u$  must be constant. This completes the proof.  $\square$

### 4.3 Maximum Principles for Bundles

To state the maximum principle for sections in vector bundles, we must introduce an appropriate notion of convexity for subsets of vector bundles. Let us recall at first the following definition of Hamilton [45]:

**Definition 16 (Hamilton)** Suppose that  $E$  is a vector bundle over the manifold  $M$  and let  $K$  be a closed subset of  $E$ .

- (a) The set  $K$  is called *fiber-convex* or *convex in the fiber*, if for each  $x \in M$ , the set  $K_x = K \cap E_x$  is a convex subset of the fiber  $E_x$ .
- (b) The set  $K$  is called *invariant under parallel transport*, if for every smooth curve  $\gamma : [0, b] \rightarrow M$  and any vector  $v \in K_{\gamma(0)}$ , the unique parallel section  $v(t) \in E_{\gamma(t)}$ ,  $t \in [0, b]$ , along  $\gamma(t)$  with  $v(0) = v$ , is contained in  $K$ .
- (c) A *fiberwise map*  $\Psi : E \rightarrow E$  is a map such that  $\pi \circ \Psi = \pi$ , where  $\pi$  denotes the bundle projection. We say that a fiberwise map  $\Psi$  points into  $K$  (or is inward pointing), if for any  $x \in M$  and any  $\vartheta \in \partial K_x$  the vector  $\Psi(\vartheta)$  belongs to the tangent cone  $C_\vartheta K_x$  of  $K_x$  at  $\vartheta$ .

Let  $E$  be a Riemannian vector bundle over a manifold  $M$  equipped with a metric compatible connection. We consider uniformly elliptic operators  $\mathcal{L}$  that are given locally by

$$\mathcal{L} = \sum_{i,j=1}^m a_{ij} \nabla_{e_i, e_j}^2 + \sum_{j=1}^m b_j \nabla_{e_j}, \tag{E}$$

where  $\{e_1, \dots, e_m\}$  is a local orthonormal frame of  $M$ ,  $\mathcal{A} = (a_{ij})$  a symmetric, uniformly positive definite tensor and  $b = \sum_{i=1}^m b_i e_i$  is a smooth vector field.

For the proof of the maximum principle, we will use a result due to Böhm and Wilking [9].

**Lemma 1** *Let  $M$  be a Riemannian manifold and  $E$  a Riemannian vector bundle over  $M$  equipped with a metric compatible connection. Let  $K \subset E$  be a closed and fiber-convex subset which is invariant under parallel transport. If  $\phi$  is a smooth section with values in  $K$  then, for any  $x \in M$  and  $v \in T_x M$ , the Hessian*

$$\nabla_{v,v}^2 \phi = \nabla_v \nabla_v \phi - \nabla_{\nabla_v v} \phi$$

*belongs into the tangent cone of  $K_x$  at the point  $\phi|_x$ .*

**Proof** It suffices to prove the result in the case where there exists a point  $x_0$  which is mapped via  $\phi$  in the boundary of  $K$ , since otherwise the result is trivially true. Consider a unit vector  $v \in T_{x_0} M$  and an normal coordinate system  $(x_1, \dots, x_m)$  in an open neighbourhood  $U$  around a point  $x_0$  such that  $\partial_{x_1}|_{x_0} = v$ . Moreover, pick a basis  $\{\phi_1|_{x_0}, \dots, \phi_k|_{x_0}\}$  of  $E_{x_0}$  and extend it into a local geodesic orthonormal frame field. Then,

$$\phi = u_1 \phi_1 + \dots + u_k \phi_k,$$

where the components  $u_i : U \rightarrow \mathbb{R}, i \in \{1, \dots, k\}$ , are smooth functions. A simple computation shows

$$\begin{aligned} \nabla_{v,v}^2 \phi|_{x_0} &= \nabla_{\partial_{x_1}} \nabla_{\partial_{x_1}} \phi|_{x_0} - \nabla_{\nabla_{\partial_{x_1}} \partial_{x_1}} \phi|_{x_0} = \sum_{i=1}^k (\partial_{x_1} \partial_{x_1} u_i)(x_0) \phi_i|_{x_0} \\ &= \sum_{i=1}^k (u_i \circ \gamma)''(0) \phi_i|_{x_0}, \end{aligned}$$

where  $\gamma : (-\varepsilon, \varepsilon) \rightarrow U$  is a length minimizing geodesic such that

$$\gamma(0) = x_0 \quad \text{and} \quad \gamma'(0) = \partial_{x_1}|_{x_0}.$$

Define now the set

$$\mathcal{K} = \left\{ (y_1, \dots, y_k) \in \mathbb{R}^k : \sum_{i=1}^k y_i \phi_i|_{x_0} \in K_{x_0} \right\}.$$

Clearly  $\mathcal{K}$  is a closed and convex subset of  $\mathbb{R}^k$ . Since  $\phi \in K$  and  $K$  is invariant under parallel transport, we deduce that the curve  $\sigma : (-\varepsilon, \varepsilon) \rightarrow \mathbb{R}^k$ , given by

$$\sigma = (u_1 \circ \gamma, \dots, u_k \circ \gamma),$$

lies in  $\mathcal{K}$ . It suffices to prove that  $\sigma''(0)$  points into  $\mathcal{K}$ . Indeed, because  $\mathcal{K}$  is convex, for any unit inward pointing normal  $\xi$  of  $\mathcal{K}$  at  $\sigma(0)$ , we have

$$g(t) = \langle \xi, \sigma(t) - \sigma(0) \rangle \geq 0,$$

for any  $t \in (-\varepsilon, \varepsilon)$ . Because  $g$  attains its minimum at  $t = 0$ , from standard calculus we get that  $g''(0) \geq 0$ , which implies  $\langle \sigma''(0), \xi \rangle \geq 0$ . This completes the proof.  $\square$

*Remark 2* According to the above result, it follows that if  $\phi$  is a section lying in a set satisfying the conditions of Lemma 1 and  $\mathcal{L}$  is a uniformly elliptic operator of second order, then section  $\mathcal{L}\phi$  always points into  $K$ .

**Theorem 8 (Strong Elliptic Maximum Principle)** *Suppose that  $M$  is a Riemannian manifold (without boundary) and  $E$  a vector bundle of rank  $k$  over  $M$  equipped with a Riemannian metric  $g_E$  and a metric compatible connection. Let  $K$  be a closed fiber-convex subset of the bundle  $E$  that is invariant under parallel transport and  $\phi \in \Gamma(E)$ ,  $\phi : M \rightarrow K$ , a smooth section such that*

$$\mathcal{L}\phi + \Psi(\phi) = 0,$$

where  $\mathcal{L}$  is a uniformly elliptic operator of second order of the form given in (E) and  $\Psi$  is a smooth fiberwise map that points into  $K$ . If there exists a point  $x_0 \in M$  such that  $\phi|_{x_0} \in \partial K_{x_0}$ , then  $\phi|_x \in \partial K_x$  for any point  $x \in M$ . If, additionally, in a neighbourhood of  $\phi|_{x_0}$  the set  $K_{x_0}$  is strictly convex and the boundary  $\partial K_{x_0}$  is  $C^2$ -smooth, then  $\phi$  is a parallel section.

**Proof** We follow the exposition in [76]. Let  $\{\phi_1, \dots, \phi_k\}$  be a geodesic orthonormal frame field defined in a neighbourhood  $U$  around  $x_0 \in M$ . Hence,  $\phi = u_1\phi_1 + \dots + u_k\phi_k$ , where  $u_i : U \rightarrow \mathbb{R}$ ,  $i \in \{1, \dots, k\}$ , are smooth functions. With respect to this frame we have

$$\begin{aligned} \mathcal{L}\phi &= \sum_{i=1}^k (\mathcal{L}u_i + (\text{gradient terms of } u_i)) + \sum_{j=1}^k u_j g_E(\mathcal{L}\phi_j, \phi_i)\phi_i \\ &= -\sum_{i=1}^k g_E(\Psi(\phi), \phi_i)\phi_i. \end{aligned}$$

Therefore, the map  $u : U \rightarrow \mathbb{R}^k$ ,  $u = (u_1, \dots, u_k)$ , satisfies a uniformly elliptic system of second order of the form

$$\tilde{\mathcal{L}}u + \Phi(u) = 0, \tag{17}$$

where  $\Phi : \mathbb{R}^k \rightarrow \mathbb{R}^k$ ,  $\Phi = (\Phi_1, \dots, \Phi_k)$ , is given by

$$\Phi_i(u) = \mathfrak{g}_E(\Psi(\sum_{j=1}^k u_j \phi_j) + \sum_{j=1}^k u_j \mathcal{L}\phi_j, \phi_i), \tag{18}$$

for any  $i \in \{1, \dots, k\}$ . Consider now the convex set

$$\mathcal{K} = \{(y_1, \dots, y_k) \in \mathbb{R}^k : \sum_{i=1}^k y_i \phi_i|_{x_0} \in K_{x_0}\}.$$

**Claim 1:** For any point  $x \in U$  we have  $u(x) \in \mathcal{K}$ .

Indeed, fix a point  $x \in U$  and let  $\gamma : [0, 1] \rightarrow U$  be the geodesic curve joining the points  $x$  and  $x_0$ . Denote by  $\theta$  the parallel section which is obtained by the parallel transport of  $\phi|_x$  along the geodesic  $\gamma$ . Then,

$$\theta|_{\gamma(t)} = \sum_{i=1}^k y_i \phi_i|_{\gamma(t)},$$

where  $y_i : [0, 1] \rightarrow \mathbb{R}$ ,  $i \in \{1, \dots, k\}$ , are smooth functions. Because,  $\theta$  and  $\phi_i$ ,  $i \in \{1, \dots, k\}$ , are parallel along  $\gamma$ , it follows that

$$0 = \nabla_{\partial_t}^{\gamma^* E} \theta = \sum_{i=1}^k y_i'(t) \phi_i|_{\gamma(t)}.$$

Hence,  $y_i(t) = y_i(0) = u_i(x)$ , for any  $t \in [0, 1]$  and  $i \in \{1, \dots, k\}$ . Therefore,

$$\theta|_{\gamma(1)} = \theta|_{x_0} = \sum_{i=1}^k u_i(x) \phi_i|_{x_0}.$$

Since by our assumptions  $K$  is invariant under parallel transport, it follows that  $\theta|_{x_0} \in K_{x_0}$ . Hence,  $u(U) \subset \mathcal{K}$  and this proves Claim 1.

**Claim 2:** For any  $y \in \partial\mathcal{K}$  the vector  $\Phi(y)$  defined in (18) points into  $\mathcal{K}$  at  $y$ .

First note that the boundary of each slice  $K_x$  is invariant under parallel transport. From (18) we deduce that it suffices to prove that both terms appearing on the right hand side of (18) point into  $\mathcal{K}$ . The first term points into  $\mathcal{K}$  by assumption on  $\Psi$ . The second term is inward pointing due to Lemma 1. This completes the proof of Claim 2.

Observe now that the solution of the uniformly second order elliptic partial differential system (17) satisfies all the assumptions of Theorem 7. Therefore, because  $u(x_0) \in \partial\mathcal{K}$  it follows that  $u(U)$  is contained in the boundary  $\partial\mathcal{K}$  of  $\mathcal{K}$ . Consequently,  $\phi|_x \in \partial K$  for any  $x \in U$ . Since  $M$  is connected, we deduce that

$\phi(M) \subset \partial K$ . Note, that if  $\mathcal{K}$  is additionally strictly convex at  $u(x_0)$ , then the map  $u$  is constant. This implies that

$$\phi|_x = \sum_{i=1}^k u_i(x_0)\phi_i|_x$$

for any  $x \in U$ . Thus,  $\phi$  is a parallel section taking all its values in  $\partial K$ . □

### 4.3.1 Maximum Principles for Symmetric Tensors

Let  $(E, g_E)$  be a Riemannian vector bundle over a manifold  $M$ . For any  $\phi \in \text{Sym}(E^* \otimes E^*)$ , a real number  $\lambda$  is called *eigenvalue* of  $\phi$  with respect to  $g_E$  at the point  $x \in M$ , if there exists a non-zero vector  $v \in E_x$ , such that

$$\phi(v, w) = \lambda g_E(v, w),$$

for any  $w \in E_x$ . The linear subspace  $\text{Eig}_{\lambda, \phi}(x)$  of  $E_x$  given by

$$\text{Eig}_{\lambda, \phi}(x) = \{v \in E_x : \phi(v, w) = \lambda g_E(v, w), \text{ for any } w \in E_x\},$$

is called the *eigenspace* of  $\lambda$  at  $x$ . Since  $\phi$  is symmetric, it admits  $k$  real eigenvalues  $\lambda_1(x), \dots, \lambda_k(x)$  at each point  $x \in M$ . We will always arrange the eigenvalues such that  $\lambda_1(x) \leq \dots \leq \lambda_k(x)$ . If  $\lambda_1(x) \geq 0$  (resp.  $> 0$ ) we say that  $\phi$  is *non-negatively* (resp. *positively*) *definite* at  $x$ .

Before stating the main results, let us recall the following definition due to Hamilton [44].

**Definition 17** A fiberwise map  $\Psi : \text{Sym}(E^* \otimes E^*) \rightarrow \text{Sym}(E^* \otimes E^*)$  is said to satisfy the *null-eigenvector condition*, if whenever  $\vartheta$  is a non-negative symmetric 2-tensor at a point  $x \in M$  and if  $v \in T_x M$  is a null-eigenvector of  $\vartheta$ , then  $\Psi(\vartheta)(v, v) \geq 0$ .

The next theorem consists the elliptic analogue of the maximum principle of Hamilton [45]. More precisely:

**Theorem 9** *Let  $(M, g)$  be a Riemannian manifold (without boundary) and suppose that  $(E, g_E)$  is a Riemannian vector bundle over  $M$  equipped with a metric connection. Assume that  $\phi \in \text{Sym}(E^* \otimes E^*)$  is non-negative definite and satisfies*

$$\mathcal{L}\phi + \Psi(\phi) = 0,$$

*where  $\Psi$  is a smooth fiberwise map satisfying the null-eigenvector condition. If there is a point of  $M$  where  $\phi$  has a zero eigenvalue, then  $\phi$  must have a zero eigenvalue everywhere.*

**Proof** Denote by  $K$  the set of non-negative definite symmetric 2-tensors, i.e.,

$$K = \{\vartheta \in \text{Sym}(E^* \otimes E^*) : \vartheta \geq 0\}.$$

Each set  $K_x$  is a closed and convex. Then,

$$\partial K_x = \{\vartheta \in K_x : \text{exists nonzero } v \in T_x M \text{ such that } \vartheta(v, \cdot) = 0\}.$$

The tangent cone of  $K_x$  at a point  $\vartheta \in \partial K$  is given by

$$C_{\vartheta} K_x = \{\phi \in \text{Sym}(E_x^* \otimes E_x^*) : \phi(v, v) \geq 0, \text{ for all } v \in \text{Eig}_{0, \vartheta}(x)\}.$$

**Claim 1.** *The set  $K$  is invariant under parallel translation.*

Let  $\gamma : [0, 1] \rightarrow M$  be a geodesic,  $P_t$  the parallel transport operator of vectors along  $\gamma$  and  $\Pi_t$  the parallel transport operator of 2-tensors along the curve  $\gamma$ . Consider  $\vartheta \in K_{\gamma(0)}$ . Then, for any  $v \in T_{\gamma(0)}M$ , we have

$$\partial_t \{(\Pi_t \vartheta)(P_t v, P_t v)\} = (\nabla_{\partial_t} \Pi_t \vartheta)(P_t v, P_t v) + 2\Pi_t \theta (\nabla_{\partial_t} P_t v, P_t v) = 0.$$

Therefore, for any vector  $v \in T_{\gamma(0)}M$ , it holds  $(\Pi_t \vartheta)(P_t v, P_t v) = \vartheta(v, v)$ . Consequently, for any  $w \in T_{\gamma(t)}M$ , we obtain that

$$(\Pi_t \vartheta)(w, w) = \vartheta(P_t^{-1} w, P_t^{-1} w) \geq 0.$$

This proves the claim.

**Claim 2.** *Let  $\Psi : \text{Sym}(E^* \otimes E^*) \rightarrow \text{Sym}(E^* \otimes E^*)$  be a smooth fiberwise map satisfying the null-eigenvector condition. Then, for any  $x \in M$  and  $\vartheta \in \partial K$ , the vector  $\Psi(x, \vartheta)$  points into  $K$ .*

Indeed, let  $\vartheta \in \partial K_{x_0}$ . Then  $\Psi$  points inwards of  $K_{x_0}$  at  $\vartheta$  if and only if

$$\langle v^* \otimes v^*, \Psi(x, \vartheta) \rangle = \Psi(x, \vartheta)(v, v) \geq 0,$$

for any  $x$  in  $M$  and null-eigenvector  $v \in T_{x_0}M$  of  $\vartheta$ .

This complete the proof. □

### 4.3.2 A Second Derivative Test for Symmetric 2-tensors

**Theorem 10** *Let  $(M, g)$  be a Riemannian manifold (without boundary) and  $(E, g_E)$  a Riemannian vector bundle of rank  $k$  over the manifold  $M$  equipped with a metric connection  $\nabla$ . Suppose that  $\phi \in \text{Sym}(E^* \otimes E^*)$  is a smooth symmetric 2-tensor. If the biggest eigenvalue  $\lambda_k$  of  $\phi$  admits a local maximum  $\lambda$  at an interior point  $x_0 \in M$ , then*

$$(\nabla \phi)(v, v) = 0 \quad \text{and} \quad (\mathcal{L}\phi)(v, v) \leq 0,$$



for all  $v \in \text{Eig}_{\lambda, \phi}(x_0)$  and for all uniformly elliptic second order operators  $\mathcal{L}$ .

*Remark 3* The above theorem is due to Hamilton [44]. Replacing  $\phi$  by  $-\phi$  in Theorem 10, we get a similar result for the smallest eigenvalue  $\lambda_1$  of  $\phi$ .

**Proof** Let  $v \in \text{Eig}_{\lambda, \phi}(x_0)$  be a unit vector and  $V \in \Gamma(E)$  such that  $V|_{x_0} = v$  and  $\nabla V|_{x_0} = 0$ . Define the symmetric 2-tensor  $S$  given by  $S = \phi - \lambda g_E$ . From our assumptions,  $S$  is non-positive definite in a small neighbourhood of  $x_0$ . Moreover, the biggest eigenvalue of  $S$  at  $x_0$  equals 0. Consider the smooth function  $f : M \rightarrow \mathbb{R}$ , given by  $f(x) = S(V|_x, V|_x)$ . The function  $f$  is non-positive in the same neighbourhood around  $x_0$  and attains a local maximum at  $x_0$ . In particular,  $f(x_0) = 0$ ,  $df(x_0) = 0$  and  $(\mathcal{L}f)(x_0) \leq 0$ . Consider a local orthonormal frame  $\{e_1, \dots, e_m\}$  with respect to  $g$  defined in a neighbourhood of the point  $x_0$ . A simple calculation yields

$$df(e_i) = (\nabla_{e_i} S)(V, V) + 2S(\nabla_{e_i} V, V).$$

Taking into account that  $g_E$  is parallel, we deduce that

$$0 = (\nabla f)(x_0) = (\nabla S)(v, v) = (\nabla \phi)(v, v).$$

Furthermore,

$$\begin{aligned} \nabla_{e_i, e_j}^2 f &= (\nabla_{e_i, e_j}^2 S)(V, V) + 2S(V, \nabla_{e_i, e_j}^2 V) \\ &\quad + 2(\nabla_{e_i} S)(\nabla_{e_j} V, V) + 2(\nabla_{e_j} S)(\nabla_{e_i} V, V) \\ &\quad + 2S(\nabla_{e_i} V, \nabla_{e_j} V). \end{aligned}$$

Bearing in mind the definition of  $S$  and using the fact that  $g_E$  is parallel with respect to  $\nabla$ , we obtain

$$\begin{aligned} \mathcal{L}f &= (\mathcal{L}\phi)(V, V) + 2S(V, \mathcal{L}V) \\ &\quad + 2\sum_{i,j=1}^m a_{ij} \{S(\nabla_{e_i} V, \nabla_{e_j} V) + 2(\nabla_{e_i} S)(\nabla_{e_j} V, V)\} \\ &= (\mathcal{L}\phi)(V, V) + 2S(V, \mathcal{L}V) \\ &\quad + 2\sum_{i,j=1}^m a_{ij} \{S(\nabla_{e_i} V, \nabla_{e_j} V) + 2(\nabla_{e_i} S)(\nabla_{e_j} V, V)\}. \end{aligned}$$

Estimating at  $x_0$  and taking into account that  $V|_{x_0} = v$  is a null eigenvector of  $S$  at  $x_0$ , we get

$$0 \geq (\mathcal{L}f)(x_0) = (\mathcal{L}\phi)(v, v).$$

This completes the proof. □

## 5 Graphical Submanifolds

### 5.1 Definitions

Let  $(M, g_M)$  and  $(N, g_N)$  be Riemannian manifolds of dimension  $m$  and  $n$ , respectively. The induced metric on  $M \times N$  will be denoted by  $g_{M \times N} = g_M \times g_N$ . We often denote the product metric also by  $\langle \cdot, \cdot \rangle$ . A smooth map  $f : M \rightarrow N$  defines an embedding  $F : M \rightarrow M \times N$ , given by  $F(x) = (x, f(x))$ , for any  $x \in M$ . The *graph* of  $f$  is defined to be the submanifold

$$\Gamma(f) = F(M) = \{(x, f(x)) \in M \times N : x \in M\}.$$

Since  $F$  is an embedding, it induces another Riemannian metric  $g = F^*g_{M \times N}$  on  $M$ . The two natural projections  $\pi_M : M \times N \rightarrow M$  and  $\pi_N : M \times N \rightarrow N$  are submersions, that is they are smooth and have maximal rank. Note that the tangent bundle of the product manifold  $M \times N$ , splits as a direct sum

$$T(M \times N) = TM \oplus TN.$$

The four metrics  $g_M, g_N, g_{M \times N}$  and  $g$  are related by

$$g_{M \times N} = \pi_M^*g_M + \pi_N^*g_N \quad \text{and} \quad g = F^*g_{M \times N} = g_M + f^*g_N. \quad (19)$$

The Levi-Civita connection  $\nabla^{g_{M \times N}}$  associated to the Riemannian metric  $g_{M \times N}$  on  $M \times N$  is related to the Levi-Civita connections  $\nabla^{g_M}$  on  $(M, g_M)$  and  $\nabla^{g_N}$  on  $(N, g_N)$  by

$$\nabla^{g_{M \times N}} = \pi_M^*\nabla^{g_M} \oplus \pi_N^*\nabla^{g_N}.$$

The corresponding curvature operator  $\tilde{R}$  on the product  $M \times N$  is related to the curvature operators on  $(M, g_M)$  and  $R_N$  on  $(N, g_N)$  by

$$\tilde{R} = \pi_M^*R_M \oplus \pi_N^*R_N.$$

The map  $f : M \rightarrow N$  is called *minimal* if  $\Gamma(f) \subset M \times N$  is minimal.

### 5.2 Singular Value Decomposition

For any fixed point  $x \in M$ , let  $\lambda_1^2(x) \leq \dots \leq \lambda_m^2(x)$  be the eigenvalues of  $f^*g_N$  with respect to  $g_M$ . The corresponding values  $\lambda_i \geq 0$ ,  $i \in \{1, \dots, m\}$ , are usually called *singular values* of the differential  $df$  of  $f$  at the point  $x$ . Let  $r = r(x) = \text{rank } df(x)$ . Then,  $r \leq \min\{m, n\}$  and  $\lambda_1(x) = \dots = \lambda_{m-r}(x) = 0$ .

It is well-known that the singular values can be used to define the so-called *singular decomposition* of  $df$ . At the point  $x$ , consider an orthonormal basis  $\{\alpha_1, \dots, \alpha_{m-r}; \alpha_{m-r+1}, \dots, \alpha_m\}$  with respect to  $g_M$  which diagonalizes  $f^*g_N$ . Moreover, at  $f(x)$  consider an orthonormal basis  $\{\beta_1, \dots, \beta_{n-r}; \beta_{n-r+1}, \dots, \beta_n\}$  with respect to  $g_N$  such that, for any  $i \in \{m-r+1, \dots, m\}$ ,

$$df(\alpha_i) = \lambda_i(x)\beta_{n-m+i}.$$

It is well-known fact that, with the above ordering, the singular values give rise to continuous functions. In the matter of fact, they are even smooth on an open and dense subset of  $M$ . In particular, they are smooth on open subsets where the corresponding multiplicities are constant and the corresponding eigenspaces are smooth distributions; see [83].

We may define a special basis for the tangent and the normal space of the graph in terms of the singular values. The vectors

$$e_i = \begin{cases} \alpha_i & , 1 \leq i \leq m-r, \\ \frac{\alpha_i}{\sqrt{1 + \lambda_i^2(x)}} & , m-r+1 \leq i \leq m, \end{cases} \tag{20}$$

form an orthonormal basis with respect to the metric  $g_{M \times N}$  of the tangent space  $dF(T_x M)$  of the graph  $\Gamma(f)$  at  $x$ . Moreover, the vectors

$$\xi_i = \begin{cases} \beta_i & , 1 \leq i \leq n-r, \\ \frac{-\lambda_{i+m-n}(x)\alpha_{i+m-n} \oplus \beta_i}{\sqrt{1 + \lambda_{i+m-n}^2(x)}} & , n-r+1 \leq i \leq n, \end{cases} \tag{21}$$

give an orthonormal basis with respect to  $g_{M \times N}$  of the normal space  $N_x M$ .

### 5.3 Length and Area Decreasing Maps

Let  $(M, g_M)$  and  $(N, g_N)$  be two Riemannian manifolds of dimensions  $m$  and  $n$  respectively. For any smooth map  $f : M \rightarrow N$  its differential  $df$  induces a map  $\Lambda^k df : \Lambda^k TM \rightarrow \Lambda^k TN$  given by

$$\left(\Lambda^k df\right)(v_1, \dots, v_k) = df(v_1) \wedge \dots \wedge df(v_k),$$

for any smooth vector fields  $v_1, \dots, v_k \in TM$ . The map  $\Lambda^k df$  is called the *k-Jacobian* of  $f$ . The *supremum norm* or the *k-dilation*  $|\Lambda^k df|(x)$  of the map  $f$  at a point  $x \in M$  is defined as the supremum of

$$\sqrt{\det ((f^* g_N(v_i, v_j))_{1 \leq i, j \leq k})}$$

when  $\{v_1, \dots, v_m\}$  runs over all orthonormal bases of  $T_x M$ .

**Definition 18** A smooth map  $f : (M, g_M) \rightarrow (N, g_N)$  between Riemannian manifolds is called (weakly)  $k$ -volume decreasing if  $|\Lambda^k df| \leq 1$ , strictly  $k$ -volume decreasing if  $|\Lambda^k df| < 1$  and  $k$ -volume preserving if  $|\Lambda^k df| = 1$ . For  $k = 1$  we use the term *length* instead of 1-volume and if  $k = 2$  we use the term *area* instead of 2-volume.

There is a way to express the length and area decreasing property of a map in terms of positivity of symmetric tensors. Define on  $M$  the symmetric 2-tensors  $S_{M \times N}$  and  $S$  given by

$$S_{M \times N} = \pi_M^* g_M - \pi_N^* g_N \quad \text{and} \quad S = F^* S_{M \times N} = g_M - f^* g_N.$$

With respect to the basis of the singular value decomposition, we have

$$S_{M \times N}(e_i, e_j) = \frac{1 - \lambda_i^2}{1 + \lambda_i^2} \delta_{ij}, \quad 1 \leq i, j \leq m. \tag{22}$$

Hence, the eigenvalues  $\mu_1, \mu_2, \dots, \mu_m$  of  $S$  with respect to  $g$ , are

$$\mu_1 = \frac{1 - \lambda_m^2}{1 + \lambda_m^2} \leq \dots \leq \mu_m = \frac{1 - \lambda_1^2}{1 + \lambda_1^2}.$$

Hence,  $f$  is length decreasing if  $S \geq 0$ . Additionally let us mention that

$$S_{M \times N}(\xi_i, \xi_j) = \begin{cases} -\delta_{ij} & , 1 \leq i \leq n - r, \\ -\frac{1 - \lambda_{i+m-n}^2}{1 + \lambda_{i+m-n}^2} \delta_{ij} & , n - r + 1 \leq i \leq n. \end{cases} \tag{23}$$

and

$$S_{M \times N}(e_{m-r+i}, \xi_{n-r+j}) = -\frac{2\lambda_{m-r+i}}{1 + \lambda_{m-r+i}^2} \delta_{ij}, \quad 1 \leq i, j \leq r. \tag{24}$$

Observe now that, for any pair of indices  $i, j \in \{1, \dots, m\}$ , we have

$$\mu_i + \mu_j = \frac{1 - \lambda_i^2}{1 + \lambda_i^2} + \frac{1 - \lambda_j^2}{1 + \lambda_j^2} = \frac{2(1 - \lambda_i^2 \lambda_j^2)}{(1 + \lambda_i^2)(1 + \lambda_j^2)}.$$

Hence, the map is strictly area decreasing, if and only if the tensor  $S$  is *strictly 2-positive*, i.e., the sum of the two smallest eigenvalues is positive. The 2-positivity of a tensor  $T \in \text{Sym}(T^*M \otimes T^*M)$  can be expressed as the positivity of another tensor  $T^{[2]} \in \text{Sym}(\Lambda^2 T^*M \otimes \Lambda^2 T^*M)$ . Indeed, let  $P$  and  $Q$  be two symmetric 2-tensors. Then, the *Kulkarni-Nomizu product*  $P \oslash Q$  given by

$$(P \oslash Q)(v_1 \wedge w_1, v_2 \wedge w_2) = P(v_1, v_2)Q(w_1, w_2) + P(w_1, w_2)Q(v_1, v_2) \\ - P(w_1, v_2)Q(v_1, w_2) - P(v_1, w_2)Q(w_1, v_2)$$

is an element of the vector bundle  $\text{Sym}(\Lambda^2 T^*M \otimes \Lambda^2 T^*M)$ . Now, to every element  $T \in \text{Sym}(T^*M \otimes T^*M)$  let us assign an element  $T^{[2]}$  of the bundle  $\text{Sym}(\Lambda^2 T^*M \otimes \Lambda^2 T^*M)$ , by setting

$$T^{[2]} = T \oslash g.$$

We point out that the Riemannian metric  $G$  of  $\Lambda^2 TM$  is given by

$$G = \frac{1}{2}g \oslash g = \frac{1}{2}g^{[2]}.$$

The relation between the eigenvalues of the tensor  $T$  and the eigenvalues of  $T^{[2]}$  is explained in the following lemma:

**Lemma 2** *Let  $T$  be a symmetric 2-tensor with eigenvalues  $\mu_1 \leq \dots \leq \mu_m$  and corresponding eigenvectors  $\{v_1, \dots, v_m\}$  with respect to the metric  $g$ . Then the eigenvalues of the symmetric 2-tensor  $T^{[2]}$  with respect to  $G$  are  $\mu_i + \mu_j$ , for any  $1 \leq i < j \leq m$ , with corresponding eigenvectors  $v_i \wedge v_j$ , for any  $1 \leq i < j \leq m$ .*

### 5.4 Minimal Graphs in the Euclidean Space

Let us discuss the case of graphs generated by smooth maps  $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ . The induced metric  $g$  on the graph is given in local coordinates in the form

$$g_{ij} = \delta_{ij} + \sum_{i,j=1}^m \sum_{\alpha=1}^n f_{x_i}^\alpha f_{x_j}^\alpha.$$

As usual, the components of the inverse matrix of the induced metric  $g$  are denoted by  $g^{ij}$ . It is not difficult to show that  $\Gamma(f)$  is minimal if and only if the components of

$$f = (f^1, \dots, f^m)$$

satisfy the following system of differential equations

$$\sum_{i,j=1}^m g^{ij} f_{x_i x_j}^\alpha = 0. \tag{MSE}$$

The equation is known in the literature as the *minimal surface equation*. For graphical hypersurfaces, that is for graphs generated by functions smooth  $f : U \subset \mathbb{R}^m \rightarrow \mathbb{R}$ , minimality is expressed by the equation

$$\operatorname{div} \left( \frac{\operatorname{grad} f}{\sqrt{1 + |\operatorname{grad} f|^2}} \right) = 0.$$

There is a long history of attempts to study entire solutions of the minimal surface equation. Bernstein [8] proved that the only entire minimal graphs in the  $\mathbb{R}^3$  are planes. However, there was a gap in the original proof of Bernstein which was fixed 40 years later; see [51, 65]. In the meantime, several complex analysis proofs have been obtained; for more details see the surveys of Osserman [70, 71].

It was conjectured for a long time that the theorem of Bernstein holds in any dimension for graphical hypersurfaces. For  $m = 3$ , its validity was proved by De Giorgi [30], for  $n = 4$  by Almgren [3] and for  $m = 5, 6, 7$  by Simons [80]. It was a big surprise when Bombieri, De Giorgi and Giusti [10] proved that, for  $m \geq 8$ , there are entire solutions of the minimal surface equation other than the affine ones.

In higher codimensions, the situation is more complicated. There are plenty of non-flat entire minimal graphs. For example, the graph of an entire holomorphic map  $f : \mathbb{C}^m \rightarrow \mathbb{C}^n$  is minimal. Moreover, Osserman [70] has constructed examples of complete minimal two-dimensional graphs in  $\mathbb{R}^4$ , which are not complex analytic with respect to any orthogonal complex structure on  $\mathbb{R}^4$ . For instance, the graph  $\Gamma(f)$  over the map  $f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ , given by

$$f(x, y) = \left( e^{\frac{x}{2}} - 3e^{-\frac{x}{2}} \right) \left( \cos \frac{y}{2}, -\sin \frac{y}{2} \right)$$

for any  $(x, y) \in \mathbb{R}^2$ , is such an example. Now the obvious questions became:

*Question 1* If entire solutions of the minimal surface equation need not be linear, do they have any other distinguishing characteristics? What additional restrictions on entire solutions would guarantee linearity in all dimensions?

The first result in this direction was obtained by Osserman [69] for two-dimensional graphs, generated by maps  $f : \mathbb{R}^2 \rightarrow \mathbb{R}^n$ . He proved that if the differential  $df$  of the map  $f$  is bounded, then must be a plane. In fact, he proved the following more general theorem:

**Theorem 11** *Suppose that  $\Sigma$  is a complete, oriented minimal surface (not necessarily graphical) in the euclidean space  $\mathbb{R}^n$ . Assume that the Gauss map of  $\Sigma$  omits an open neighbourhood in the Grassmannian. Then,  $\Sigma$  is flat.*

Let us restrict now in two dimensional graphs in  $\mathbb{R}^4$ , i.e., minimal graphs generated by maps  $f = (f_1, f_2)$ . For such graphs, Simon [81] proved that if one component of  $f$  have bounded gradient, then  $f$  is affine. Later on, Schoen [79] obtained a Bernstein-type result by imposing the assumption that  $f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  is a diffeomorphism. Moreover, Ni [68] has derived a result of Bernstein type under the assumption that  $f$  is an area-preserving map. In this case, area preserving is equivalent with the condition  $|\det(df)| = 1$ . The function  $\text{Jac}(f) = \det(df)$  is called the *Jacobian determinant* of  $f$ . All these result were generalized by Hasanis, Savas-Halilaj and Vlachos in [47, 48], just by assuming that  $\text{Jac}(f)$  is bounded. In fact, the following result is shown:

**Theorem 12** *Let  $f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  be an entire solution of the minimal surface equation. Assume that  $\Gamma(f)$  is not a plane. Then,  $\Gamma(f)$  is a complex analytic curve if and only if the Jacobian determinant  $\text{Jac}(f)$  of  $f$  does not take every real value. In particular if  $\Gamma(f)$  is a complex analytic curve, then:*

- (a) *The Jacobian determinant  $\text{Jac}(f)$  takes every real value in  $(0, +\infty)$  or in  $[0, +\infty)$  if  $f$  is holomorphic.*
- (b) *The Jacobian determinant  $\text{Jac}(f)$  takes every real value in  $(-\infty, 0)$  or in  $(-\infty, 0]$ , if  $f$  is anti-holomorphic.*

All these proofs use strongly the fact that the Gauss map of a minimal surface in the euclidean space is anti-holomorphic.

The first Bernstein-type theorem which was valid for arbitrary dimension and codimension is due to Hildebrandt, Jost, and Widman [50]. They obtained such a result under the assumption of a certain quantitative bound for the slope, that is a bound on the norm of the differential of the generating map.

Let us describe here briefly their technique. Note at first that a bound on the differential of the map forces the Gauss map of the graph to lie in a bounded region of the Grassmannian manifold. In particular, the first step is to determine which bounds on the differential will force the Gauss map to have its range in a sufficiently small convex subset of the Grassmannian. The second step is to find a convex function defined on the convex set, which contains the Gauss image of the graph, and to compose it with the Gauss map. By Theorem 2 of Ruh and Vilms, the Gauss map is harmonic. Consequently, the composition of the Gauss map with the convex function will give rise to a subharmonic function defined on the graph. The third step is to show that this particular subharmonic function is constant and the Gauss map is parallel. Of course, there are many difficulties to overcome to run this program. The first problem is the complexity of the Grassmannian manifolds. For example, it is not so easy to identify which are the convex subsets of the Grassmannian and their corresponding convex supporting functions. One way is to consider distance balls. In fact, Hildebrandt, Jost, and Widman [50] identified the largest ball in the Grassmannian manifold on which the square of the distance function is convex. Another major difficulty is that an entire euclidean minimal graph is complete and non-compact. Consequently, the standard maximum principle cannot be applied directly. Let us mention here that the original assumption on the slope was obtained by Hildebrandt, Jost, and Widman in [50] was

$$E(f) = \sqrt{\det(I + df^t df)} \leq \beta_0 < \cos^{-p} \left( \frac{\pi}{2\sqrt{2p}} \right)$$

where  $\beta_0$  is a constant and  $p = \min\{m, n\}$ . Over the years, the bound on  $E(f)$  was improved. Recently, Jost, Xin, and Yang [56] proved the following:

**Theorem 13** *Let  $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$  be an entire solution of the minimal surface equation. Suppose that there exists a number  $\beta_0$  such that*

$$\beta_0 < \begin{cases} 3, & \text{if } n \geq 2, \\ \infty, & \text{if } n = 1, \end{cases}$$

and

$$E(f) = \sqrt{\det(I + df^t df)} \leq \beta_0.$$

Then  $\Gamma(f)$  is an affine subspace of  $\mathbb{R}^m \times \mathbb{R}^n$ .

*Remark 4* For codimension one graphs, the above theorem was first obtained by Moser [67].

*Question 2* Let  $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$  be an entire solution of the minimal surface equation such that

$$E(f) = \sqrt{\det(I + df^t df)} < 9.$$

Is it true that  $\Gamma(f)$  is an affine subspace of  $\mathbb{R}^m \times \mathbb{R}^n$ ?

*Remark 5* The number 9 in the above conjecture should be the sharp bound. The reason is that there are examples of Lipschitz minimal maps constructed by Lawson and Osserman [58] with  $E(f) = 9$ ; see also [37]. These examples are generated from the map  $f : \mathbb{C}^2 - \{0\} = \mathbb{R}^4 - \{0\} \rightarrow \mathbb{R} \times \mathbb{C} = \mathbb{R}^3$  given by

$$f(x) = \frac{\sqrt{5}}{2} |x| \mathcal{H} \left( \frac{x}{|x|} \right),$$

where  $\mathcal{H} : \mathbb{C}^2 \rightarrow \mathbb{R} \times \mathbb{C}$  is the Hopf-map  $\mathcal{H}(z, w) = (|z|^2 - |w|^2, 2z\bar{w})$ .

Let us conclude this section by mentioning some results in special situations. The first one is due to Fischer-Colbrie [37] and it says that a 3-dimensional complete minimal graph with bounded differential is totally geodesic. In the matter of fact, the following holds:

**Theorem 14** *Let  $f : \mathbb{R}^3 \rightarrow \mathbb{R}^n$  be an entire solution of the minimal surface equation. If  $|df|$  is uniformly bounded, then  $\Gamma(f)$  is flat.*



In codimension two, no specific bound on the differential of  $f$  is needed. Recently, Assimos and Jost [6] obtained the following interesting theorem:

**Theorem 15** *Let  $f : \mathbb{R}^m \rightarrow \mathbb{R}^2$  be an entire solution of the minimal surface equation. Suppose that there exists a number  $\beta_0$  such that*

$$E(f) = \sqrt{\det(I + df^t df)} \leq \beta_0.$$

*Then  $\Gamma(f)$  is an affine subspace of  $\mathbb{R}^m \times \mathbb{R}^2$ .*

The next result we would like to mention is due to Wang [97]. He obtained the following theorem for strictly area decreasing minimal graphs.

**Theorem 16** *Let  $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$  be an entire solution of the minimal surface equation. Suppose that there exists numbers  $\delta_1 \in (0, 1)$  and  $\delta_2 > 0$  such that  $|\Lambda^2 df| \leq 1 - \delta_1$  and  $E(f) \leq \delta_2$ . Then  $\Gamma(f)$  is an affine subspace of  $\mathbb{R}^m \times \mathbb{R}^n$ .*

*Remark 6* The above cannot be extended for  $k$ -volume decreasing minimal maps with  $k > 2$ . For example, consider  $f : \mathbb{C}^2 = \mathbb{R}^4 \rightarrow \mathbb{C}^2 = \mathbb{R}^4$ , given by

$$f(z, w) = (\beta_0 z + h(w), w),$$

where  $z, w \in \mathbb{C}$ ,  $h : \mathbb{C} \rightarrow \mathbb{C}$  is a non-affine holomorphic map and  $\beta_0$  a real number. Observe that the graph  $\Gamma(f)$  is a non-flat minimal submanifold of  $\mathbb{R}^8$  and  $|\Lambda^4 df| = |\beta_0|$ . Consequently, there exists an abundance of non-flat minimal graphs in the euclidean space with arbitrary small 4-Jacobian.

## 6 Mean Curvature Flow

In this section, we introduce the notion of the mean curvature flow. Later, we will examine how various geometric quantities evolve under the mean curvature flow. Suppose that  $M$  is a manifold of dimension  $m$ , let  $T > 0$  be a real number and  $F : M \times [0, T) \rightarrow N$  a smooth time-dependent family of immersions of  $M$  into a Riemannian manifold  $N$  of dimension  $n$ . We follow the exposition in [5, 32, 63, 86].

**Definition 19** Let  $N$  be a Riemannian manifold. We say that a family of immersions  $F : M \times [0, T) \rightarrow N$  evolves by *mean curvature flow* (MCF for short) with initial data the immersion  $F_0 : M \rightarrow N$  if it satisfies the initial value problem

$$\begin{cases} dF_{(x,t)}(\partial_t) = H(F(x, t)) \\ F(x, 0) = F_0(x) \end{cases},$$

for any  $(x, t) \in M \times [0, T)$ , where  $H(F(x, t))$  denotes the mean curvature vector of the immersion  $F(\cdot, t) : M \rightarrow N$  at the point  $x \in M$ .

### 6.1 Basic Facts for Systems of Parabolic PDEs

In this section we recall basic facts about solvability of Cauchy problems; for more details see [57].

#### 6.1.1 Differential Operators

Let  $M$  be a smooth manifold equipped with a Riemannian metric  $g$  whose associated Levi-Civita connection is  $\nabla^M$ . Suppose that  $E_1$  and  $E_2$  are two vector bundles over  $M$  and assume that  $E_1$  is equipped with a Riemannian metric  $h$  and a compatible connection  $\nabla^{E_1}$ . As in Definition 4, from the connections  $\nabla^M$  and  $\nabla^{E_1}$ , one can form the  $k$ -th derivative  $\nabla^k$  of a section  $\phi \in \Gamma(E_1)$ .

**Definition 20** A map  $P: \Gamma(E_1) \rightarrow \Gamma(E_2)$  of the form

$$(P\phi)(x) = Q(x, \nabla^1\phi(x), \dots, \nabla^k\phi(x)) \in (E_2)_x,$$

where  $Q$  is smooth in all its variables, will be called *differential operator of order  $k$* . In the case where  $P$  is  $\mathbb{R}$ -linear, we say that  $P$  is a *linear differential operator of order  $k$* . Otherwise, we say that  $P$  is non-linear.

Suppose that  $P: \Gamma(E_1) \rightarrow \Gamma(E_2)$  is a linear differential operator of degree  $k$ . Then, in index notation, it can be written in the form

$$P\phi = \sum_{i_1, \dots, i_k} A^{i_1 \dots i_k} \nabla_{\partial_{x_{i_1}} \dots \partial_{x_{i_k}}}^k \phi + \dots + \sum_{i_1} A^{i_1} \nabla_{\partial_{x_{i_1}}}^1 \phi + A^0 \phi,$$

where for each  $x \in M$ ,  $A^{i_1 \dots i_k}(x): (E_1)_x \rightarrow (E_2)_x$  is linear map. These maps are called the *coefficients* of the linear operator  $P$ .

**Definition 21** Let  $P: \Gamma(E_1) \rightarrow \Gamma(E_2)$  be a linear differential operator of order  $k$ , let  $x$  be a point in  $M$  and  $\zeta = \sum_i \zeta_i dx_i \in T_x^*M$ . The linear map  $\sigma_\zeta(P; x): (E_1)_x \rightarrow (E_2)_x$ , given by

$$\sigma_\zeta(P; x)\phi = \sum_{i_1, \dots, i_k} \zeta_{i_1} \dots \zeta_{i_k} A^{i_1 \dots i_k} \phi|_x,$$

is called the *principal symbol of the operator  $P$  at the point  $x$  and in the direction  $\zeta$* . In particular, the operator  $P$  is called *elliptic* if its principal symbol is an isomorphism, for every point  $x$  and every non-zero direction  $\zeta$ .

**Definition 22** The *differential* or the *linearization* of  $P$  at  $\phi_0$ , if it exists, is defined to be the linear map  $DP|_{\phi_0}: \Gamma(E_1) \rightarrow \Gamma(E_2)$ , given by the expression

$$DP|_{\phi_0}(\psi) = \lim_{s \rightarrow 0} \frac{P(\phi_0 + s\psi) - P(\phi_0)}{s},$$

for any  $\psi \in \Gamma(E_1)$ .

**Definition 23** Let  $P: \Gamma(E_1) \rightarrow \Gamma(E_2)$  be a differential operator of order  $k$ . We say that  $P$  is *elliptic, undetermined elliptic or overdetermined elliptic* if its linearization is so.

*Example 1* Let  $f: M \rightarrow N$  be a smooth map between manifolds endowed with Riemannian metrics  $g_M$  and  $g_N$ , respectively, and consider the operator  $\Delta_{g_M, g_N}: C^\infty(M) \rightarrow C^\infty(M)$ , given by

$$\Delta_{g_M, g_N} f = \text{tr}_{g_M} B,$$

where  $B$  stands for the Hessian of  $f$ . In local coordinates, we have

$$\Delta_{g_M, g_N} f = \sum_{i, j, \alpha} g_M^{ij} (f_{x_i x_j}^\alpha - \sum_k \Gamma_{ij}^k f_{x_k}^\alpha + \sum_{\gamma, \delta} \Gamma_{\gamma \delta}^\alpha f_{x_i}^\gamma f_{x_j}^\delta) \partial_{y_\alpha}.$$

The linearization of  $\Delta_{g_M, g_N} f$  is

$$\begin{aligned} D\Delta_{g_M, g_N}|_f(G) &= \lim_{s \rightarrow \infty} \frac{\Delta_{g_M, g_N}(f + sG) - \Delta_{g_M, g_N}(f)}{s} \\ &= \sum_{i, j} g_M^{ij} G_{x_i x_j}^\alpha \partial_{y_\alpha} + \text{lower order terms.} \end{aligned}$$

Hence, for any

$$\zeta = (\zeta_1, \dots, \zeta_m) \quad \text{and} \quad \phi = (\phi_1, \dots, \phi_n)$$

we have

$$\sigma_\zeta(D\Delta_{g_M, g_N}, x)\phi = \sum_{i, j} g_M^{ij} \zeta_i \zeta_j \phi|_x = |\zeta|_g^2 \phi|_x,$$

Consequently, the Laplacian operator  $\Delta_{g_M, g_N}$  is elliptic.

### 6.1.2 Time-Dependent Vector Bundles

Suppose that  $I \subset \mathbb{R}$  is an open interval and let  $\{g(t)\}_{t \in I}$  be a smooth family of Riemannian metrics on a manifold  $M$ . This means that for any  $(x, t) \in M \times I$  we have an inner product  $g_{(x, t)}$  on  $T_x M$ . We can regard  $\{g(t)\}_{t \in I}$  as a metric  $g$  acting on the *spatial tangent bundle*  $\mathcal{H}$ , defined by

$$\mathcal{H} = \{v \in T(M \times \mathbb{R}) : d\pi_2(v) = 0\},$$

where  $\pi_2 : M \times I \rightarrow I$  is given by  $\pi_2(x, t) = t$ . Note that each  $g(t)$  is a metric on  $\mathcal{H}$  since  $\mathcal{H}_{(x,t)}$  is isomorphic to  $T_x M$  via  $\pi_2$ . We can even extend  $g$  into a metric on  $M \times I$ , with respect to which we have the orthogonal decomposition

$$T(M \times I) = \mathcal{H} \oplus \mathbb{R}\partial_t.$$

Since  $\mathcal{H}$  is a subbundle of  $T(M \times I)$ , any section of  $\mathcal{H}$  is also a section of  $T(M \times I)$ . Sections of  $\Gamma(\mathcal{H})$  are called *spatial vector fields*. There is a natural connection  $\nabla$  on  $M \times I$ . Namely, define  $\nabla$  by

$$\nabla_v w = \nabla_v^{g(t)} w, \nabla_v \partial_t = 0, \nabla_{\partial_t} \partial_t = 0 \text{ and } \nabla_{\partial_t} v = [\partial_t, v], \tag{25}$$

for any  $v, w \in \Gamma(\mathcal{H})$ , where  $\nabla^{g(t)}$  stand for the Levi-Civita connection of  $g(t)$ . One can readily check that  $\nabla$  is compatible with  $g$ , i.e.,

$$v g(w_1, w_2) = g(\nabla_v w_1, w_2) + g(w_1, \nabla_v w_2),$$

for any  $v \in \mathfrak{X}(M \times \mathbb{R})$  and  $w_1, w_2 \in \Gamma(\mathcal{H})$ . Moreover, for any  $w_1, w_2 \in \Gamma(\mathcal{H})$ ,

$$\nabla_{w_1} w_2 - \nabla_{w_2} w_1 = [w_1, w_2].$$

The situation we discussed above occurs, when we have a family of immersions  $F : M \times I \rightarrow N$ . In this case,  $F^*h$  gives a family of metrics on  $M$ . Endowing  $M \times I$  with the connection  $\nabla$ , we have for any  $v \in \Gamma(\mathcal{H})$  that

$$\nabla_{\partial_t}^{F^*TN} dF(v) - \nabla_v^{F^*TN} dF(\partial_t) = dF([\partial_t, v]) = dF(\nabla_{\partial_t} v).$$

### 6.1.3 Parabolic Differential Equations

Let  $M$  be a manifold equipped with a family of metrics  $\{g(t)\}_{[0,T]}$ . Denote by  $\{\nabla^{g(t)}\}_{t \in [0,T]}$  the corresponding Levi-Civita connections. Let  $E_1$  and  $E_2$  be vector bundles over  $M$  and assume that  $E_1$  is equipped with a fixed time independent metric  $h$  and connections  $\{\nabla(t)\}_{t \in [0,T]}$  that are compatible with  $h$ , i.e.,

$$vh(\phi_1, \phi_2) = h(\nabla(t)_v \phi_1, \phi_2) + h(\phi_1, \nabla(t)_v \phi_2),$$

for any tangent vector  $v$ , sections  $\phi_1, \phi_2 \in \Gamma(E)$  and any time  $t \in [0, T)$ .

As in Definition 4, by coupling  $\nabla(t)$  with  $\nabla^{g(t)}$  we obtain repeated covariant derivatives  $\nabla^k(t)$  acting on sections of  $E_1$ . Suppose now that  $\{\phi(t)\}_{t \in [0,T]}$  is a smooth time-dependent family of sections of  $E_1$ , where smooth means that for any fixed  $(x, t) \in M \times [0, T)$ , the time-derivative

$$(\nabla_{\partial_t} \phi)(x, t) = \lim_{h \rightarrow 0} \frac{\phi(x, t+h) - \phi(x, t)}{h}$$

exists. Hence,  $\{\nabla_{\partial_t}\phi\}_{t \in [0, T]}$  is another one parameter family of sections on  $E_1$ . We are interested now in expressions of the form:

$$(\nabla_{\partial_t}\phi)(x, t) = (P\phi)(x, t) = \mathcal{Q}(x, t, \nabla^1(t)\phi(x, t), \dots, \nabla^k(t)\phi(x, t)), \tag{26}$$

where now  $P : \Gamma(E_1) \rightarrow \Gamma(E_2)$  is a time-dependent differentiable operator of order  $k$ . If for each fixed  $t$  the operator  $P$  is linear elliptic, we say that (26) is a *linear parabolic differential equation*. We say that (26) represents a *non-linear parabolic differential equation* if and only if, for any  $\phi \in \Gamma(E_1)$ , its linearization is parabolic.

**Theorem 17** *If the differential operator  $P$  is parabolic at  $\phi_0 \in \Gamma(E_1)$ , then there exist a  $T > 0$  and a smooth family  $\phi(t) \in \Gamma(E_1)$ , for  $t \in [0, T]$ , such that there exists a unique smooth solution for the initial value problem*

$$\begin{cases} \nabla_{\partial_t}\phi = P\phi, \\ \phi(0) = \phi_0. \end{cases}$$

for  $t \in [0, T]$ , where  $T$  depends on the initial data  $\phi_0$ .

We close this section with an application of this general theory.

**Definition 24** Let  $(M, g_M)$  and  $(N, g_N)$  be Riemannian manifolds. We say that a family of smooth maps  $F : M \times [0, T) \rightarrow N$  evolves by (*harmonic*) *heat flow*, with initial data  $F_0 : M \rightarrow N$ , if it satisfies the initial value problem

$$\begin{cases} \nabla_{\partial_t}dF = dF(\partial_t) = \Delta_{g_M, g_N} F, \\ F(\cdot, 0) = F_0. \end{cases} \tag{27}$$

**Theorem 18** *Let  $(M, g_M)$  be a compact Riemannian manifold and suppose that  $F_0 : (M, g_M) \rightarrow (N, g_N)$  is a smooth map into a Riemannian manifold  $(N, g_N)$ . Then, (27) admits a unique, smooth solution on a maximal time interval  $[0, T_{\max})$ , where  $0 < T_{\max} \leq \infty$ .*

**Proof** We already computed that for  $\zeta \in T^*M$ , we have

$$\sigma_{\zeta}(D\Delta_{g_M, g_N}, x) = |\zeta|_g^2 I.$$

Hence, the parabolic theory can be used to ensure short-time existence. □

## 6.2 Short-time Existence of the Mean Curvature Flow

A supposed solution  $F$  of MCF can be represented in local coordinates as

$$F(x_1, \dots, x_m, t) = (F^1(x_1, \dots, x_m, t), \dots, F^n(x_1, \dots, x_m, t)).$$

Then, from (2) we have

$$H = \sum_{i,j,\alpha} g^{ij} (F_{x_i x_j}^\alpha - \sum_k \Gamma_{ij}^k F_{x_k}^\alpha + \sum_{\gamma,\delta} \Gamma_{\gamma\delta}^\alpha F_{x_i}^\gamma F_{x_j}^\delta) \partial_{y_\alpha},$$

where

$$g_{ij} = \sum_{\alpha,\beta} h_{\alpha\beta} F_{x_i}^\alpha F_{x_j}^\beta \quad \text{and} \quad \Gamma_{ij}^k = \frac{1}{2} \sum_l g^{kl} (\partial_{x_i} g_{jl} + \partial_{x_j} g_{il} - \partial_{x_l} g_{ij}).$$

Note that  $g$  is the induced metric and it depends on  $F$ . Hence,

$$\partial_{x_i} g_{jl} = \sum_{\beta,\gamma} (h_{\beta\gamma} F_{x_i x_j}^\gamma F_{x_l}^\beta + h_{\beta\gamma} F_{x_j}^\gamma F_{x_i x_l}^\beta) + \text{lower order terms}$$

and consequently

$$\Gamma_{ij}^k = \sum_{l,\beta,\gamma} g^{kl} h_{\beta\gamma} F_{x_i}^\gamma F_{x_j x_l}^\beta + \text{lower order terms.} \tag{28}$$

Combining the formula (28) with equation (2), we obtain

$$H = \sum_{i,j,\alpha,\beta} g^{ij} (\delta_{\alpha\beta} - \sum_{k,l,\gamma} g^{kl} h_{\beta\gamma} F_{x_k}^\alpha F_{x_l}^\gamma) F_{x_i x_j}^\beta \partial_{y_\alpha} + \text{lower order terms.}$$

By a straightforward computation, we get

$$\begin{aligned} DH|_F(G) &= \lim_{s \rightarrow 0} \frac{H(F + sG) - H(F)}{s} \\ &= \sum_{i,j,\alpha,\beta} g^{ij} (\delta_{\alpha\beta} - \sum_{k,l,\gamma} g^{kl} h_{\beta\gamma} F_{x_k}^\alpha F_{x_l}^\gamma) G_{x_i x_j}^\beta \partial_{y_\alpha} + \text{lower order terms.} \end{aligned}$$

Denote by  $\pi_{TM}$  and  $\pi_{NM}$  the projections of  $F^*TN$  onto  $dF(T_xM)$  and  $NM$ , respectively. Then, for any  $\phi = \sum_\alpha \phi_\alpha \partial_{y_\alpha} \in \Gamma(F^*TN)$ , we have

$$\pi_{NM}(\phi) = \phi - \pi_{TM}(\phi) = \sum_{\alpha,\beta} (\delta_{\alpha\beta} - \sum_{k,l,\gamma} g^{kl} h_{\beta\gamma} F_{x_k}^\alpha F_{x_l}^\gamma) \phi_\beta \partial_{y_\alpha}.$$

Therefore, the principal symbol is given by

$$\begin{aligned} \sigma_\zeta(DH; x)\phi &= \sum_{i,j} g^{ij} \zeta_i \zeta_j \sum_{\alpha,\beta} (\delta_{\alpha\beta} - \sum_{k,l,\gamma} g^{kl} h_{\beta\gamma} F_{x_k}^\alpha F_{x_l}^\gamma) \phi_\beta \partial_{y_\alpha} \\ &= |\zeta|_g^2 \pi_{NM}(\phi|_x). \end{aligned}$$

Observe that the principal symbol is zero for tangent directions. Thus, MCF is degenerate and we cannot obtain information from the standard theory about short-time existence. Short-time existence and uniqueness of MCF was originally proven using results of Hamilton [43, 44] based on the Nash-Moser iteration method. We present a proof adapting a variant of the DeTurck’s trick which was first used in Ricci flow [31]; see also [7, 63, 92].

**Theorem 19 (Invariance Under Tangential Variations)** *Suppose that  $F: M \times [0, T] \rightarrow N$  is a family of immersions satisfying the system of PDEs*

$$\begin{cases} dF_{(x,t)}(\partial_t) = H(F(x, t)) + dF_{(x,t)}(V(x, t)), \\ F(x, 0) = F_0(x), \end{cases} \tag{29}$$

where  $(x, t) \in M \times [0, T]$ , the manifold  $M$  is compact and  $V$  is a time-dependent family of smooth vector fields. Then, there exists a unique family of diffeomorphisms  $\psi: M \times [0, T] \rightarrow M$ , such that the map  $\widehat{F}: M \times [0, T] \rightarrow N$  given by  $\widehat{F}(x, t) = F(\psi(x, t), t)$ , is a solution of

$$\begin{cases} d\widehat{F}_{(x,t)}(\partial_t) = H(\widehat{F}(x, t)), \\ \widehat{F}(x, 0) = F_0(\psi(x, 0)). \end{cases}$$

Conversely, if  $F: M \times [0, T] \rightarrow N$  is a solution of the mean curvature flow and  $\psi: M \times [0, T] \rightarrow M$  is a family of diffeomorphisms, then  $\widehat{F}: M \times [0, T] \rightarrow N$  satisfies a system of the form (29).

**Proof** Consider for the moment an arbitrary family a time-dependent of diffeomorphisms  $\psi: M \times [0, T] \rightarrow M$  and define  $\widehat{F}: M \times [0, T] \rightarrow N$  given by  $\widehat{F}(x, t) = F(\psi(x, t), t)$ , for  $(x, t) \in M \times [0, T]$ . From the chain rule, we have

$$d\widehat{F}_{(x,t)}(\partial_t) = H(\widehat{F}(x, t)) + dF_{(\psi(x,t),t)}(V(\psi(x, t), t) + d\psi_{(x,t)}(\partial_t)),$$

for any  $(x, t) \in M \times [0, T]$ . Hence, it suffices to find a one-parameter family of diffeomorphisms  $\psi: M \times [0, T] \rightarrow M$  solving the initial value problem

$$\begin{cases} d\psi_{(x,t)}(\partial_t) = -V(\psi(x, t), t), \\ \psi(x, 0) = I, \end{cases}$$

for any  $(x, t) \in M \times [0, T]$ , where  $I: M \rightarrow M$  is the identity map. By Picard-Lindelöf theorem there exists a unique smooth solution of the above initial value problem. Moreover, because the initial data is the identity, taking  $T > 0$  small enough we can assume that for any  $t \in [0, T]$  the map  $\psi(\cdot, t): M \rightarrow M$  is a diffeomorphism. The converse is straightforward.  $\square$

**Theorem 20 (Short-time Existence)** *Let  $M$  be a compact Riemannian manifold and  $F_0: M \rightarrow N$  an immersion into a Riemannian manifold  $N$ . Then, the mean curvature flow with initial data the immersion  $F_0$  admits a smooth solution on a maximal time interval  $[0, T_{\max})$ , where  $0 < T_{\max} \leq \infty$ .*

**Proof** The idea is to modify MCF by adding some tangential component in order to make it parabolic. Suppose that  $F : M \times [0, T_{\max}) \rightarrow N$  solves MCF. Fix a Riemannian metric  $\widehat{g}$  on  $M$ , denote its Levi-Civita connection by  $\widehat{\nabla}$  and consider the vector field  $V_{DT}$  on  $M$  given by

$$V_{DT} = \text{tr}_g(\nabla - \widehat{\nabla}). \tag{30}$$

Note that in local coordinates,  $V_{DT}$  has the form

$$V_{DT} = \sum_{i,j,k} g^{ij}(\Gamma_{ij}^k - \widehat{\Gamma}_{ij}^k)\partial_{x_k},$$

where  $\Gamma_{ij}^k$  and  $\widehat{\Gamma}_{ij}^k$  are the Christoffel symbols of the connections  $\nabla$  and  $\widehat{\nabla}$ , respectively. Consider now the initial value problem,

$$\begin{cases} dF(\partial_t) = H + dF(V_{DT}) \\ F(\cdot, 0) = F_0 \end{cases}, \tag{31}$$

The first equation of (31) in local coordinates takes the form

$$F_t = \sum_{i,j,\alpha} g^{ij}(F_{x_i x_j}^\alpha - \sum_k \widehat{\Gamma}_{ij}^k F_{x_k}^\alpha + \sum_{\gamma,\delta} \Gamma_{\gamma\delta}^\alpha F_{x_i}^\gamma F_{x_j}^\delta)\partial_{y_\alpha}.$$

Since  $\widehat{\Gamma}_{ij}^k$  does not depend on time, the principal symbol of (31) is

$$\sigma_\zeta(D(H + V_{DT}), \cdot) = |\zeta|^2 I.$$

Hence (31) is parabolic and has a unique solution. According to Theorem 19, from a solution of (31) we obtain a solution of the mean curvature flow.  $\square$

**Definition 25** Let  $F : M \times [0, T) \rightarrow N$  be a solution of MCF. Fix a metric  $\widehat{g}$  and consider the vector field  $V_{DT}$ . The modified flow (31) is called *DeTurck mean curvature flow*.

**Lemma 3** *The vector field  $V_{DT}$  defined in (30) is minus the Laplacian of the identity map  $I : (M, g) \rightarrow (M, \widehat{g})$ .*

**Proof** The Hessian  $B$  of the map  $I$  is given by

$$B(v_1, v_2) = \widehat{\nabla}_{dI(v_1)} dI(v_2) - dI(\nabla_{v_1} v_2) = \widehat{\nabla}_{v_1} v_2 - \nabla_{v_1} v_2,$$

for any  $v_1, v_2 \in \mathfrak{X}(M)$ . Hence,  $\Delta_{g,\widehat{g}} I = -V_{DT}$ . This completes the proof.  $\square$



**Theorem 21 (Uniqueness)** *Let  $M$  be a compact Riemannian manifold and  $F_0: M \rightarrow N$  an immersion into a Riemannian manifold  $N$ . Then, the solution of MCF, with initial data the immersion  $F_0: M \rightarrow N$ , is unique up to diffeomorphisms.*

**Proof** Suppose that  $\tilde{F}: M \times [0, T_{\max}) \rightarrow N$  is the maximal solution of MCF, with initial data the given immersion  $F_0$ , and denote the induced metrics by  $\tilde{g}$ . As in the existence part, fix a metric  $\hat{g}$  and denote by  $\tilde{\nabla}$  its associated Levi-Civita connection. Consider the initial value problem

$$\begin{cases} d\phi(\partial_t) = \Delta_{\tilde{g}, \hat{g}}\phi \\ \phi(\cdot, 0) = I \end{cases}.$$

Observe that the above problem is a parabolic and thus its solution gives rise to a unique one parameter family of diffeomorphisms  $\phi: M \times [0, \varepsilon) \rightarrow M$ , for at least some short time  $\varepsilon > 0$ . Denote by  $\psi: M \times [0, \varepsilon) \rightarrow M$  the one parameter family of diffeomorphisms with the property that, for each  $t$ , the map  $\psi(\cdot, t)$  is the inverse of  $\phi(\cdot, t)$ , i.e.,

$$\psi(\phi(x, t), t) = x = \phi(\psi(x, t), t)$$

for any  $(x, t)$  in space-time. From the chain rule, we have

$$d\psi_{(\phi(x,t),t)}(\partial_t) = -d\psi_{(\phi(x,t),t)}((\Delta_{\tilde{g}, \hat{g}}\phi)(x)). \tag{32}$$

Define the map  $F: M \times [0, \varepsilon) \rightarrow N$  given by  $F(x, t) = \tilde{F}(\psi(x, t), t)$ , for any  $(x, t) \in M \times [0, T_{\max})$ . The induced time-dependent metric on  $M$  is  $g = \psi^*\tilde{g}$ . Moreover, the map  $F$  satisfies the evolution equation

$$F_t = H + d\tilde{F}(W), \tag{33}$$

where for any point  $(x, t)$  in space-time, we have

$$W(\psi(x, t), t) = d\psi_{(x,t)}(\partial_t).$$

Taking into account (32) and the composition formula for the Laplacian (see for example [24, page 116, equation (2.56)]), we have

$$W(\psi(x, t), t) = d\psi_{(x,t)}(V_{DT}(x)), \tag{34}$$

for any  $(x, t) \in M \times [0, \varepsilon)$ . From (33) and (34), we see that  $F$  satisfies the DeTurck mean curvature flow

$$dF(\partial_t) = H + dF(V_{DT}),$$

with initial data the immersion  $F_0: M \rightarrow N$ .

Suppose now that  $\tilde{F}_1, \tilde{F}_2: M \times [0, T_{\max}) \rightarrow N$  are two solutions of the mean curvature flow, with the same initial condition  $F_0: M \rightarrow N$ . As before fix a metric  $\hat{g}$  on  $M$  and denote by  $\tilde{g}_1$  and  $\tilde{g}_2$  the induced time-dependent metrics on  $M$  by  $\tilde{F}_1$  and  $\tilde{F}_2$ , respectively. Denote by

$$\phi^1: M \times [0, \varepsilon) \rightarrow N \quad \text{and} \quad \phi^2: M \times [0, \varepsilon) \rightarrow N$$

the one-parameter family of diffeomorphisms solving the initial value problem

$$\begin{cases} d\eta(\partial_t) = \Delta_{\tilde{g}_i, \hat{g}}\eta, \\ \eta(\cdot, 0) = I. \end{cases}$$

Then, as we verified above, the maps

$$F_i: M \times [0, \varepsilon) \rightarrow N, \quad i \in \{1, 2\},$$

satisfy

$$\tilde{F}_i(x, t) = F_i(\phi^i(x, t), t),$$

for any  $(x, t) \in M \times [0, \varepsilon)$ , form solutions of the DeTurck mean curvature flow, with common initial data the immersion  $F_0: M \rightarrow N$ . Since the DeTurck mean curvature flow is parabolic, it follows that its solution is unique.  $\square$

### 6.3 Parabolic Maximum Principles

In this subsection, we state the weak and strong version of the parabolic maximum principle for scalar functions obeying a diffusion-reaction equation on a manifold equipped with a smooth time-dependent family of Riemannian metrics. Then we also present Hamilton’s version [44, 45] of the parabolic maximum principle for arbitrary sections of a vector bundle; for detailed proofs see also the excellent monograph [26].

#### 6.3.1 Scalar Parabolic Maximum Principle

Suppose that  $M$  is a smooth manifold, possibly with boundary  $\partial M$ , and  $\{g(t)\}_{t \in [0, T)}$  a smooth family of Riemannian metrics. We will consider the second order time-dependent operator  $\mathcal{L}$  given by

$$\mathcal{L}u = \Delta_{g(t)}u + g(t)(X, \nabla^{g(t)}u) \tag{P}$$

where

$$u \in C^2(M \times (0, T)) \cup C^0(\overline{M} \times [0, T]).$$

Note that, for each fixed time,  $\mathcal{L}$  is an elliptic operator.

**Theorem 22 (Comparison Principle)** *Suppose that  $M$  is a compact, without boundary, manifold equipped with a smooth family  $\{g(t)\}_{t \in [0, T]}$  of Riemannian metrics and  $u: M \times [0, T) \rightarrow \mathbb{R}$  a  $C^2$ -smooth function, which satisfies the differential inequality*

$$\partial_t u - \mathcal{L}u \leq \Psi(u, t),$$

where  $\mathcal{L}$  is the (time-dependent) operator defined in (P) and  $\Psi: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  a smooth map. Let  $\varphi$  be the solution of the associated ODE

$$\begin{cases} \varphi'(t) = \Psi(\varphi(t), t), \\ \varphi(0) = \max_{x \in M} u(x, 0). \end{cases}$$

Then, the solution  $u$  of the differential inequality is bounded from above by the solution  $\varphi$  of the ODE, that is  $u(x, t) \leq \varphi(t)$ , for every  $(x, t) \in M \times [0, T)$ .

As in the elliptic case, there exists a criterion which forces a solution of a parabolic differential inequality to be constant.

**Theorem 23 (Strong Maximum Principle)** *Suppose that  $M$  is a smooth manifold, possibly with boundary, equipped with a smooth family  $\{g(t)\}_{t \in [0, T]}$  of Riemannian metrics. Let  $u \in C^2(M \times (0, T)) \cup C^0(\overline{M} \times [0, T])$  be a solution of*

$$\partial_t u - \mathcal{L}u + cu \leq 0$$

where  $c$  is a non-negative constant.

- (a) *If  $c = 0$  and  $u$  attains a maximum at point  $(x_0, t_0) \in M \times (0, T)$  then  $u$  is constant on  $M \times [0, t_0]$ .*
- (b) *If  $c < 0$  and the function  $u$  attains a non-negative maximum at a point  $(x_0, t_0) \in M \times (0, T)$ , then  $u$  is constant on  $M \times [0, t_0]$ .*

By reversing both inequalities we obtain the corresponding minimum version of the comparison and strong principle; for the proofs see [5, 35] or [72].

### 6.3.2 Vectorial Parabolic Maximum Principle

Let  $M$  be a smooth manifold, possibly with boundary  $\partial M$ , equipped with a smooth family of metrics  $\{g(t)\}_{t \in [0, T]}$  and associated Levi-Civita connections  $\nabla^{g(t)}$ . Let  $E$  be a vector bundle over  $M$  equipped with a time-independent metric  $h$  and a

family  $\{\nabla(t)\}_{t \in [0, T]}$  of connections that are compatible with  $h$ . The *time-dependent Laplacian* acting on smooth sections of  $E$  is defined by

$$\Delta(t)\phi = \sum_{i=1}^m (\nabla(t)_{v_i} \nabla(t)_{v_i} \phi - \nabla(t)_{\nabla_{v_i}^{g(t)} v_i} \phi)$$

where  $\{v_1, \dots, v_m\}$  is an orthonormal basis of  $g(t)$ .

Following the same lines as in the elliptic case, we can derive Weinberger-Hamilton's versions of the parabolic maximum principle.

**Theorem 24 (Weak Vectorial Maximum Principle)** *Suppose that  $M$  is a compact manifold, possibly with boundary  $\partial K$ , equipped with a smooth family  $\{g(t)\}_{t \in [0, T]}$  of Riemannian metrics. Let  $E$  be a vector bundle over  $M$  endowed with time independent bundle metric  $h$  and a family  $\{\nabla(t)\}_{t \in [0, T]}$  of connections that are compatible with  $h$ . Let  $K$  be a closed fiber-convex subset of  $E$  that is invariant under parallel transport with respect to each connection  $\nabla(t)$ ,  $t \in [0, T]$ , and let  $\{\phi(t)\}_{t \in [0, T]}$  be a smooth family of sections such that*

$$\nabla_{\partial_t} \phi - \Delta(t)\phi = \nabla(t)_X \phi + \Psi(\phi)$$

where  $X$  is a smooth time dependent vector field and  $\Psi$  is a smooth fiberwise map that points into  $K$ . If  $\phi_{(x,t)} \in K$  for any  $(x, t)$  in the parabolic boundary of  $M \times [0, T]$ , i.e., for any  $(x, t) \in (M \times \{0\}) \cup (\partial M \times [0, T])$ , then  $\phi_{(x,t)} \in K$  for any  $(x, t) \in M \times [0, T]$ .

**Theorem 25 (Strong Vectorial Maximum Principle)** *Suppose that  $M$  is a smooth, not necessarily compact, manifold equipped with a smooth family  $\{g(t)\}_{t \in [0, T]}$  of Riemannian metrics. Moreover, let  $E$  be a vector bundle over  $M$  endowed with time independent metric  $h$  and a family  $\{\nabla(t)\}_{t \in [0, T]}$  of connections that are compatible with  $h$ . Assume that  $K$  is a closed fiber-convex subset of the vector bundle  $E$  that is invariant under parallel transport with respect to each connection  $\nabla(t)$ ,  $t \in [0, T]$ , and let  $\{\phi(t)\}_{t \in [0, T]}$  be a smooth family of sections such that*

$$\nabla_{\partial_t} \phi - \Delta(t)\phi = \nabla(t)_X \phi + \Psi(\phi)$$

where  $X$  is a smooth time dependent vector field and  $\Psi$  is a smooth fiberwise map that points into  $K$ . If there exists a point  $(x_0, t_0) \in M \times (0, T)$  such that  $\phi_{(x_0, t_0)} \in \partial K$ , then  $\phi_{(x,t)} \in \partial K$  for any  $(x, t) \in M \times [0, t_0]$ .

Let us describe now the parabolic maximum principle in the special case where as vector bundle we consider the space of symmetric 2-tensors.

**Theorem 26** *Let  $M$  be a compact manifold equipped with a smooth family  $\{g(t)\}_{t \in [0, T]}$  of Riemannian metrics. Suppose that  $\{\phi(t)\}_{t \in [0, T]}$  is smooth family of symmetric 2-tensors on  $M$  such that*

$$\nabla_{\partial_t} \phi - \Delta(t)\phi = \nabla(t)_X \phi + \Psi(\phi)$$

where  $\Psi : \text{Sym}(T^*M \otimes T^*M) \rightarrow \text{Sym}(T^*M \otimes T^*M)$  is a smooth fiberwise map satisfying the null-eigenvector condition and  $X$  is a smooth time dependent vector field. If  $\phi(0) \geq 0$ , then  $\phi(t) \geq 0$  for all  $t \in [0, T)$ . Additionally, if there is a point  $(x_0, t_0) \in M \times (0, T)$  where  $\phi(t_0)$  has a zero eigenvalue then  $\phi(t)$  has a zero eigenvalue for any  $t \in (0, t_0)$ .

### 6.4 Evolution Equations

We will compute the evolution of some important quantities. In order to simplify the notation, we omit upper or lower indices on connections and Laplacians which identify the corresponding bundles where they are defined. Most of these computations can be found in [4, 75–78, 86, 94, 96].

**Lemma 4** *Suppose that  $F : M \times [0, T) \rightarrow N$  is a solution of the mean curvature flow. Then, the following facts are true:*

- (a) *The induced metrics  $g$  evolve in time under the equation*

$$(\nabla_{\partial_t} g)(v_1, v_2) = -2\langle H, A(v_1, v_2) \rangle = -2A^H(v_1, v_2),$$

*for any  $v_1, v_2 \in \mathfrak{X}(M)$ .*

- (b) *The induced volume form  $\Omega$  on  $(M, g)$  evolves according to the equation*

$$\nabla_{\partial_t} \Omega = -|H|^2 \Omega.$$

*Moreover, the volume of the evolved submanifolds satisfy*

$$\partial_t \text{Vol} = - \int_M |H|^2 \Omega.$$

- (c) *There exists a local smooth time-dependent tangent orthonormal frame field and a local smooth time-dependent orthonormal frame field along the normal bundle of the evolving submanifolds.*

**Proof**

- (a) Let  $v_1, \dots, v_m$  be time-independent tangent vector fields. Keeping in mind the notation introduced in Section 6.1.2, we have

$$\nabla_{\partial_t} dF(v_i) = \nabla_{v_i} dF(\partial_t) + dF([\partial_t, v_i]) = \nabla_{v_i} H,$$

for any  $i \in \{1, \dots, m\}$ . Therefore, for any  $i, j \in \{1, \dots, m\}$ , we deduce that

$$\begin{aligned}
 (\nabla_{\partial_t} g)(v_i, v_j) &= \partial_t(g(v_i, v_j)) - g(\nabla_{\partial_t} v_i, v_j) - g(v_i, \nabla_{\partial_t} v_j) \\
 &= \partial_t \langle dF(v_i), dF(v_j) \rangle = \langle \nabla_{v_i} H, dF(v_j) \rangle + \langle \nabla_{v_j} H, dF(v_i) \rangle \\
 &= -\langle H, \nabla_{v_i} dF(v_j) \rangle - \langle H, \nabla_{v_j} dF(v_i) \rangle \\
 &= -2\langle H, A(v_i, v_j) \rangle.
 \end{aligned}$$

(b) We compute

$$\begin{aligned}
 \partial_t \sqrt{\det g_{ij}} &= \sum_{k,l} \frac{(g^{kl} \partial_t g_{kl}) \det g_{ij}}{2\sqrt{\det g_{ij}}} = - \sum_{k,l} \langle H, g^{kl} A_{kl} \rangle \sqrt{\det g_{ij}} \\
 &= -|H|^2 \sqrt{\det g_{ij}}.
 \end{aligned}$$

(c) The associated adjoint operator  $P : (TM, g) \rightarrow (TM, g)$  of  $A^H$  satisfies

$$A^H(v_1, v_2) = g(Pv_1, v_2) = g(v_1, Pv_2), \tag{35}$$

for any  $v_1, v_2 \in \mathfrak{X}(M)$ . Consider now the family of bundle isomorphism  $U(t) : (TM, g(0)) \rightarrow (TM, g(t))$ , given as the solution of the initial value problem

$$\begin{cases} \nabla_{\partial_t} U(t) = P \circ U(t), \\ U(0) = I. \end{cases} \tag{36}$$

By a straightforward computation, we can show that  $U^*(t)g(t) = g(0)$ . Hence, if  $\{e_1(0), \dots, e_m(0)\}$  is a local orthonormal frame with respect to  $g(0)$ , then  $\{e_1(t) = U(t)e_1(0), \dots, e_m(t) = U(t)e_m(0)\}$  is a local orthonormal frame of  $g(t)$ . By taking the complement of  $\{e_1, \dots, e_m\}$ , we get a time-dependent frame field on the normal bundles of the evolving submanifolds.

□

**Lemma 5** *The time-derivative of the second fundamental form is given by*

$$(\nabla_{\partial_t}^\perp A)_{ij}^\alpha = (\nabla^{\perp 2} H)_{ij}^\alpha - \sum_{k,\beta} H^\beta A_{jk}^\beta A_{ik}^\alpha - \sum_\beta H^\beta \tilde{R}_{\beta ij\alpha},$$

where the indices are with respect to a local orthonormal frame.

**Proof** Suppose that  $\{e_1, \dots, e_m; \xi_{m+1}, \dots, \xi_n\}$  is a local adapted orthonormal frame field around a fixed point  $(x_0, t_0)$ . Recall that

$$\nabla_{\partial_t} \partial_t = 0, \quad \nabla_{e_i} \partial_t = 0 \quad \text{and} \quad [\partial_t, e_i] = \nabla_{\partial_t} e_i = \sum_{j,\beta} H^\beta A_{ij}^\beta e_j. \tag{37}$$

In order to simplify the computations, we may assume that  $\{e_1, \dots, e_m\}$  is normal frame at  $(x_0, t_0)$ . Under these considerations, we have that at  $(x_0, t_0)$

$$\begin{aligned}
(\nabla_{\partial_t} A)_{ij} &= \nabla_{\partial_t} \nabla_{e_i} dF(e_j) - \nabla_{\partial_t} dF(\nabla_{e_i} e_j) - A(\nabla_{\partial_t} e_i, e_j) - A(e_i, \nabla_{\partial_t} e_j) \\
&= \nabla_{e_i} \nabla_{\partial_t} dF(e_j) + \tilde{R}(H, dF(e_i), dF(e_j)) + \nabla_{\nabla_{\partial_t} e_i} dF(e_j) \\
&\quad - dF(\nabla_{\partial_t} \nabla_{e_i} e_j) - A(\nabla_{\partial_t} e_i, e_j) - A(e_i, \nabla_{\partial_t} e_j).
\end{aligned}$$

Hence,

$$\begin{aligned}
(\nabla_{\partial_t} A)_{ij} &= \nabla_{e_i} (\nabla_{e_j} H + dF(\nabla_{\partial_t} e_j)) + \tilde{R}(H, dF(e_i), dF(e_j)) \\
&\quad + \nabla_{\nabla_{\partial_t} e_i} dF(e_j) - dF(\nabla_{\partial_t} \nabla_{e_i} e_j) - A(\nabla_{\partial_t} e_i, e_j) - A(e_i, \nabla_{\partial_t} e_j) \\
&= \nabla_{e_i, e_j}^2 H + \tilde{R}(H, dF(e_i), dF(e_j)) + \nabla_{e_i} dF(\nabla_{\partial_t} e_j) \\
&\quad + \nabla_{\nabla_{\partial_t} e_i} dF(e_j) - dF(\nabla_{\partial_t} \nabla_{e_i} e_j) - A(\nabla_{\partial_t} e_i, e_j) - A(e_i, \nabla_{\partial_t} e_j) \\
&= \nabla_{e_i, e_j}^2 H + \tilde{R}(H, dF(e_i), dF(e_j)) + \nabla_{e_i} dF(\nabla_{\partial_t} e_j) \\
&\quad + \nabla_{\nabla_{\partial_t} e_i} dF(e_j) - dF(\nabla_{\partial_t} \nabla_{e_i} e_j) - A(\nabla_{\partial_t} e_i, e_j) - A(e_i, \nabla_{\partial_t} e_j)
\end{aligned}$$

and so

$$(\nabla_{\partial_t} A)_{ij} = \nabla_{e_i, e_j}^2 H + \tilde{R}(H, dF(e_i), dF(e_j)) - dF(R^\nabla(\partial_t, e_i, e_j))$$

where  $R^\nabla$  is the curvature operator of  $\nabla$  on  $T(M \times (0, T))$ . Consequently, at  $(x_0, t_0)$  we have

$$\begin{aligned}
(\nabla_{\partial_t}^\perp A)_{ij} &= \sum_\alpha \langle (\nabla_{\partial_t}^\perp A)_{ij}, \xi_\alpha \rangle \xi_\alpha = \sum_\alpha \langle (\nabla_{\partial_t} A)_{ij}, \xi_\alpha \rangle \xi_\alpha \\
&= \sum_\alpha \langle \nabla_{e_i} \nabla_{e_j} H, \xi_\alpha \rangle \xi_\alpha + \sum_{\alpha, \beta} H^\beta \tilde{R}_{\beta i j \alpha} \xi_\alpha.
\end{aligned}$$

On the other hand,

$$\begin{aligned}
\langle \nabla_{e_i} \nabla_{e_j} H, \xi_\alpha \rangle &= \langle \nabla_{e_i}^\perp (\nabla_{e_j}^\perp H + \sum_k \langle \nabla_{e_j} H, dF(e_k) \rangle dF(e_k)), \xi_\alpha \rangle \\
&= (\nabla^{2\perp} H)_{ij}^\alpha - \sum_{k, \beta} H^\beta A_{jk}^\beta A_{ik}^\alpha.
\end{aligned}$$

Combining the last two equalities we obtain the result.  $\square$

**Lemma 6** *The mean curvature  $H$  evolves in time under the equation*

$$(\nabla_{\partial_t}^\perp H)^\alpha = (\Delta^\perp H)^\alpha - \sum_{i, \beta} H^\beta \tilde{R}_{\beta i i \alpha} + \sum_{i, j, \beta} H^\beta A_{ij}^\beta A_{ij}^\alpha.$$

Moreover,

$$\partial_t |H|^2 = \Delta |H|^2 - 2|\nabla^\perp H|^2 + 2|A^H|^2 - 2 \sum_{i,\alpha,\beta} H^\alpha H^\beta \tilde{R}_{\alpha i i \beta},$$

where the indices are with respect to a local orthonormal frame.

**Proof** Let  $(x_0, t_0) \in M \times (0, T)$  and  $\{e_1, \dots, e_m; \xi_{m+1}, \dots, \xi_n\}$  be a local orthonormal frame field around of  $(x_0, t_0)$ . From (37) and Lemma 5, we have

$$\begin{aligned} (\nabla_{\partial_t}^\perp H)^\alpha &= \sum_i (\nabla_{\partial_t}^\perp A_{ii})^\alpha = \sum_i (\nabla_{\partial_t}^\perp A)_{ii}^\alpha + 2 \sum_i A^\alpha (\nabla_{\partial_t} e_i, e_i) \\ &= (\Delta^\perp H)^\alpha + \sum_{i,\beta} H^\beta \tilde{R}_{\beta i i \alpha} - \sum_{i,j,\beta} H^\beta A_{ij}^\beta A_{ij}^\alpha + 2 \sum_{i,j,\beta} H^\beta A_{ij}^\beta A_{ij}^\alpha, \end{aligned}$$

from where we deduce the evolution equation for  $H$ . Moreover

$$\begin{aligned} \partial_t |H|^2 &= \partial_t \langle H, H \rangle = 2 \langle \nabla_{\partial_t}^\perp H, H \rangle = \sum_\alpha (\nabla_{\partial_t}^\perp H)^\alpha H^\alpha \\ &= 2 \sum_\alpha (\Delta H)^\alpha H^\alpha - 2 \sum_{i,\alpha,\beta} H^\alpha H^\beta \tilde{R}_{\alpha i i \beta} + 2 \sum_{i,j,\alpha,\beta} H^\alpha H^\beta A_{ij}^\alpha A_{ij}^\beta. \end{aligned}$$

On the other hand

$$\sum_\alpha \Delta (H^\alpha)^2 = 2 \sum_\alpha (\Delta H)^\alpha H^\alpha + 2 \sum_\alpha |\nabla H^\alpha|^2.$$

Combining the last two identities we obtain the desired identity. □

### 6.5 Evolution Equations of Parallel Forms

Let  $F : M \times [0, T) \rightarrow N$  be a solution of the mean curvature flow and suppose that  $\Phi$  is a parallel  $k$ -tensor on  $N$ . Then, the pullback via  $F$  of  $\Phi$  gives rise to a time-dependent  $k$ -form on  $M$ . For example, the volume form of  $N$  is such a tensor. As we will see in the next section, interesting situations occurs when  $N$  is a Riemannian product  $N_1 \times N_2$  and we consider the volume forms  $\Omega_1$  and  $\Omega_2$  of  $N_1$  and  $N_2$ , respectively.

In the next lemmata, we will compute how these pullback tensors evolve under the mean curvature flow.

**Lemma 7** *The covariant derivative of the tensor  $F^*\Phi$  is given by*

$$(\nabla_{e_s} F^*\Phi)_{i_1 \dots i_k} = \sum_\alpha (A_{s i_1}^\alpha \Phi_{\alpha i_2 \dots i_k} + \dots + A_{s i_m}^\alpha \Phi_{i_1 \dots i_{m-1} \alpha}),$$

for any adapted orthonormal frame field  $\{e_1, \dots, e_m; \xi_{m+1}, \dots, \xi_{n-m}\}$ .



**Proof** As usually let us suppose that  $\{e_1, \dots, e_m\}$  is a normal frame at a fixed point  $(x_0, t_0)$  in space-time. By a direct computation, we get that at  $(x_0, t_0)$  we have

$$\begin{aligned} (\nabla_{e_s} F^* \Phi)_{i_1 \dots i_k} &= e_s \Phi(dF(e_{i_1}), \dots, dF(e_{i_m})) \\ &= \Phi(\nabla_{e_s} dF(e_{i_1}), \dots, dF(e_{i_m})) + \dots + \Phi(dF(e_{i_1}), \dots, \nabla_{e_s} dF(e_{i_m})) \\ &= \Phi(A(e_s, e_{i_1}), \dots, dF(e_{i_m})) + \dots + \Phi(dF(e_{i_1}), \dots, A(e_s, e_{i_m})) \\ &= \sum_{\alpha} (A_{s i_1}^{\alpha} \Phi_{\alpha i_2 \dots i_k} + \dots + A_{s i_m}^{\alpha} \Phi_{i_1 \dots i_{m-1} \alpha}). \end{aligned}$$

This completes the proof. □

By a direct computation we can derive the expression for the Laplacian of the pullback of a parallel  $k$ -tensor on  $N$ .

**Lemma 8** *The Laplacian of the  $k$ -tensor  $F^* \Phi$  is given by*

$$\begin{aligned} (\Delta F^* \Phi)_{i_1 \dots i_m} &= \sum_{\alpha} (\nabla_{e_1}^{\perp} H)^{\alpha} \Phi_{\alpha i_2 \dots i_m} + \dots + \sum_{\alpha} (\nabla_{e_m}^{\perp} H)^{\alpha} \Phi_{i_1 \dots i_{m-1} \alpha} \\ &\quad + 2 \sum_{k, \alpha, \beta} A_{k i_1}^{\alpha} A_{k i_2}^{\beta} \Phi_{\alpha \beta i_2 \dots i_m} + \dots + 2 \sum_{k, \alpha, \beta} A_{k i_{m-1}}^{\alpha} A_{k i_m}^{\beta} \Phi_{i_1 \dots \alpha \beta} \\ &\quad - \sum_{k, l, \alpha} (A_{k i_1}^{\alpha} A_{k l}^{\alpha} \Phi_{l i_2 \dots i_m} + \dots + A_{k i_m}^{\alpha} A_{k l}^{\alpha} \Phi_{i_1 \dots i_{m-1} l}) \\ &\quad - \sum_{k, \alpha} (\tilde{R}_{k \alpha k i_1} \Phi_{\alpha i_2 \dots i_m} + \dots + \tilde{R}_{k \alpha k i_m} \Phi_{i_1 \dots i_{m-1} \alpha}), \end{aligned}$$

for any adapted orthonormal frame field  $\{e_1, \dots, e_m; \xi_{m+1}, \dots, \xi_{n-m}\}$ .

**Proof** Let  $\{e_1, \dots, e_m; \xi_{m+1}, \dots, \xi_{n-m}\}$  be an adapted normal frame at the point  $(x_0, t_0)$  in space-time. We compute,

$$\begin{aligned} (\nabla_{e_k} \nabla_{e_k} F^* \Phi)_{i_1 \dots i_m} &= e_k (\Phi(A_{k i_1}, \dots, dF(e_{i_m})) + \dots + \Phi(dF(e_{i_1}), \dots, A_{k i_m})) \\ &= \Phi((\nabla_{e_k} A)_{k i_1}, \dots, dF(e_{i_m})) + \dots + \Phi(dF(e_{i_1}), \dots, (\nabla_{e_k} A)_{k i_m}) \\ &\quad + 2\Phi(A_{k i_1}, A_{k i_2}, \dots, dF(e_{i_m})) + \dots + 2\Phi(dF(e_{i_1}), \dots, A_{k i_{m-1}}, A_{k i_m}) \\ &= \Phi((\nabla_{e_k}^{\perp} A)_{k i_1}, \dots, dF(e_{i_m})) + \dots + \Phi(dF(e_{i_1}), \dots, (\nabla_{e_k}^{\perp} A)_{k i_m}) \\ &\quad + 2\Phi(A_{k i_1}, A_{k i_2}, \dots, dF(e_{i_m})) + \dots + 2\Phi(dF(e_{i_1}), \dots, A_{k i_{m-1}}, A_{k i_m}) \\ &\quad - \sum_l \langle A_{k i_1}, A_{k l} \rangle F^* \Phi(e_l, \dots, e_{i_m}) - \dots - \sum_l \langle A_{k i_m}, A_{k l} \rangle F^* \Phi(e_{i_1}, \dots, e_l). \end{aligned}$$

Summing over  $k$  and using the Codazzi equation (4), we get the result. □

**Lemma 9** *Suppose that  $F : M \times [0, T) \rightarrow N$  is a solution of the mean curvature flow and let  $\Phi$  be a parallel  $m$ -form on  $N$ . Then,  $u = *(F^*\Phi)$ , where  $*$  is the Hodge star operator with respect to the induced Riemannian metric  $g$ , evolves in time under the equation*

$$\begin{aligned} \partial_t u - \Delta u &= -2 \sum_{k,\alpha,\beta} A_{k1}^\alpha A_{k2}^\beta \Phi_{\alpha\beta 2\dots m} - \dots - 2 \sum_{k,\alpha,\beta} A_{km-1}^\alpha A_{km}^\beta \Phi_{1\dots\alpha\beta} \\ &+ \sum_{k,l,\alpha} (A_{k1}^\alpha A_{kl}^\alpha \Phi_{l2\dots m} + \dots + A_{km}^\alpha A_{kl}^\alpha \Phi_{1\dots m-l}) \\ &+ \sum_{k,\alpha} (\tilde{R}_{k\alpha k1} \Phi_{\alpha 2\dots m} + \dots + \tilde{R}_{k\alpha km} \Phi_{1\dots m-l\alpha}), \end{aligned}$$

for any adapted orthonormal frame field  $\{e_1, \dots, e_m; \xi_{m+1}, \dots, \xi_{n-m}\}$ .

**Proof** Let us make our computations again, with respect to a time-dependent orthonormal frame field as in Lemma 4. We compute,

$$\begin{aligned} \partial_t u &= \partial_t ((F^*\Phi)(e_1, \dots, e_m)) \\ &= \Phi(\nabla_{\partial_t} dF(e_1), \dots, dF(e_m)) + \dots + \Phi(dF(e_1), \dots, \nabla_{\partial_t} dF(e_m)). \end{aligned}$$

Taking into account the formulas (37), we have

$$\begin{aligned} \nabla_{\partial_t} dF(e_i) &= \nabla_{e_i} dF(\partial_t) + dF(\nabla_{\partial_t} e_i) = \nabla_{e_i} H + \sum_{k,\beta} H^\beta A_{ik}^\beta dF(e_k) \\ &= \nabla_{e_i}^\perp H, \end{aligned}$$

for any  $i \in \{1, \dots, m\}$ . Hence, putting everything together, we deduce that

$$\partial_t u = \Phi(\nabla_{e_1}^\perp H, \dots, dF(e_m)) + \dots + \Phi(dF(e_1), \dots, \nabla_{e_m}^\perp H).$$

Combining with Lemma 8 we obtain the result. □

## 7 Formation of Singularities Under Mean Curvature Flow

In this section, we present how one can build smooth singularity models for the mean curvature flow by rescaling properly around points, where the second fundamental form attains its maximum. The proof relies heavily on a compactness theorem of Cheeger-Gromov-Taylor [14] for pointed Riemannian manifolds and on the standard compactness theorem for immersions; see for example [27].

### 7.1 Characterization of the Maximal Time of Existence

In the following theorem, we give a characterization of the maximal time of solutions of the mean curvature flow. Its proof has been done by Huisken in [52, 53] and is based on the parabolic maximum principle. The key observation is that all higher derivatives  $\nabla^k A$ ,  $k \in \mathbb{N}$ , of the second fundamental tensor are uniformly bounded, once  $A$  is uniformly bounded. More precisely, the following result holds:

**Theorem 27** *Let  $M$  be a compact manifold and let  $F_0: M \rightarrow N$  a smooth immersion into a complete Riemannian manifold  $N$ . Then, the maximal time  $T_{\max}$  of the solution of the mean curvature flow, with initial data  $F_0$ , is finite if and only if*

$$\limsup_{t \rightarrow T_{\max}} (\max_{M \times [0,t]} |A|) = \infty.$$

An immediate consequence of the above result is the following theorem.

**Theorem 28** *Let  $M$  be a compact manifold and  $F : M \rightarrow [0, T_{\max}) \rightarrow N$  a solution of the mean curvature flow on a maximal time interval in a complete Riemannian manifold  $N$ . If the norm  $|A|$  of the second fundamental form is uniformly bounded, then the maximal time of solution of the flow is infinite.*

*Remark 7* When the target space  $N$  is compact and the maximal time of solution of the flow is infinite, due to a deep result of Simon [82], it follows that the flow converges smoothly and uniformly to a minimal submanifold. However, long-time existence does not automatically imply convergence. For instance, start with a latitude circle  $\mathbb{S}^1$  on a complete surface of revolution that does not admit closed embedded curves as geodesics. Then the flow with initial that particular circle will run forever, but it will not converge.

*Remark 8* Due to a recent result of Cooper [27], it is not necessary to have boundedness on the full norm of the second fundamental form in order to get long-time existence of the flow. In the matter of fact, he showed that uniform boundedness of the second fundamental form only in the direction of the mean curvature also leads to long-time existence.

### 7.2 Cheeger-Gromov Compactness for Metrics

Let us recall here the basic notions and definitions. For more details, see [5, 25] and [66]. We closely follow the exposition in [78].

**Definition 26** Let  $(E, \pi, \Sigma)$  be a vector bundle endowed with a Riemannian metric  $g$  and a metric connection  $\nabla$  and suppose that  $\{\xi_k\}_{k \in \mathbb{N}}$  is a sequence of sections of  $E$ . Let  $U$  be an open subset of  $\Sigma$  with compact closure  $\bar{U}$  in  $\Sigma$ . Fix a natural number  $p \geq 0$ . We say that  $\{\xi_k\}_{k \in \mathbb{N}}$  converges  $C^p$ -smoothly to  $\xi_\infty \in \Gamma(E|_{\bar{U}})$ , if for every  $\varepsilon > 0$ , there exists  $k_0 = k_0(\varepsilon)$ , such that

$$\sup_{0 \leq \alpha \leq p} \sup_{x \in \bar{U}} |\nabla^\alpha (\xi_k - \xi_\infty)| < \varepsilon$$

where  $k \geq k_0$ . We say that  $\{\xi_k\}_{k \in \mathbb{N}}$   $C^\infty$ -smoothly converges to  $\xi_\infty \in \Gamma(E|_{\bar{U}})$  if  $\{\xi_k\}_{k \in \mathbb{N}}$  converges in  $C^p$  to  $\xi_\infty \in \Gamma(E|_{\bar{U}})$ , for any  $p \geq 0$ .

**Definition 27** Let  $(E, \pi, \Sigma)$  be a vector bundle endowed with a Riemannian metric  $g$  and a metric connection  $\nabla$ . Let  $\{U_n\}_{n \in \mathbb{N}}$  be an exhaustion of  $\Sigma$  and  $\{\xi_k\}_{k \in \mathbb{N}}$  be a sequence of sections of  $E$  defined on open sets  $A_k$  of  $\Sigma$ . We say that  $\{\xi_k\}_{k \in \mathbb{N}}$  converges smoothly on compact sets to  $\xi_\infty \in \Gamma(E)$  if:

- (a) For every  $n \in \mathbb{N}$  there exists  $k_0$  such that  $\bar{U}_n \subset A_k$ , for all natural numbers  $k \geq k_0$ .
- (b) The sequence  $\{\xi|_{\bar{U}_k}\}_{k \geq k_0}$  converges in  $C^\infty$  to the restriction of the section  $\xi_\infty$  on  $\bar{U}_n$ .

In the next definitions, we recall the notion of the smooth Cheeger-Gromov convergence of sequences of Riemannian manifolds.

**Definition 28** A pointed Riemannian manifold  $(\Sigma, g, x)$  is a Riemannian manifold  $(\Sigma, g)$  with a choice of origin or base point  $x \in \Sigma$ . If the metric  $g$  is complete, we say that  $(\Sigma, g, x)$  is a complete pointed Riemannian manifold.

**Definition 29** We will say that a sequence  $\{(\Sigma_k, g_k, x_k)\}_{k \in \mathbb{N}}$  of complete, pointed Riemannian manifolds smoothly converges in the sense of Cheeger-Gromov to a complete pointed Riemannian manifold  $(\Sigma_\infty, g_\infty, x_\infty)$ , if there exists:

- (a) An exhaustion  $\{U_k\}_{k \in \mathbb{N}}$  of  $\Sigma_\infty$  with  $x_\infty \in U_k$ , for all  $k \in \mathbb{N}$ .
- (b) A sequence of diffeomorphisms  $\Phi_k : U_k \rightarrow \Phi_k(U_k) \subset \Sigma_k$ , with

$$\Phi_k(x_\infty) = x_k$$

and such that the sequence  $\{\Phi_k^* g_k\}_{k \in \mathbb{N}}$  smoothly converges in  $C^\infty$  to  $g_\infty$  on compact sets in  $\Sigma_\infty$ .

The sequence  $\{(U_k, \Phi_k)\}_{k \in \mathbb{N}}$  is called a *family of convergence pairs* of the sequence  $\{(\Sigma_k, g_k, x_k)\}_{k \in \mathbb{N}}$ , with respect to the limit  $(\Sigma_\infty, g_\infty, x_\infty)$ .

When we say *smooth convergence*, we always mean smooth convergence in the sense of Cheeger-Gromov. The family of convergence pairs is not unique. Two such families  $\{(U_k, \Phi_k)\}_{k \in \mathbb{N}}, \{(W_k, \Psi_k)\}_{k \in \mathbb{N}}$  are equivalent in the sense that there exists an isometry  $I$  of the limit  $(\Sigma_\infty, g_\infty, x_\infty)$ , such that for every compact subset  $K$  of  $\Sigma_\infty$ , there exists a natural number  $k_0$ , such that for any natural  $k \geq k_0$ :

- (a) The mapping  $\Phi_k^{-1} \circ \Psi_k$  is well defined over  $K$ .
- (b) The sequence  $\{\Phi_k^{-1} \circ \Psi_k\}_{k \geq k_0}$  smoothly converges to  $I$  on  $K$ .

The limiting pointed Riemannian manifold  $(\Sigma_\infty, g_\infty, x_\infty)$  of the Definition 29 is unique up to isometries.

**Definition 30** Let  $M$  be a Riemannian manifold. The injectivity radius at  $x \in M$  is the supremum of all values  $r$ , such that the exponential map from the unit ball  $B_r(x)$  in  $T_x M$ , to the manifold  $M$ , is injective.

**Definition 31** A complete Riemannian manifold  $(\Sigma, g)$  is said to have *bounded geometry*, if the following conditions are satisfied:

- (a) For any integer  $j \geq 0$ , there exists a uniform positive constant  $C_j$ , such that  $|\nabla^j R| \leq C_j$ .
- (b) The injectivity radius satisfies  $inj_g(\Sigma) > 0$ .

The following proposition is standard and will be useful in the proof of the long-time existence of the mean curvature flow.

**Proposition 3** Suppose  $(\Sigma, g)$  is a complete Riemannian manifold with bounded geometry. Suppose that  $\{\alpha_k\}_{k \in \mathbb{N}}$  is an increasing sequence of real numbers that tends to  $+\infty$  and let  $\{x_k\}_{k \in \mathbb{N}}$  be a sequence of points on  $\Sigma$ . Then, the sequence  $\{(\Sigma, \alpha_k^2 g, x_k)\}_{k \in \mathbb{N}}$  smoothly subconverges to the euclidean space  $(\mathbb{R}^m, g_{euc}, 0)$ .

We will use the following definition of uniformly bounded geometry for a sequence of pointed Riemannian manifolds.

**Definition 32** We say that a sequence  $\{(\Sigma_k, g_k, x_k)\}_{k \in \mathbb{N}}$  of complete pointed Riemannian manifolds has *uniformly bounded geometry*, if the following two conditions are satisfied:

- (a) For any integer  $j \geq 0$ , there exists a uniform constant  $C_j$ , such that for each  $k \in \mathbb{N}$  it holds  $|\nabla^j R_k| \leq C_j$ , where  $R_k$  is the curvature operator of the metric  $g_k$ .
- (b) There exists a uniform constant  $c_0$ , such that  $inj_{g_k}(\Sigma_k) \geq c_0 > 0$ .

In the next result, we state the Cheeger-Gromov compactness theorem for sequences of complete pointed Riemannian manifolds. The version that we present here is due to Hamilton [46].

**Theorem 29** Let  $\{(\Sigma_k, g_k, x_k)\}_{k \in \mathbb{N}}$  be a sequence of complete pointed Riemannian manifolds with uniformly bounded geometry. Then, the sequence  $\{(\Sigma_k, g_k, x_k)\}_{k \in \mathbb{N}}$  subconverges smoothly to a complete pointed Riemannian manifold  $(\Sigma_\infty, g_\infty, x_\infty)$ .

*Remark 9* We would like to mention here that due to an estimate from Cheeger, Gromov and Taylor [14], the above compactness theorem still holds under the weaker assumption that the injectivity radius is uniformly bounded from below by a positive constant, only along the base points  $\{x_k\}_{k \in \mathbb{N}}$ , thereby avoiding the assumption of the uniform lower bound for  $inj_{g_k}(\Sigma_k)$ .

### 7.3 Convergence of Immersions

**Definition 33** Let  $F_k: (\Sigma_k, g_k, x_k) \rightarrow (P_k, h_k, y_k)$  be a sequence of isometric immersions, such that  $F(x_k) = y_k$ , for any  $k \in \mathbb{N}$ . We say that the sequence  $\{F_k\}_{k \in \mathbb{N}}$  converges smoothly to an isometric immersion

$$F_\infty: (\Sigma_\infty, g_\infty, x_\infty) \rightarrow (P_\infty, h_\infty, y_\infty)$$

if the following conditions are satisfied:

- (a) The sequence  $\{(\Sigma_k, g_k, x_k)\}_{k \in \mathbb{N}}$  smoothly converges to  $(\Sigma_\infty, g_\infty, x_\infty)$ .
- (b) The sequence  $\{(P_k, h_k, y_k)\}_{k \in \mathbb{N}}$  smoothly converges to  $(P_\infty, h_\infty, y_\infty)$ .
- (c) If  $\{(U_k, \Phi_k)\}_{k \in \mathbb{N}}$  is a family of convergence pairs of  $\{(\Sigma_k, g_k, x_k)\}_{k \in \mathbb{N}}$  and  $\{(W_k, \Psi_k)\}_{k \in \mathbb{N}}$  is a family of convergence pairs of  $\{(P_k, h_k, y_k)\}_{k \in \mathbb{N}}$ , then for each  $k \in \mathbb{N}$ , we have  $F_k \circ \Phi_k(U_k) \subset \Psi_k(W_k)$  and  $\Psi_k^{-1} \circ F_k \circ \Phi_k$  smoothly converges to  $F_\infty$  on compact sets.

**Lemma 10** Suppose that  $(P, h)$  is a complete Riemannian manifold with bounded geometry. Then, for any  $C > 0$ , there exists a positive constant  $r > 0$ , such that  $\text{inj}_g(\Sigma) > r$ , for any isometric immersion  $F: (\Sigma, g) \rightarrow (P, h)$  such that the norm  $|A_F|$  of its second fundamental form satisfies  $|A_F| \leq C$ .

The last lemma and the Cheeger–Gromov compactness theorem allow us to deduce a compactness theorem in the category of sequences of immersions; see for example [27].

**Theorem 30** Let  $\{(\Sigma_k, g_k, x_k)\}_{k \in \mathbb{N}}$  and  $\{(P_k, h_k, y_k)\}_{k \in \mathbb{N}}$  be two sequences of complete Riemannian manifolds with dimensions  $m$  and  $l$ , respectively. Suppose that  $F_k: (\Sigma_k, g_k, x_k) \rightarrow (P_k, h_k, y_k)$  is a family of isometric immersions, where  $F_k(x_k) = y_k$ . Assume that:

- (a) Each  $\Sigma_k$  is compact.
- (b) The sequence  $\{(P_k, h_k, y_k)\}_{k \in \mathbb{N}}$  has uniformly bounded geometry.
- (c) For any integer  $j \geq 0$ , there exists a uniform constant  $C_j$ , such that

$$|(\nabla^{F_k})^j A_{F_k}| \leq C_j,$$

for any  $k \in \mathbb{N}$ . Here,  $A_{F_k}$  stands for the second fundamental form of  $F_k$ .

Then, the sequence  $\{F_k\}_{k \in \mathbb{N}}$  subconverges smoothly to a complete isometric immersion  $F_\infty: (\Sigma_\infty, g_\infty, x_\infty) \rightarrow (P_\infty, h_\infty, y_\infty)$ .

### 7.4 Modeling the Singularities

In the following theorem, we describe a method of rescaling around points, where the second fundamental form attains its maximum.

**Theorem 31** *Let  $\Sigma$  be a compact manifold and let  $F : \Sigma \times [0, T_{\max}) \rightarrow (P, h)$  be a solution of mean curvature flow, where  $P$  is a Riemannian manifold with bounded geometry and  $T_{\max} \leq \infty$  is the maximal time of existence of a smooth solution. Suppose that there exists a point  $x_\infty \in \Sigma$  and a sequence of points  $\{(x_k, t_k)\}_{k \in \mathbb{N}}$  in  $\Sigma \times [0, T)$  with  $\lim x_k = x_\infty, \lim t_k = T_{\max}$  such that*

$$a_k = \max_{M \times [0, t_k]} |A(x, t)| = |A(x_k, t_k)| \rightarrow \infty.$$

Then:

- (a) *The family of maps  $F_k : \Sigma \times [-a_k^2 t_k, 0] \rightarrow (P, a_k^2 h), k \in \mathbb{N}$ , given by*

$$F_k(x, s) = F_{k,s}(x) = F(x, s/a_k^2 + t_k),$$

*form a sequence of mean curvature flow solutions. The mean curvature  $H_{F_k}$  and the norm  $|A_{F_k}|$  of the second fundamental form of  $F_k$  satisfy the equation*

$$H_{F_k} = \frac{1}{a_k^2} H(x, s/a_k^2 + t_k) \quad \text{and} \quad |A_{F_k}(x, s)| = \frac{1}{a_k} |A(x, s/a_k^2 + t_k)|.$$

*Moreover, for any  $s \leq 0$  we have*

$$|A_{F_k}(x, s)| \leq 1 \quad \text{and} \quad |A_{F_k}(x_k, 0)| = 1,$$

*for any  $k \in \mathbb{N}$ .*

- (b) *For any fixed  $s \leq 0$ , the sequence  $\{(\Sigma, F_{k,s}^*(a_k^2 h), x_k)\}_{k \in \mathbb{N}}$  smoothly subconverges in the Cheeger-Gromov sense to a connected complete pointed Riemannian manifold  $(\Sigma_\infty, g_\infty(s), x_\infty)$ , where  $\Sigma_\infty$  does not depend on the choice of  $s$ . Moreover, the sequence  $\{(P, a_k^2 h, F_k(x_k, s))\}_{k \in \mathbb{N}}$  smoothly subconverges in the Cheeger-Gromov sense to the standard Euclidean space  $(\mathbb{R}^l, g_{\text{euc}}, 0)$ .*
- (c) *There is an ancient smooth solution  $F_\infty : \Sigma_\infty \times (-\infty, 0] \rightarrow \mathbb{R}^l$  of the mean curvature flow, such that for each fixed time  $s \leq 0$ , the sequence  $\{F_{k,s}\}_{k \in \mathbb{N}}$  smoothly subconverges in the Cheeger-Gromov sense to  $F_{\infty,s}$ . Additionally,*

$$|A_{F_\infty}| \leq 1 \quad \text{and} \quad |A_{F_\infty}(x_\infty, 0)| = 1.$$

- (d) *If  $\dim \Sigma = 2$  and  $H_{F_\infty} = 0$ , then the limiting Riemann surface  $\Sigma_\infty$  has finite total curvature. In the matter of fact, the limiting surface  $\Sigma_\infty$  is conformally*

*diffeomorphic to a compact Riemann surface minus a finite number of points and is of parabolic type.*

For the proof see [15] and [66].

## 8 Graphical MCF of Surfaces in Four Manifolds

Let  $(M, g_M)$  and  $(N, g_N)$  be compact Riemann surfaces. Recall that a smooth map  $f : M \rightarrow N$  is called *area decreasing* if  $|\Lambda^2 df| \leq 1$ , where  $\Lambda^2 df$  is the 2-Jacobian of  $f$ . Being area decreasing means that the map  $f$  contracts 2-dimensional regions of  $M$ . If  $|\Lambda^2 df| < 1$  the map is called *strictly area decreasing* and if  $|\Lambda^2 df| \equiv 1$  the map is said *area preserving*.

We will deform area decreasing maps  $f$  by evolving their corresponding graphs

$$\Gamma(f) = \{(x, f(x)) \in M \times N : x \in M\},$$

under the mean curvature flow in the Riemannian product 4-manifold

$$(M \times N, g_{M \times N} = \pi_M^* g_M + \pi_N^* g_N),$$

where  $\pi_M : M \times N \rightarrow M$  and  $\pi_N : M \times N \rightarrow N$  are the natural projection maps.

Our goal is to give a detailed, unified proof of the following theorem, which was shown in [78, 85, 95, 100]. For the strictly area decreasing case, we closely follow the presentation in [78].

**Theorem 32** *Let  $(M, g_M)$  and  $(N, g_N)$  be compact Riemann surfaces and  $f : M \rightarrow N$  be a smooth area decreasing map. Suppose that the curvatures  $\sigma_M$  of  $g_M$  and  $\sigma_N$  of  $g_N$  are related by*

$$\min \sigma_M \geq \max \sigma_N.$$

*Then there exists a family of smooth area decreasing maps  $f_t : M \rightarrow N$ ,  $t \in [0, \infty)$ ,  $f_0 = f$ , such that the graphs  $\Gamma(f_t)$  of  $f_t$  move by mean curvature flow in  $(M \times N, g_{M \times N})$ . Furthermore, there exist only two possible categories of initial data sets and corresponding solutions:*

- (I) *The curvatures  $\sigma_M$  and  $\sigma_N$  are constant and equal and the map  $f_0$  is area preserving. In this category, each  $f_t$  is area preserving and  $\Gamma(f_t)$  smoothly converges to a minimal Lagrangian graph  $\Gamma(f_\infty)$  in  $M \times N$ , with respect to the symplectic form*

$$\Omega_{M \times N} = \pi_M^* \Omega_M \mp \pi_N^* \Omega_N,$$



depending on whether the map  $f_0$  is orientation preserving or reversing, respectively. Here  $\Omega_M$  and  $\Omega_N$  are the positively oriented volume forms of  $M$  and  $N$ , respectively.

- (II) All other possible cases. In this category, for  $t > 0$  each map  $f_t$  is strictly area decreasing. Moreover, depending on the sign of  $\sigma = \min \sigma_M$  we have the following behavior:
  - (a) If  $\sigma > 0$ , then the family  $\Gamma(f_t)$  smoothly converges to the graph of a constant map.
  - (b) If  $\sigma = 0$ , then  $\Gamma(f_t)$  smoothly converges to a totally geodesic graph  $\Gamma(f_\infty)$  of  $M \times N$ .
  - (c) If  $\sigma < 0$ , then  $\Gamma(f_t)$  smoothly converges to a minimal surface  $M_\infty$  of the product manifold  $M \times N$ .

### 8.1 Jacobians of the Projection Maps

Let  $\Omega_M$  denote the Kähler form of the Riemann surface  $(M, g_M)$  and  $\Omega_N$  the Kähler form of  $(N, g_N)$ . We can extend  $\Omega_M$  and  $\Omega_N$  to two parallel 2-forms on the product manifold  $M \times N$  by pulling them back via the projection maps  $\pi_M$  and  $\pi_N$ . That is we may define the parallel forms  $\Omega_1 = \pi_M^* \Omega_M$  and  $\Omega_2 = \pi_N^* \Omega_N$ . Define now two smooth functions  $u_1$  and  $u_2$  given by

$$u_1 = *(F^* \Omega_1) = *{(\pi_M \circ F)^* \Omega_M} = *(I^* \Omega_M)$$

and

$$u_2 = *(F^* \Omega_2) = *{(\pi_N \circ F)^* \Omega_N} = *(f^* \Omega_N)$$

where here  $*$  stands for the Hodge star operator with respect to the metric  $g$ . Note that  $u_1$  is the Jacobian of the projection map from  $\Gamma(f)$  to the first factor of  $M \times N$  and  $u_2$  is the Jacobian of the projection map of  $\Gamma(f)$  to the second factor of  $M \times N$ . With respect to the basis  $\{e_1, e_2; \xi_3, \xi_4\}$  of the singular decomposition, we can write

$$u_1 = \frac{1}{\sqrt{(1 + \lambda^2)(1 + \mu^2)}} \quad \text{and} \quad |u_2| = \frac{\lambda \mu}{\sqrt{(1 + \lambda^2)(1 + \mu^2)}}.$$

Another important quantity that plays a crucial role in the case of maps between equi-dimensional manifolds is the *Jacobian determinant*, i.e., the map given by

$$\text{Jac}(f) = \frac{*(f^* \Omega_N)}{*(I^* \Omega_M)} = \frac{u_2}{u_1}.$$

Moreover, the difference between  $u_1$  and  $|u_2|$  measures how far  $f$  is from being area preserving. In particular:

$$\begin{aligned}
 u_1 - |u_2| \geq 0 & \quad \text{if and only if} \quad f \text{ is area decreasing,} \\
 u_1 - |u_2| > 0 & \quad \text{if and only if} \quad f \text{ is strictly area decreasing,} \\
 u_1 - |u_2| = 0 & \quad \text{if and only if} \quad f \text{ is area preserving.}
 \end{aligned}$$

### 8.2 The Kähler Angles

There are two natural complex structures associated to the product space  $(M \times N, \mathfrak{g}_{M \times N})$ , namely  $J_1 = \pi_M^* J_M - \pi_N^* J_N$  and  $J_2 = \pi_M^* J_M + \pi_N^* J_N$ , where  $J_M$  and  $J_N$  are the complex structures on  $M$  and  $N$  defined by

$$\Omega_M(\cdot, \cdot) = g_M(J_M \cdot, \cdot) \quad \text{and} \quad \Omega_N(\cdot, \cdot) = g_N(J_N \cdot, \cdot).$$

Chern and Wolfson [23] introduced a function which measures the deviation of the tangent plane  $dF(T_x M)$  from a complex line of the space  $T_{F(x)}(M \times N)$ . More precisely, if we consider  $(M \times N, \mathfrak{g}_{M \times N})$  as a complex manifold with respect to  $J_1$  then its corresponding Kähler angle  $a_1$  is given by the formula

$$\cos a_1 = \varphi = g_{M \times N}(J_1 dF(v_1), dF(v_2)) = u_1 - u_2.$$

For our convenience we require that  $a_1 \in [0, \pi]$ . Note that in general  $a_1$  is not smooth at points where  $\varphi = \pm 1$ . If there exists a point  $x \in M$  such that  $a_1(x) = 0$  then  $dF(T_x M)$  is a complex line of  $T_{F(x)}(M \times N)$  and  $x$  is called a *complex point* of  $F$ . If  $a_1(x) = \pi$  then  $dF(T_x M)$  is an anti-complex line of  $T_{F(x)}(M \times N)$  and  $x$  is said *anti-complex point* of  $F$ . In the case where  $a_1(x) = \pi/2$ , the point  $x$  is called *Lagrangian point* of the map  $F$ . In this case  $u_1 = u_2$ . Similarly, if we regard the product manifold  $(M \times N, \mathfrak{g}_{M \times N})$  as a Kähler manifold with respect to the complex structure  $J_2$ , then its corresponding Kähler angle  $a_2$  is defined by the formula

$$\cos a_2 = \vartheta = g_{M \times N}(J_2 dF(v_1), dF(v_2)) = u_1 + u_2.$$

The graph  $\Gamma(f)$  in the product Kähler manifold  $(M \times N, \mathfrak{g}_{M \times N}, J_i)$  is called *symplectic* with respect to the Kähler form related to  $J_i$ , if the corresponding Kähler angle satisfies  $\cos a_i > 0$ . Therefore a map  $f$  is strictly area decreasing if and only if its graph  $\Gamma(f)$  is symplectic with respect to both Kähler forms.

### 8.3 Structure Equations

Around each point  $x \in \Gamma(f)$  we choose an adapted local orthonormal frame  $\{e_1, e_2; \xi_3, \xi_4\}$  along the graph. In this special case the *Gauss equation* reads

$$2\sigma_g = 2u_1^2\sigma_M + 2u_2^2\sigma_N + |H|^2 - |A|^2,$$

where here  $\sigma_g$  is the Gauss curvature of the induced metric. From the *Ricci equation* we see that the curvature  $\sigma_n$  of the normal bundle of  $\Gamma(f)$  is given by the formula

$$\sigma_n = R_{1234}^\perp = u_1u_2(\sigma_M + \sigma_N) + A_{11}^3A_{12}^4 - A_{12}^3A_{11}^4 + A_{12}^3A_{22}^4 - A_{22}^3A_{12}^4.$$

The sum of the last four terms in the above formula is equal to minus the commutator  $\sigma^\perp$  of the matrices  $A^3 = (A_{ij}^3)$  and  $A^4 = (A_{ij}^4)$ , that is

$$\sigma^\perp = \langle [A^3, A^4]e_1, e_2 \rangle = -A_{11}^3A_{12}^4 + A_{12}^3A_{11}^4 - A_{12}^3A_{22}^4 + A_{22}^3A_{12}^4. \tag{38}$$

In the case where  $u_1 = u_2$  and  $\sigma_M = \sigma = \sigma_n$ , it turns out that the immersion  $F$  is Lagrangian and  $\sigma_g = \sigma_n$ . In this case, the following algebraic equality holds

$$\sigma^\perp = \frac{|A|^2 - |H|^2}{2}. \tag{39}$$

### 8.4 Estimates for the Jacobians and the Kähler Angles

Let us evolve now by mean curvature flow the graph  $\Gamma(f)$ . Denote by  $T_{\max}$  the maximal time of solution of the flow and by  $T_\Gamma$  the time until graphical property is preserved. Of course,  $0 < T_\Gamma \leq T_{\max}$ . We will give here several a priori estimates for the Jacobians  $u_1$  and  $u_2$  and the Kähler angles. The proofs are straightforward and follow directly as special cases of the general formulas of Section 6.5.

**Lemma 11** *The gradients of the functions  $\varphi, \vartheta$  at a point  $x \in M$  satisfy the equations*

$$\begin{aligned} |\nabla\varphi|^2 &= (1 - \varphi^2)((A_{11}^3 + A_{12}^4)^2 + (A_{12}^3 + A_{22}^4)^2), \\ |\nabla\vartheta|^2 &= (1 - \vartheta^2)((A_{11}^3 - A_{12}^4)^2 + (A_{12}^3 - A_{22}^4)^2), \end{aligned}$$

*As long the mean curvature flow remains graphical, the Jacobians  $u_1$  and  $u_2$  satisfy the following coupled system of parabolic equations*

$$\partial_t u_1 - \Delta u_1 = |A|^2 u_1 + 2\sigma^\perp u_2 + \sigma_M(1 - u_1^2 - u_2^2)u_1 - 2\sigma_N u_1 u_2^2,$$

$$\partial_t u_2 - \Delta u_2 = |A|^2 u_2 + 2\sigma^\perp u_1 + \sigma_N(1 - u_1^2 - u_2^2)u_2 - 2\sigma_M u_1^2 u_2.$$

Moreover,  $\varphi$  and  $\vartheta$  satisfy the following system of equations

$$\partial_t \varphi - \Delta \varphi = (|A|^2 - 2\sigma^\perp)\varphi + \frac{1}{2}(\sigma_M(\varphi + \vartheta) + \sigma_N(\varphi - \vartheta))(1 - \varphi^2),$$

$$\partial_t \vartheta - \Delta \vartheta = (|A|^2 + 2\sigma^\perp)\vartheta + \frac{1}{2}(\sigma_M(\varphi + \vartheta) - \sigma_N(\varphi - \vartheta))(1 - \vartheta^2).$$

**Lemma 12** *Let  $f : (M, g_M) \rightarrow (N, g_N)$  be an area decreasing map between compact Riemann surfaces. Suppose that the curvatures of  $g_M$  and  $g_N$  satisfy  $\sigma = \min \sigma_M \geq \max \sigma_N$ . Then the following statements hold.*

- (a) *The conditions  $\text{Jac}(f) \leq 1$  or  $\text{Jac}(f) \geq -1$  are both preserved as long as the flow remains graphical. In particular, the area decreasing property is preserved as long as the flow remains graphical.*
- (b) *If there is a point  $(x_0, t_0) \in M \times (0, T_\Gamma)$  where  $\text{Jac}^2(f)(x_0, t_0) = 1$ , then  $\text{Jac}^2(f) \equiv 1$  in space and time and  $\sigma_M \equiv \sigma \equiv \sigma_N$ .*
- (c) *The flow remains graphical as long as it exists, that is  $T_\Gamma = T_{\max}$ .*

**Proof**

- (a) From Lemma 11, we deduce that

$$\partial_t \varphi - \Delta \varphi = (|A|^2 - 2\sigma^\perp + \sigma_N(1 - \varphi^2))\varphi + \frac{1}{2}(\sigma_M - \sigma_N)(\varphi + \vartheta)(1 - \varphi^2).$$

Note that the quantities  $1 - \varphi^2$  and  $\varphi + \vartheta$  are non-negative. Hence, because of our curvature assumptions, the last line of the above equality is positive. Thus, there exists a time dependent function  $h$  such that

$$\partial_t \varphi - \Delta \varphi \geq h \varphi.$$

From the maximum principle we deduce that  $\varphi$  stays non-negative in time.

- (b) From the strong maximum principle it follows that if  $\varphi$  vanishes at a point  $(x_0, t_0) \in M \times (0, T_\Gamma)$ , then it vanishes identically in space and time. In this case,  $\vartheta$  is positive. Going back to the evolution equation of  $\varphi$ , we see that  $\sigma_M$  and  $\sigma_N$  must be constant equal to  $\sigma$ . Similarly, we prove the results concerning  $\vartheta$ .
- (c) By compactness, initially, we have that  $\min_{x \in M} u_1(x, 0) = \varepsilon > 0$ . By continuity, the minimum of  $u_1$  stays positive for small values of  $t$ . However, we will show that the flow remains graphical as long as it exists. As a matter of fact, we will show that

$$\min_{x \in M} u_1(x, t) > 0,$$

as long as the flow exists. Suppose to the contrary, that there exists a first time where the graphical property does not hold. This means that there exists a point  $(x_0, t_0)$  in space-time with  $t_0 < T$ , such that  $u_1(x_0, t_0) = 0$  and  $u_1(x, t) > 0$ , for all  $(x, t) \in M \times [0, t_0)$ . Since the area decreasing property is preserved by the flow and  $|A|^2$  is bounded on  $M \times [0, t_0]$ , there exists a constant  $c(t_0) \in \mathbb{R}$ , such that

$$\partial_t u_1 - \Delta u_1 \geq c(t_0)u_1,$$

for all  $(x, t) \in M \times [0, t_0)$ . From the parabolic maximum principle, we get  $u_1(x, t) \geq e^{c(t_0)t}$ , for all  $(x, t) \in M \times [0, t_0)$ . Passing to the limit as  $t$  approaches  $t_0$ , we obtain

$$u_1(x_0, t_0) = \lim_{t \rightarrow t_0} u_1(x_0, t) \geq e^{c(t_0)t_0} > 0,$$

which leads to a contradiction.

This completes the proof. □

From the Lemma 12 we see that, under our assumptions, the evolved maps  $\{f_t\}_{t \in (0, T_{\max})}$  are either strictly area decreasing or area preserving. This fact leads us to investigate these two cases separately.

### 8.4.1 Strictly Area Decreasing Case

We will explore the behaviour of  $\rho : M \times [0, T_{\max}) \rightarrow \mathbb{R}$  given by  $\rho = \varphi \vartheta$  under the graphical mean curvature flow.

**Lemma 13** *Let  $(M, g_M)$  and  $(N, g_N)$  be compact Riemann surfaces such that their curvatures  $\sigma_M$  and  $\sigma_N$  are related by  $\sigma = \min \sigma_M \geq \max \sigma_N$ . The following hold true:*

(a) *If  $\sigma \geq 0$ , then there exists a positive constant  $c_0$  such that*

$$\rho \geq \frac{c_0 e^{\sigma t}}{\sqrt{1 + c_0^2 e^{2\sigma t}}},$$

*for any  $(x, t)$  in space-time.*

(b) *If  $\sigma < 0$ , then there exists a positive constant  $c_0$  such that*

$$\rho \geq \frac{c_0 e^{2\sigma t}}{\sqrt{1 + c_0^2 e^{4\sigma t}}},$$

*for any  $(x, t)$  in space-time.*

**Proof** From Lemma 11 we get,

$$\partial_t \rho - \Delta \rho = 2\rho |A|^2 - 2\langle \nabla \varphi, \nabla \vartheta \rangle + 2(1 - \rho)\sigma_M u_1^2 - 2(1 + \rho)\sigma_N u_2^2.$$

Note that

$$\begin{aligned} -2\rho \langle \nabla \varphi, \nabla \vartheta \rangle + \frac{1}{2} |\nabla \rho|^2 &= \frac{1}{2} (|\nabla(\varphi \vartheta)|^2 - 4\varphi \vartheta \langle \nabla \varphi, \nabla \vartheta \rangle) \\ &= \frac{1}{2} (\varphi^2 |\nabla \vartheta|^2 + \vartheta^2 |\nabla \varphi|^2 - 2\varphi \vartheta \langle \nabla \varphi, \nabla \vartheta \rangle) \\ &\geq \frac{1}{2} (|\varphi \nabla \vartheta| - |\vartheta \nabla \varphi|)^2. \end{aligned}$$

Since by assumption  $\sigma_M \geq \sigma \geq \sigma_N$ , we deduce that

$$\partial_t \rho - \Delta \rho \geq -\frac{1}{2\rho} |\nabla \rho|^2 + 2\sigma \rho (1 - u_1^2 - u_2^2).$$

One can algebraically check that

$$1 - \rho^2 \leq 2(1 - u_1^2 - u_2^2) \leq 2(1 - \rho^2). \tag{40}$$

(a) Suppose at first that  $\sigma \geq 0$ . Then

$$\partial_t \rho - \Delta \rho \geq -\frac{1}{2\rho} |\nabla \rho|^2 + \sigma \rho (1 - \rho^2).$$

From the comparison maximum principle we obtain

$$\rho \geq \frac{c_0 e^{\sigma t}}{\sqrt{1 + c_0^2 e^{2\sigma t}}},$$

where  $c_0$  is a positive constant.

(b) In the case where  $\sigma < 0$ , from the Equation (40) we deduce that

$$\partial_t \rho - \Delta \rho \geq -\frac{1}{2\rho} |\nabla \rho|^2 + 2\sigma \rho (1 - \rho^2),$$

from where we get the desired estimate.

This completes the proof. □

### 8.4.2 Area Preserving Case

Suppose that the family of the graphs is generated orientation preserving by area decreasing maps. This means that  $\varphi$  is identically zero. In the next lemma we derive an estimate for the Kähler angle  $\vartheta$ .

**Lemma 14** *Suppose that  $M$  and  $N$  are compact with the same constant sectional curvature  $\sigma$  and that  $f : M \rightarrow N$  is an area preserving map. Then, there exists a positive real number  $c_0$  such that*

$$1 \geq \vartheta(x, t) \geq \frac{c_0 e^{\sigma t}}{\sqrt{1 + c_0^2 e^{2\sigma t}}},$$

for any point  $(x, t)$  in space-time.

**Proof** Since  $|A|^2 + 2\sigma^\perp \geq 0$ , from the evolution equation of  $\vartheta$ , we get

$$\partial_t \vartheta - \Delta \vartheta \geq \sigma \vartheta (1 - \vartheta^2).$$

According to the parabolic maximum principle, there exist a positive real number  $c_0$  such that

$$\vartheta(x, t) \geq \frac{c_0 e^{\sigma t}}{\sqrt{1 + c_0^2 e^{2\sigma t}}},$$

for any  $(x, t)$  in space-time. This completes the proof. □

## 8.5 Curvature Decay Estimates

### 8.5.1 Strictly Area Decreasing Case

**Lemma 15** *Let  $f : (M, g_M) \rightarrow (N, g_N)$  be a strictly area decreasing map. Suppose that the curvatures of  $M$  and  $N$  satisfy  $\sigma = \min \sigma_M \geq \max \sigma_N$ . Let  $\delta : [0, T) \rightarrow \mathbb{R}$  be a positive increasing real function and  $\tau$  the time dependent function given by  $\tau = \log(\delta|H|^2 + \varepsilon)$ , where  $\varepsilon$  is a non-negative number. Then,*

$$\begin{aligned} \partial_t \tau - \Delta \tau \leq & \frac{2\delta}{\delta|H|^2 + \varepsilon} |H|^2 |A|^2 + \frac{\delta'}{\delta|H|^2 + \varepsilon} |H|^2 \\ & + \frac{2\delta}{\delta|H|^2 + \varepsilon} |H|^2 \sigma_M (1 - u_1^2 - u_2^2) + \frac{1}{2} |\nabla \tau|^2. \end{aligned}$$

**Proof** Recall from Lemma 6 that  $|H|^2$  evolves in time under the equation

$$\begin{aligned} \partial_t |H|^2 - \Delta |H|^2 &= 2|A^H|^2 - 2|\nabla^\perp H|^2 \\ &\quad + 2\tilde{R}(H, dF(e_1), H, dF(e_1)) + 2\tilde{R}(H, dF(e_2), H, dF(e_2)), \end{aligned}$$

where  $\{e_1, e_2\}$  is a local orthonormal frame with respect to  $g$ . Using the special frames of the singular value decomposition we see that

$$\begin{aligned} &\tilde{R}(H, dF(e_1), H, dF(e_1)) + \tilde{R}(H, dF(e_2), H, dF(e_2)) \\ &= \sigma_M u_1^2 (\lambda^2 + \mu^2) |H|^2 - (\sigma_M - \sigma_N) u_1^2 (\lambda^2 (H^4)^2 + \mu^2 (H^3)^2) \\ &\leq \sigma_M (1 - u_1^2 - u_2^2) |H|^2. \end{aligned}$$

Note that from Cauchy–Schwarz inequality  $|A^H| \leq |A| \cdot |H|$ . Moreover, observe that at points where the mean curvature vector is non-zero, from Kato’s inequality, we have that

$$|\nabla^\perp H|^2 \geq |\nabla |H||^2.$$

Consequently, at points where the norm  $|H|$  of the mean curvature is not zero the following inequality holds

$$\partial_t |H|^2 - \Delta |H|^2 \leq -2|\nabla |H||^2 + 2|A|^2 |H|^2 + 2\sigma_M (1 - u_1^2 - u_2^2) |H|^2.$$

Now let us compute the evolution equation of the function  $\tau$ . We have,

$$\begin{aligned} \partial_t \tau - \Delta \tau &= \frac{\delta(\partial_t |H|^2 - \Delta |H|^2)}{\delta |H|^2 + \varepsilon} + \frac{\delta^2 |\nabla |H||^2}{(\delta |H|^2 + \varepsilon)^2} + \frac{\delta' |H|^2}{\delta |H|^2 + \varepsilon} \\ &\leq -\frac{2\delta}{\delta |H|^2 + \varepsilon} |\nabla |H||^2 + \frac{\delta^2}{(\delta |H|^2 + \varepsilon)^2} |\nabla |H||^2 \\ &\quad + \frac{2\delta}{\delta |H|^2 + \varepsilon} |H|^2 |A|^2 + \frac{\delta'}{\delta |H|^2 + \varepsilon} |H|^2 \\ &\quad + \frac{2\delta}{2\delta |H|^2 + \varepsilon} |H|^2 \sigma_M (1 - u_1^2 - u_2^2). \end{aligned}$$

Note that

$$-\frac{2\delta}{\delta |H|^2 + \varepsilon} |\nabla |H||^2 + \frac{1}{2} \frac{\delta^2}{(\delta |H|^2 + \varepsilon)^2} |\nabla |H||^2 \leq 0.$$



Therefore,

$$\begin{aligned} \partial_t \tau - \Delta \tau \leq & \frac{1}{2} |\nabla \tau|^2 + \frac{2\delta}{\delta |H|^2 + \varepsilon} |H|^2 |A|^2 \\ & + \frac{\delta'}{\delta |H|^2 + \varepsilon} |H|^2 + \frac{2\delta}{\delta |H|^2 + \varepsilon} |H|^2 \sigma_M (1 - u_1^2 - u_2^2), \end{aligned}$$

and this completes the proof. □

**Theorem 33** *Let  $f : (M, g_M) \rightarrow (N, g_N)$  be an area decreasing map, where  $M$  and  $N$  are compact Riemann surfaces. Suppose that the curvatures of  $M$  and  $N$  satisfy  $\sigma = \min \sigma_M \geq \sup \sigma_N$ . Then the following statements hold:*

- (a) *There exist a positive time independent constant  $C$  such that  $|H|^2 \leq C$ .*
- (b) *If  $\sigma \geq 0$ , there exist a time independent constant  $C$  so that  $|H|^2 \leq Ct^{-1}$ .*

**Proof** Consider the time dependent function  $\Theta = \log(\delta |H|^2 + \varepsilon) - \log \rho$ , where  $\delta$  is a positive increasing function. From Lemmas 6 and 13 and  $|H|^2 \leq 2|A|^2$ , we deduce that

$$\partial_t \Theta - \Delta \Theta \leq \frac{1}{2} \langle \nabla \Theta, \nabla \tau + \nabla \rho \rangle + \frac{\delta' |H|^2 - \varepsilon |H|^2 - 2\varepsilon \sigma (1 - u_1^2 - u_2^2)}{\delta |H|^2 + \varepsilon}.$$

Choosing  $\delta = 1$  and  $\varepsilon = 0$ , we obtain that

$$\partial_t \Theta - \Delta \Theta \leq \frac{1}{2} \langle \nabla \Theta, \nabla \tau + \nabla \rho \rangle.$$

From the maximum principle the norm  $|H|$  remains uniformly bounded in time regardless of the sign of the constant  $\sigma$ . In the case where  $\sigma \geq 0$ , choosing  $\varepsilon = 1$  and  $\delta = t$ , we deduce that  $\Theta$  remains uniformly bounded in time which gives the desired decay estimate for  $H$ . □

### 8.5.2 Area Preserving Case

In the sequel, we provide a very important decay estimate due to Wang [91] for the mean curvature in the area preserving case.

**Theorem 34** *Suppose that  $M$  and  $N$  are compact Riemannian manifolds with the same constant sectional curvature  $\sigma$  and that  $f : M \rightarrow N$  is an area preserving map. Then, the following decay estimate holds:*

$$\int \frac{|H|^2}{\vartheta} \Omega \leq e^{\sigma t},$$

where  $\Omega$  is the volume element of the induced metric.

**Proof** The idea is to compare  $|H|$  with  $\vartheta$ . We compute

$$\begin{aligned} \partial_t \left( \vartheta^{-1} |H|^2 \right) - \Delta \left( \vartheta^{-1} |H|^2 \right) &= \vartheta^{-1} (\partial_t |H|^2 - \Delta |H|^2) - \vartheta^{-2} |H|^2 (\partial_t \vartheta - \Delta \vartheta) \\ &\quad + 2\vartheta^{-2} \langle \nabla |H|^2, \nabla \vartheta \rangle - 2\vartheta^{-3} |H|^2 |\nabla \vartheta|^2. \end{aligned}$$

But from the evolution equation of  $\vartheta$  and  $|H|^2$ , we obtain

$$\begin{aligned} \partial_t \left( \vartheta^{-1} |H|^2 \right) - \Delta \left( \vartheta^{-1} |H|^2 \right) & \tag{41} \\ &= \vartheta^{-1} \left( -2|\nabla^\perp H|^2 + 2 \sum_{k,\alpha,\beta} H^\alpha H^\beta \tilde{R}_{\alpha k \beta k} + 2 \sum_{i,j} (A_{ij}^H)^2 \right) \\ &\quad - \vartheta^{-2} |H|^2 \left( (|A|^2 + 2\sigma^\perp) \vartheta + \sigma \vartheta (1 - \vartheta^2) \right) + 2\vartheta^{-2} \langle \nabla |H|^2, \nabla \vartheta \rangle - 2\vartheta^{-3} |H|^2 |\nabla \vartheta|^2. \end{aligned}$$

Using the Equation (39) and the formula

$$\sum_{k,\alpha,\beta} H^\alpha H^\beta \tilde{R}_{\alpha k \beta k} = \sigma \left( 1 - \frac{\vartheta^2}{2} \right) |H|^2 \tag{42}$$

the identity (41) becomes

$$\begin{aligned} \partial_t \left( \vartheta^{-1} |H|^2 \right) - \Delta \left( \vartheta^{-1} |H|^2 \right) & \\ &= \vartheta^{-3} \left( 4\vartheta |H| \langle \nabla \vartheta, \nabla |H| \rangle - 2|\nabla \vartheta|^2 |H|^2 - 2\vartheta^2 |\nabla^\perp H|^2 \right) \\ &\quad + \vartheta^{-1} \left( 2 \sum_{i,j} (A_{ij}^H)^2 - 2|H|^2 |A|^2 + |H|^4 \right) + \sigma \vartheta^{-1} |H|^2. \end{aligned}$$

Integrating and using Stokes' theorem, we have

$$\begin{aligned} \partial_t \left( \int \vartheta^{-1} |H|^2 \Omega \right) &= \int \vartheta^{-1} |H|^2 \nabla_{\partial_t} \Omega \\ &\quad + 2 \int \vartheta^{-3} \left( 2\vartheta |H| \langle \nabla \vartheta, \nabla |H| \rangle - |\nabla \vartheta|^2 |H|^2 - \vartheta^2 |\nabla^\perp H|^2 \right) \Omega \\ &\quad + \int \vartheta^{-1} \left( 2 \sum_{i,j} (A_{ij}^H)^2 - 2|H|^2 |A|^2 + |H|^4 \right) \Omega + \sigma \int \vartheta^{-1} |H|^2 \Omega. \end{aligned}$$

Using

$$|\nabla |H|| \leq |\nabla^\perp H|$$

in the first term on the right hand side of the above equation and completing the square, we have

$$2\vartheta|H|\langle \nabla\vartheta, \nabla|H| \rangle - |\nabla\vartheta|^2|H|^2 - \vartheta^2|\nabla|H||^2 = -||H|\nabla\vartheta - \vartheta\nabla|H||^2 \leq 0.$$

Moreover, from Lemma 4, we have  $\nabla_{\partial_t}\Omega = -|H|^2\Omega$ . Also, by Cauchy–Schwarz inequality, we get

$$\sum_{i,j} (A_{ij}^H)^2 \leq \sum_{i,j} |A_{ij}|^2|H|^2 = |A|^2|H|^2.$$

Therefore, putting everything together, we get

$$\partial_t \left( \int \vartheta^{-1}|H|^2\Omega \right) \leq \sigma \int \vartheta^{-1}|H|^2\Omega$$

and by integration, we obtain the result. □

### 8.6 Long-time Existence

We will show now that the graphical MCF exists for all times.

**Theorem 35** *Let  $(M, g_M)$  and  $(N, g_N)$  be compact Riemann surfaces such that their curvatures  $\sigma_M$  and  $\sigma_N$  are related by  $\sigma = \min \sigma_M \geq \max \sigma_N$ . Also, let  $f : M \rightarrow N$  be an area preserving map. Evolve the graph off under the mean curvature flow. Then, the norm of the second fundamental form of the evolved graphs stays uniformly bounded in time and so the graphical mean curvature flow exists for all times.*

**Proof** Suppose that  $|A|$  is not uniformly bounded. Then, there exists a sequence  $\{(x_k, t_k)\}_{k \in \mathbb{N}}$  in  $M \times [0, T_{\max})$  with  $\lim t_k = T_{\max} \leq \infty$ , and such that

$$a_k = \max_{(x,t) \in M \times [0,t_k]} |A(x, t)| = |A(x_k, t_k)| \rightarrow \infty.$$

Let  $F_k : M \times [-a_k^2 t_k, 0] \rightarrow (N, a_k^2 g_N)$  be the graph of the “rescaled map”

$$f : (M, a_k^2 g_M) \rightarrow (N, a_k^2 g_N).$$

**Claim:** *The singular values are invariant under parabolic rescalings.*

Let  $\{\alpha_1, \alpha_2\}$  and  $\{\beta_1, \beta_2\}$  orthonormal frames of the singular value decomposition of  $f$ . Then  $\{\tilde{\alpha}_1 = \alpha_1/a_k, \tilde{\alpha}_2 = \alpha_2/a_k\}$  is an orthonormal frame with respect to  $a_k^2 g_M$  and  $\{\tilde{\beta}_1 = \beta_1/\alpha_k, \tilde{\beta}_2 = \beta_2/\alpha_k\}$  is orthonormal with respect to  $g_N$ . Therefore, the singular values of the rescaled map  $f$  are given by

$$df(\tilde{\alpha}_1) = \frac{1}{a_k} df(\alpha_1) = \lambda \frac{\beta_1}{a_k} = \lambda \tilde{\beta}_1$$

and

$$df(\tilde{\alpha}_2) = \frac{1}{a_k} df(\alpha_2) = \mu \frac{\beta_2}{a_k} = \mu \tilde{\beta}_2.$$

This completes the proof of the claim.

Thus,  $\varphi_{F_k} = \varphi$  and  $\vartheta_{F_k} = \vartheta$ . Also, from Theorem 31(a) we have

$$H_{F_k}(x, s) = \frac{1}{a_k^2} H(x, s/a_k^2 + t_k),$$

for any  $(x, s) \in M \times [-a_k^2 t_k, 0]$ .

**CASE 1** Suppose that the evolved graphs are generated by strictly area decreasing maps. Since from the estimate of Lemma 15 the norm  $|H|$  of the mean curvature is uniformly bounded and the convergence is smooth, it follows that  $F_\infty : \Sigma_\infty \rightarrow \mathbb{R}^4$  is a complete minimal immersion of parabolic type. Hence, any non-negative superharmonic function must be constant. Since the convergence is smooth, the corresponding Kähler angles  $\varphi_\infty, \vartheta_\infty$  of  $F_\infty$  with respect to the complex structures  $J = (J_{\mathbb{R}^2}, -J_{\mathbb{R}^2})$  and  $J_2 = (J_{\mathbb{R}^2}, J_{\mathbb{R}^2})$  of  $\mathbb{R}^4$  are non-negative. As in Lemma 11 we get that

$$\Delta\varphi_\infty + (|A_{F_\infty}|^2 - 2\sigma_{F_\infty}^\perp)\varphi_\infty = 0, \tag{43}$$

$$\Delta\vartheta_\infty + (|A_{F_\infty}|^2 + 2\sigma_{F_\infty}^\perp)\vartheta_\infty = 0, \tag{44}$$

where  $-\sigma_{F_\infty}^\perp$  is the normal curvature of  $F_\infty$ . Moreover,

$$|\nabla\varphi_\infty|^2 = (1 - \varphi_\infty^2) \left( ((A_{F_\infty})_{11}^3 + (A_{F_\infty})_{12}^4)^2 + ((A_{F_\infty})_{12}^3 - (A_{F_\infty})_{11}^4)^2 \right), \tag{45}$$

$$|\nabla\vartheta_\infty|^2 = (1 - \vartheta_\infty^2) \left( ((A_{F_\infty})_{11}^3 - (A_{F_\infty})_{12}^4)^2 + ((A_{F_\infty})_{12}^3 + (A_{F_\infty})_{11}^4)^2 \right). \tag{46}$$

Note that from (38) one can easily derive the inequalities

$$|A_{F_\infty}|^2 \pm 2\sigma_{F_\infty}^\perp \geq 0.$$

From (43) and (44) we deduce that  $\varphi_\infty$  and  $\vartheta_\infty$  are superharmonic and consequently they must be constants. Thus, the functions  $(u_1)_\infty$  and  $(u_2)_\infty$  are also constants. We will distinguish three subcases:

**Sub-case A** Suppose at first that  $\varphi_\infty > 0$  and  $\vartheta_\infty > 0$ . Then from (43) and (44) we deduce that

$$|A_{F_\infty}|^2 \pm 2\sigma_{F_\infty}^\perp = 0$$

which implies that  $|A_{F_\infty}| = 0$ . This contradicts the fact that there is a point where  $|A_{F_\infty}| = 1$ .

**Sub-case B** Suppose that both constants  $\varphi_\infty$  and  $\vartheta_\infty$  are zero. Then from the Equations (45) and (46) we obtain that  $A_{F_\infty}$  vanishes identically, which is a again a contradiction.

**Sub-case C** Suppose now that only one of the constants  $\varphi_\infty, \vartheta_\infty$  is zero. Let us assume that  $\varphi_\infty = 0$  and  $\vartheta_\infty > 0$ . The case  $\varphi_\infty > 0$  and  $\vartheta_\infty = 0$  is treated in a similar way. Since  $\varphi_\infty = 0$ ,  $F_\infty : \Sigma_\infty \rightarrow \mathbb{R}^4$  must be a minimal Lagrangian immersion. Note that in this case necessarily  $(u_1)_\infty = (u_2)_\infty = \text{const} > 0$ . Recall from Theorem 5 that the minimal Lagrangian  $F_\infty$  can be locally reparametrized in the form

$$F_\infty = \frac{1}{\sqrt{2}} e^{i\beta/2} (\mathcal{F}_1 - i\overline{\mathcal{F}}_2, \mathcal{F}_2 + i\overline{\mathcal{F}}_1),$$

where  $\beta$  is a constant and  $\mathcal{F}_1, \mathcal{F}_2 : \mathbb{D} \subset \mathbb{C} \rightarrow \mathbb{C}$  are holomorphic functions defined in a simply connected domain  $\mathbb{D}$  such that

$$|(\mathcal{F}_1)_z|^2 + |(\mathcal{F}_2)_z|^2 > 0.$$

The Gauss map of  $F_\infty$  is described by  $\mathcal{G} : \mathbb{D} \rightarrow \mathbb{S}^2 = \mathbb{C} \cup \{\infty\}$  given by

$$\mathcal{G} = (\mathcal{F}_1)_z / (\mathcal{F}_2)_z.$$

Because  $(u_1)_\infty = \text{const} > 0$  we get that  $F_\infty$  is the graph of an area preserving map  $h$ . Then

$$\mathcal{F}_1 = (z + i\overline{h})/2, \quad \mathcal{F}_2 = (-i\overline{z} + h)/2 \quad \text{and} \quad |h_z|^2 - |h_{\overline{z}}|^2 = 1.$$

Therefore

$$\mathcal{G} = (\mathcal{F}_1)_z / (\mathcal{F}_2)_z = (1 - ih_{\overline{z}}) / h_z.$$

A straightforward computation shows that

$$|\mathcal{G}|^2 = \frac{|1 + i\overline{h_{\overline{z}}}|^2}{|h_z|^2} = \frac{1 + |h_{\overline{z}}|^2 + i(\overline{h_{\overline{z}}} - h_z)}{1 + |h_{\overline{z}}|^2} = 1 + \frac{2 \text{Im}(h_{\overline{z}})}{1 + |h_{\overline{z}}|^2} \leq 2.$$

Hence, the image of  $\mathcal{G}$  is contained in a bounded subset of  $\mathbb{C} \cup \{\infty\}$ . But then, due to Theorem 11 the immersion  $F_\infty$  must be flat, which is a contradiction.

**CASE 2** Let us treat now the area preserving case. In this situation, we have that

$$\frac{|H_{F_k}|^2}{\vartheta_{F_k}} = \frac{1}{a_k^2} \frac{|H|^2}{\vartheta}.$$

We distinguish two subcases:

**Sub-case A** Let us suppose that  $\sigma \leq 0$ . Using Lemma 34, we have

$$\int \frac{|H_{F_k}|^2}{\vartheta_{F_k}} \Omega_k = \frac{1}{a_k^2} \int \frac{|H|^2}{\vartheta} \Omega \leq \frac{1}{a_k^2} e^{\sigma(s/a_k^2 + tk)} \leq \frac{1}{a_k^2} c,$$

where  $c > 0$ . Since the convergence is smooth, we have

$$0 = \lim_{k \rightarrow \infty} \int \frac{|H_{F_k}|^2}{\vartheta_{F_k}} \Omega = \int \lim_{k \rightarrow \infty} \frac{|H_{F_k}|^2}{\vartheta_{F_k}} \Omega = \int \frac{|H_{F_\infty}|^2}{\vartheta_\infty} \Omega.$$

Therefore,  $H_{F_\infty} = 0$ . Proceeding exactly in the same way as in CASE 1 we can prove that  $F_\infty$  is flat, something which leads to a contradiction.

**Sub-case B** Let us treat now the case  $\sigma > 0$ . We will show at first that  $T_{\max} = \infty$ . To show this, assume in contrary that  $T_{\max} < +\infty$ . Then,

$$\int \frac{|H|^2}{\vartheta} \Omega \leq e^{\sigma t} \leq e^{\sigma T_{\max}} < +\infty.$$

As in the previous case, we deduce that  $H_{F_\infty} = 0$ . Performing exactly the same procedure as above, we get a contradiction. Therefore, there is no finite time singularity and the flow exists for all times. It remains to show that  $|A|^2 \leq C$ , where  $C$  is time independent. Indeed, since  $\lambda\mu = 1$ , we obtain

$$\vartheta = \frac{2\lambda}{1 + \lambda^2} \leq 1.$$

On the other hand, from Lemma 14, we have

$$1 \geq \vartheta \geq \frac{c_0 e^t}{\sqrt{1 + c_0^2 e^{2t}}},$$

which tends to 1 as  $t \rightarrow \infty$ . Therefore,  $\vartheta_\infty = 1$  and  $\lambda_\infty = 1$ . Therefore,  $f_\infty$  is an isometry and, thus,  $F_\infty$  must be totally geodesic. The latter implies  $|A_{F_\infty}| = 0$  and this is again a contradiction.

This completes the proof. □

### 8.7 Convergence and Proof of Theorem 32

We are ready to prove the main theorem stated in the introduction of this section. We will show that the graphical mean curvature flow of an area preserving map converges to an isometry in the positive case, to an affine map in the zero case, and to a minimal surface in the negative case. Recall that from Theorem 35, we already know that the norm of the second fundamental form stays uniformly bounded in time. Since

$$\nabla_{\partial_t} \Omega = - \int_M |H|^2 \Omega$$

and since the graphical flow exists for all time we have that there exists a time-independent constant  $C$ , such that

$$\int_0^\infty \left( \int_M |H|^2 \Omega \right) dt \leq C.$$

Therefore, there exists a sequence  $\{t_k\}_{k \in \mathbb{N}}$ , such that

$$\lim_{k \rightarrow \infty} \int_M |H|^2 \Omega = 0. \tag{47}$$

From Theorem 35, the norms of the second fundamental forms of the evolving submanifolds and their derivatives are uniformly bounded in time. Since the product manifold  $M \times N$  is compact, after passing to a subsequence of  $\{t_k\}_{k \in \mathbb{N}}$  if necessary, we deduce that the flow subconverges smoothly to a smooth surface  $M_\infty$  of  $M \times N$ ; see for example [11, Theorem 1.1]. From (47)  $M_\infty$  should be minimal. Due to a deep result of Simon [82], it follows that the flow converges smoothly and uniformly to a minimal surface  $M_\infty \subset M \times N$ . Additionally, we have the following situations:

**Area Preserving Case** Let us treat the case where the evolving maps are area preserving.

- (a) If  $\sigma > 0$ , then from Lemma 14(c), we have  $\vartheta \rightarrow 1$ , as  $t \rightarrow \infty$ . Therefore,  $M_\infty$  is the graph of an isometry  $f_\infty : M \rightarrow N$ .
- (b) If  $\sigma = 0$ , then from Lemma 14(c), we have that  $\vartheta \geq c_0 > 0$ . Hence, the surface  $M_\infty$  is the graph of a map  $f_\infty : M \rightarrow N$ . From Lemma 11 and the fact that  $2\sigma_\infty^\perp = |A_\infty|^2$ , we have

$$-\Delta \vartheta_\infty = 2|A_\infty|^2 \vartheta_\infty \geq 0.$$

By the strong maximum principle, we get  $|A_\infty|^2 = 0$ . Hence,  $M_\infty$  is totally geodesic.

**Strictly Area Decreasing Case** Assume that our maps are area decreasing.

- (a) Suppose that  $\sigma > 0$ . In this case the flow is smoothly converging to a graphical minimal surface  $M_\infty = \Gamma(f_\infty)$  of  $M \times N$ . Due to Theorem 13(a), the biggest singular value tends to zero as time goes to infinity. Hence,  $M_\infty$  must be totally geodesic and  $f_\infty$  is a constant map.
- (b) Assume that  $\sigma = 0$ . As in the previous case, we have smooth convergence of the flow to a minimal graphical surface  $M_\infty = \Gamma(f_\infty)$  of  $M \times N$ , where  $f_\infty$  is a strictly area decreasing map. Because of the formula

$$\partial_t \int_M \Omega = - \int_M |H|^2 \Omega \leq 0,$$

we obtain that

$$\int_M \Omega \leq \int_M \Omega_M = \text{constant}.$$

From Theorem 33(b), there is a non-negative constant  $C$  such that

$$\int_M |H|^2 \Omega \leq \frac{C}{t} \int_M \Omega \leq \frac{C}{t} \int_M \Omega.$$

Due to our assumptions we have  $u_2^2 \leq u_1^2 \leq 1$  and  $\min \sigma_M \geq 0 \geq \sup \sigma_N$ . Moreover, recall that

$$\Omega = \sqrt{(1 + \lambda^2)(1 + \mu^2)} \Omega_M = u_1^{-1} \Omega_M.$$

From the Gauss equation (8.3) and the Gauss-Bonnet formula we get

$$\begin{aligned} \int_M |A|^2 \Omega &= \int_M |H|^2 \Omega + 2 \int_M (\sigma_M u_1^2 + \sigma_N u_2^2) \Omega - 2 \int_M \sigma_{g(t)} \Omega \\ &\leq 2 \int_M \sigma_M u_1^2 \Omega - 2 \int_M \sigma_{g(t)} \Omega + \int_M |H|^2 \Omega \\ &\leq 2 \int_M \sigma_M u_1 \Omega - 2 \int_M \sigma_{g(t)} \Omega + \int_M |H|^2 \Omega \\ &\leq 2 \int_M \sigma_M \Omega_M - 2 \int_M \sigma_{g(t)} \Omega + \int_M |H|^2 \Omega = \int_M |H|^2 \Omega \\ &\leq Ct^{-1}, \end{aligned}$$

where  $C$  is a non-negative constant. Passing to the limit, we deduce that

$$\int_M |A_\infty|^2 \Omega_\infty = 0.$$



Thus,  $M_\infty = F_\infty(M)$  must be a totally geodesic graphical surface.

This completes the proof.  $\square$

## References

1. R. Aiyama, Totally real surfaces in the complex 2-space, in *Steps in Differential Geometry (Debrecen, 2000)*. Institute of Mathematical Information, Debrecen (2001), pp. 15–22
2. R. Aiyama, Lagrangian surfaces with circle symmetry in the complex two-space. *Michigan Math. J.* **52**, 491–506 (2004)
3. F. Almgren, Some interior regularity theorems for minimal surfaces and an extension of Bernstein's theorem. *Ann. Math.* **84**, 277–292 (1966)
4. B. Andrews, C. Baker, Mean curvature flow of pinched submanifolds to spheres. *J. Differ. Geom.* **85**, 357–395 (2010)
5. B. Andrews, C. Hopper, *The Ricci flow in Riemannian geometry, A complete proof of the differentiable 1/4-pinching sphere theorem*. Lecture Notes in Mathematics, vol. 2011 (Springer, Heidelberg, 2011)
6. R. Assimos, J. Jost, The geometry of maximum principles and a Bernstein theorem in codimension 2. arXiv:1811.09869 (2018)
7. C. Baker, *The mean curvature flow of submanifolds of high codimension*, PhD Thesis (Australian National University, Mathematical Sciences Institute, Australia, 2010). <https://arxiv.org/abs/1104.4409>
8. S. Bernstein, Über ein geometrisches Theorem und seine Anwendung auf die partiellen Differentialgleichungen vom elliptischen Typus. *Math. Z.* **26**, 551–558 (1927). (translation of the original version in *Comm. Soc. Math. Kharkov 2-'eme sér.* **15**, 38–45 (1915-1917))
9. C. Böhm, B. Wilking, Nonnegatively curved manifolds with finite fundamental groups admit metrics with positive Ricci curvature. *Geom. Funct. Anal.* **17**, 665–681 (2007)
10. E. Bombieri, E. De Giorgi, E. Giusti, Minimal cones and the Bernstein conjecture. *Invent. Math.* **7**, 243–268 (1969)
11. P. Breuning, Immersions with bounded second fundamental form. *J. Geom. Anal.* **25**, 1344–1386 (2015)
12. A. Chau, J. Chen, W. He, Lagrangian mean curvature flow for entire Lipschitz graphs. *Calc. Var. Partial Differ. Equations* **44**, 199–220 (2012)
13. A. Chau, J. Chen, Y. Yuan, Lagrangian mean curvature flow for entire Lipschitz graphs II. *Math. Ann.*, **357**, 165–183 (2013)
14. J. Cheeger, M. Gromov, M. Taylor, Finite propagation speed, kernel estimates for functions of the Laplace operator, and the geometry of complete Riemannian manifolds. *J. Differ. Geom.* **17**, 15–53 (1982)
15. J. Chen, W. He, A note on singular time of mean curvature flow. *Math. Z.*, **266**, 921–931 (2010)
16. J. Chen, J. Li, Mean curvature flow of surface in 4-manifolds. *Adv. Math.* **163**, 287–309 (2001)
17. J. Chen, J. Li, Singularity of mean curvature flow of Lagrangian submanifolds. *Invent. Math.* **156**, 25–51 (2004)
18. J. Chen, G. Tian, Moving symplectic curves in Kähler-Einstein surfaces. *Acta Math. Sin. (Engl. Ser.)* **16**, 541–548 (2000)
19. B.-Y. Chen, J.-M. Morvan, Géométrie des surfaces lagrangiennes de  $\mathbb{C}^2$ . *J. Math. Pures Appl.* **66**, 321–325 (1987)
20. S.-S. Chern, Minimal surfaces in an euclidean space of  $N$  dimensions, in *Differential and Combinatorial Topology* (Princeton University, Princeton, 1965), pp. 187–198

21. S.-S. Chern, Minimal submanifolds in a Riemannian manifold, in *University of Kansas, Department of Mathematics Technical Report*, vol. 19 (New Series) (University of Kansas Lawrence, Lawrence, 1968)
22. S.S. Chern, R. Osserman, Complete minimal surfaces in Euclidean  $n$ -space. *J. Anal. Math.* **19**, 15–34 (1967)
23. S.-S. Chern, J.G. Wolfson, Minimal surfaces by moving frames. *Am. J. Math.* **105**, 59–83 (1983)
24. B. Chow, P. Lu, L. Ni, *Hamilton's Ricci Flow* (American Mathematical Society/Science Press Beijing, Providence/New York, 2006)
25. B. Chow, S.-C. Chu, D. Glickenstein, C. Guenther, J. Isenberg, T. Ivey, D. Knopf, P. Lu, F. Luo, L. Ni, The Ricci flow: techniques and applications. Geometric aspects. Part I, in *Mathematical Surveys and Monographs*, vol. 135 (American Mathematical Society, Providence, RI, 2007)
26. B. Chow, S.-C. Chu, D. Glickenstein, C. Guenther, J. Isenberg, T. Ivey, D. Knopf, P. Lu, F. Luo, L. Ni, The Ricci flow: techniques and applications. Part II, in *Mathematical Surveys and Monographs*, vol. 144 (American Mathematical Society, Providence RI, 2008)
27. A. Cooper, *Mean Curvature Flow in Higher Codimension*, Thesis (Ph.D.) (Michigan State University, Michigan, 2011)
28. D. Crowley, T. Schick, The Gromoll filtration, KO-characteristic classes and metrics of positive scalar curvature. *Geom. Topol.* **17**, 1773–1789 (2013)
29. M. Dajczer, R. Tojeiro, *Submanifold Theory Beyond an Introduction* (Springer, New York, 2019)
30. E. De Giorgi, Una estensione del teorema di Bernstein. *Ann. Sci. Norm. Super. Pisa* **19**, 79–85 (1965)
31. D. DeTurck, Deforming metrics in the direction of their Ricci tensors. *J. Differ. Geom.* **18**, 157–162 (1983)
32. K. Ecker, *Regularity Theory for Mean Curvature Flow* (Birkhäuser Boston Inc., Boston, 2004)
33. K. Ecker, G. Huisken, Mean curvature evolution of entire graphs. *Ann. Math.* **130**, 453–471 (1989)
34. J. Eells, J. Sampson, Harmonic mappings of Riemannian manifolds. *Am. J. Math.* **86**, 109–160 (1964)
35. L.C. Evans, Partial differential equations, in *Graduate Studies in Mathematics*, vol. 19 (2010)
36. L.C. Evans, A strong maximum principle for parabolic systems in a convex set with arbitrary boundary. *Proc. Am. Math. Soc.* **138**, 3179–3185 (2010)
37. D. Fischer-Colbrie, Some rigidity theorems for minimal submanifolds of the sphere. *Acta Math.* **145**, 29–46 (1980)
38. M. Gromov, Homotopical effects of dilatation. *J. Differ. Geom.* **13**, 303–310 (1978)
39. M. Gromov, Pseudoholomorphic curves in symplectic manifolds. *Invent. Math.* **82**, 307–347 (1985)
40. M. Gromov, Carnot-Carathéodory spaces seen from within, in *Sub-Riemannian Geometry*. Program of Mathematical, vol. 144 (Birkhäuser, Basel, 1996), pp. 79–323
41. M. Gromov, *Quantitative homotopy theory*, in *Prospects in Mathematics* (American Mathematical Society, Providence, 1999), pp. 45–49
42. L. Guth, Contraction of areas vs. topology of mappings. *Geom. Funct. Anal.* **23**, 1804–1902 (2013)
43. R. Hamilton, The inverse function theorem of Nash and Moser. *Bull. Am. Math. Soc. (N.S.)* **7**, 65–222 (1982)
44. R. Hamilton, Three-manifolds with positive Ricci curvature. *J. Differ. Geom.* **17**, 255–306 (1982)
45. R. Hamilton, Four-manifolds with positive curvature operator. *J. Differ. Geom.* **24**, 153–179 (1986)
46. R. Hamilton, A compactness property for solutions of the Ricci flow. *Am. J. Math.* **117**, 545–572 (1995)

47. T. Hasanis, A. Savas-Halilaj, T. Vlachos, Minimal graphs in  $\mathbb{R}^4$  with bounded Jacobians. Proc. Am. Math. Soc. **137**, 3463–3471 (2009)
48. T. Hasanis, A. Savas-Halilaj, T. Vlachos, On the Jacobian of minimal graphs in  $\mathbb{R}^4$ . Bull. Lond. Math. Soc. **43**, 321–327 (2011)
49. A. Hatcher, A proof of the Smale conjecture,  $\text{Diff}(\mathbb{S}^3) \simeq \mathbb{O}(4)$ . Ann. Math. **117**, 553–607 (1983)
50. S. Hildebrandt, J. Jost, K.O. Widman, Harmonic mappings and minimal submanifolds. Invent. Math. **62**, 269–298 (1980)
51. E. Hopf, On S. Bernstein's theorem on surfaces  $z(x, y)$  of nonpositive curvature. Proc. Am. Math. Soc. **1**, 80–85 (1950)
52. G. Huisken, Flow by mean curvature of convex surfaces into spheres. J. Differ. Geom. **20**, 237–266 (1984)
53. G. Huisken, Asymptotic behavior for singularities of the mean curvature flow. J. Differ. Geom. **31**, 285–299 (1990)
54. J. Jost, *Compact Riemann Surfaces: An Introduction to Contemporary Mathematics* (Springer, Berlin, 2006)
55. J. Jost, *Riemannian Geometry and Geometric Analysis* Universitext edn., vol. 7 (Springer, Cham, 2017)
56. J. Jost, Y.-L. Xin, L. Yang, The Geometry of Grassmannian manifolds and Bernstein type theorems for higher codimension. Ann. Sci. Norm. Super. Pisa Cl. Sci. **16**, 1–39 (2016)
57. J. Kazdan, *Applications of Partial Differential Equations to Problems in Differential Geometry* (1983). <http://www.math.upenn.edu/~kazdan/japan/japan.pdf>
58. H.-B. Lawson, R. Osserman, Non-existence, non-uniqueness and irregularity of solutions to the minimal surface system. Acta Math. **139**, 1–17 (1977)
59. J.M. Lee, Riemannian manifolds, in *Graduate Texts in Mathematical*, vol. 176 (Springer, Berlin, 1997)
60. K.-W. Lee, Y.-I. Lee, Mean curvature flow of the graphs of maps between compact manifolds. Trans. Am. Math. Soc. **363**, 5745–5759 (2011)
61. J. Li, L. Yang, Symplectic mean curvature flows in Kähler surfaces with positive holomorphic sectional curvatures. Geom. Dedicata **170**, 63–69 (2014)
62. F. Lubbe, Curvature estimates for graphical mean curvature flow in higher codimension, in *Technische Informationsbibliothek und Universitätsbibliothek Hannover (TIB), Hannover* (2015)
63. C. Mantegazza, *Lecture Notes on Mean Curvature Flow* (Birkhäuser/Springer, Basel, 2011)
64. I. Medoš, M.-T. Wang, Deforming symplectomorphisms of complex projective spaces by the mean curvature flow. J. Differ. Geom. **87**, 309–341 (2011)
65. E.J. Mickle, A remark on a theorem of Serge Bernstein. Proc. Am. Math. Soc. **1**, 86–89 (1950)
66. J. Morgan, G. Tian, Ricci flow and the Poincaré conjecture, in *Clay Mathematics Monographs*, vol. 3 (American Mathematical Society/Clay Mathematics Institute, Providence/Cambridge, 2007)
67. J. Moser, On Harnack's theorem for elliptic differential equations. Comm. Pure Appl. Math. **14**, 577–591 (1961)
68. L. Ni, A Bernstein type theorem for minimal volume preserving maps. Proc. Am. Math. Soc. **130**, 1207–1210 (2002)
69. R. Osserman, Global properties of minimal surfaces in  $E^3$  and  $E^n$ . Ann. Math. (2) **80**, 340–364 (1964)
70. R. Osserman, *A Survey of Minimal Surfaces* (Van Nostrand-Reinhold, New York, 1969)
71. R. Osserman, The minimal surface equation, in *Seminar on Nonlinear Partial Differential Equations, Berkeley, California* (1983), pp. 237–259
72. M.H. Protter, H.F. Weinberger, *Maximum Principles in Differential Equations* (Springer, New York, 1984)
73. K. Reidemeister, Über die singulären Randpunkte eines konvexen Körpers. Math. Ann. **83**, 116–118 (1921)

74. E.A. Ruh, J. Vilms, The tension field of the Gauss map. *Trans. Am. Math. Soc.* **149**, 569–573 (1970)
75. A. Savas-Halilaj, K. Smoczyk, Homotopy of area decreasing maps by mean curvature flow. *Adv. Math.* **255**, 455–473 (2014)
76. A. Savas-Halilaj, K. Smoczyk, Bernstein theorems for length and area decreasing minimal maps. *Calc. Var. Partial Differ. Equations* **50**, 549–577 (2014)
77. A. Savas-Halilaj, K. Smoczyk, Evolution of contractions by mean curvature flow. *Math. Ann.* **361**, 725–740 (2015)
78. A. Savas-Halilaj, K. Smoczyk, Mean curvature flow of area decreasing maps between Riemann surfaces. *Ann. Global Anal. Geom.* **53**, 11–37 (2018)
79. R. Schoen, *The role of harmonic mappings in rigidity and deformation problems*. Lecture Notes in Pure and Application Mathematical, vol. 143 (Dekker, New York, 1993), pp. 179–200
80. J. Simons, Minimal varieties in Riemannian manifolds. *Ann. Math. (2)* **88**, 62–105 (1968)
81. L. Simon, A Hölder estimate for quasiconformal maps between surfaces in euclidean space. *Acta Math.* **139**, 19–51 (1977)
82. L. Simon, Asymptotics for a class of nonlinear evolution equations, with applications to geometric problems. *Ann. Math. (2)* **118**, 525–571 (1983)
83. D. Singley, Smoothness theorems for the principal curvatures and principal vectors of a hypersurface. *Rocky Mountain J. Math.* **5**, 135–144 (1975)
84. S. Smale, Diffeomorphisms of the 2-sphere. *Proc. Am. Math. Soc.* **10**, 621–626 (1959)
85. K. Smoczyk, Angle theorems for the Lagrangian mean curvature flow. *Math. Z.* **240**, 849–883 (2002)
86. K. Smoczyk, Mean curvature flow in higher codimension: introduction and survey, in *Global Differential Geometry, Springer Proceedings of the Mathematical*, vol. 17 (Springer, Heidelberg, 2012), pp. 231–274
87. K. Smoczyk, M.-T. Wang, Mean curvature flows of Lagrangians submanifolds with convex potentials. *J. Differ. Geom.* **62**, 243–257 (2002)
88. K. Smoczyk, M.-T. Wang, Generalized Lagrangian mean curvature flows in symplectic manifolds. *Asian J. Math.* **15**, 129–140 (2011)
89. K. Smoczyk, M.-P. Tsui, M.-T. Wang, Curvature decay estimates for mean curvature flow in higher codimensions. *Trans. Am. Math. Soc.* **368**, 7763–7775 (2016)
90. K. Smoczyk, M.-P. Tsui, M.-T. Wang, Generalized Lagrangian mean curvature flows: the cotangent bundle case. *J. Reine Angew. Math.* **750**, 97–121 (2019)
91. M.-P. Tsui, M.-T. Wang, Mean curvature flows and isotopy of maps between spheres. *Comm. Pure Appl. Math.* **57**, 1110–1126 (2004)
92. A. Vogiatzi, *Mean Curvature Flow and Isotopy of Problems*, M.S. Thesis (University of Ioannina, Ioannina, 2020), pp. 1–91
93. X. Wang, A remark on strong maximum principle for parabolic and elliptic systems. *Proc. Am. Math. Soc.* **109**, 343–348 (1990)
94. M.-T. Wang, Mean curvature flow of surfaces in Einstein four-manifolds. *J. Differ. Geom.* **57**, 301–338 (2001)
95. M.-T. Wang, Deforming area preserving diffeomorphism of surfaces by mean curvature flow. *Math. Res. Lett.* **8**, 651–661 (2001)
96. M.-T. Wang, Long-time existence and convergence of graphic mean curvature flow in arbitrary codimension. *Invent. Math.* **148**, 525–543 (2002)
97. M.-T. Wang, On graphic Bernstein type results in higher codimension. *Trans. Am. Math. Soc.* **355**, 265–271 (2003)
98. M.-T. Wang, Subsets of Grassmannians preserved by mean curvature flow. *Comm. Anal. Geom.* **13**, 981–998 (2005)
99. M.-T. Wang, *Lectures on mean curvature flows in higher codimensions*. Handbook of Geometric Analysis, vol. 1 (2008), pp. 525–543
100. M.-T. Wang, Mean curvature flows and isotopy problems, in *Surveys in Differential Geometry. Geometry and Topology* (International Press, Somerville, 2013), pp. 227–235

101. H.F. Weinberger, Invariant sets for weakly coupled parabolic and elliptic systems. *Rend. Mat.* (6) **8**, 295–310 (1975)
102. Y.-L. Xin, Geometry of harmonic maps, in *Progress in Nonlinear Differential Equations and their Applications*, vol. 23 (Birkhäuser Boston, Inc., Boston, 1996)

# Critical Point Theory in Infinite Dimensional Spaces Using the Leray–Schauder Index



Martin Schechter

**Abstract** Many problems arising in science and engineering call for the solving of the Euler–Lagrange equations of functionals. Thus, solving the Euler–Lagrange equations is tantamount to finding critical points of the corresponding functional. An idea that has been very successful is to find appropriate sets that sandwich the functional. This means that the functional is bounded from above on one of the sets and bounded from below on the other. Two sets of the space are said to form a sandwich if they produce a critical sequence whenever they sandwich a functional. If the critical sequence has a convergent subsequence, then that produces a critical point. Finding sets that sandwich a functional is quite easy, but determining whether or not the sets form a sandwich is quite another story. It appears that the only way we can check to see if two sets form a sandwich is to require that one of them be contained in a finite-dimensional subspace. The reason is that in order to verify the definition, we need to invoke the Brouwer fixed point theorem. Our aim is to find a counterpart that holds true when both sets are infinite dimensional. We adjust our definitions to accommodate infinite dimensions. These definitions reduce to the usual when one set is finite dimensional. In order to prove the corresponding theorems, we make adjustments to the topology of the space and introduce infinite dimensional splitting. This allows us to use a form of compactness on infinite dimensional subspaces that does not exist under the usual topology. We lose the Brouwer index, but we are able to replace it with the Leray–Schauder index. We carry out the details in Sections 5, 7, and 8. In Section 6 we solve a system of equations which require infinite dimensional splitting.

**AMS Subject Classification** Primary 35J35, 47J30, 49J35, 49J40, 58K05

---

M. Schechter (✉)

Department of Mathematics, University of California, Irvine, CA, USA  
e-mail: [mschecht@math.uci.edu](mailto:mschecht@math.uci.edu)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,  
Springer Optimization and Its Applications 173,  
[https://doi.org/10.1007/978-3-030-72563-1\\_21](https://doi.org/10.1007/978-3-030-72563-1_21)

579

## 1 Introduction

Many problems arising in science and engineering call for the solving of the Euler–Lagrange equations of functionals, i.e., equations equivalent to

$$G'(u) = 0, \tag{1}$$

where  $G(u)$  is a  $C^1$ -functional (usually representing the energy) arising from the given data. By this we mean that functions are solutions of the Euler–Lagrange equations of  $G$  iff they satisfy (1). Solutions of (1) are called **critical points** of  $G$ . Thus, solving the Euler–Lagrange equations is tantamount to finding critical points of the corresponding functional.

The variational approach to solving differential equations and systems has its roots in the calculus of variations. The original problem was to minimize or maximize a given functional. The approach was to obtain the Euler–Lagrange equations of the functional, solve them, and show that the solutions provided the required minimum or maximum. This worked well for one-dimensional problems. However, when it came to higher dimensions, it was recognized quite early that it was more difficult to solve the Euler–Lagrange equations than it was to find minima or maxima of the corresponding functional. Consequently, the approach was abandoned for many years.

Eventually, when nonlinear partial differential equations and systems arose in applications and people were searching for solutions, they began to check if the equations and systems were the Euler–Lagrange equations of functionals. If so, a natural approach would be to find critical points of the corresponding functionals. The problem is that there is no uniform way of finding them.

The initial approach to finding critical points is to look for maxima or minima. Global extrema are the easiest to find, but they can exist only if the functional is semi bounded. For instance, if the continuously differential functional  $G$  is bounded from below, then we can find a minimizing sequence  $\{u_k\}$  such that

$$G(u_k) \rightarrow a = \inf G > -\infty. \tag{2}$$

If this series converges or has a convergent subsequence, we have a minimum.

However, if the functional  $G$  is bounded from below, it can be shown that there is a sequence satisfying

$$G(u_k) \rightarrow a, \quad G'(u_k) \rightarrow 0. \tag{3}$$

If the sequence has a convergent subsequence, this will produce a minimum. The gain is that a sequence satisfying (3) has a better chance of having a convergent subsequence than a sequence satisfying only (2).

When the functional is not semibounded, there is no clear way of obtaining critical points. In general, one would like to determine when a functional has a

critical sequence, i.e., a sequence satisfying

$$G(u_k) \rightarrow a, \quad G'(u_k) \rightarrow 0. \quad (4)$$

This would give one the same advantages that one has in the case of semi-bounded functionals.

## 2 Sandwich Sets

An idea that has been very successful is to find appropriate sets that **sandwich the functional**. By this we mean the following: Two sets  $A, B$  **sandwich** the functional  $G(u)$  if  $G(u)$  is bounded from below on one of them and bounded from above on the other, e.g., if

$$a_0 := \sup_A G < \infty, \quad b_0 := \inf_B G > -\infty. \quad (5)$$

We would like to find sets  $A$  and  $B$  such that (5) will imply

$$\exists u : G(u) \geq b_0, \quad G'(u) = 0. \quad (6)$$

This is too much to expect, since even semi-boundedness alone does not imply the existence of a critical point. Consequently, we weaken our requirements and look for sets  $A, B$  such that (5) implies the existence of a critical sequence (4) with  $a \geq b_0$ . This leads to

**Definition 1** We shall say that the set  $A$  **forms a sandwich** with the set  $B$  if (5) implies (4) with  $a \geq b_0$  for every  $C^1$ -functional  $G(u)$ .

Of course, (4) is a far cry from (6), but if, e.g., the sequence (4) has a convergent subsequence, then (4) implies (6). Whether or not this is true depends on the functional  $G(u)$ .

## 3 The Finite Dimensional Case

It appears that the only way we can check to see if two sets form a sandwich, is to require that one of them is contained in a finite-dimensional subspace. The reason is that in order to verify the definition, we need to invoke the Brouwer fixed point theorem.

The following three results hold when a subspace  $N$  is finite dimensional.

**Theorem 2** *Let  $N$  be a finite dimensional subspace of a Banach space  $E$ , and for each  $R > R_0$ , let  $\Omega_R(p)$  be an open bounded set in  $N$  containing a point  $p$  such*



that  $d(\partial\Omega_R(p), p) \rightarrow \infty$  as  $R \rightarrow \infty$ . Let  $F$  be a continuous map of  $E$  onto  $N$  such that  $F = I$  on  $N$ . Assume also that

$$d(A_R, F^{-1}(p)) \rightarrow \infty, \quad R \rightarrow \infty,$$

where  $A_R = \partial\Omega_R(p)$ . Let  $G$  be a  $C^1$ -functional on  $E$  such that

$$-\infty < b_0 = \inf_B G, \quad \sup_{A^R} G \leq a_0 < \infty, \tag{7}$$

for  $R > R_0$ , where  $A^R = N \setminus \Omega_R(p)$  and  $B = F^{-1}(p)$ . Then for each  $\rho \in \mathcal{Q}$  there is a sequence  $\{u_k\} \subset E$  such that

$$G(u_k) \rightarrow c, \quad b_0 \leq c \leq a_0, \quad \rho(d(u_k, B))\|G'(u_k)\| \rightarrow 0. \tag{8}$$

**Corollary 3** Let  $N$  be a finite dimensional subspace of a Hilbert space  $E$  and let  $M = N^\perp$ . For  $G \in C^1(E, \mathbb{R})$ , assume

$$a_0 = \sup_N G < \infty, \quad b_0 = \inf_M G > -\infty. \tag{9}$$

Then for each  $\rho \in \mathcal{Q}$  there is a sequence  $\{u_k\} \subset E$  such that

$$G(u_k) \rightarrow c, \quad b_0 \leq c \leq a_0, \quad \rho(d(u_k, M))\|G'(u_k)\| \rightarrow 0. \tag{10}$$

**Theorem 4** Let  $N$  be a finite dimensional subspace of a Hilbert space  $E$  with complement  $M \oplus \{v_0\}$ , where  $v_0$  is an element in  $E$  having unit norm, and let  $\delta$  be any positive number. Let  $\varphi(t) \in C^1(\mathbb{R})$  be such that

$$0 \leq \varphi(t) \leq 1, \quad \varphi(0) = 1,$$

and

$$\varphi(t) = 0, \quad |t| \geq 1.$$

Let

$$F(v+w+sv_0) = v+[s+\delta-\delta\varphi(\|w\|^2/\delta^2)]v_0, \quad w \in M, \quad v \in N, \quad s \in \mathbb{R}. \tag{11}$$

Let  $G$  be a  $C^1$ -functional on  $E$  such that (5) holds with  $A = [N \oplus \{v_0\}] \setminus \mathcal{B}_{R_0}$  and  $B = F^{-1}(\delta v_0) = \{w + rv_0 : w \in M, r = \delta\varphi(\|w\|^2/\delta^2)\}$ . Then for each  $\rho \in \mathcal{Q}$  there is a sequence  $\{u_k\} \subset E$  such that

$$G(u_k) \rightarrow c, \quad b_0 \leq c \leq a_0, \quad \rho(d(u_k, B))\|G'(u_k)\| \rightarrow 0. \tag{12}$$

However, there are many applications for which we would like to obtain critical points in infinite-dimensional situations. It is not obvious how to proceed. It is not clear that we can obtain similar results for such situations. We now describe one method that works in the infinite-dimensional case. It involves adjusting the topology of the underlying space. In order to deal with the infinite dimensional situation, we are forced to make several adjustments. We must restrict the functional  $G$  and the mapping  $F$  and use a more general index for verification. The restrictions that we have chosen do not apply when one of the subspaces is finite dimensional. While we could not use the Brouwer index, we were able to use the Leray–Schauder index in its place. Thus the new theorems that we prove reverts to the old theorems when one of the subsets is finite dimensional. For this purpose we use flows described in the next section.

### 4 Flows

Let  $\mathcal{Q}$  be a set of positive functions  $\rho(t)$  on  $[0, \infty)$ , which are

- (a) locally Lipschitz continuous,
- (b) nondecreasing
- (c) satisfy

$$\int_0^\infty \frac{dt}{\rho(t)} = \infty. \tag{13}$$

Moreover,  $\mathcal{Q}$  is to satisfy

$$\rho_1, \rho_2 \in \mathcal{Q} \implies \max(\rho_1, \rho_2) \in \mathcal{Q},$$

and contain functions of the form

$$(1 + |t|)^\beta, \quad \beta \in \mathbb{R}.$$

Let  $Q \neq \emptyset$  be a subset of a Banach space  $E$ , and let  $\Sigma_Q$  be the set of all continuous maps  $\sigma = \sigma(t)$  from  $E \times [0, 1]$  to  $E$  such that

- 1.  $\sigma(0)$  is the identity map,
- 2. for each  $t \in [0, 1]$ ,  $\sigma(t)$  is a homeomorphism of  $E$  onto  $E$ ,
- 3.  $\sigma'(t)$  is piecewise continuous and satisfies

$$\|\sigma'(t)u\| \leq C\rho(d(\sigma(t)u, Q)), \quad u \in E, \tag{14}$$

for some  $\rho \in \mathcal{Q}$ . If  $Q = \{0\}$ , we write  $\Sigma = \Sigma_Q$ . The mappings in  $\Sigma_Q$  are called *flows*. We note the following.

*Remark 5* If  $\sigma_1, \sigma_2$  are in  $\Sigma_Q$ , define  $\sigma_3 = \sigma_1 \circ \sigma_2$  by

$$\sigma_3(s) = \begin{cases} \sigma_1(2s), & 0 \leq s \leq \frac{1}{2}, \\ \sigma_2(2s - 1)\sigma_1(1), & \frac{1}{2} < s \leq 1. \end{cases}$$

Then  $\sigma_3 \in \Sigma_Q$ , and  $\sigma_3(1) = \sigma_2(1)\sigma_1(1)$ .

**Proof** The first two properties are obvious. To check the third, note that

$$\sigma'_3(s) = \begin{cases} 2\sigma'_1(2s), & 0 \leq s \leq (\frac{1}{2})_-, \\ 2\sigma'_2(2s - 1)\sigma_1(1), & (\frac{1}{2})_+ \leq s \leq 1. \end{cases}$$

Thus, if

$$\|\sigma'_i(t)u\| \leq C_i \rho_i(d(\sigma_i(t)u, Q)), \quad u \in E, \quad i = 1, 2, \tag{15}$$

then

$$\|\sigma'_3(s)u\| \leq \begin{cases} 2\|\sigma'_1(2s)u\|, & 0 \leq s \leq (\frac{1}{2})_-, \\ 2\|\sigma'_2(2s - 1)\sigma_1(1)u\|, & (\frac{1}{2})_+ \leq s \leq 1, \end{cases}$$

or

$$\|\sigma'_3(s)u\| \leq \begin{cases} 2C_1\rho(d(\sigma_3(s)u, Q)), & 0 \leq s \leq (\frac{1}{2})_-, \\ 2C_2\rho(d(\sigma_3(s)u, Q)), & (\frac{1}{2})_+ \leq s \leq 1, \end{cases}$$

where  $\rho = \max(\rho_1, \rho_2)$ . We can now take  $C_3 = 2 \max(C_1, C_2)$ . □

The following theorem can be found in [21, 22].

**Theorem 6** *Let  $g(t, x)$  be a continuous map from  $\mathbb{R} \times X$  to  $X$ , where  $X$  is a Banach space. Assume that for each point  $(\hat{t}, \hat{x}) \in \mathbb{R} \times X$ , there are constants  $K, b > 0$  such that*

$$\|g(t, x) - g(t, y)\| \leq K\|x - y\|, \quad |t - \hat{t}| < b, \quad \|x - \hat{x}\| < b, \quad \|y - \hat{x}\| < b, \tag{16}$$

and

$$\|g(t, x)\| \leq \rho(d(x, Q)), \quad x \in X, \quad t \in [t_0, \infty), \tag{17}$$

where  $Q$  is a subset of  $X$ , with  $\rho$  nondecreasing or bounded. Assume

$$\int_0^{u_0} \frac{d\tau}{\rho(\tau)} = \int_{u_0}^\infty \frac{d\tau}{\rho(\tau)} = \infty \tag{18}$$

where  $u_0 > 0$ . Then for each  $x_0 \in X$  and  $t_0 \in \mathbb{R}$  there is a unique solution  $x(t)$  of the equation

$$\frac{dx(t)}{dt} = g(t, x(t)), \quad t \in [t_0, \infty), \quad x(t_0) = x_0. \tag{19}$$

Moreover,  $x(t)$  depends continuously on  $x_0$  and satisfies

$$u_1(t) \leq d(x(t), Q) \leq u_2(t), \quad t \in [t_0, \infty), \tag{20}$$

where  $u_1(t)$  is the solution of

$$u'(t) = -\rho(u(t)) \tag{21}$$

in  $[t_0, \infty)$  satisfying  $u(t_0) = u_0 = d(x_0, Q)$ , and  $u_2(t)$  is the solution of

$$u'(t) = \rho(u(t)) \tag{22}$$

in  $[t_0, \infty)$  satisfying  $u(t_0) = u_0 = d(x_0, Q)$ .

## 5 The Infinite Dimensional Case

We now describe a sandwich theory that works in the infinite-dimensional case.

Let  $N$  be a closed, separable subspace of a Hilbert space  $E$ . We can define a new norm  $|v|_w$  satisfying  $|v|_w \leq \|v\| \ \forall v \in E$  and such that the topology induced by this norm is equivalent to the weak topology of  $N$  on bounded subsets of  $N$ . We construct the norm so that  $v_j \rightarrow v$  weakly in  $N$  implies  $|v_j - v|_w \rightarrow 0$ . Conversely, if  $\|v_j\|, \|v\| \leq C$  for all  $j > 0$  and  $|v_j - v|_w \rightarrow 0$ , then  $v_j \rightarrow v$  weakly in  $N$ .

We adjust our assumptions on  $G$  and  $F$  for the infinite dimensional case of  $\dim N = \infty$ , but they reduce to the same assumptions that are made in the finite dimensional case  $\dim N < \infty$ .

Our requirements on  $G$  are given by

**Definition 7** Let  $N$  be a closed separable subspace of a Hilbert space  $E$ . A functional  $G(u)$  on  $E$  will be called an  $N$ -weak-to-weak continuously differentiable functional on  $E$  if

$$|v_n - v|_w \rightarrow 0 \tag{23}$$

implies that there is a renamed subsequence satisfying

$$G(v_n) \rightarrow G(v), \ |G'(v_n) - G'(v)|_w \rightarrow 0. \tag{24}$$

This means that  $G$  is a continuous functional on  $E_w$ , continuously differentiable on  $E$  and such that

$$v_n = Pu_n \rightarrow v \text{ weakly in } E, \quad w_n = (I - P)u_n \rightarrow w \text{ strongly in } E \tag{25}$$

implies that there is a renamed subsequence satisfying

$$G'(v_n + w_n) \rightarrow G'(v + w) \text{ weakly in } E, \tag{26}$$

where  $P$  is the projection of  $E$  onto  $N$ .

Note that every  $C'$  functional is  $N$ -weak-to-weak continuously differentiable when  $\dim N < \infty$ .

Concerning the mapping  $F$  we define

**Definition 8** Let  $N$  be a closed separable subspace of a Hilbert space  $E$ . We shall call a map  $F$  of  $E$  onto  $N$  an  $N$ -weakly continuous mapping if  $F$  is a  $|\cdot|_w$ -continuous map from  $E$  onto  $N$  satisfying

- $F_N = I$ ;  $F$  maps any finite dimensional subspace of  $N$  containing  $p$  into itself; it maps bounded sets into bounded sets;
- There exists a fixed finite-dimensional subspace  $E_0$  of  $E$  such that  $F(u - v) - (F(u) - F(v)) \in E_0, \forall u, v \in E$ ;
- $F$  maps finite-dimensional subspaces of  $E$  to finite-dimensional subspaces of  $E$ ;

Note that every continuous map  $F$  of  $E$  onto  $N$  satisfying  $F_N = I$  is  $N$ -weakly continuous when  $N$  is finite dimensional.

We have

**Theorem 9** Let  $N$  be a closed separable subspace of a Hilbert space  $E$ . For each  $R > R_0$ , let  $\Omega_R(p)$  be an open, convex, bounded set in  $N$  containing a point  $p$  such that  $d(\partial\Omega_R(p), p) \rightarrow \infty$  as  $R \rightarrow \infty$ . Let  $F$  be a  $N$ -weakly continuous mapping of  $E$  onto  $N$  such that  $F = I$  on  $N$ . Assume also that

$$d(A_R, F^{-1}(p)) \rightarrow \infty, \quad R \rightarrow \infty,$$

where  $A_R = \partial\Omega_R(p)$ . Let  $G$  be a  $N$ -weak-to-weak continuously differentiable functional on  $E$  such that

$$-\infty < b_0 = \inf_B G, \quad \sup_{A^R} G \leq a_0 < \infty, \tag{27}$$

for  $R > R_0$ , where  $A^R = N \setminus \Omega_R(p)$  and  $B = F^{-1}(p)$ . Then for each  $\rho \in \mathcal{Q}$  there is a sequence  $\{u_k\} \subset E$  such that

$$G(u_k) \rightarrow c, \quad b_0 \leq c \leq a_0, \quad \rho(d_w(u_k, B)) \|G'(u_k)\| \rightarrow 0. \tag{28}$$

Note that when  $\dim N < \infty$ , Theorem 9 reduces to Theorem 2.

**Corollary 10** *Let  $N$  be a closed separable subspace of a Hilbert space  $E$ , and let  $M = N^\perp$ . For  $G$  a  $N$ -weak-to-weak continuously differentiable functional on  $E$  assume*

$$a_0 = \sup_N G < \infty, \quad b_0 = \inf_M G > -\infty. \tag{29}$$

Then for each  $\rho \in \mathcal{Q}$  there is a sequence  $\{u_k\} \subset E$  such that

$$G(u_k) \rightarrow c, \quad b_0 \leq c \leq a_0, \quad \rho(d_w(u_k, M))\|G'(u_k)\| \rightarrow 0. \tag{30}$$

Note that when  $\dim N < \infty$ , Corollary 10 reduces to Corollary 3.

**Corollary 11** *Let  $N$  be a closed, separable subspace of a Banach space  $E$ , and for each  $R > R_0$  let  $\Omega_R(p)$  be an open, convex, bounded set in  $N$  containing a point  $p$ . Let  $G$  be a  $N$ -weak-to-weak continuously differentiable functional on  $E$ , and let  $F$  be an  $N$ -weakly continuous mapping. Assume*

$$d_R = d_w(A^R, F^{-1}(p)) \rightarrow \infty, \quad R \rightarrow \infty,$$

where  $A^R = N \setminus \Omega_R(p)$ . Assume

$$-\infty < b_0 = \inf_B G, \quad \sup_{A^R} G \leq a_0 < \infty, \tag{31}$$

for  $R > R_0$ . Then for each sequence  $v_k \geq 2d_k - \varepsilon$  there is a  $\beta > 0$  and a sequence  $\{u_k\} \subset E$  such that

$$G(u_k) \rightarrow c, \quad b_0 \leq c \leq b_1, \quad (v_k + |u_k|_w)\|G'(u_k)\| \leq \beta. \tag{32}$$

**Theorem 12** *Let  $N$  be a closed separable subspace of a Hilbert space  $E$  with complement  $M \oplus \{v_0\}$ , where  $v_0$  is an element in  $E$  having unit norm, and let  $\delta$  be any positive number. Let  $\varphi(t) \in C^1(\mathbb{R})$  be such that*

$$0 \leq \varphi(t) \leq 1, \quad \varphi(0) = 1,$$

and

$$\varphi(t) = 0, \quad |t| \geq 1.$$

Let

$$F(v+w+sv_0) = v + [s + \delta - \delta\varphi(\|w\|^2/\delta^2)]v_0, \quad w \in M, \quad v \in N, \quad s \in \mathbb{R}. \tag{33}$$

Let  $G$  be a  $N$ -weak-to-weak continuously differentiable functional on  $E$  such that (5) holds with  $A = [N \oplus \{v_0\}] \setminus \mathcal{B}_{R_0}$  and  $B = F^{-1}(\delta v_0) = \{w + rv_0 : w \in M, r = \delta\varphi(\|w\|^2/\delta^2)\}$ . Then for each  $\rho \in \mathcal{Q}$  there is a sequence  $\{u_k\} \subset E$  such that

$$G(u_k) \rightarrow c, \quad b_0 \leq c \leq a_0, \quad \rho(d_w(u_k, B))\|G'(u_k)\| \rightarrow 0. \tag{34}$$

Note that when  $\dim N < \infty$ , Theorem 12 reduces to Theorem 4.

**Corollary 13** *Let  $N$  be a closed, separable subspace of a Hilbert space  $E$  with orthogonal complement  $M \oplus \{v_0\}$ , where  $v_0$  is an element in  $E$  having unit norm and orthogonal to both  $M$  and  $N$ , and let  $\delta < R$  be positive numbers. Let  $\varphi(t) \in C^1(\mathbb{R})$  be such that*

$$0 \leq \varphi(t) \leq 1, \quad \varphi(0) = 1,$$

and

$$\varphi(t) = 0, \quad |t| \geq 1.$$

Let

$$F(v+w+sv_0) = v + [s + \delta - \delta\varphi(\|w\|^2/\delta^2)]v_0, \quad w \in M, \quad v \in N, \quad s \in \mathbb{R}. \tag{35}$$

Let  $G$  be a an  $N$ -weak-to-weak continuously differentiable functional on  $E$ . Assume

$$-\infty < b_0 = \inf_B G, \quad \sup_A G = b_1 < \infty, \tag{36}$$

holds with  $A = [N \oplus \{v_0\}]$ ,  $A_R = A \cap \partial\mathcal{B}_R$  and  $B = F^{-1}(\delta v_0) = \{w + rv_0 : w \in M, r = \delta\varphi(\|w\|^2/\delta^2)\}$ . Then for each sequence  $v_k \rightarrow \infty$  there is a  $\beta > 0$  and a sequence  $\{u_k\} \subset E$  such that

$$G(u_k) \rightarrow c, \quad b_0 \leq c \leq b_1, \quad (v_k + |u_k|_w)\|G'(u_k)\| \leq \beta. \tag{37}$$

## 6 Applications

Let  $\mathcal{A}, \mathcal{B}$  be positive, self-adjoint operators on  $L^2(\Omega)$  with compact resolvents, where  $\Omega \subset \mathbb{R}^n$  ( $\Omega$  may be unbounded). Let  $F(x, v, w)$  be a Caratheodory function on  $\Omega \times \mathbb{R}^2$  such that

$$f(x, v, w) = \partial F/\partial v, \quad g(x, v, w) = \partial F/\partial w \tag{38}$$

are also Caratheodory functions satisfying

$$|f(x, v, w)| + |g(x, v, w)| \leq C_0(|v| + |w| + 1), \quad v, w \in \mathbb{R}. \tag{39}$$

We wish to solve the system

$$\mathcal{A}v = -f(x, v, w) \tag{40}$$

$$\mathcal{B}w = g(x, v, w). \tag{41}$$

The reason for the minus sign is because it leads to a variational problem involving infinite dimensional subspaces. The plus sign leads to a minimization problem.

Let  $\lambda_0(\mu_0)$  be the lowest eigenvalue of  $\mathcal{A}(\mathcal{B})$ .

**Theorem 14** *Assume*

$$2F(x, s, 0) \geq -\lambda_0 s^2 - W_1(x), \quad x \in \Omega, s \in \mathbb{R}, \tag{42}$$

where  $W_1(x) \in L^1(\Omega)$ . In addition, assume that the eigenfunctions of  $\lambda_0$  and  $\mu_0$  are bounded and  $\neq 0$  a.e. in  $\Omega$ , and there is a  $q > 2$  such that

$$\|w\|_q^2 \leq Cb(w), \quad w \in M. \tag{43}$$

*Assume*

$$2F(x, 0, t) \leq \mu(x)t^2, \quad x \in \Omega, \quad t \in \mathbb{R} \tag{44}$$

where

$$\mu(x) \leq \neq \mu_0, \quad x \in \Omega, \tag{45}$$

and for some  $\delta > 0$ ,

$$2F(x, s, t) \leq \mu_0 t^2 - \lambda_0 s^2, \quad |t| + |s| \leq \delta. \tag{46}$$

Also

$$H(x, s, t) \leq W(x) \tag{47}$$

and

$$H(x, s, t) \rightarrow -\infty \text{ as } |s| + |t| \rightarrow \infty, \tag{48}$$

where  $W(x) \in L^2(\mathbb{R}^n)$  and

$$H(x, s, t) = f(x, s, t)s + g(x, s, t)t - 2F(x, s, t). \tag{49}$$



Then the system (40) (41) has a nontrivial solution.

**Proof** Let  $D = D(\mathcal{A}^{1/2}) \times D(\mathcal{B}^{1/2})$ . Then  $D$  becomes a Hilbert space with norm given by

$$\|u\|_D^2 = (\mathcal{A}v, v) + (\mathcal{B}w, w), \quad u = (v, w) \in D. \tag{50}$$

We define

$$G(u) = b(w) - a(v) - 2 \int_{\Omega} F(x, v, w) dx, \quad u \in D \tag{51}$$

where

$$a(v) = (\mathcal{A}v, v), \quad b(w) = (\mathcal{B}w, w). \tag{52}$$

Then  $G \in C^1(D, \mathbb{R})$  and

$$(G'(u), h)/2 = b(w, h_2) - a(v, h_1) - (f(u), h_1) - (g(u), h_2), \tag{53}$$

where we write  $f(u), g(u)$  in place of  $f(x, v, w), g(x, v, w)$ , respectively. It is readily seen that the system (40), (41) is equivalent to

$$G'(u) = 0. \tag{54}$$

We let  $N$  be the set of those  $(v, 0) \in D$  and  $M$  the set of those  $(0, w) \in D$ . Then  $M, N$  are orthogonal closed subspaces such that

$$D = M \oplus N. \tag{55}$$

If we define

$$Lu = 2(-v, w), \quad u = (v, w) \in D \tag{56}$$

then  $L$  is a selfadjoint bounded operator on  $D$ . Also

$$G'(u) = Lu + c_0(u) \tag{57}$$

where

$$c_0(u) = -(\mathcal{A}^{-1/2} f(u), \mathcal{B}^{-1/2} g(u)) \tag{58}$$

is compact on  $D$ . This follows from (39) and the fact that  $\mathcal{A}$  and  $\mathcal{B}$  have compact resolvents. It also follows that  $G'$  has  $D - weak - to - weak$  continuity. For if  $u_k \rightarrow u$  weakly, then  $Lu_k \rightarrow Lu$  weakly and  $c_0(u_k)$  has a convergent subsequence.

Let  $N'$  be the orthogonal complement of  $N_0 = \{\varphi_0\}$  in  $N$ , where  $\varphi_0$  is the eigenfunction of  $\mathcal{A}$  corresponding to  $\lambda_0$ . Then  $N = N' \oplus N_0$ . Let  $M_0$  be the subspace of  $M$  spanned by the eigenfunctions of  $\mathcal{B}$  corresponding to  $\mu_0$ , and let  $M'$  be its orthogonal complement in  $M$ . Since  $N_0$  and  $M_0$  are contained in  $L^\infty(\Omega)$ , there is a positive constant  $\rho$  such that

$$a(y) \leq \rho^2 \Rightarrow \|y\|_\infty \leq \delta/4, \quad y \in N_0 \tag{59}$$

$$b(h) \leq \rho^2 \Rightarrow \|h\|_\infty \leq \delta/4, \quad h \in M_0 \tag{60}$$

where  $\delta$  is the number given in (46). If

$$a(y) \leq \rho^2, b(w) \leq \rho^2, |y(x)| + |w(x)| \geq \delta \tag{61}$$

we write  $w = h + w', h \in M_0, w' \in M'$  and

$$\delta \leq |y(x)| + |w(x)| \leq |y(x)| + |h(x)| + |w'(x)| \leq (\delta/2) + |w'(x)|. \tag{62}$$

Thus

$$|y(x)| + |h(x)| \leq \delta/2 \leq |w'(x)| \tag{63}$$

and

$$|y(x)| + |w(x)| \leq 2|w'(x)|. \tag{64}$$

Now by (46) and (64)

$$\begin{aligned} G(y, w) &= b(w) - a(y) - 2 \int_{\Omega} F(x, y, w) dx \tag{65} \\ &\geq b(w) - a(y) - \int_{|y|+|w|<\delta} \{\mu_0 w^2 - \lambda_0 y^2\} dx \\ &\quad - c_0 \int_{|y|+|w|>\delta} (|y| + |w| + 1)^2 dx \\ &\geq b(w) - a(y) - \mu_0 \|w\|^2 + \lambda_0 \|y\|^2 - c_1 \int_{2|w'|>\delta} |w'|^q dx \\ &\geq b(w') - \mu_0 \|w'\|^2 - c_2 b(w')^{q/2} \\ &\geq \left(1 - \frac{\mu_0}{\mu_1} - c_2 b(w')^{(q/2)-1}\right) b(w'), \quad a(y) \leq \rho^2, b(w) \leq \rho^2 \end{aligned}$$

where  $\mu_1$  is the next eigenvalue of  $\mathcal{B}$  after  $\mu_0$ . If we reduce  $\rho$  accordingly, we can find a positive constant  $\nu$  such that

$$G(y, w) \geq \nu b(w'), \quad a(y) \leq \rho^2, \quad b(w) \leq \rho^2. \tag{66}$$

I claim that either (40) (41) has a nontrivial solution or there is an  $\epsilon > 0$  such that

$$G(y, w) \geq \epsilon, \quad a(y) + b(w) = \rho^2. \tag{67}$$

For suppose (67) did not hold. Then there would be a sequence  $\{y_k, w_k\}$  such that  $a(y_k) + b(w_k) = \rho^2$  and  $G(y_k, w_k) \rightarrow 0$ . If we write  $w_k = w'_k + h_k$ ,  $w'_k \in M'$ ,  $h_k \in M_0$ , then (66) tells us that  $b(w'_k) \rightarrow 0$ . Thus  $a(y_k) + b(h_k) \rightarrow \rho^2$ . Since  $N_0, M_0$  are finite dimensional, there is a renamed subsequence such that  $y_k \rightarrow y$  in  $N_0$  and  $h_k \rightarrow h$  in  $M_0$ . Thus  $G(y, h) = 0$ . By (59) and (60),  $\|y\|_\infty \leq \delta/4$  and  $\|h\|_\infty \leq \delta/4$ . Consequently (46) implies

$$2F(x, y, h) \leq \mu_0 h^2 - \lambda_0 y^2. \tag{68}$$

Since

$$G(y, h) = b(h) - a(y) - 2 \int_\Omega F(x, y, h) dx = 0, \tag{69}$$

we have

$$\int_\Omega \{2F(x, y, h) + \lambda_0 y^2 - \mu_0 h^2\} dx = 0. \tag{70}$$

In view of (68), this implies

$$2F(x, y, h) \equiv \mu_0 h^2 - \lambda_0 y^2. \tag{71}$$

For  $\zeta \in C_0^\infty(\Omega)$  and  $t > 0$  small we have

$$2[F(x, y + t\zeta, h) - F(x, y, h)]/t \leq -\lambda_0[(y + t\zeta)^2 - y^2]/t. \tag{72}$$

Taking  $t \rightarrow 0$ , we have

$$f(x, y, h)\zeta \leq -\lambda_0 y\zeta. \tag{73}$$

Since this is true for all  $\zeta \in C_0^\infty(\Omega)$ , we have

$$f(x, y, h) = -\lambda_0 y = -\mathcal{A}y. \tag{74}$$

Similarly,

$$2[F(x, y, h + t\zeta) - F(x, y, h)]/t \leq \mu_0[(h + t\zeta)^2 - h^2]/t \tag{75}$$

and consequently

$$g(x, y, h)\zeta \leq \mu_0 h\zeta \tag{76}$$

and

$$g(x, y, h) = \mu_0 h = \mathcal{B}h \tag{77}$$

We see from (74) and (77) that (40) (41) has a nontrivial solution. Thus, we may assume that (67) holds.

Next, we note that there is an  $\varepsilon > 0$  depending on  $\rho$  such that

$$G(0, w) \geq \varepsilon, \quad b(w) \geq \rho > 0.$$

To see this, suppose that  $\{w_k\} \subset M$  is a sequence such that

$$G(0, w_k) \rightarrow 0, \quad b(w_k) \geq \rho.$$

If

$$b_k = b(w_k) \leq C,$$

this implies

$$b(w_k) - \mu_0 \|w_k\|^2 \rightarrow 0$$

and

$$\int [\mu_0 - \mu(x)] w_k^2 dx \rightarrow 0,$$

since

$$G(0, w) \geq b(w) - \mu_0 \|w\|^2 + \int [\mu_0 - \mu(x)] w^2 dx, \quad w \in M.$$

If we write  $w_k = w'_k + h_k$ ,  $w'_k \in M'$ ,  $h_k \in M_0$  as before, then this tells us that  $b(w'_k) \rightarrow 0$ . Since  $M_0$  is finite dimensional, there is a renamed subsequence such that  $h_k \rightarrow h$ . But the two conclusions above tell us that  $h = 0$ . Since  $b(h) \geq \rho$ , we see that  $\varepsilon > 0$  exists for any constant  $C$ . If the sequence  $\{b_k\}$  is not bounded, we take  $\tilde{w}_k = w_k/b_k^{1/2}$ . Then

$$G(0, w_k)/b_k \geq b(\tilde{w}_k) - \mu_0 \|\tilde{w}_k\|^2 + \int [\mu_0 - \mu(x)](\tilde{w}_k)^2 dx.$$

Next we note that there is a  $\nu > 0$  such that

$$G(0, w) \geq \nu b(w), \quad w \in M. \tag{78}$$

Assuming this for the moment, we see that

$$b_0 := \inf_B G \geq \varepsilon_1 > 0 \tag{79}$$

where

$$B = \{w \in M : b(w) \geq \rho^2\} \cup \{u = (s\varphi_0, w) : s \geq 0, w \in M, \|u\|_D = \rho\}, \tag{80}$$

and  $\varepsilon_1 = \min\{\varepsilon, \nu\rho^2\}$ . By (14) there is an  $R > \rho$  such that

$$\sup_{A^R} G \leq a_0 < \infty, \tag{81}$$

where  $A^R = N \setminus \mathcal{B}_R$ . By Theorem 9 there is a sequence  $\{u_k\} \subset D$  such that (34) holds with  $c \geq \varepsilon_1$ . To complete the proof, we show that the sequence  $u_k$  is bounded in  $D$  and has a convergent subsequence. To see this, assume that  $r_k = \|u_k\|_D \rightarrow \infty$ , and let  $\tilde{u}_k = u_k/r_k$ . Then  $\|\tilde{u}_k\|_D = 1$ , and there is a renamed subsequence such that  $\tilde{u}_k \rightarrow \tilde{u}$ , weakly in  $D$ , strongly in  $L^2(\Omega)$ , and a.e. in  $\Omega$ . Since,

$$b(w_k) + a(v_k) = (g(u_k), w_k) - (f(u_k), v_k),$$

we have

$$b(\tilde{w}_k) + a(\tilde{v}_k) \leq C \|\tilde{u}_k\|^2.$$

Hence,

$$1 = \|\tilde{u}_k\|_D^2 \leq C \|\tilde{u}_k\|^2 \rightarrow C \|\tilde{u}\|^2.$$

Consequently,  $\tilde{u} \neq 0$ .

Let  $\Omega_0$  be the subset of  $\Omega$  on which  $\tilde{u} \neq 0$ . Then

$$|u_k(x)| = r_k |\tilde{u}_k(x)| \rightarrow \infty, \quad x \in \Omega_0. \tag{82}$$

If  $\Omega_1 = \Omega \setminus \Omega_0$ , then we have

$$\int_{\Omega} H(x, u_k) dx = \int_{\Omega_0} + \int_{\Omega_1} \leq \int_{\Omega_0} H(x, u_k) dx + \int_{\Omega_1} W_1(x) dx \rightarrow -\infty. \tag{83}$$

But

$$\int_{\Omega} H(x, u_k) dx = G(u_k) - (G'(u_k), u_k)$$

and

$$|(G'(u_k), u_k)| \leq (\nu_k + \|u_k\|_w) \|G'(u_k)\| \leq \beta.$$

Thus

$$\left| \int_{\Omega} H(x, u_k) dx \right| \leq K.$$

This contradicts (83), and we see that  $r_k = \|u_k\|_D$  is bounded. Hence, there is a renamed subsequence converging weakly to a function  $u \in D$ . Since

$$(G'(u_k), h)/2 = b(w_k, h_2) - a(\nu_k, h_1) - (f(u_k), h_1) - (g(u_k), h_2) \rightarrow 0, \quad h \in D, \tag{84}$$

we have

$$(G'(u), h)/2 = b(w, h_2) - a(\nu, h_1) - (f(u), h_1) - (g(u), h_2) = 0, \quad h \in D. \tag{85}$$

This shows that (40) and (41) hold. Since  $c \neq 0$  and  $G(0) = 0$ , we see that  $u \neq 0$ , and we have a nontrivial solution of the system (40) (41).

It therefore remains only to prove (78). Clearly  $\nu \geq 0$ . If  $\nu = 0$ , then there is a sequence  $\{w_k\} \subset M$  such that

$$G(0, w_k) \rightarrow 0, \quad b(w_k) = 1. \tag{86}$$

Thus there is a renamed subsequence such that  $w_k \rightarrow w$  weakly in  $M$ , strongly in  $L^2(\Omega)$  and a.e. in  $\Omega$ . Consequently

$$\int_{\Omega} [\mu_0 - \mu(x)] w_k^2 dx \leq 1 - \int_{\Omega} \mu(x) w_k^2 dx \leq G(0, w_k) \rightarrow 0 \tag{87}$$

and

$$1 = \int_{\Omega} \mu(x) w^2 dx \leq \mu_0 \|w\|^2 \leq b(w) \leq 1 \tag{88}$$

which means that we have equality throughout. It follows that we must have  $w \in E(\mu_0)$ , the eigenspace of  $\mu_0$ . Since  $w \neq 0$ , we have  $w \neq 0$  a.e. But

$$\int_{\Omega} [\mu_0 - \mu(x)] w^2 dx = 0 \quad (89)$$

implies that the integrand vanishes identically on  $\Omega$ , and consequently  $\mu(x) \equiv \mu_0$ , violating (45). This establishes (78) and completes the proof of the theorem.  $\square$

## 7 Infinite Dimensional Splitting

Let  $N$  be a closed, separable subspace of a Hilbert space  $E$ . We can define a new norm  $|v|_w$  satisfying  $|v|_w \leq \|v\| \forall v \in N$  and such that the topology induced by this norm is equivalent to the weak topology of  $N$  on bounded subsets of  $N$ . This can be done as follows: Let  $\{e_k\}$  be an orthonormal basis for  $N$ . Define

$$(u, v)_w = \sum_{k=1}^{\infty} \frac{(u, e_k)(v, e_k)}{2^k}, \quad u, v \in N.$$

This is a scalar product. The corresponding norm squared is

$$|v|_w^2 = \sum_{k=1}^{\infty} \frac{|(v, e_k)|^2}{2^k}, \quad v \in N.$$

Then  $|v|_w$  satisfies  $|v|_w \leq \|v\|$ ,  $v \in N$ . If  $v_j \rightarrow v$  weakly in  $N$ , then there is a  $C > 0$  such that

$$\|v_j\|, \|v\| \leq C, \quad \forall j > 0.$$

For any  $\varepsilon > 0$ , there exist  $K > 0, M > 0$ , such that  $1/2^K < \varepsilon^2/(8C^2)$  and  $|(v_j - v, e_k)| < \varepsilon/2$  for  $1 \leq k \leq K, j > M$ . Therefore,

$$\begin{aligned} |v_j - v|_w^2 &= \sum_{k=1}^{\infty} \frac{|(v_j - v, e_k)|^2}{2^k} \\ &\leq \sum_{k=1}^K \frac{\varepsilon^2/4}{2^k} + \sum_{k=K+1}^{\infty} \frac{4C^2}{2^k} \\ &\leq \frac{\varepsilon^2}{4} \sum_{k=1}^{\infty} \frac{1}{2^k} + \frac{4C^2}{2^K} \sum_{k=1}^{\infty} \frac{1}{2^k} \\ &\leq \frac{\varepsilon^2}{2} + \frac{\varepsilon^2}{2}. \end{aligned}$$

Therefore,  $v_j \rightharpoonup v$  weakly in  $N$  implies  $|v_j - v|_w \rightarrow 0$ .

Conversely, let  $\|v_j\|, \|v\| \leq C$  for all  $j > 0$  and  $|v_j - v|_w \rightarrow 0$ . Let  $\varepsilon > 0$  be given. If  $h = \sum_{k=1}^{\infty} \alpha_k e_k \in N$ , take  $K$  so large that  $\|h_K\| < \varepsilon/(4C)$ , where  $h_K = \sum_{k=K+1}^{\infty} \alpha_k e_k$ . Take  $M$  so large that  $|v_j - v|_w^2 < \varepsilon^2/(4 \sum_{k=1}^K 2^k |\alpha_k|^2)$  for all  $j > M$ . Then

$$\begin{aligned} |(v_j - v, h - h_K)|^2 &= \left| \sum_{k=1}^K \alpha_k (v_j - v, e_k) \right|^2 \\ &\leq \sum_{k=1}^K 2^k |\alpha_k|^2 \sum_{k=1}^{\infty} \frac{|(v_j - v, e_k)|^2}{2^k} \\ &< \varepsilon^2/4 \end{aligned}$$

for  $j > M$ . Also,  $|(v_j - v, h_K)| \leq 2C\|h_K\| < \varepsilon/2$ . Therefore,

$$|(v_j - v, h)| < \varepsilon, \quad \forall j > M,$$

that is,  $v_j \rightharpoonup v$  weakly in  $N$ .

For  $u = v + h$ ,  $u_1 = v_1 + h_1 \in E = N \oplus N^\perp$  with  $v, v_1 \in N, h, h_1 \in N^\perp$ , we define the scalar product  $(u, u_1)_w = (v, v_1)_w + (h, h_1)$ . Thus, the corresponding norm satisfies  $|u|_w \leq \|u\| \forall u \in E$ . Clearly, when  $\dim N < \infty$ , the norms  $\|\cdot\|$  and  $|\cdot|_w$  are equivalent.

We denote  $E$  equipped with this scalar product and norm by  $E_w$ . It is a scalar product space with the same elements as  $E$ . In particular, if  $(u_n = v_n + w_n)$  is  $\|\cdot\|$ -bounded and  $u_n \xrightarrow{|\cdot|_w} u$ , then  $v_n \rightharpoonup v$  weakly in  $N, w_n \rightarrow w$  strongly in  $N^\perp, u_n \rightharpoonup v + w$  weakly in  $E$ .

For  $u \in E$  and  $Q \subset E$ , we define

$$d_w(u, Q) = \inf_{v \in Q} |u - v|_w.$$

Let  $L$  be a bounded, convex, closed subset of  $N$ . Then  $L$  is  $|\cdot|_w$ -compact. In fact, since  $L$  is bounded with respect to both norms  $|\cdot|_w$  and  $\|\cdot\|$ , for any  $v_n \in L$ , there is a renamed subsequence such that  $v_n \rightharpoonup v_0$  weakly in  $E$ . Then  $v_0 \in L$  since  $L$  is convex, and on the bounded set  $L$  the  $|\cdot|_w$ -topology is equivalent to the weak topology. Thus,  $v_n \xrightarrow{|\cdot|_w} v_0$  and  $L$  is  $|\cdot|_w$ -compact.

Let  $L$  be a compact subset of  $E_w$ . We define  $\Sigma_w(L)$  to be the set of all  $\sigma(t) \in \Sigma : [0, 1] \times E \mapsto E$  such that



1.  $\sigma(t)$  is  $|\cdot|_w$ -continuous.
2.  $\sigma(0)u = u, u \in E$ .
3. There is a finite dimensional subspace  $E_f$  of  $E$  such that  $\dim E_f > 0$  and  $\sigma(t)u - u \in E_f, (t, u) \in I \times L$ .

Here we use  $E_f$  to denote various finite-dimensional subspaces of  $E$  when exact dimensions are irrelevant. Note that  $\Sigma_w(L)$  is not empty since  $\sigma(t) \equiv 1$  is a member.

We let  $\Sigma_{wQ}$  denote the set of those  $\sigma \in \Sigma_w$  which satisfy

$$|\sigma'(t)u|_w \leq C\rho(d_w(\sigma(t)u, Q)), \quad u \in E, \tag{90}$$

where  $Q \subset E$ .

We have

**Lemma 15** *If  $L$  is compact in  $E_w$  and  $\sigma \in \Sigma_w(L)$ , then*

$$\tilde{L} = \{\sigma(t)L : t \in I\}$$

*is compact in  $E_w$ .*

**Proof** Suppose  $\{t_k\} \subset I, \{u_k\} \subset L$  are sequences. Then there are renamed subsequences such that

$$t_j \rightarrow t_0, \quad |u_k - u_0|_w \rightarrow 0.$$

Thus  $I \times L$  is a compact subset of  $I \times E_w$ . By definition, there is a finite dimensional subspace  $E_f$  containing the set  $\{\sigma(t)u - u, t \in I, u \in L\}$ . Since this set is bounded, every sequence has a convergent subsequence. Since every sequence in  $L$  has a convergent subsequence, the same must be true of  $\tilde{L}$ . □

**Lemma 16** *If  $\sigma_1, \sigma_2 \in \Sigma_w(L)$ , then  $\sigma_3 = \sigma_1 \circ \sigma_2 \in \Sigma_w(L)$ .*

**Proof** By the definition of  $\Sigma_w(L)$ , for any  $(s_0, u_0) \in I \times L$ , there is a  $|\cdot|_w$ -neighborhood  $U_{(s_0, u_0)}$  such that  $\{u - \sigma_1(t)u : (t, u) \in U_{(s_0, u_0)} \cap L\} \subset E_f$ . Note that,

$$L \subset \bigcup_{(s, u) \in L} U_{(s, u)}. \text{ Since } L \text{ is } |\cdot|_w\text{-compact, } L \subset \bigcup_{i=1}^{j_0} U_{(s_i, u_i)} \text{ where } (s_i, u_i) \in L.$$

Consequently,  $\{u - \sigma_1(t)u : (t, u) \in L\} \subset E_f$ . The same is true of  $\sigma_2$ . Since

$$\sigma_3(s) = \begin{cases} \sigma_1(2s), & 0 \leq s \leq \frac{1}{2}, \\ \sigma_2(2s - 1)\sigma_1(1), & \frac{1}{2} < s \leq 1, \end{cases}$$

$u - \sigma_3(t)u \in E_f$  as well. □

## 8 Some Lemmas

Before giving the proof of Theorem 9, we prove a few lemmas.

**Lemma 17** *Let  $N$  be a closed, separable subspace of a Banach space  $E$ , and let  $\Omega$  be a bounded, convex, open subset of  $N$  containing a point  $p$ . Let  $F$  be an  $N$ -weakly continuous mapping. Assume*

$$\sigma(t)\partial\Omega \cap F^{-1}(p) = \phi, \quad 0 \leq t \leq 1,$$

for some  $\sigma \in \Sigma_w(\overline{\Omega})$ . Then

$$\sigma(t)\Omega \cap F^{-1}(p) \neq \phi, \quad 0 \leq t \leq 1.$$

**Proof** Assume that there is a  $\sigma \in \Sigma_w(\overline{\Omega})$  such that

$$\sigma(t)\partial\Omega \cap F^{-1}(p) = \phi, \quad 0 \leq t \leq 1, \tag{91}$$

and

$$\sigma(t)\Omega \cap F^{-1}(p) = \phi, \quad 0 \leq t \leq 1,$$

or, equivalently,

$$F(\sigma(t)\Omega) \cap \{p\} = \phi, \quad 0 < t \leq 1. \tag{92}$$

Let

$$\gamma(t)x = F\sigma(t)x, \quad (t, x) \in I \times \overline{\Omega}.$$

Then  $\gamma(t) \in C(I \times \overline{\Omega}, E_w \cap N)$  and

$$\gamma(t)x \neq p, \quad x \in \partial\Omega, \quad t \in [0, 1]. \tag{93}$$

Also

$$\gamma(0)x = Fx = x, \quad x \in \overline{\Omega}. \tag{94}$$

By hypothesis, there exists a fixed finite-dimensional subspace  $E_0$  of  $E$  such that  $F(u - v) - (F(u) - F(v)) \in E_0, \forall u, v \in E$ . Take  $u = \sigma(t)x, v = x$ . Since  $\overline{\Omega}$  is compact in  $E_w$  and  $\sigma \in \Sigma_w(\overline{\Omega})$ , there is a finite dimensional subspace  $E_1$  of  $E$  such that  $\dim E_1 > 0$  and  $\sigma(t)u - u \in E_1, (t, u) \in I \times \overline{\Omega}$ . Hence

$$\gamma(t)x = P_0(F\sigma(t)x - Fx - F[\sigma(t)x - x])$$

$$\begin{aligned}
 &+ FP_1[\sigma(t)x - x] + x \\
 &= x - \varphi(t)x, \quad (t, x) \in I \times \overline{\Omega},
 \end{aligned}$$

where  $\varphi(t)x = -P_0(F\sigma(t)x - Fx - F[\sigma(t)x - x]) - FP_1[\sigma(t)x - x]$ , and the  $P_0, P_1$  are projections onto the finite dimensional subspaces  $E_0, E_1$ . Thus,  $\varphi(t)$  is a compact map from  $I \times \overline{\Omega}$  to  $I \times E_f$ . In view of (91), the Leray–Schauder degree  $i$  satisfies

$$i(\gamma(t), \Omega, p) = i(\gamma(0), \Omega, p) = 1$$

for all  $t \in [0, 1]$ . But this contradicts (92). Hence

$$\sigma(t)\Omega \cap F^{-1}(p) \neq \emptyset, \quad 0 \leq t \leq 1.$$

□

**Lemma 18** *Let  $N$  be a closed separable subspace of a Hilbert space  $E$ , and let  $\Omega$  be a bounded, convex, open subset of  $N$  containing a point  $p$ . Let  $G$  be an  $N$ -weak-to-weak continuously differentiable functional on  $E$ , and let  $F$  be an  $N$ -weakly continuous mapping. Assume that  $d = d(A, B) > 0$ , and*

$$-\infty < b_0 := \inf_B G, \quad b_1 := \sup_{\overline{\Omega}} G < \infty,$$

where  $A = \partial\Omega$  and  $B = F^{-1}(p)$ . If  $\overline{\Omega} \subset \mathcal{B}_R$ , let  $\tilde{B} = B \cap \mathcal{B}_v$ , where

$$\beta \int_R^v \frac{dt}{\rho(t)} > b_1 - b_0, \quad \beta \int_0^d \frac{dt}{\rho(t)} > b_1 - b_0$$

for some  $\rho \in \mathcal{Q}$  and  $\beta > 0$ . Then there is a sequence  $\{u_k\} \subset E$  such that

$$G(u_k) \rightarrow c, \quad b_0 \leq c \leq b_1, \quad \rho(d_w(u_k, \tilde{B}))\|G'(u_k)\| \leq \beta. \tag{95}$$

**Proof** If the lemma were false, then there would be a  $\delta > 0$  such that

$$\rho(d_w(u, \tilde{B}))\|G'(u)\| > \beta \tag{96}$$

when

$$u \in U = \{u \in E : b_0 - 3\delta \leq G(u) \leq b_1 + 3\delta\}. \tag{97}$$

For  $u \in \hat{E} = \{u \in E : G'(u) \neq 0\}$ , let  $h(u) = G'(u)/\|G'(u)\|$ . Then by (96)

$$(G'(u), h(u)) > \beta/\rho(d_w(u, \tilde{B})), \quad u \in U. \tag{98}$$

For each  $u \in U$  there is an  $\tilde{E}$  neighborhood  $W(u)$  of  $u$  such that

$$(G'(v), h(u)) > \beta/\rho(d_w(v, \tilde{B})), \quad v \in W(u) \cap U. \quad (99)$$

For otherwise there would be a sequence  $\{v_k\} \subset U$  such that

$$|v_k - u|_w \rightarrow 0 \text{ and } (G'(v_k), h(u)) \leq \beta/\rho(d_w(v_k, \tilde{B})). \quad (100)$$

Since  $G$  is an  $N$ -weak-to-weak continuously differentiable functional on  $E$ , we would have

$$(G'(v_k), h(u)) \rightarrow (G'(u), h(u)) \leq \beta/\rho(d_w(u, \tilde{B})), \quad (101)$$

by (24) in view of (100). This contradicts (98). Thus (99) holds.

Let  $\tilde{U}$  be the set  $U$  with the inherited topology of  $\tilde{E}$ . It is a metric space, and  $W(u) \cap \tilde{U}$  is an open set in this space. Thus,  $\{W(u) \cap \tilde{U}, u \in \tilde{U}\}$  is an open covering of the paracompact space  $\tilde{U}$  (cf., e.g., [7]). Consequently, there is a locally finite refinement  $\{W_\tau\}$  of this cover. For each  $\tau$  there is an element  $u_\tau$  such that  $W_\tau \subset W(u_\tau)$ . Let  $\{\psi_\tau\}$  be a partition of unity subordinate to this covering. Each  $\psi_\tau$  is locally Lipschitz continuous with respect to the norm  $|u|_w$  and consequently with respect to the norm of  $E$ . Let

$$Y(u) = \sum \psi_\tau(u)h(u_\tau), \quad u \in \tilde{U}. \quad (102)$$

Then  $Y(u)$  is locally Lipschitz continuous with respect to both norms. Moreover,

$$\|Y(u)\| \leq \sum \psi_\tau(u)\|h(u_\tau)\| \leq 1 \quad (103)$$

and

$$(G'(u), Y(u)) = \sum \psi_\tau(u)(G'(u), h(u_\tau)) \geq \beta/\rho(d_w(u, \tilde{B})), \quad u \in \tilde{U}. \quad (104)$$

Reduce  $\delta$  to satisfy

$$\beta \int_\delta^d \frac{dt}{\rho(t)} \geq b_1 - b_0 + \delta.$$

Let

$$Q_0 = \{u \in E : b_0 - 2\delta \leq G(u) \leq b_1 + 2\delta\},$$

$$Q_1 = \{u \in E : b_0 - \delta \leq G(u) \leq b_1 + \delta\},$$

$$Q_2 = E \setminus Q_0,$$

$$\eta(u) = d_w(u, Q_2) / [d_w(u, Q_1) + d_w(u, Q_2)].$$

It is easily checked that  $\eta(u)$  is locally Lipschitz continuous (with respect to the  $E_w$  norm) on  $E$  and satisfies

$$\begin{cases} \eta(u) = 1, & u \in Q_1, \\ \eta(u) = 0, & u \in \bar{Q}_2, \\ \eta(u) \in (0, 1), & \text{otherwise.} \end{cases} \tag{105}$$

Let

$$\tilde{W}(u) = -\eta(u)Y(u)\rho(d_w(u, \tilde{B})).$$

Then

$$\|\tilde{W}(u)\| \leq \rho(d_w(u, \tilde{B})) \leq \rho(d(u, \tilde{B})), \quad u \in \tilde{U}.$$

By Theorem 6, for each  $v \in U$  there is a unique solution  $\sigma(t)v$  of

$$\sigma'(t) = \tilde{W}(\sigma(t)), \quad t \in \mathbb{R}^+, \quad \sigma(0) = v. \tag{106}$$

Take

$$T = \int_{\delta}^d \frac{dt}{\rho(t)} \geq (b_1 - b_0 + \delta) / \beta. \tag{107}$$

Let

$$K = \{(u, t) : u = \sigma(t)v, v \in \bar{\Omega}, t \in [0, T]\}.$$

Then  $K$  is a compact subset of  $\tilde{E} \times \mathbb{R}$ . To see this, let  $(u_k, t_k)$  be any sequence in  $K$ . Then  $u_k = \sigma(t_k)v_k$ , where  $v_k \in \bar{\Omega}$ . Since  $\Omega$  is bounded, there is a subsequence such that  $v_k \rightharpoonup v_0$  weakly in  $E$  and  $t_k \rightarrow t_0$  in  $[0, T]$ . Since  $\bar{\Omega}$  is convex and bounded,  $v_0$  is in  $\bar{\Omega}$  and  $|v_k - v_0|_w \rightarrow 0$ . Since  $\sigma(t)$  is continuous in  $\tilde{E} \times \mathbb{R}$ , we have

$$u_k = \sigma(t_k)v_k \rightharpoonup \sigma(t_0)v_0 \in K.$$

Each  $u_0 \in U$  has a neighborhood  $W(u_0)$  in  $\tilde{E}$  and a finite dimensional subspace  $S(u_0)$  such that  $Y(u) \in S(u_0)$  for  $u \in W(u_0) \cap U$ . Since  $\sigma(t)u$  is continuous in  $\tilde{E} \times \mathbb{R}$ , for each  $(u_0, t_0) \in K$  there is a neighborhood  $W(u_0, t_0) \subset \tilde{E} \times \mathbb{R}$  and a finite dimensional subspace  $S(u_0, t_0) \subset E$  such that  $z_t(u) \in S(u_0, t_0)$  for  $(u, t) \in W(u_0, t_0)$ , where

$$z_t(u) := u - \sigma(t)u = \begin{cases} \int_0^t Y(\sigma(s)u)\rho(d_w(\sigma(s)u, \tilde{B}))ds, & u \in U, \\ 0, & u \notin U. \end{cases} \tag{108}$$

Since  $K$  is compact, there is a finite number of points  $(u_j, t_j) \subset K$  such that  $K \subset W = \cup W(u_j, t_j)$ . Let  $S$  be a finite dimensional subspace of  $E$  containing  $p$  and all the  $S(u_j, t_j)$  and such that  $FS \neq \{0\}$ . Then for  $v \in \bar{\Omega}$  and  $t \in [0, T]$  we have  $z_t(v) \in S$ . Thus  $\sigma \in \Sigma_w(\bar{\Omega})$ .

We also have

$$dG(\sigma(t)v)/dt = -\eta(\sigma(t)v)(G'(\sigma(t)v), Y(\sigma(t)v))\rho(d_w(\sigma(t)v, \tilde{B})) \tag{109}$$

$$\leq -\beta\eta(\sigma).$$

Let  $v \in \bar{\Omega}$ . If there is a  $t_1 \leq T$  such that  $\sigma(t_1)v \notin Q_1$ , then

$$G(\sigma(T)v) \leq G(\sigma(t_1)v) \leq b_0 - \delta. \tag{110}$$

On the other hand, if  $\sigma(t)v \in Q_1$  for all  $t \in [0, T]$ , then we have by (109)

$$G(\sigma(T)v) \leq b_1 - \beta T \leq b_0 - \delta.$$

Hence

$$G(\sigma(T)v) \leq b_0 - \delta, \quad v \in \bar{\Omega}. \tag{111}$$

Let  $u_1(t)$  be the solution of

$$u'(t) = -\rho(u(t)), \quad t \in [0, T], \quad u(0) = d = d(A, \tilde{B}).$$

By Theorem 6,

$$d(\sigma(t)v, \tilde{B}) \geq u_1(t), \quad t \in [0, T], \quad v \in A.$$

But

$$\int_{u_1(t)}^d \frac{d\tau}{\rho(\tau)} = t, \quad t \in [0, T].$$

Consequently,

$$u_1(t) \geq u_1(T) \geq \delta, \quad t \in [0, T],$$

since

$$T = \int_{\delta}^d \frac{dt}{\rho(t)} \geq (b_1 - b_0 + \delta)/\beta.$$

Thus,

$$d(\sigma(t)v, \tilde{B}) \geq \delta, \quad t \in [0, T], \quad v \in A.$$

Consequently,  $\sigma(t)v \cap \tilde{B} = \phi$ ,  $t \in (0, T]$ . This means that

$$\sigma(t)v \cap \tilde{B} = \phi, \quad v \in A, \quad t \in (0, T].$$

Hence,

$$\sigma(t)A \cap \tilde{B} = \phi, \quad t \in (0, T], \tag{112}$$

and

$$\sup_{\sigma(T)A} G \leq b_0 - \delta.$$

Let  $u_2(t)$  be the solution of

$$u'(t) = \rho(u(t)), \quad t \in [0, T], \quad u(0) = R.$$

By Theorem 6,

$$\|\sigma(t)v\| \leq u_2(t), \quad t \in [0, T], \quad v \in A.$$

Now it follows from the choice of  $v$  and the fact that  $A \subset \mathcal{B}_R$ , that

$$\int_R^{u_2(T)} \frac{dt}{\rho(t)} = T \leq \int_R^v \frac{dt}{\rho(t)}.$$

Thus,

$$\|\sigma(t)v\| \leq v, \quad v \in A, \quad t \in [0, T].$$

Hence,

$$\sigma(t)A \cap [B \setminus \tilde{B}] = \phi, \quad t \in [0, T]. \tag{113}$$

If we combine this with (112), we obtain

$$\sigma(t)A \cap B = \phi, \quad t \in [0, T]. \tag{114}$$

But  $\sigma \in \Sigma_w(\overline{\Omega})$ . By Lemma 17, this implies

$$\sigma(t)\Omega \cap B \neq \phi, \quad 0 < t \leq T.$$

Thus, there is a  $u \in \Omega$  such that  $\sigma(T)u \in B$ . But that would mean that  $G(\sigma(T)u) \geq b_0$ , contradicting (111). This completes the proof.  $\square$

**Lemma 19** *Let  $N$  be a closed separable subspace of a Hilbert space  $E$ , and let  $\Omega$  be a bounded, convex, open subset of  $N$  containing a point  $p$ . Let  $G$  be an  $N$ -weak-to-weak continuously differentiable functional on  $E$ . Let  $F$  be an  $N$ -weakly continuous mapping. Assume  $d = d(A, B) > 0$ , and*

$$-\infty < b_0 := \inf_B G, \quad b_1 := \sup_{\overline{\Omega}} G < \infty,$$

where  $A = \partial\Omega$  and  $B = F^{-1}(p)$ . If  $\overline{\Omega} \subset \mathcal{B}_R$ , let  $\nu > 0$ ,  $\beta > 0$  be such that

$$\beta \ln \frac{2\nu + d}{2\nu} > b_1 - b_0, \quad \beta \ln \frac{3\nu}{2\nu + R} > b_1 - b_0. \tag{115}$$

Then there is a sequence  $\{u_k\} \subset E$  such that

$$G(u_k) \rightarrow c, \quad b_0 \leq c \leq b_1, \quad (\nu + |u_k|_w) \|G'(u_k)\| \leq \beta. \tag{116}$$

**Proof** We take  $\rho(s) = 2\nu + s$ . Then Lemma 18 requires (115) to hold. Since

$$|u|_w \leq \nu + d_w(u, \tilde{B}), \quad u \in E,$$

(116) follows from (95).  $\square$

**Lemma 20** *Let  $N$  be a closed separable subspace of a Hilbert space  $E$ , and let  $\Omega_n$  be a sequence of bounded, convex, open subsets of  $N$  containing a point  $p$ . Let  $G$  be an  $N$ -weak-to-weak continuously differentiable functional on  $E$ . Let  $F$  be an  $N$ -weakly continuous mapping. Assume  $d_n = d(A_n, B) \rightarrow \infty$ , and*

$$-\infty < b_0 := \inf_B G, \quad b_1 := \sup_n \sup_{\overline{\Omega}_n} G < \infty,$$

where  $A_n = \partial\Omega_n$  and  $B = F^{-1}(p)$ . Assume that  $\overline{\Omega}_n \subset \mathcal{B}_{R_n}$ , and there are  $\nu_n > 0$ ,  $\beta > 0$  be such that

$$\beta \ln \frac{2\nu_n + d_n}{2\nu_n} > b_1 - b_0, \quad \beta \ln \frac{3\nu_n}{2\nu_n + R_n} > b_1 - b_0. \tag{117}$$

Then there is a sequence  $\{u_n\} \subset E$  such that



$$G(u_n) \rightarrow c, \quad b_0 \leq c \leq b_1, \quad (v_n + |u_n|_w)\|G'(u_n)\| \leq \beta. \tag{118}$$

**Proof** By Lemma 19, for each  $v_n$  there is a sequence satisfying

$$G(u_k) \rightarrow c, \quad b_0 \leq c \leq b_1, \quad (v_n + |u_k|_w)\|G'(u_k)\| \leq \beta. \tag{119}$$

Pick one member. □

**Lemma 21** *Let  $N$  be a closed, separable subspace of a Banach space  $E$ , and for each  $R > R_0$  let  $\Omega_R(p)$  be an open, convex, bounded set in  $N$  containing a point  $p$ . Let  $G$  be an  $N$ -weak-to-weak continuously differentiable functional on  $E$ , and let  $F$  be an  $N$ -weakly continuous mapping. Assume*

$$d_R = d_w(A^R, F^{-1}(p)) \rightarrow \infty, \quad R \rightarrow \infty,$$

where  $A^R = N \setminus \Omega_R(p)$ . Assume

$$-\infty < b_0 = \inf_B G, \quad \sup_{A^R} G \leq a_0 < \infty, \tag{120}$$

for  $R > R_0$ . Then for each sequence  $v_k \geq 2d_k - \varepsilon$  there is a  $\beta > 0$  and a sequence  $\{u_k\} \subset E$  such that

$$G(u_k) \rightarrow c, \quad b_0 \leq c \leq b_1, \quad (v_k + |u_k|_w)\|G'(u_k)\| \leq \beta. \tag{121}$$

**Proof** Take  $v_n = 2R_n = 2(d_n + \varepsilon)$ . Apply Lemma 20. □

**Lemma 22** *Let  $N$  be a closed, separable subspace of a Hilbert space  $E$  with orthogonal complement  $M \oplus \{v_0\}$ , where  $v_0$  is an element in  $E$  having unit norm and orthogonal to both  $M$  and  $N$ , and let  $\delta < R$  be positive numbers. Let  $\varphi(t) \in C^1(\mathbb{R})$  be such that*

$$0 \leq \varphi(t) \leq 1, \quad \varphi(0) = 1,$$

and

$$\varphi(t) = 0, \quad |t| \geq 1.$$

Let

$$F(v+w+s v_0) = v + [s + \delta - \delta\varphi(\|w\|^2/\delta^2)]v_0, \quad w \in M, \quad v \in N, \quad s \in \mathbb{R}. \tag{122}$$

Let  $G$  be an  $N$ -weak-to-weak continuously differentiable functional on  $E$ . Assume

$$-\infty < b_0 = \inf_B G, \quad \sup_A G = b_1 < \infty, \tag{123}$$

holds with  $A = [N \oplus \{v_0\}]$ ,  $A_R = A \cap \partial B_R$  and  $B = F^{-1}(\delta v_0) = \{w + r v_0 : w \in M, r = \delta\varphi(\|w\|^2/\delta^2)\}$ . Then for each sequence  $v_k \rightarrow \infty$  there is a  $\beta > 0$  and a sequence  $\{u_k\} \subset E$  such that

$$G(u_k) \rightarrow c, \quad b_0 \leq c \leq b_1, \quad (v_k + |u_k|_w)\|G'(u_k)\| \leq \beta. \tag{124}$$

**Proof** Note that  $F$  is an  $N$ -weakly continuous mapping. Let

$$\Omega_R = \{v + s v_0 : |v|_w^2 + s^2 < R^2, \quad v \in N\},$$

$$A_R = \{v + s v_0 : |v|_w^2 + s^2 = R^2, \quad v \in N\}$$

and

$$B = F^{-1}(\delta v_0) = \{w + r v_0 : w \in M, \quad r = \delta\varphi(\|w\|^2/\delta^2)\}.$$

Then,

$$\begin{aligned} d_w(A_R, B)^2 &= \inf |v + s v_0 - w - r v_0|_w^2 \\ &= |v|_w^2 + \|w\|^2 + (s - r)^2 \\ &= R^2 - s^2 + \|w\|^2 + s^2 - 2s\delta\varphi(\|w\|^2/\delta^2) + \delta^2\varphi(\|w\|^2/\delta^2)^2 \\ &\geq R^2 - 2R\delta\varphi(\|w\|^2/\delta^2) + \delta^2\varphi(\|w\|^2/\delta^2)^2 \\ &= [R - \delta\varphi(\|w\|^2/\delta^2)]^2 \\ &\geq (R - \delta)^2 \rightarrow \infty, \quad R \rightarrow \infty. \end{aligned}$$

The hypotheses of Lemma 21 are satisfied. □

**Proof of Theorem 9** Apply Lemma 18. □

**Proof of Theorem 11** Apply Lemma 21. □

**Proof of Corollary 13** Apply Lemma 22. □

**Proof of Theorem 2** Apply Theorem 9 □

**Proof of Theorem 4** Apply Theorem 12. □

## References

1. D.G. Costa, On a class of elliptic systems in  $\mathbb{R}^N$ . Electron. J. Differ. Equ. 7 (1994), 14 pp. (electronic)

2. D.G. Costa, C.A. Magalhes, A variational approach to subquadratic perturbations of elliptic systems. *J. Differ. Equ.* **111**(1), 103–122 (1994)
3. D.G. de Figueiredo, P.L. Felmer, On superquadratic elliptic systems. *Trans. Am. Math. Soc.* **343**(1), 99–116 (1994)
4. M.F. Furtado and E.A.B. Silva, Double resonant problems which are locally non-quadratic at infinity, in *Proceedings of the USA–Chile Workshop on Nonlinear Analysis* (Vi-a del Mar-Valparaiso, 2000), 155–171 (electronic), *Electron. J. Differ. Equ. Conf.*, vol. 6, Southwest Texas State Univ., San Marcos, TX, 2001
5. M.F. Furtado, L.A. Maia, E.A.B. Silva, On a double resonant problem in  $\mathbb{R}^N$ . *Differ. Integr. Equ.* **15**(11), 1335–1344 (2002)
6. M.F. Furtado, L.A. Maia, E.A.B. Silva, Solutions for a resonant elliptic system with coupling in  $\mathbb{R}^N$ . *Commun. Partial Differ. Equ.* **27**(7–8), 1515–1536 (2002)
7. J.L. Kelley, *General Topology* (Van Nostrand Reinhold, 1955)
8. W. Kryszewski, A. Szulkin, Generalized linking theorems with an application to semilinear Schrödinger equation. *Adv. Differ. Equ.* **3**, 441–472 (1998)
9. G. Li, J. Yang, Asymptotically linear elliptic systems (English summary). *Commun. Partial Differ. Equ.* **29**(5–6), 925–954 (2004)
10. J. Mawhin, Nonlinear functional analysis and periodic solution of semilinear wave equation, in *Nonlinear Phenomena in Mathematical Sciences*, ed. by Lakshmikantham (Academic Press, London, 1982), pp. 671–681
11. M. Schechter, New saddle point theorems, Generalized functions and their applications (Varanasi, 1991) (Plenum, New York, 1993), pp. 213–219
12. M. Schechter, A generalization of the saddle point method with applications. *Ann. Polon. Math.* **57**(3), 269–281 (1992)
13. M. Schechter, New linking theorems. *Rend. Sem. Mat. Univ. Padova* **99**, 255–269 (1998)
14. M. Schechter, Infinite-dimensional linking. *Duke Math. J.* **94**(3), 573–595 (1998)
15. M. Schechter, Critical point theory with weak-to-weak linking. *Commun. Pure Appl. Math.* **51**(11–12), 1247–1254 (1998)
16. M. Schechter, Rotationally invariant periodic solutions of semilinear wave equations. *Abstr. Appl. Anal.* **3**(1–2), 171–180 (1998)
17. M. Schechter, *Linking Methods in Critical Point Theory* (Birkhauser Boston, 1999)
18. M. Schechter, Periodic solutions of semilinear higher dimensional wave equations. *Chaos Solitons Fractals* **12**(6), 1029–1034 (2001)
19. M. Schechter, Sandwich pairs in critical point theory. *Trans. Am. Math. Soc.* **360**(6), 2811–2823 (2008)
20. M. Schechter, Strong sandwich pairs. *Indiana Univ. Math. J.* **57**(3), 1105–1131 (2008)
21. M. Schechter, *Minimax Systems and Critical Point Theory* (Birkhauser Boston, 2009)
22. M. Schechter, Critical point theory: sandwich and linking systems (2020)
23. M. Schechter, W. Zou, Weak linking. *Nonlinear Anal.* **55**(6), 695–706 (2003)
24. E.A. de B.e. Silva, Linking theorems and applications to semilinear elliptic problems at resonance. *Nonlinear Anal. TMA* **16**, 455–477 (1991)
25. E.A. B. Silva, Nontrivial solutions for noncooperative elliptic systems at resonance, in *Proceedings of the USA–Chile Workshop on Nonlinear Analysis* (Vi-a del Mar-Valparaiso, 2000), 267–283 (electronic), *Electron. J. Differ. Equ. Conf.*, vol. 6, Southwest Texas State Univ., San Marcos, TX, 2001
26. K. Tintarev, Solutions to elliptic systems of Hamiltonian type in  $\mathbb{R}^N$ . *Electron. J. Differ. Equ.* **29** (1999), 11 pp.

# Canonical Systems of Partial Differential Equations



Martin Schechter

**Abstract** We use critical point theory to find solutions of the nonlinear steady state Schrödinger equations arising in the study of photonic lattices.

## 1 Introduction

Systems of partial differential equations arise in many investigations in the physical sciences. Depending on the application and on the questions asked, different types of systems emerge. Usually, if one is interested in finding steady states solutions, the resulting system is elliptic in nature. Such systems may display severe difficulties when one tries to solve them. Most of the time they admit a trivial solution, where all of the unknown functions are identically zero. However, the physical application requires a solution which is not identically zero. In such cases, the methods of solution may be very difficult. In particular, one has to show that the solution obtained is not trivial. The system that we study is not only deceptive, but it is almost impossible to tell if one has solved the whole system or only parts of the system. I call it “canonical.” I shall elaborate on this later.

Many general systems are the form

$$\mathcal{A}v = f(x, v, w), \quad x \in Q \subset \mathbb{R}^n, \quad (1)$$

$$\mathcal{B}w = g(x, v, w), \quad x \in Q \subset \mathbb{R}^n, \quad (2)$$

where  $\mathcal{A}$ ,  $\mathcal{B}$  are linear partial differential operators. I call this system “deceptive” if  $(v, 0)$  is a solution of (1) and (2) whenever  $v$  satisfies

$$\mathcal{A}v = f(x, v, 0), \quad x \in Q \subset \mathbb{R}^n, \quad (3)$$

---

M. Schechter (✉)

Department of Mathematics, University of California, Irvine, CA, USA

e-mail: [mschecht@math.uci.edu](mailto:mschecht@math.uci.edu)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,

Springer Optimization and Its Applications 173,

[https://doi.org/10.1007/978-3-030-72563-1\\_22](https://doi.org/10.1007/978-3-030-72563-1_22)

609

or  $(0, w)$  is a solution whenever  $w$  satisfies

$$\mathcal{B}w = g(x, 0, w), \quad x \in Q \subset \mathbb{R}^n. \quad (4)$$

In this case it is very difficult to determine if both components of a solution are nontrivial.

The particular system I have chosen consist of nonlinear Schrödinger equations arising in optics (cf. [16]) describing the propagation of a light wave in induced photonic lattices. They can be written in the form

$$iV_t + \Delta V = \frac{PV}{1 + |V|^2 + |W|^2}$$

$$iW_t + \Delta W = \frac{PW}{1 + |V|^2 + |W|^2}$$

for the periodic wave functions  $V(x, t), W(x, t)$  over a periodic bounded spacial domain  $\Omega \subset \mathbb{R}^2$ , where  $P, Q$  are parameters (cf. [2, 21]). To find a steady state solution, we look for solutions of the form

$$V(x, t) = e^{i\lambda t} v(x), \quad W(x, t) = e^{i\lambda t} w(x),$$

where  $\lambda$  is a real constant. This leads to the following system of equations over a periodic domain  $\Omega \subset \mathbb{R}^2$  :

$$\Delta v = \frac{Pv}{1 + v^2 + w^2} + \lambda v, \quad (5)$$

$$\Delta w = \frac{Qw}{1 + v^2 + w^2} + \lambda w, \quad (6)$$

where  $P, Q, \lambda$  are parameters. The solutions  $v, w$  are to be periodic in  $\Omega$  with the same periods. One wishes to obtain intervals of the parameter  $\lambda$  for which there are nontrivial solutions. This will provide continuous energy spectrum that allows the existence of steady state solutions. This system was studied in [2], where it was shown that

1. If  $P, Q, \lambda$  are all positive, then the only solution is trivial.
2. If  $P < 0$  and  $0 < \lambda < -P$ , then the system (5) and (6) has a nontrivial solution.
3. If  $P, Q > 0$ , there is a constant  $\delta > 0$  such that the system (5) and (6) has a nontrivial solution provided  $0 < -\lambda < \delta$ .
4. All of these statements are true if we replace  $P$  by  $Q$ .

Wave propagation in nonlinear periodic lattices has been studied by many reseachers (cf., e.g., [1–11, 17, 20–23] and their bibliographies.)

In the present paper, we wish to cover some remaining situations not mentioned in [2] as well as extending their results to higher dimensions. We shall show that there are many intervals of the parameters in which nontrivial solutions exist. Our results are true in any dimension.

In stating our results, we shall make use of the following considerations. Let  $\Omega$  be a bounded periodic domain in  $\mathbb{R}^n$ ,  $n \geq 1$ . Consider the operator  $-\Delta$  on functions in  $L^2(\Omega)$  having the same periods as  $\Omega$ . The spectrum of  $-\Delta$  consists of isolated eigenvalues of finite multiplicity:

$$0 = \lambda_0 < \lambda_1 < \dots < \lambda_\ell < \dots ,$$

with eigenfunctions in  $L^\infty(\Omega)$ . Let  $\lambda_\ell$ ,  $\ell \geq 0$ , be one of these eigenvalues, and define

$$N = \bigoplus_{\lambda \leq \lambda_\ell} E(\lambda), \quad M = N^\perp.$$

As noted in [2], to prove the existence of a nontrivial solution of system (5) and (6), it suffices to obtain a nontrivial solution of either

$$\Delta v = \frac{Pv}{1 + v^2} + \lambda v, \tag{7}$$

or

$$\Delta w = \frac{Qw}{1 + w^2} + \lambda w. \tag{8}$$

This stems from the fact that  $(v, 0)$  is a solution of (5) and (6) if  $v$  is a solution of (7) and  $(0, w)$  is a solution of (5) and (6) if  $w$  is a solution of (8). The author is unaware if such solutions are desirable from the physical point of view. However, we have been able to find values of  $P, Q, \lambda$  for which the system (5) and (6) has a solution  $(v, w)$  where  $v \neq 0, w \neq 0$ .

We shall prove

**Theorem 1** *If  $0 < \lambda < -P$  or  $0 < \lambda < -Q$ , then (5) and (6) has a nontrivial solution.*

**Theorem 2** *If  $0 < -\lambda < P$  or  $0 < -\lambda < Q$ , then (5) and (6) has a nontrivial solution.*

**Theorem 3** *If  $P > 0, Q > 0, \sigma = -\lambda > 0$ , and either  $0 \leq \sigma - P < \lambda_1 < \sigma$  or  $0 \leq \sigma - Q < \lambda_1 < \sigma$ , then (5) and (6) has a nontrivial solution.*

**Theorem 4** *If  $P > 0, Q > 0, \sigma = -\lambda > 0$ , and either  $\lambda_\ell \leq \sigma - P < \lambda_{\ell+1} < \sigma$  or  $\lambda_\ell \leq \sigma - Q < \lambda_{\ell+1} < \sigma$  then (5) and (6) has a nontrivial solution.*

**Theorem 5** *If  $P > 0, Q = 0, \lambda = -\lambda_\ell < 0, -\Delta w = \lambda_\ell w$  and*

$$\lambda_\ell < P \int_\Omega \frac{1}{1+w^2} / |\Omega|,$$

then (5) and (6) has a nontrivial solution. If  $w \neq 0$ , then the solution has both components nonzero.

**Theorem 6** If  $Q > 0$ ,  $P = 0$ ,  $\lambda = -\lambda_\ell < 0$ ,  $-\Delta v = \lambda_\ell v$  and

$$\lambda_\ell < Q \int_\Omega \frac{1}{1+v^2} / |\Omega|,$$

then (5) and (6) has a nontrivial solution. If  $v \neq 0$ , then the solution has both components nonzero.

## 2 Some Lemmas

In proving our results we shall make use of the following lemmas (cf., e.g., [12, 14, 15, 18]). For the definition of linking, cf. [12].

**Lemma 1** Let  $M, N$  be closed subspaces of a Hilbert space  $E$  such that one of them is finite dimensional and  $E = M \oplus N$ . Take  $B = \partial \mathcal{B}_\delta \cap M$ , and let  $w_0$  be any element in  $\partial \mathcal{B}_1 \cap M$ . Take  $A$  to be the set of all  $u$  of the form

$$u = v + sw_0, \quad v \in N, \quad s \in \mathbb{R},$$

satisfying the following

- (a)  $\|v\|_E \leq R, \quad s = 0$
- (b)  $\|v\|_E \leq R, \quad s = 2R_0$
- (c)  $\|v\|_E = R, \quad 0 \leq s \leq 2R_0,$

where  $0 < \delta < \min(R, R_0)$ . Then  $A$  and  $B$  link each other.

**Lemma 2** The sets  $\|u\|_E = R > 0$  and  $\{e_1, e_2\}$  link each other provided  $\|e_1\|_E < R$  and  $\|e_2\|_E > R$ .

**Lemma 3** If  $G(u) \in C^1(E, \mathbb{R})$  satisfies

$$\alpha = \inf_E G > -\infty, \tag{9}$$

then there is a sequence  $\{u_k\}$  such that

$$G(u_k) \rightarrow \alpha, \quad (1 + \|u_k\|_E) \|G'(u_k)\| \rightarrow 0. \tag{10}$$

**Lemma 4** If  $A$  links  $B$ , and  $G(u) \in C^1(E, \mathbb{R})$  satisfies

$$a_0 = \sup_A G \leq b_0 = \inf_B G, \tag{11}$$

then there is a sequence  $\{u_k\}$  such that

$$G(u_k) \rightarrow c \geq b_0, \quad (1 + \|u_k\|_E)\|G'(u_k)\| \rightarrow 0. \tag{12}$$

We let  $E$  be the subspace of  $H^{1,2}(\Omega)$  consisting of those functions having the same periodicity as  $\Omega$  with norm given by

$$\|w\|_E^2 = \|\nabla w\|^2 + \|w\|^2.$$

Assume  $P \neq 0, Q \neq 0, \lambda \neq 0$ . Let

$$a(u) = \frac{1}{P} [ \|\nabla v\|^2 + \lambda \|v\|^2 ] + \frac{1}{Q} [ \|\nabla w\|^2 + \lambda \|w\|^2 ], \quad v, w \in E \tag{13}$$

and

$$G(u) = a(u) + \int_{\Omega} \ln(1 + u^2) dx. \tag{14}$$

We have

**Lemma 5** *If  $G(u)$  is given by (14), then every sequence satisfying (10) has a subsequence converging in  $E$ . Consequently, there is a  $u \in E$  such that  $G(u)=c$  and  $G'(u) = 0$ .*

**Proof** The sequence satisfies

$$G(u_k) = \frac{1}{P} \|\nabla v_k\|^2 + \frac{\lambda}{P} \|v_k\|^2 + \frac{1}{Q} \|\nabla w_k\|^2 + \frac{\lambda}{Q} \|w_k\|^2 + \int_{\Omega} \ln\{1 + |u_k|^2\} dx \rightarrow c, \tag{15}$$

$$\begin{aligned} (G'(u_k), q)/2 &= \frac{1}{P} (\nabla v_k, \nabla q) + \frac{\lambda}{P} (v_k, q) \\ &+ \frac{1}{Q} (\nabla w_k, \nabla q) + \frac{\lambda}{Q} (w_k, q) \\ &+ \int_{\Omega} \frac{u_k q}{1 + u_k^2} dx \rightarrow 0, \quad q = (g, h), \end{aligned} \tag{16}$$

$$(G'(u_k), v_k)/2 = \frac{1}{P} (\nabla v_k, \nabla v_k) + \frac{\lambda}{P} (v_k, v_k) \tag{17}$$



$$+ \int_{\Omega} \frac{u_k v_k}{1 + u_k^2} dx \rightarrow 0.$$

and

$$\begin{aligned} (G'(u_k), w_k)/2 &= \frac{1}{Q} (\nabla w_k, \nabla w_k) + \frac{\lambda}{Q} (w_k, w_k) \\ &+ \int_{\Omega} \frac{u_k w_k}{1 + u_k^2} dx \rightarrow 0. \end{aligned} \tag{18}$$

Thus,

$$\int_{\Omega} H(x, u_k) dx \rightarrow c, \tag{19}$$

where

$$H(x, t) = \ln(1 + t^2) - \frac{t^2}{1 + t^2}. \tag{20}$$

Let  $\rho_k = \|u_k\|_H$ , where

$$\begin{aligned} \|u\|_H^2 &= \frac{1}{|P|} [\|\nabla v\|^2 + |\lambda| \|v\|^2] \\ &+ \frac{1}{|Q|} [\|\nabla w\|^2 + |\lambda| \|w\|^2], \quad u = (v, w) \in E. \end{aligned} \tag{21}$$

Assume first that  $\rho_k \rightarrow \infty$ . Let  $\tilde{u}_k = u_k/\rho_k$ . Then  $\|\tilde{u}_k\|_H = 1$ . Hence, there is a renamed subsequence such that  $\tilde{u}_k \rightarrow \tilde{u}$  in  $E$ , and  $\tilde{u}_k \rightarrow \tilde{u}$  in  $L^2(\Omega)$  and a.e. Now

$$\|u_k\|_H^2 = \frac{1}{|P|} [\|\nabla v_k\|^2 + |\lambda| \|v_k\|^2] + \frac{1}{|Q|} [\|\nabla w_k\|^2 + |\lambda| \|w_k\|^2]. \tag{22}$$

By (17) and (18),

$$\begin{aligned} \|u_k\|_H^2 &\leq |(G'(u_k), v_k)|/2 + |(G'(u_k), w_k)|/2 \\ &+ \frac{|\lambda| - \lambda}{|P|} \|v_k\|^2 + \frac{|\lambda| - \lambda}{|Q|} \|w_k\|^2 \\ &+ \int_{\Omega} \frac{u_k^2}{1 + u_k^2} dx. \end{aligned}$$

Hence,

$$1 = \|\tilde{u}_k\|_H^2 \leq [|(G'(u_k), v_k)|/2 + |(G'(u_k), w_k)|/2]/\rho_k^2 + C\|\tilde{u}_k\|^2. \tag{23}$$

In the limit we have,

$$1 \leq C\|\tilde{u}\|^2.$$

This shows that  $\tilde{u} \neq 0$ . Let  $\Omega_0$  be the subset of  $\Omega$  where  $\tilde{u}(x) \neq 0$ . Then  $|\Omega_0| \neq 0$ . Thus

$$\begin{aligned} \int_{\Omega} H(x, u_k) dx &= \int_{\Omega_0} H(x, u_k) dx + \int_{\Omega \setminus \Omega_0} H(x, u_k) dx \\ &\geq \int_{\Omega_0} H(x, u_k) dx \rightarrow \infty. \end{aligned}$$

This contradicts (19). Thus, the sequence satisfying (10) is bounded in  $E$ . Hence, there is a renamed subsequence such that  $u_k \rightharpoonup$  in  $E$ , and  $u_k \rightarrow u_0$  in  $L^2(\Omega)$  and a.e. Taking the limit in (17), we obtain

$$\begin{aligned} (G'(u_0), q)/2 &= \frac{1}{P}(\nabla v_0, \nabla g) + \frac{\lambda}{P}(v_0, g) \\ &+ \frac{1}{Q}(\nabla w_0, \nabla h) + \frac{\lambda}{Q}(w_0, h) \\ &+ \int_{\Omega} \frac{u_0 q}{1 + u_0^2} dx = 0, \quad q = (g, h), \end{aligned} \tag{24}$$

Thus,  $u_0$  satisfies  $G'(u_0) = 0$ . Since  $u_0 \in E$ , it satisfies

$$\begin{aligned} (G'(u_0), u_0)/2 &= \frac{1}{P}(\nabla v_0, \nabla v_0) + \frac{\lambda}{P}(v_0, v_0) \\ &+ \frac{1}{Q}(\nabla w_0, \nabla w_0) + \frac{\lambda}{Q}(w_0, w_0) \\ &+ \int_{\Omega} \frac{u_0^2}{1 + u_0^2} dx = 0 \end{aligned} \tag{25}$$

Also, from the limit in (17), we have

$$\begin{aligned} \lim \frac{1}{P} \|\nabla v_k\|^2 &= \lim (G'(u_k), v_k)/2 \\ &- \lim \left[ \frac{\lambda}{P} \|v_k\|^2 + \int_{\Omega} \frac{v_k^2}{1 + u_k^2} dx \right] \end{aligned}$$

$$\begin{aligned}
 &= - \left[ \frac{\lambda}{P} \|v\|^2 + \int_{\Omega} \frac{v^2}{1+u^2} dx \right] \\
 &= \frac{1}{P} \|\nabla v\|^2,
 \end{aligned}$$

with a similar statement for  $\|\nabla w\|^2$ . Consequently,  $\nabla u_k \rightarrow \nabla u$  in  $L^2(\Omega)$ . This shows that  $G(u_k) \rightarrow G(u_0)$ . Hence,  $G(u_0) = c$ .

**Lemma 6** *If  $G'(u) = 0$ , then  $(v,w)$  is a solution of (5) and (6).*

**Proof** From (24) we see that

$$|(\nabla u, \nabla q)| \leq C\|q\|, \quad q \in E.$$

From the fact that the functions and  $\Omega$  are periodic with the same period, it follows that  $u \in H^{2,2}(\Omega)$  and satisfies (5) and (6) (cf., e.g., [13]).

**Lemma 7**

$$\int_{\Omega} \ln(1+u^2) dx / \|u\|_H^2 \rightarrow 0, \quad \|u\|_H \rightarrow \infty. \tag{26}$$

**Proof** Suppose  $u_k \in H$  is a sequence such that  $\rho_k = \|u_k\|_H \rightarrow \infty$ . Let  $\tilde{u}_k = u_k/\rho_k$ . Then  $\|\tilde{u}_k\|_H = 1$ . Hence, there is a renamed subsequence such that  $\tilde{u}_k \rightarrow \tilde{u}$  in  $H$ , and  $\tilde{u}_k \rightarrow \tilde{u}$  in  $L^2(\Omega)$  and a.e. Now

$$\frac{\ln(1+u_k^2)}{\rho_k^2} = \frac{\ln(1+u_k^2)}{u_k^2} \tilde{u}_k^2 \rightarrow 0 \text{ a.e.}$$

and it is dominated a.e. by  $\tilde{u}_k^2 \rightarrow \tilde{u}^2$  in  $L^1(\Omega)$ . Thus

$$\int_{\Omega} \frac{\ln(1+u_k^2)}{\rho_k^2} dx \rightarrow 0.$$

Since this is true for any sequence satisfying  $\|u_k\|_H \rightarrow \infty$ , we see that (26) holds.

**Corollary 1** *If*

$$I(u) = \|u\|_H^2 - \int_{\Omega} \ln(1+u^2) dx,$$

then

$$I(v) \rightarrow \infty \text{ as } \|v\|_H \rightarrow \infty. \tag{27}$$

**Proof** We have

$$I(u)/\|u\|_H^2 = 1 - \int_{\Omega} \ln(1 + u^2)dx/\|u\|_H^2 \rightarrow 1, \quad \|u\|_H \rightarrow \infty$$

by Lemma 7. This gives (27).

**Lemma 8**

$$\int_{\Omega} [u^2 - \ln(1 + u^2)]dx/\|u\|_H^2 \rightarrow 0, \quad \|u\|_H \rightarrow 0. \tag{28}$$

**Proof** Suppose  $u_k \in H$  is a sequence such that  $\rho_k = \|u_k\|_H \rightarrow 0$ . In particular, there is a renamed subsequence such that  $u_k \rightarrow 0$  a.e. Let  $\tilde{u}_k = u_k/\rho_k$ . Then  $\|\tilde{u}_k\|_H = 1$ . Hence, there is a renamed subsequence such that  $\tilde{u}_k \rightarrow \tilde{u} \in H$ , and  $\tilde{u}_k \rightarrow \tilde{u}$  in  $L^2(\Omega)$  and a.e. Now

$$\frac{u_k^2 - \ln(1 + u_k^2)}{\rho_k^2} \leq \frac{u_k^2}{1 + u_k^2} \tilde{u}_k^2 \rightarrow 0 \text{ a.e.}$$

and it is dominated a.e. by  $\tilde{u}_k^2 \rightarrow \tilde{u}^2$  in  $L^1(\Omega)$ . Thus

$$\int_{\Omega} \frac{u_k^2 - \ln(1 + u_k^2)}{\rho_k^2} dx \rightarrow 0.$$

Since this is true for any sequence satisfying  $\|u_k\|_H \rightarrow 0$ , we see that (28) holds.

### 3 Proofs of the Theorems

**Proof of Theorem 1** We let  $E$  be the subspace of  $H^{1,2}(\Omega)$  consisting of those functions having the same periodicity as  $\Omega$  with norm given by

$$\|w\|_E^2 = \|\nabla w\|^2 + \|w\|^2.$$

Let  $u = (v, w)$ , where  $v, w \in E$  and  $u^2 = v^2 + w^2$ . If  $q = (g, h)$ , we write  $uq = vg + wh$ . Define

$$\begin{aligned} \|u\|_H^2 &= \frac{1}{|P|} [\|\nabla v\|^2 + |\lambda| \|v\|^2] \\ &+ \frac{1}{|Q|} [\|\nabla w\|^2 + |\lambda| \|w\|^2], \quad v, w \in E. \end{aligned} \tag{29}$$

Assume that  $P, Q, \lambda$  do not vanish. Then  $\|u\|_H^2$  is a norm on  $H = E \times E$  having a scalar product  $(u, h)_H$ .

Let

$$I(u) = \|u\|_H^2 - \int_{\Omega} \ln(1 + u^2) \, dx. \tag{30}$$

Then,

$$(I'(u), q)/2 = (u, q)_H - \int_{\Omega} \frac{uq}{1 + u^2} \, dx, \quad q \in H. \tag{31}$$

If  $I'(u) = 0$ , then

$$\Delta v = \frac{-|P|v}{1 + |v|^2 + |w|^2} + |\lambda|v, \tag{32}$$

$$\Delta w = \frac{-|Q|w}{1 + |v|^2 + |w|^2} + |\lambda|w. \tag{33}$$

This is equivalent to (5) and (6) if  $P < 0, Q < 0, \lambda > 0$ . To prove the theorem, we must show that there is a nontrivial solution of  $I'(u) = 0$  when either  $0 < \lambda < -P$  or  $0 < \lambda < -Q$ .

Assume  $0 < \lambda < -P$ . We show that  $I(u)$  has a minimum  $u \neq 0$ .

Let the sequence  $u_k \in H$  satisfy

$$I(u_k) \searrow \alpha = \inf_H I$$

(which may be  $-\infty$ ). By (27),  $\rho_k = \|u_k\|_H$  is bounded. Hence, there is a renamed subsequence such that  $u_k \rightharpoonup u_0$  in  $H$ , and  $u_k \rightarrow u_0$  in  $L^2(\Omega)$  and a.e. Since

$$\|u_k\|_H^2 - 2([u_k - u_0], u_0)_H = \|u_0\|_H^2 + \|u_k - u_0\|_H^2,$$

we have

$$\begin{aligned} I(u_0) &\leq \|u_k\|_H^2 - 2([u_k - u_0], u_0)_H \\ &\quad - \int_Q \ln(1 + u_0^2) \, dx \\ &= I(u_k) - 2([u_k - u_0], u_0)_H \\ &\quad - \int_Q [\ln(1 + u_0^2) - \ln(1 + u_k^2)] \, dx \\ &\rightarrow \alpha. \end{aligned}$$

Thus,

$$\alpha \leq I(u_0) \leq \alpha,$$

showing that  $\alpha$  is finite and that  $u_0$  is a minimum. Thus,  $I'(u_0) = 0$  and  $u_0$  is a solution of

$$\Delta v = \frac{-|P|v}{1 + |v|^2 + |w|^2} + \lambda v, \tag{34}$$

$$\Delta w = \frac{-|Q|w}{1 + |v|^2 + |w|^2} + \lambda w. \tag{35}$$

Next, we show that  $u_0 \neq 0$ . We do this by showing that  $\alpha < 0$ . Consider a constant function  $u = (s, 0)$ . Then,

$$I(u) = \left[ \frac{\lambda}{|P|} s^2 - \ln(1 + s^2) \right] |\Omega|, \quad s \in \mathbb{R}.$$

This has a negative minimum if  $\lambda < |P|$ . Thus  $I(u_0) = \alpha < 0$ . Since  $I(0, 0) = 0$ , we see that  $u_0 \neq 0$ . However,  $u_0$  satisfies (34) and (35), not (5) and (6). To rectify the situation, we merely note that the same method produces a negative minimum  $v_0$  for  $I(v, 0)$ , and  $(v_0, 0)$  is a nontrivial solution of (5) and (6). This completes the proof for the case  $0 < \lambda < -P$ . The case  $0 < \lambda < -Q$  is treated similarly.

**Proof of Theorem 2** Assume  $0 < \sigma < P$ ,  $0 < \sigma < Q$ , and let  $a(u)$  and  $G(u)$  be given by (13) and (14), respectively. Then  $G'(u) = 0$  iff  $u = (v, w)$  is a solution of (5) and (6). We search for a nontrivial solution.

Let  $\rho_k = \|u_k\|_H$ , where

$$\begin{aligned} \|u\|_H^2 &= \frac{1}{|P|} [\|\nabla v\|^2 + |\lambda| \|v\|^2] \\ &+ \frac{1}{|Q|} [\|\nabla w\|^2 + |\lambda| \|w\|^2], \quad u = (v, w) \in E. \end{aligned} \tag{36}$$

Assume that  $\rho_k \rightarrow 0$ . Let  $\tilde{u}_k = u_k/\rho_k$ . Then  $\|\tilde{u}_k\|_H = 1$ . Hence, there is a renamed subsequence such that  $\tilde{u}_k \rightharpoonup \tilde{u}$  in  $E$ , and  $\tilde{u}_k \rightarrow \tilde{u}$  in  $L^2(\Omega)$  and a.e. We have

$$\begin{aligned} G(u_k)/\rho_k^2 &= \frac{1}{P} \|\nabla \tilde{v}_k\|^2 + \frac{\lambda + P}{P} \|\tilde{v}_k\|^2 \\ &+ \frac{1}{Q} \|\nabla \tilde{w}_k\|^2 + \frac{\lambda + Q}{Q} \|\tilde{w}_k\|^2 \\ &+ \int_{\Omega} [\ln\{1 + |u_k|^2\} - u_k^2] dx / \rho_k^2. \end{aligned}$$

Since  $P > \sigma$  and  $Q > \sigma$ , we see in view of Lemma 7 that there are positive constants  $\varepsilon, \eta$  such that

$$G(u)/\|u\|_H^2 \geq \varepsilon, \quad \|u\|_H \leq \eta.$$

Let  $A$  be the set of those  $u \in H$  such that  $\|u\|_H = \eta$ . Consider a constant function  $u = (s, 0)$ . Then,

$$G(u)/s^2 = \left[ \frac{\lambda}{P} + s^{-2} \ln(1 + s^2) \right] |\Omega| \rightarrow \frac{\lambda}{P} |\Omega| < 0, \quad s \rightarrow \infty.$$

Hence, there is a  $u \in H$  such that  $\|u\|_H > \eta$  and  $G(u) < \varepsilon \eta^2$ . Since  $G(0, 0) = 0$ , there is a  $u \in H$  such that  $\|u\|_H < \eta$  and  $G(u) < \varepsilon \eta^2$ . The theorem now follows from Lemmas 2, 4, and 5.

**Proof of Theorem 3** Assume  $P > 0$ ,  $Q > 0$ ,  $\lambda < 0$  and  $\sigma = -\lambda > \max[P, \lambda_1]$ . Let

$$a(u) = \frac{1}{P} [\|\nabla v\|^2 - \sigma \|v\|^2] + \frac{1}{Q} [\|\nabla w\|^2 - \sigma \|w\|^2], \quad v, w \in E \quad (37)$$

and

$$G(u) = a(u) + \int_{\Omega} \ln(1 + u^2) dx. \quad (38)$$

Then  $G'(u) = 0$  iff  $u$  satisfies (5) and (6).

First, we note that

$$G(u) \leq 0, \quad u \in N,$$

if  $\sigma \geq P$ ,  $\sigma \geq Q$ . To see this, let  $u = (c, d) \in N$ . Then

$$a(c, d) = -\frac{\sigma}{P} c^2 |\Omega| - \frac{\sigma}{Q} d^2 |\Omega|$$

and

$$\int_{\Omega} \ln(1 + c^2 + d^2) dx \leq (c^2 + d^2) |\Omega|.$$

Thus,

$$G(u) \leq \left[1 - \frac{\sigma}{P}\right] c^2 |\Omega| + \left[1 - \frac{\sigma}{Q}\right] d^2 |\Omega|.$$

This means that

$$G(u) \leq 0, \quad u \in N, \quad (39)$$

provided  $\sigma \geq P$ ,  $\sigma \geq Q$ .

Next, let  $\psi$  be an eigenfunction of  $-\Delta$  corresponding to the eigenvalue  $\lambda_1$ . If we take  $u = (\psi + c, \psi + d)$ , we have

$$a(u) = \frac{1}{P} [ \|\nabla(\psi + c)\|^2 - \sigma \|(\psi + c)\|^2 ] + \frac{1}{Q} [ \|\nabla(\psi + d)\|^2 - \sigma \|(\psi + d)\|^2 ],$$

and this gives

$$a(u) = \frac{1}{P} [ (\lambda_1 - \sigma) \|\psi\|^2 - \sigma c^2 ] + \frac{1}{Q} [ (\lambda_1 - \sigma) \|\psi\|^2 - \sigma d^2 ],$$

which will be negative if  $\sigma > \lambda_1$ . Moreover,

$$\int_{\Omega} \ln(1 + 2\psi^2 + c^2 + d^2) dx / \|u\|_H^2 \rightarrow 0, \quad \|u\|_H \rightarrow \infty.$$

This follows from the fact that

$$\int_{\Omega} \ln(1 + u^2) dx / \|u\|_H^2 \rightarrow 0, \quad \|u\|_H \rightarrow \infty \tag{40}$$

(Lemma 7). Consequently,

$$\limsup_{\|(\psi+c, \psi+d)\|_H \rightarrow \infty} G(\psi + c, \psi + d) < 0 \tag{41}$$

provided  $\sigma > \lambda_1$ .

Next, let  $u = (v, w)$  be any function in  $M$ . Then  $\|\nabla u\|^2 = \|\nabla v\|^2 + \|\nabla w\|^2 \geq \lambda_1 \|v\|^2 + \lambda_1 \|w\|^2 = \lambda_1 \|u\|^2$ . Then

$$a(u) + \|u\|^2 \geq \frac{1}{P} [ 1 - \frac{\sigma - P}{\lambda_1} ] \|\nabla v\|^2 + \frac{1}{Q} [ 1 - \frac{\sigma - Q}{\lambda_1} ] \|\nabla w\|^2.$$

Thus, there is an  $\varepsilon > 0$  such that

$$a(u) + \|u\|^2 \geq 2\varepsilon \|\nabla u\|^2, \quad u \in M, \tag{42}$$

when  $\sigma - \lambda_1 < \min[P, Q]$ .

Now

$$\int_{\Omega} [u^2 - \ln(1 + u^2)] dx / \|u\|_H^2 \rightarrow 0, \quad \|u\|_H \rightarrow 0 \tag{43}$$

by Lemma 8. If we combine (42) and (43), we see that there is an  $\varepsilon > 0$  such that



$$G(u) \geq \varepsilon \|\nabla u\|^2, \quad u \in M, \tag{44}$$

when  $\|\nabla u\|^2$  is small and  $\sigma - \lambda_1 < \min[P, Q]$ .

Take  $A = \partial(N \oplus \{\psi\})$ ,  $B = \partial\mathcal{B}_\delta \cap M$ . By (39), (41), and (44) one can apply Lemma 1 to obtain (10) and then Lemma 5 to conclude that (5) and (6) has a nontrivial solution. To see this, note that  $a_0 = 0 < \varepsilon \leq b_0$ , showing that the solution  $u_0$  satisfies  $G(u_0) \geq \varepsilon > 0$ . Since  $G(0) = 0$ , we see that  $u_0 \neq 0$ . If  $\max[P, \lambda_1] < -\lambda < P + \lambda_1$  is true, but  $\max[Q, \lambda_1] < -\lambda < Q + \lambda_1$ , is not, we can apply the argument used in the proof of Theorem 1. The same is true in the other direction. This completes the proof.

**Proof of Theorem 4** First, we note that

$$G(u) \leq 0, \quad u \in N, \tag{45}$$

if  $\sigma \geq \lambda_\ell + \max[P, Q]$ . To see this, let  $u = (v, w) \in N$ . Then  $\|\nabla u\|^2 = \|\nabla v\|^2 + \|\nabla w\|^2 \leq \lambda_\ell \|v\|^2 + \lambda_\ell \|w\|^2 = \lambda_\ell \|u\|^2$ . Then

$$G(u) \leq \frac{1}{P}[\lambda_\ell - \sigma + P] \|v\|^2 + \frac{1}{Q}[\lambda_\ell - \sigma + Q] \|w\|^2 \leq 0.$$

Next, let  $g$  be an eigenfunction of  $-\Delta$  corresponding to the eigenvalue  $\lambda_{\ell+1}$ . If we take  $u = (g + v, g + w)$ , we have

$$\begin{aligned} a(u) &= \frac{1}{P} [\|\nabla(g + v)\|^2 - \sigma \|(g + v)\|^2] \\ &\quad + \frac{1}{Q} [\|\nabla(g + w)\|^2 - \sigma \|(g + w)\|^2], \end{aligned}$$

and this gives

$$\begin{aligned} a(u) &= \frac{1}{P} [(\lambda_{\ell+1} - \sigma) \|g\|^2 + (\lambda_\ell - \sigma)\|v\|^2] \\ &\quad + \frac{1}{Q} [(\lambda_{\ell+1} - \sigma) \|g\|^2 + (\lambda_\ell - \sigma^2)\|w\|^2], \end{aligned}$$

which will be negative if  $\sigma > \lambda_{\ell+1}$ . Moreover, by Lemma 7,

$$\int_\Omega \ln(1 + 2g^2 + v^2 + w^2)dx / \|u\|_H^2 \rightarrow 0, \quad \|u\|_H \rightarrow \infty. \tag{46}$$

Consequently, (46) holds provided  $\sigma > \lambda_{\ell+1}$ .

Next, let  $u = (v, w)$  be any function in  $M$ . Then  $\|\nabla u\|^2 = \|\nabla v\|^2 + \|\nabla w\|^2 \geq \lambda_{\ell+1} \|v\|^2 + \lambda_{\ell+1} \|w\|^2 = \lambda_{\ell+1} \|u\|^2$ . Then

$$a(u) + \|u\|^2 \geq \frac{1}{P} \left[1 - \frac{\sigma - P}{\lambda_{\ell+1}}\right] \|\nabla v\|^2 + \frac{1}{Q} \left[1 - \frac{\sigma - Q}{\lambda_{\ell+1}}\right] \|\nabla w\|^2.$$

Thus, there is an  $\varepsilon > 0$  such that

$$a(u) + \|u\|^2 \geq 2\varepsilon \|\nabla u\|^2, \quad u \in M, \tag{47}$$

when  $\sigma - \lambda_{\ell+1} < \min[P, Q]$ .

Now by Lemma 8,

$$\int_{\Omega} [u^2 - \ln(1 + u^2)] dx / \|u\|_H^2 \rightarrow 0, \quad \|u\|_H \rightarrow 0. \tag{48}$$

If we combine (47) and (48), we see that there is an  $\varepsilon > 0$  such that

$$G(u) \geq \varepsilon \|\nabla u\|^2, \quad u \in M, \tag{49}$$

when  $\|\nabla u\|^2$  is small and  $\sigma - \lambda_{\ell+1} < \min[P, Q]$ .

By (45), (46), and (49) one can apply Lemma 1 to obtain (10) and then Lemma 5 to conclude that (5) and (6) has a nontrivial solution  $u_0$  taking  $A = \partial(N \oplus \{g\})$ ,  $B = \partial \mathcal{B}_\delta \cap M$ . Then  $a_0 = 0 < \varepsilon \leq b_0$ , showing that  $G(u_0) \geq \varepsilon > 0$ . Since  $G(0, 0) = 0$ , we see that  $u_0 \neq 0$ . If  $\lambda_\ell < \sigma - P < \lambda_{\ell+1} < \sigma$  is true, but  $\lambda_\ell < \sigma - Q < \lambda_{\ell+1} < \sigma$  is not, we can apply the argument used in the proof of Theorem 1. The same is true in the other direction. This completes the proof.

**Proof of Theorem 5** If  $w = 0$ , this follows from Theorem 1 since  $0 < \lambda_\ell < -P$ . Otherwise, let

$$I_w(v) = \frac{1}{P} \|\nabla v\|^2 - \frac{\lambda_\ell}{P} \|v\|^2 + \int_{\Omega} \ln\{1 + v^2 + w^2\} dx, \quad v \in H. \tag{50}$$

Then,

$$(I'_w(v), g)/2 = \frac{1}{P} (\nabla v, \nabla g) - \frac{\lambda_\ell}{P} (v, g) + \int_{\Omega} \frac{vg}{1 + v^2 + w^2} dx. \tag{51}$$

If  $I'_w(v) = 0$ , then  $u = (v, w)$  satisfies

$$\Delta v = \frac{Pv}{1 + v^2 + w^2} - \lambda_\ell v, \tag{52}$$

$$\Delta w = -\lambda_\ell w, \tag{53}$$

which is (5) and (6) for the case  $Q = 0$ ,  $\lambda = -\lambda_\ell$ . If we can find a solution  $v \neq 0$  of  $I'_w(v) = 0$ , then we shall have a solution  $u = (v, w)$  of (5) and (6) with  $v \neq 0$ ,  $w \neq 0$ . This was done in Theorem 3 of [19].

The proof of Theorem 6 is similar to that of Theorem 5 and is omitted.

## References

1. G. Bartal, O. Manela, O. Cohen, J.W. Fleischer, M. Segev, Observation of second-band vortex solitons in 2D photonic lattices. *Phys. Rev. Lett.* **95**, 053904 (2005)
2. S. Chen, Y. Lei, Existence of steady-state solutions in a nonlinear photonic lattice model. *J. Math. Phys.* **52**(6), 063508 (2011)
3. W. Chen, D.L. Mills, Gap solitons and the nonlinear optical response of superlattices. *Phys. Rev. Lett.* **62**, 1746–1749 (1989)
4. N.K. Efremidis, S. Sears, D.N. Christodoulides, Discrete solitons in photorefractive optically-induced photonic lattices. *Phys.Rev.Lett.* **85**, 1863–1866 (2000)
5. W.J.W. Fleischer, M. Segev, N.K. Efremidis, D.N. Christodoulides, Observation of two-dimensional discrete solitons in optically induced nonlinear photonic lattices. *Nature* **422**(6928), 147–149 (2003)
6. J.W. Fleischer, G. Bartal, O. Cohen, O. Manela, M. Segev, J. Hudock, D.N. Christodoulides, Observation of vortex-ring discrete solitons in photonic lattices. *Phys. Rev. Lett.* **92**, 123904 (2004)
7. P. Kuchment, The mathematics of photonic crystals, in *Mathematical Modeling in Optical Science*. Frontiers Application of Mathematical, vol. 22 (SIAM, Philadelphia, 2001), pp. 207–272
8. C. Liu, Q. Ren, On the steady-state solutions of a nonlinear photonic lattice model. *J. Math. Phys.* **56**, 031501, 1–12 (2015). <https://doi.org/10.1063/1.4914333>
9. H. Martin, E.D. Eugeniya, Z. Chen, Discrete solitons and soliton-induced dislocations in partially coherent photonic lattices. Martin et al. *Phys. Rev. Lett.* **92**, 123902 (2004)
10. D.N. Neshev, T.J. Alexander, E.A. Ostrovskaya, Y.S. Kivshar, H. Martin, I. Makasyuk, Z. Chen, Observation of discrete vortex solitons in optically induced photonic lattices. *Phys. Rev. Lett.* **92**, 123903 (2004)
11. A. Pankov, Periodic nonlinear Schrödinger equation with application to photonic crystals. *Milan J. Math.* **73**, 259–287 (2005)
12. M. Schechter, *Linking Methods in Critical Point Theory* (Birkhauser, Boston, 1999)
13. M. Schechter, An introduction to nonlinear analysis, in *Cambridge Studies in Advanced Mathematics*, vol. 95 (Cambridge University, Cambridge, 2004)
14. M. Schechter, The use of Cerami sequences in critical point theory theory. *Abstr. Appl. Anal.* **2007**, 28 (2007). Art. ID 58948
15. M. Schechter, *Minimax Systems and Critical Point Theory* (Birkhauser, Boston, 2009)
16. M. Schechter, Steady state solutions for Schrödinger equations governing nonlinear optics. *J. Math. Phys.* **53**, 043504, 8 pp. (2012)
17. M. Schechter, Photonic lattices. *J. Math. Phys.* **54**, 061502, 7 pp. (2013)
18. M. Schechter, *Critical Point Theory, Sandwich and Linking Systems* (Birkhauser, Boston, 2020)
19. M. Schechter, Schrödinger equations in nonlinear optics, in *Nonlinear Analysis and Global Optimization*, ed. by Th. M. Rassias, P.M. Pardalos (Springer, 2021), pp. 449–459
20. Y. Yang, *Soliton in Field Theory and Nonlinear Analysis* (Springer, New York, 2001)
21. Y. Yang, R. Zhang, Steady state solutions for nonlinear Schrödinger equation arising in optics. *J. Math. Phys.* **50**, 053501–053509 (2009)
22. J. Yang, A. Bezryadina, Z. Chen, I. Makasyuk, Observation of two-dimensional lattice vector solitons. *Opt. Lett.* **29**, 1656 (2004)
23. J. Yang, I. Makasyuk, A. Bezryadina, Z. Chen, Dipole and quadrupole solitons in optically induced two-dimensional photonic lattices: theory and experiment. *Studies Appl. Math.* **113**, 389–412 (2004)

# The Semi-discrete Method for the Approximation of the Solution of Stochastic Differential Equations



Ioannis S. Stamatiou

**Abstract** We study the numerical approximation of the solution of stochastic differential equations (SDEs) that do not follow the standard smoothness assumptions. In particular, we focus on SDEs that admit solutions which take values in a certain domain; examples of these equations appear in various fields of application such as mathematical finance and natural sciences among others, where the quantity of interest may be the interest rate, which takes non-negative values, or the population dynamics which takes values between zero and one. We review the Semi-Discrete method (SD), a numerical method that has the qualitative feature of domain preservation among other desirable properties.

## 1 Introduction

We are interested in the numerical approximation of stochastic differential equations (SDEs) that admit solutions in a certain domain and do not satisfy the usual assumptions. Such equations appear in mathematical finance, e.g. interest rate models, but also in other fields of applications such as natural and social sciences. Generally speaking, explicit solutions of these SDEs are unknown, so numerical methods have to be used to simulate them. While numerical methods exist that converge strongly to the true solution of SDEs with non-standard coefficients, few of them are able to maintain the solution process domain. Implicit methods can in some cases succeed in that direction, but they are usually more time-consuming. Let us state the problem in mathematical terms.

Throughout, let  $T > 0$  and  $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_{0 \leq t \leq T}, \mathbb{P})$  be a complete probability space, meaning that the filtration  $\{\mathcal{F}_t\}_{0 \leq t \leq T}$  is right continuous and  $\mathcal{F}_0$  includes all  $\mathbb{P}$ -null sets. We are interested in the following SDE in integral form

---

I. S. Stamatiou (✉)

Department of Biomedical Sciences, University of West Attica, Athens, Greece

e-mail: [istamatiou@uniwa.gr](mailto:istamatiou@uniwa.gr)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,

Springer Optimization and Its Applications 173,

[https://doi.org/10.1007/978-3-030-72563-1\\_23](https://doi.org/10.1007/978-3-030-72563-1_23)

625

$$X_t = X_0 + \int_0^t a(s, X_s)ds + \int_0^t b(s, X_s)dW_s, \quad t \in [0, T], \tag{1}$$

where  $W_{t,\omega} : [0, T] \times \Omega \rightarrow \mathbb{R}^m$  is an  $m$ -dimensional Wiener process adapted to the filtration  $\{\mathcal{F}_t\}_{0 \leq t \leq T}$ , the drift coefficient  $a : [0, T] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$  and the diffusion coefficient  $b : [0, T] \times \mathbb{R}^d \rightarrow \mathbb{R}^{d \times m}$  are measurable functions such that (1) has a unique strong solution and  $X_0$  is independent of all  $\{W_t\}_{0 \leq t \leq T}$ . SDE (1) has non-autonomous coefficients, i.e.  $a(t, x), b(t, x)$  depend explicitly on  $t$ . More precisely, we assume the existence of a predictable stochastic process  $X : [0, T] \times \Omega \rightarrow \mathbb{R}^d$  such that, c.f. [22, Def. 5.2.1], [24, Def. 2.1],

$$\{a(t, X_t)\} \in \mathcal{L}^1([0, T]; \mathbb{R}^d), \quad \{b(t, X_t)\} \in \mathcal{L}^2([0, T]; \mathbb{R}^{d \times m})$$

and

$$\mathbb{P} \left[ X_t = X_0 + \int_0^t a(s, X_s)ds + \int_0^t b(s, X_s)dW_s \right] = 1, \quad \text{for every } t \in [0, T].$$

The drift coefficient  $a$  is the infinitesimal mean of the process  $(X_t)$  and the diffusion coefficient  $b$  is the infinitesimal standard deviation of the process  $(X_t)$ . SDEs like (1) have rarely explicit solutions so numerical approximations are required for path simulations of the solution process  $X_t(\omega)$ .

We are interested in strong approximations (mean-square) of (1), in the case of nonlinear drift and diffusion coefficients. Strongly converging numerical schemes have applications in many areas, such as simulating scenarios, filtering or visualizing stochastic dynamics (c.f [20, Sec. 4] and references therein), they are of theoretical interest (they provide basic insight into weak-sense schemes) and usually do not require simulations over long-time periods or of a significant number of trajectories. In the same time we aim for numerical methods that preserve the domain of the original process, or as we say possess an *eternal life time*.

**Definition 1 (Eternal Life Time of Numerical Solution)** Let  $D \subseteq \mathbb{R}^d$  and consider a process  $(X_t)$  well defined on the domain  $\overline{D}$ , with initial condition  $X_0 \in \overline{D}$  and such that

$$\mathbb{P}(\{\omega \in \Omega : X(t, \omega) \in \overline{D}\}) = 1,$$

for all  $t > 0$ . A numerical solution  $(Y_n)_{n \in \mathbb{N}}$  has an *eternal life time* if

$$\mathbb{P}(Y_{t_{n+1}} \in \overline{D} \mid Y_{t_n} \in \overline{D}) = 1.$$

Let us consider the following nonlinear model both in the drift and diffusion coefficient:

$$x_t = x_0 + \int_0^t (\alpha x_s - \beta x_s^2)ds + \int_0^t \sigma x_s^{3/2}dW_s, \quad t \in [0, T], \tag{2}$$

where  $x_0$  is independent of  $\{W_t\}_{0 \leq t \leq T}$ ,  $x_0 > 0$  a.s. and  $\sigma \in \mathbb{R}$ . SDE (2) is referred to as the 3/2-model [18] or the inverse square root process [1] and is used for modeling stochastic volatility. The conditions  $\alpha > 0$  and  $\beta > 0$  are necessary and sufficient for the stationarity of the process  $(x_t)$  and such that neither zero nor infinity is attainable in finite time [1, App. A].

A “good” numerical scheme for the approximation of the solution of an SDE that takes positive values, as (2), should preserve positivity, c.f. [2, 21]. The explicit Euler scheme does not have that property, since its increments are conditionally Gaussian and therefore there is a positive probability of producing negative values. We refer, among other papers, to [23] that considers Euler type schemes, modifications of them to overcome the above drawback, and the importance of positivity.

SDE (2) is a special case of super-linear models of the form (1) where one of the coefficients  $a(\cdot)$ ,  $b(\cdot)$  is super-linear, i.e. when we have that

$$a(x) \geq \frac{|x|^\beta}{C}, \quad b(x) \leq C|x|^\alpha, \quad \text{for every } |x| \geq C, \tag{3}$$

or

$$b(x) \geq \frac{|x|^\beta}{C}, \quad a(x) \leq C|x|^\alpha, \quad \text{for every } |x| \geq C, \tag{4}$$

where  $\beta > 1$ ,  $\beta > \alpha \geq 0$ ,  $C > 0$ .

Another issue that arises at the numerical approximation of super-linear problems like (3) or (4), is that the moments of the scheme may explode, see [19, Th. 1]. A method that overcomes this drawback is the tamed Euler method, which reads in a general form

$$Y_{n+1}^N(\omega) := Y_n^N(\omega) + a_\Delta(Y_n^N(\omega)) \cdot \Delta + b_\Delta(Y_n^N(\omega))\Delta W_n(\omega), \tag{5}$$

for every  $n \in \{0, 1, \dots, N - 1\}$ ,  $N \in \mathbb{N}$  and all  $\omega \in \Omega$  where  $\Delta W_n(\omega) := W_{\frac{(n+1)T}{N}}(\omega) - W_{\frac{nT}{N}}(\omega)$  are the increments of the Wiener process,  $Y_0^N(\omega) := x_0(\omega)$  and the control functions are such that  $a_\Delta \rightarrow a$  and  $b_\Delta \rightarrow b$  as  $\Delta \rightarrow 0$ , c.f. [20, (4)], [31, Rel. (3.1)], [27], for various choices of  $a_\Delta$  and  $b_\Delta$ . These balanced type schemes are explicit, do not explode in finite time and converge strongly to the exact solution. Nevertheless, in general they do not preserve positivity. We should also mention here other interesting implicit methods, c.f. [26] and [25], which are unfortunately time-consuming.

We study SDEs of the general type (1) with solutions in a certain domain and our aim is to construct explicit numerical schemes which on the one hand, converge strongly to the solution process and on the other, preserve the domain of the original SDE.

The semi-discrete (SD) method, originally proposed in [7], has all the above properties and more, that is:

- it is explicit in general and therefore does not require a lot of computational time,
- it does not explode in non-linear problems, see [8, Sec. 3], [15, Sec. 4], [11]

- it strongly converges to the exact solution of the original SDE, [7, Sec. 3], [10–15, 28, 29]
- has the qualitative property of domain preservation, [7, Sec. 3.2], [10, 12–14], [15, Sec. 4], [11, 28, 29]
- preserves monotonicity, [7, Sec. 3.1]
- preserves the a.s. asymptotic stability of the underlying SDE, [16].

## 2 The Semi-discrete Method: Setting and General Results

We address first the scalar differential equation (1), that is the one-dimensional case ( $d = 1$ ), which we rewrite here

$$x_t = x_0 + \int_0^t a(s, x_s) ds + \int_0^t b(s, x_s) dW_s, \quad t \in [0, T]. \quad (6)$$

Consider the equidistant partition  $0 = t_0 < t_1 < \dots < t_N = T$  with step-size  $\Delta = T/N$ . We assume that there is a unique strong solution a.s. to the following SDE

$$y_t = y_{t_n} + \int_{t_n}^t f(t_n, s, y_{t_n}, y_s) ds + \int_{t_n}^t g(t_n, s, y_{t_n}, y_s) dW_s, \quad t \in (t_n, t_{n+1}], \quad (7)$$

for every  $n \in \mathbb{N}, n \leq N - 1$ , with  $y_0 = x_0$ . Here, the auxiliary functions  $f$  and  $g$  satisfy the following assumption.

**Assumption 2.1** *Let  $f(s, r, x, y), g(s, r, x, y) : [0, T]^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$  be such that  $f(s, s, x, x) = a(s, x), g(s, s, x, x) = b(s, x)$ , where  $f, g$  satisfy the following conditions:*

$$\begin{aligned} |f(s_1, r_1, x_1, y_1) - f(s_2, r_2, x_2, y_2)| &\leq C_R (|s_1 - s_2| + |r_1 - r_2| + |x_1 - x_2| + |y_1 - y_2|) \\ |g(s_1, r_1, x_1, y_1) - g(s_2, r_2, x_2, y_2)| &\leq C_R (|s_1 - s_2| + |r_1 - r_2| + |x_1 - x_2| + |y_1 - y_2| \\ &\quad + \sqrt{|x_1 - x_2|}), \end{aligned}$$

for any  $R > 0$  such that  $|x_1| \vee |x_2| \vee |y_1| \vee |y_2| \leq R$ , where the constant  $C_R$  depends on  $R$  and  $x \vee y$  denotes the maximum of  $x, y$ .

We consider the following interpolation process of the semi-discrete approximation, in a compact form,

$$y_t = y_0 + \int_0^t f(\hat{s}, s, y_{\hat{s}}, y_s) ds + \int_0^t g(\hat{s}, s, y_{\hat{s}}, y_s) dW_s, \quad (8)$$

where  $\hat{s} = t_n$  when  $s \in [t_n, t_{n+1})$ . In that way we may compare with the exact solution  $x_t$ , which is a continuous time process. The first and third variable in  $f, g$  denote the discretized part of the original SDE. We observe from (8) that in order to solve for  $y_t$ , we have to solve, in general, an SDE and not an algebraic equation. We can reproduce the Euler scheme if we choose  $f(s, r, x, y) = a(s, x)$  and  $g(s, r, x, y) = b(s, x)$ . The semi-discrete method (8) can be appropriately modified to produce an implicit scheme that is explicitly and easily solved if necessary (see [11, 14, 29]).

In the case of superlinear coefficients the numerical scheme (8) converges to the true solution  $x_t$  of SDE (6) and this is stated in the following, see [15, Th. 2.1].

**Theorem 1 (Strong Convergence)** *Suppose Assumption 2.1 holds and (7) has a unique strong solution for every  $n \leq N - 1$ , where  $x_0 \in \mathcal{L}^p(\Omega, \mathbb{R})$ . Let also*

$$\mathbb{E}(\sup_{0 \leq t \leq T} |x_t|^p) \vee \mathbb{E}(\sup_{0 \leq t \leq T} |y_t|^p) < A,$$

for some  $p > 2$  and  $A > 0$ . Then the semi-discrete numerical scheme (8) converges to the true solution of (6) in the  $\mathcal{L}^2$ -sense, that is

$$\lim_{\Delta \rightarrow 0} \mathbb{E} \sup_{0 \leq t \leq T} |y_t - x_t|^2 = 0. \tag{9}$$

Theorem 1 is an extension of [8, Th. 1] to time-dependent coefficients which covers super-linear diffusion coefficients, like for example of the form  $b(t, x) = \beta(t) \cdot x^{3/2}$ . In all other cases we may assume the usual local Lipschitz assumption for both  $f$  and  $g$ .

We understand by the general form of decomposition (7) that we may produce many different semi-discrete numerical schemes. In a sense the method is problem dependent, since the form of the drift and diffusion coefficients,  $a$  and  $b$ , of the original SDE suggest the way of discretization. We will see in the following Sections 3 and 4 applications of the semi-discrete method which all have in common the qualitative property of domain preservation.

Relation (9) does not reveal the order of convergence. In order to show the order of convergence, we work with a truncated version of the SD method, see [30].

We choose a strictly increasing function  $\mu : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  such that for every  $s, r \leq T$

$$\sup_{|x| \leq u} (|f(s, r, x, y)| \vee |g(s, r, x, y)|) \leq \mu(u)(1 + |y|), \quad u \geq 1. \tag{10}$$

The inverse function of  $\mu$ , denoted by  $\mu^{-1}$ , maps  $[\mu(1), \infty)$  to  $\mathbb{R}_+$ . Moreover, we choose a strictly decreasing function  $h : (0, 1] \rightarrow [\mu(1), \infty)$  and a constant  $\hat{h} \geq 1 \vee \mu(1)$  such that

$$\lim_{\Delta \rightarrow 0} h(\Delta) = \infty \quad \text{and} \quad \Delta^{1/6} h(\Delta) \leq \hat{h} \quad \text{for every} \quad \Delta \in (0, 1]. \tag{11}$$



Let  $\Delta \in (0, 1]$  and  $f_\Delta, g_\Delta$  defined by

$$\phi_\Delta(s, r, x, y) := \phi\left(s, r, (|x| \wedge \mu^{-1}(h(\Delta))) \frac{x}{|x|}, y\right), \tag{12}$$

for  $x, y \in \mathbb{R}$  where we set  $x/|x| = 0$  when  $x = 0$ . Using the truncated auxiliary functions  $f_\Delta$  and  $g_\Delta$  we may redefine SDEs (7) and (8), which now read

$$y_t^\Delta = y_{t_n}^\Delta + \int_{t_n}^t f_\Delta(t_n, s, y_{t_n}^\Delta, y_s^\Delta) ds + \int_{t_n}^t g_\Delta(t_n, s, y_{t_n}^\Delta, y_s^\Delta) dW_s, \quad t \in (t_n, t_{n+1}], \tag{13}$$

and

$$y_t^\Delta = y_0 + \int_0^t f_\Delta(\hat{s}, s, y_s^\Delta, y_s^\Delta) ds + \int_0^t g_\Delta(\hat{s}, s, y_s^\Delta, y_s^\Delta) dW_s. \tag{14}$$

respectively, with  $y_0 = x_0$  a.s.

**Assumption 2.2** *Let the truncated versions  $f_\Delta(s, r, x, y), g_\Delta(s, r, x, y)$  of  $f, g$  satisfy the following condition ( $\phi_\Delta \equiv f_\Delta, g_\Delta$ )*

$$|\phi_\Delta(s_1, r_1, x_1, y_1) - \phi_\Delta(s_2, r_2, x_2, y_2)| \leq h(\Delta) (|s_1 - s_2| + |r_1 - r_2| + |x_1 - x_2| + |y_1 - y_2|)$$

for all  $0 < \Delta \leq 1$  and  $x_1, x_2, y_1, y_2 \in \mathbb{R}$ , where  $h(\Delta)$  is as in (11).

Let us also assume that the coefficients  $a(t, x), b(t, x)$  of the original SDE satisfy the Khasminskii-type condition.

**Assumption 2.3** *We assume the existence of constants  $p \geq 2$  and  $C_K > 0$  such that  $x_0 \in \mathcal{L}^p(\Omega, \mathbb{R})$  and*

$$xa(t, x) + \frac{p-1}{2}b(t, x)^2 \leq C_K(1 + |x|^2)$$

for all  $(t, x) \in [0, T] \times \mathbb{R}$ .

Under the local Lipschitz and the Khasminskii-type condition SDE (6) has a unique solution and finite moment bounds of order  $p$ , c.f. [24], i.e. for all  $T > 0$ , there exists a constant  $A > 0$  such that  $\sup_{0 \leq t \leq T} \mathbb{E}|x_t|^p < A$ . We rewrite the main result [30, Th. 3.1].

**Theorem 2 (Order of Strong Convergence)** *Suppose Assumption 2.2 and Assumption 2.3 hold and (13) has a unique strong solution for every  $n \leq N - 1$ , where  $x_0 \in \mathcal{L}^p(\Omega, \mathbb{R})$  for some  $p \geq 14 + 2\gamma$ . Let  $\epsilon \in (0, 1/3)$  and define for  $\gamma > 0$*

$$\mu(u) = \bar{C}u^{1+\gamma}, \quad u \geq 0 \quad \text{and} \quad h(\Delta) = \bar{C} + \sqrt{\ln \Delta^{-\epsilon}}, \quad \Delta \in (0, 1],$$

where  $\Delta \leq 1$  and  $\hat{h}$  are such that (11) holds. Then the semi-discrete numerical scheme (14) converges to the true solution of (6) in the  $\mathcal{L}^2$ -sense with order arbitrarily close to  $1/2$ , that is

$$\mathbb{E} \sup_{0 \leq t \leq T} |y_t^\Delta - x_t|^2 \leq C \Delta^{1-\epsilon}. \tag{15}$$

### 3 Applications of the Semi-discrete Method: Mathematical Finance

#### 3.1 3/2-Model

Let us first consider the more general 3/2-model (2) with super-linear drift and diffusion coefficients, see [15, Sec. 4.1],

$$x_t = x_0 + \int_0^t (k_1(s)x_s - k_2(s)x_s^2)ds + \int_0^t k_3(s)x_s^{3/2}\phi(x_s)dW_s, \quad t \in [0, T], \tag{16}$$

where  $\phi(\cdot)$  is a locally Lipschitz and bounded function with locally Lipschitz constant  $C_R^\phi$ , bounding constant  $K_\phi$ ,  $x_0$  is independent of all  $\{W_t\}_{0 \leq t \leq T}$ ,  $x_0 \in \mathcal{L}^{4p}(\Omega, \mathbb{R})$  for some  $2 < p$  and  $x_0 > 0$  a.s.,  $\mathbb{E}(x_0)^{-2} < A$ ,  $k_1(\cdot), k_2(\cdot), k_3(\cdot)$  are positive and bounded functions with  $k_{2,\min} > \frac{7}{2}(K_\phi k_{3,\max})^2$ . It holds that  $x_t > 0$  a.s. The following semi-discrete numerical scheme,

$$y_t = y_0 + \int_0^t (k_1(s) - k_2(s)y_s)y_s ds + \int_0^t k_3(s)\sqrt{y_s}\phi(y_s)y_s dW_s, \tag{17}$$

where  $\hat{s} = t_n$ , when  $s \in [t_n, t_{n+1})$ , produces a linear SDE with solution

$$y_t = x_0 \exp \left\{ \int_0^t \left( k_1(s) - k_2(s)y_{\hat{s}} - k_3^2(s) \frac{y_{\hat{s}}^2 \phi^2(y_{\hat{s}})}{2} \right) ds + \int_0^t k_3(s)\sqrt{y_{\hat{s}}}\phi(y_{\hat{s}})dW_s \right\}, \tag{18}$$

where  $y_t = y_t(t_0, x_0)$ . We call (18) an exponential semi-discrete approximation of (16). The exponential semi-discrete numerical scheme (18) converges to the true solution of (16) in the mean square sense, is positive and has finite moments  $\mathbb{E}(\sup_{0 \leq t \leq T} (y_t)^p)$  for appropriate  $p$ , see [15, Sec. 4.1]. See also the very recent work [17], a combination of the Lamperti transformation with the SD method, named LSD method.

### 3.2 CEV Process

The following SDE

$$x_t = x_0 + \int_0^t (k_1 - k_2 x_s) ds + \int_0^t k_3 (x_s)^q dW_s, \quad t \in [0, T], \tag{19}$$

where  $k_1, k_2, k_3$  are positive and  $1/2 < q < 1$  is known as a mean-reverting CEV process. Equation (19) may represent the instantaneous volatility or the instantaneous variance of the underlying financially observable. Here the diffusion coefficient is sub-linear. Feller’s test implies that there is a unique non-explosive strong solution such that  $x_t > 0$  a.s. when  $x_0 > 0$  a.s. c.f. [22, Prop. 5.22]. The steady-state level of  $x_t$  is  $k_1/k_2$  and the rate of mean-reversion is  $k_2$ .

Here we examine two versions of an implicit SD scheme that are solved explicitly. In [14], we propose

$$y_t = x_0 + \int_0^t (k_1 - k_2(1 - \theta)y_{\hat{s}} - k_2\theta y_{\tilde{s}}) ds + k_3 \int_0^t (y_{\hat{s}})^{q-\frac{1}{2}} \sqrt{y_s} dW_s + \int_t^{t_{n+1}} \left( k_1 - k_2(1 - \theta)y_{t_n} - \frac{(k_3)^2}{4(1 + k_2\theta\Delta)} (y_{t_n})^{2q-1} - k_2\theta y_t \right) ds, \tag{20}$$

for  $t \in (t_n, t_{n+1}]$  where

$$\hat{s} = t_j, s \in (t_j, t_{j+1}], j=0, \dots, n, \quad \tilde{s} = \begin{cases} t_{j+1}, & \text{for } s \in [t_j, t_{j+1}], j=0, \dots, n-1 \\ t, & \text{for } s \in [t_n, t], \end{cases}$$

and  $\theta \in [0, 1]$  represents the level of implicitness. After rearranging

$$y_t(q) = y_n + \int_{t_n}^t \frac{(k_3)^2}{4(1 + k_2\theta\Delta)^2} (y_{t_n})^{2q-1} ds + \frac{k_3}{1 + k_2\theta\Delta} (y_{t_n})^{q-\frac{1}{2}} \int_{t_n}^t \text{sgn}(z_s) \sqrt{y_s} dW_s, \tag{21}$$

with solution

$$y_t(q) = (z_t)^2, \quad z_t := \sqrt{y_n} + \frac{k_3}{2(1 + k_2\theta\Delta)} (y_{t_n})^{q-\frac{1}{2}} (W_t - W_{t_n}), \tag{22}$$

where  $y_n$  is

$$y_n := y_{t_n} \left( 1 - \frac{k_2\Delta}{1 + k_2\theta\Delta} \right) + \frac{k_1\Delta}{1 + k_2\theta\Delta} - \frac{(k_3)^2}{4(1 + k_2\theta\Delta)^2} (y_{t_n})^{2q-1} \Delta.$$

The SD method (22) is positive by construction and under some conditions on the coefficients  $k_i$ , the level of implicitness  $\theta$  and the step-size  $\Delta$ , it strongly converges to the solution of (19) with a logarithmic rate if also  $\mathbb{E}(x_0)^p < A$  for some  $p \geq 4$

and with a polynomial rate of convergence of magnitude  $\frac{1}{2}(q - \frac{1}{2})$  if  $x_0 \in \mathbb{R}$ , see [14, Th.1 and Th.2]. The other version of the implicit SD scheme, see [12], is written in each sub-interval,

$$\tilde{y}_t(q) = \tilde{y}_n + \int_{t_n}^t \frac{q(k_3)^2}{2} (\tilde{y}_s)^{2q-1} ds + k_3 \int_{t_n}^t \text{sgn}(\tilde{z}_s) (\tilde{y}_s)^q dW_s \tag{23}$$

with solution

$$\tilde{y}_t(q) = |\tilde{z}_t|^{1/(1-q)}, \quad \tilde{z}_t := (\tilde{y}_n)^{1-q} + k_3(1 - q)(W_t - W_{t_n}), \tag{24}$$

where

$$\tilde{y}_n := \tilde{y}_{t_n} (1 - k_2 \Delta) + k_1 \Delta - \frac{q(k_3)^2 \Delta}{2} (\tilde{y}_{t_n})^{2q-1}.$$

The SD method (24) is again positive by construction and under some conditions on the coefficients  $k_i$ , the level of implicitness  $\theta$  and the step-size  $\Delta$ , it strongly converges to the solution of (19) with a polynomial rate of convergence of magnitude  $q(q - \frac{1}{2})$  if  $x_0 \in \mathbb{R}$ . See also how LSD performs [17].

### 3.3 CIR/CEV Delay Models with Jump

Here we study a general model of type (19) including delay and jump terms. In particular we consider the following stochastic delay differential equation (SDDE) with jump,

$$x_t = \begin{cases} \xi_0 + \int_0^t (k_1 - k_2 x_{s-}) ds + \int_0^t k_3 b(x_{s-\tau}) x_{s-}^\alpha dW_s + \int_0^t g(x_{s-}) d\tilde{N}_s, & t \in [0, T], \\ \xi(t), & t \in [-\tau, 0], \end{cases} \tag{25}$$

where  $x_{s-} = \lim_{r \uparrow s} x_r$ , the coefficient  $b \in \mathcal{C}(\mathbb{R}_+, \mathbb{R}_+)^1$  and is assumed to be  $\gamma$ -Hölder continuous with  $\gamma > 0$ , the jump coefficient  $g : \mathbb{R} \mapsto \mathbb{R}$  is assumed deterministic for simplicity, the function  $\xi \in \mathcal{C}([-\tau, 0], (0, \infty))$  and  $\tau > 0$  is a positive constant which represents the delay. Process  $\tilde{N}(t) = N(t) - \lambda t$  a compensated Poisson process with intensity  $\lambda > 0$  independent of  $W_t$ . (25) has a unique and nonnegative solution and under some conditions on  $\|\xi\|$  and the step-size  $\Delta$  the following scheme strongly converges to the solution of (25) with polynomial or logarithmic rate, see [29],

---

<sup>1</sup> $\mathcal{C}(A, B)$  the space of continuous functions  $\phi : A \mapsto B$  with norm  $\|\phi\| = \sup_{u \in A} \phi(u)$ .

$$\begin{cases} y_{t_{k+1}^-} = (z_{t_{k+1}})^2, \\ y_{t_{k+1}} = y_{t_{k+1}^-} + g(y_{t_{k+1}^-}) \Delta \tilde{N}_k, \end{cases} \tag{26}$$

where

$$z_t = \sqrt{y_{t_k} \left( 1 - \frac{k_2 \Delta_k}{1 + k_2 \theta \Delta_k} \right) + \frac{k_1 \Delta_k}{1 + k_2 \theta \Delta_k} - \frac{(k_3)^2}{4(1 + k_2 \theta \Delta_k)^2} \frac{b^2(y_{t_{k-\tau}})}{(1 + b(y_{t_{k-\tau}}) \Delta_k^m)^2} (y_{t_k})^{2\alpha-1} \Delta_k} + \frac{k_3}{2(1 + k_2 \theta \Delta_k)} \frac{b(y_{t_{k-\tau}})}{1 + b(y_{t_{k-\tau}}) \Delta_k^m} (y_{t_k})^{\alpha-\frac{1}{2}} (W_t - W_{t_k})}$$

$y_t = \xi(t)$  when  $t \in [-\tau, 0]$  and for  $k = 0, 1, \dots, n_T - 1$ , and  $\Delta_k = t_{k+1} - t_k$ ,  $\Delta \tilde{N}_k := \tilde{N}(t_{k+1}) - \tilde{N}(t_k) = \Delta N_k - \lambda \Delta_k$  and  $\theta \in [0, 1]$  represents the level of implicitness, with  $m = 1/4$ . The SD scheme (26) combines the semi-discrete idea with a taming procedure. For the case  $\alpha = 1/2$ , known as the CIR model, where no delay and jump terms, see also [7, 11] and the application of the LSD method [17]. For extensions of the SD method to the two-factor CIR, see [10].

### 3.4 Ait-Sahalia Model

Let

$$x_t = x_0 + \int_0^t \left( \frac{a_1}{x_s} - a_2 + a_3 x_s - a_4 x_s^r \right) ds + \sigma \int_0^t x_s^\rho dW_s, \tag{27}$$

where  $x_0 > 0$ , the coefficients  $a_i$  are nonnegative and  $r > 1, \rho > 1$ . SDE (27), known as the Ait-Sahalia model, is used as an interest rate model and satisfies  $x_t > 0$  a.s. The approximation of (27), by a combination of the splitting step method and the semi-discrete method, is proposed in [13]. In fact the SD approximation for the transformed process  $z_t = x_t^2$  takes place first with dynamics given by

$$z_t = z_0 + \int_0^t (2a_1 z_s - 2a_2 \sqrt{z_s} + 2a_3 z_s - 2a_4 z_s^{(r+1)/2} + \sigma^2 z_s^\rho) ds + 2\sigma \int_0^t z_s^{(\rho+1)/2} dW_s. \tag{28}$$

Splitting (28) in each subinterval with  $t \in [t_n, t_{n+1}]$  as

$$z_1(t) = z_2(t_n) + \int_{t_n}^t (\ln(4/3) z_1(s) - 2a_2 \sqrt{z_1(s)}) ds \tag{29}$$

$$z_2(t) = z_1(t_{n+1}) + \int_{t_n}^t (2a_1 + (2a_3 - \ln(4/3)) z_2(s) - 2a_4 z_2^{(r+1)/2}(s) + \sigma^2 z_2^\rho(s)) ds + 2\sigma \int_{t_n}^t z_2^{(\rho+1)/2}(s) dW_s, \tag{30}$$

where  $z_2(0) = x_0$  suggests that we may take the solution of (29)

$$z_1(t) = \left( \frac{2a_2}{\ln(4/3)} + \left( \sqrt{z_2(t_n)} - \frac{2a_2}{\ln(4/3)} \right) \left( \frac{4}{3} \right)^{(t-t_n)/2} \right)^2 \tag{31}$$

and approximate (30) with

$$\begin{aligned} \tilde{z}_2(t) &= z_1(t_{n+1}) + 2a_1 \Delta \\ &+ \int_{t_n}^t \left( 2a_3 - \ln(4/3) - 2a_4 \tilde{z}_2^{(r-1)/2}(\hat{s}) + \sigma^2 \tilde{z}_2^{(\rho-1)/2}(\hat{s}) \right) \tilde{z}_2(s) ds \\ &+ 2\sigma \int_{t_n}^t \tilde{z}_2^{(\rho-1)/2}(\hat{s}) \tilde{z}_2(s) dW_s. \end{aligned} \tag{32}$$

We end up with the following SD numerical scheme for the transformed process  $z_t$

$$\begin{aligned} \tilde{z}_{n+1} &= \left( 2a_1 \Delta + \left( \frac{2a_2}{\ln(4/3)} + \left( \sqrt{\tilde{z}_n} - \frac{2a_2}{\ln(4/3)} \right) \left( \frac{4}{3} \right)^{\Delta/2} \right)^2 \right) \\ &\times \exp\{ (2a_3 - \ln(4/3) - 2a_4 \tilde{z}_n^{(r-1)/2} - \sigma^2 \tilde{z}_n^{\rho-2}) \Delta + 2\sigma \tilde{z}_n^{(\rho-1)/2} \Delta W_n \} \end{aligned} \tag{33}$$

and then take  $y_n = \sqrt{\tilde{z}_n}$  for the approximation of the original Ait-Sahalia model, which is positive, strongly convergent with finite moment bounds, when  $r + 1 > 2\rho$ , with  $\rho \geq 2$ , see [13]. See also the performance of LSD [17].

## 4 Applications of the Semi-discrete Method: Population Dynamics and Biology

### 4.1 Wright-Fisher Model

The next class of SDEs appears in population dynamics to describe fluctuations in gene frequency of reproducing individuals among finite populations [5] and ion channel dynamics within cardiac and neuronal cells, (cf. [3, 4, 6] and references therein),

$$x_t = x_0 + \int_0^t (k_1 - k_2 x_s) ds + k_3 \int_0^t \sqrt{x_s(1-x_s)} dW_s, \tag{34}$$

where  $k_i > 0, i = 1, 2, 3$ . If  $x_0 \in (0, 1)$  and  $(k_1 \wedge (k_2 - k_1)) \geq (k_3)^2/2$ , then  $0 < x_t < 1$  a.s. The process

$$\begin{aligned}
 y_t = & y_{t_n} + \int_{t_n}^{t_{n+1}} \left( k_1 - \frac{(k_3)^2}{4} + y_{t_n} \left( \frac{(k_3)^2}{2} - k_2 \right) \right) ds + \int_{t_n}^t \frac{(k_3)^2}{4} (1 - 2y_s) ds \\
 & + k_3 \int_{t_n}^t \sqrt{y_s(1 - y_s)} \operatorname{sgn}(z_s) dW_s,
 \end{aligned} \tag{35}$$

for  $t \in (t_n, t_{n+1}]$ , with  $y_0 = x_0$  a.s. and  $z_t = \sin(k_3 \Delta W_n^t + 2 \arcsin(\sqrt{y_n}))$ , where  $y_n := y_{t_n} + \left( k_1 - \frac{(k_3)^2}{4} + y_{t_n} \left( \frac{(k_3)^2}{2} - k_2 \right) \right) \cdot \Delta$  has the following solution

$$y_t = \sin^2 \left( \frac{k_3}{2} \Delta W_n^t + \arcsin(\sqrt{y_n}) \right), \tag{36}$$

which has the pleasant feature that  $y_t \in (0, 1)$  when  $y_0 \in (0, 1)$ . Process (36) is well defined when  $0 < y_n < 1$ , which is achieved for appropriate  $\Delta$ . To simplify conditions on the parameters and the step size  $\Delta$  we may adopt the strategy presented in [28] considering a perturbation of order  $\Delta$  in the initial condition. Here we used an additive discretization of the drift coefficient and the eternal life time SD scheme (36) strongly converges to the solution of (34), see [28]. Moreover, in [28], an application of the SD method to an extension of the Wright-Fisher model to the multidimensional case is treated, producing a strongly converging and boundary preserving scheme.

### 4.2 Predator-Prey Model

The following system of SDEs, c.f. [20],

$$\begin{aligned}
 X_t^{(1)} &= X_0^{(1)} + \int_0^t (aX_s^{(1)} - bX_s^{(1)}X_s^{(2)}) ds + \int_0^t k_1 X_s^{(1)} dW_s^{(1)}, \\
 X_t^{(2)} &= X_0^{(2)} + \int_0^t (cX_s^{(1)}X_s^{(2)} - dX_s^{(2)}) ds + \int_0^t k_2 X_s^{(2)} dW_s^{(2)},
 \end{aligned}$$

where  $a, b, c, d > 0$  and  $k_1, k_2 \in \mathbb{R}$  with independent Brownian motions  $W_t^{(1)}, W_t^{(2)}$  was studied in [9]. Under some moment bound conditions for  $(X_t^{(i)})$ ,  $i = 1, 2$  and when  $X_0^{(1)} > 0$  and  $X_0^{(2)} > 0$  then  $X_t^{(1)} > 0$  and  $X_t^{(2)} > 0$  a.s. Transforming the second equation  $Z_t^{(2)} = \ln(X_t^{(2)})$  produces the following system

$$\begin{aligned}
 X_t^{(1)} &= X_0^{(1)} + \int_0^t (a - be^{Z_s^{(2)}}) X_s^{(1)} ds + \int_0^t k_1 X_s^{(1)} dW_s^{(1)}, \\
 Z_t^{(2)} &= Z_0^{(2)} + \int_0^t (cX_s^{(1)} - d - (k_2)^2) ds + k_2 W_t^{(2)},
 \end{aligned}$$

which is approximated by the following SD scheme

$$Y_t^{(1)} = X_0^{(1)} + \int_0^t (a - be^{Y_s^{(2)}})Y_s^{(1)} ds + \int_0^t k_1 Y_s^{(1)} dW_s^{(1)},$$

$$Y_t^{(2)} = Y_0^{(2)} + \int_0^t (cY_s^{(1)} - d - (k_2)^2) ds + k_2 W_t^{(2)},$$

which reads

$$Y_{t_{n+1}}^{(1)} = Y_{t_n}^{(1)} \exp\left\{(a - be^{Y_{t_n}^{(2)}} - \frac{(k_1)^2}{2})\Delta + k_1 \Delta W_n^{(1)}\right\}$$

$$Y_{t_{n+1}}^{(2)} = Y_{t_n}^{(2)} + (cY_{t_n}^{(1)} - d - (k_2)^2)\Delta + k_2 \Delta W_n^{(2)}.$$

## References

1. D-H. Ahn, B. Gao, A parametric nonlinear model of term structure dynamics. *Rev. Financ. Stud.* **12**(4), 721–762 (1999)
2. J.A.D. Appleby, M. Guzowska, K. Cónall, A. Rodkina, Preserving positivity in solutions of discretised stochastic differential equations. *Appl. Math. Comput.* **217**(2), 763–774 (2010)
3. C.E. Dangerfield, D. Kay, S. MacNamara, K. Burrage, A boundary preserving numerical algorithm for the Wright-Fisher model with mutation. *BIT Numer. Math.* **52**(2), 283–304 (2012)
4. C.E. Dangerfield, D. Kay, K. Burrage, Modeling ion channel dynamics through reflected stochastic differential equations. *Phys. Rev. E* **85**(5), 051907 (2012)
5. W.J. Ewens, *Mathematical Population Genetics I: Theoretical Introduction*, vol. 27 (Springer Science & Business Media, New York, 2012)
6. J.H. Goldwyn, N.S. Imenov, M. Famulare, E. Shea-Brown, Stochastic differential equation models for ion channel noise in Hodgkin-Huxley neurons. *Phys. Rev. E* **83**(4), 041908 (2011)
7. N. Halidias, Semi-discrete approximations for stochastic differential equations and applications. *Int. J. Comput. Math.* **89**(6), 780–794 (2012)
8. N. Halidias, A novel approach to construct numerical methods for stochastic differential equations. *Numer. Algorithms* **66**(1), 79–87 (2014)
9. N. Halidias, Construction of positivity preserving numerical schemes for some multidimensional stochastic differential equations. *Discrete Contin. Dynam. Syst. B* **20**(1), 153–160 (2015)
10. N. Halidias, Constructing positivity preserving numerical schemes for the two-factor CIR model. *Monte Carlo Methods Appl.* **21**(4), 313–323 (2015)
11. N. Halidias, A new numerical scheme for the CIR process. *Monte Carlo Methods Appl.* **21**(3), 245–253 (2015)
12. N. Halidias, An explicit and positivity preserving numerical scheme for the mean reverting CEV model. *Japan J. Ind. Appl. Math.* **32**(2), 545–552 (2015)
13. N. Halidias, *On the construction of boundary preserving numerical schemes*. *Monte Carlo Methods Appl.* **22**, 277–289 (2016)
14. N. Halidias, I.S. Stamatiou, Approximating explicitly the mean-reverting CEV process. *J. Probab. Stat.* (2015), Article ID 513137, 20 pp.
15. N. Halidias, I.S. Stamatiou, On the numerical solution of some non-linear stochastic differential equations using the semi-discrete method. *Comput. Methods Appl. Math.* **16**(1), 105–132 (2016)



16. N. Halidias, I.S. Stamatiou, A note on the asymptotic stability of the semi-discrete method for stochastic differential equations (2020). <https://arxiv.org/abs/2008.03148>
17. N. Halidias, I.S. Stamatiou, Lamperti Semi-Discrete method (2020). <http://arxiv.org/abs/2104.06149>
18. S.L. Heston, A simple new formula for options with stochastic volatility, OLIN-97-23 (1997). <http://ssrn.com/abstract=86074>
19. M. Hutzenthaler, A. Jentzen, P.E. Kloeden, Strong and weak divergence in finite time of Euler's method for stochastic differential equations with non-globally Lipschitz continuous coefficients. *Proc. R. Soc. Lond. A Math. Phys. Eng. Sci.* **467**(2130), 1563–1576 (2011)
20. M. Hutzenthaler, A. Jentzen, Numerical approximations of stochastic differential equations with non-globally Lipschitz continuous coefficients. *Memoirs Am. Math. Soc.* **236**(1112), 1611–1641 (2015)
21. C. Kahl, M. Günther, T. Rossberg, Structure preserving stochastic integration schemes in interest rate derivative modeling. *Appl. Numer. Math.* **58**(3), 284–295 (2008)
22. I. Karatzas, S.E. Shreve, *Brownian Motion and Stochastic Calculus* (Springer, New York, 1988)
23. P. Kloeden, A. Neuenkirch, Convergence of numerical methods for stochastic differential equations in mathematical finance, in *Recent Developments in Computational Finance: Foundations, Algorithms and Applications* (2013), pp. 49–80
24. X. Mao, *Stochastic Differential Equations and Applications*, 2nd ed. (Horwood Publishing, Chichester, 2007)
25. X. Mao, L. Szpruch, Strong convergence and stability of implicit numerical methods for stochastic differential equations with non-globally Lipschitz continuous coefficients. *J. Comput. Appl. Math.* **238**, 14–28 (2013)
26. X. Mao, L. Szpruch, Strong convergence rates for backward Euler–Maruyama method for non-linear dissipative-type stochastic differential equations with super-linear diffusion coefficients. *Stochastics* **85**(1), 144–171 (2013)
27. S. Sabanis, Euler approximations with varying coefficients: the case of superlinearly growing diffusion coefficients. *Ann. Appl. Probab.* **26**(4), 2083 (2016)
28. I.S. Stamatiou, A boundary preserving numerical scheme for the Wright–Fisher model. *J. Comput. Appl. Math.* **328**, 132–150 (2018)
29. I.S. Stamatiou, An explicit positivity preserving numerical scheme for CIR/CEV type delay models with jump. *J. Comput. Appl. Math.* **360**, 78–98 (2019)
30. I.S. Stamatiou, N. Halidias, Convergence rates of the semi-discrete method for stochastic differential equations. *Theory Stoch. Process.* **24**(40), 89–100 (2019)
31. M.V. Tretyakov, Z. Zhang, A fundamental mean-square convergence theorem for SDEs with locally Lipschitz coefficients and its applications. *SIAM J. Numer. Anal.* **51**(6), 3135–3162 (2013)

# Homotopic Metric-Interval L-Contractions in Gauge Spaces



Mihai Turinici

**Abstract** A functional version—under the lines in Leader [Math. Japonica, 24 (1979), 17–24]—is given for the 1967 contraction mapping principle in Gheorghiu [Stud. Cerc. Mat., 19 (1967), 119–122]. As a by-product of this, an appropriate functional extension is proposed for the homotopic fixed point result in gauge spaces due to Frigon [L. Notes Nonlin. Anal., 16 (2017), 9–91].

## 1 Introduction

Let  $X$  be a nonempty set. Call the subset  $Y$  of  $X$ , *almost-singleton* (in short: *asingleton*), provided [ $y_1, y_2 \in Y$  implies  $y_1 = y_2$ ]; and *singleton* if, in addition,  $Y$  is nonempty; note that in this case  $Y = \{y\}$ , for some  $y \in X$ .

Further, let  $d : X \times X \rightarrow R_+ := [0, \infty[$  be a *metric* over  $X$ ; and take some  $T \in \mathcal{F}(X)$ . [Here, given the nonempty sets  $A$  and  $B$ ,  $\mathcal{F}(A, B)$  stands for the class of all functions from  $A$  to  $B$ ; when  $A = B$ , we write  $\mathcal{F}(A, A)$  as  $\mathcal{F}(A)$ ]. Denote  $\text{Fix}(T) := \{z \in X; z = Tz\}$ ; any point of it will be called *fixed* under  $T$ . These points are to be determined in the context below (cf. Rus [50, Ch 2, Sect 2.2]):

- (pic-1) We say that  $T$  is *fix-asingleton*, if  $\text{Fix}(T)$  is an asingleton; and *fix-singleton*, if  $\text{Fix}(T)$  is a singleton
- (pic-2) We say that  $x \in X$  is a *Picard point* (modulo  $(d, T)$ ) when  $(T^n x; n \geq 0)$  is  $d$ -Cauchy. If this property holds for all  $x \in X$ , we say that  $T$  is a *Picard operator* (modulo  $d$ )

---

**AMS Subject Classification:** 47H10 (Primary), 54H25 (Secondary).

---

M. Turinici (✉)

A. Myller Mathematical Seminar, A. I. Cuza University, Iași, Romania  
e-mail: [mturi@uaic.ro](mailto:mturi@uaic.ro)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,  
Springer Optimization and Its Applications 173,  
[https://doi.org/10.1007/978-3-030-72563-1\\_24](https://doi.org/10.1007/978-3-030-72563-1_24)

639

(pic-3) We say that  $x \in X$  is a *strongly Picard point* (modulo  $(d, T)$ ) when  $(T^n x; n \geq 0)$  is  $d$ -convergent with  $\lim_n(T^n x) \in \text{Fix}(T)$ . If this property holds for all  $x \in X$ , we say that  $T$  is a *strongly Picard operator* (modulo  $d$ ).

In this perspective, a basic answer to the posed question [referred to as *Banach contraction principle*; in short: (B-cp)], may be stated as follows. Given  $k \geq 0$ , let us say that  $T$  is *(Banach)  $(d; k)$ -contractive*, provided

$$(B\text{-contr}) \quad d(Tx, Ty) \leq kd(x, y), \text{ for all } x, y \in X.$$

**Theorem 1** *Suppose that  $T$  is  $(d; k)$ -contractive, for some  $k \in [0, 1[$ . In addition, let  $X$  be  $d$ -complete. Then,*

(11-a)  *$T$  is fix-singleton:  $\text{Fix}(T) = \{z\}$ , for some  $z \in X$*

(11-b)  *$T$  is a strongly Picard operator (modulo  $d$ ):  $T^n x \xrightarrow{d} z, \forall x \in X$ .*

This result, obtained in 1922 by Banach [7], found a multitude of applications in operator equations theory; so, it was the subject of many extensions. Essentially, there are two main directions of doing this.

(I) The initial metric  $d$  remains as it is, but the contractive condition is taken in the implicit way

$$(i\text{-contr}) \quad F(d(Tx, Ty), d(x, y), d(x, Tx), d(y, Ty), d(x, Ty), d(Tx, y)) \leq 0,$$

for all  $x, y \in X, x \mathcal{R} y$

where  $F : R_+^6 \rightarrow R$  is a function and  $\mathcal{R} \subseteq X \times X$  is a *relation* over  $X$ . When the function  $F$  appearing here admits the explicit form

$$F(t_1, t_2, t_3, t_4, t_5, t_6) = t_1 - G(t_2, t_3, t_4, t_5, t_6), (t_1, t_2, t_3, t_4, t_5, t_6) \in R_+^6$$

(where  $G : R_+^5 \rightarrow R_+$  is a function), one gets the explicit functional version of this (functional) contraction

$$(e\text{-contr}) \quad d(Tx, Ty) \leq G(d(x, y), d(x, Tx), d(y, Ty), d(x, Ty), d(Tx, y)),$$

for all  $x, y \in X, x \mathcal{R} y$ .

In particular, when  $\mathcal{R} = X \times X$  (the *trivial relation* over  $X$ ), some outstanding explicit results have been established in Boyd and Wong [11], Reich [48], and Matkowski [40]; see also the survey paper by Rhoades [49]. And, for the implicit setting above, certain technical aspects have been considered by Leader [39] and Turinici [54]. On the other hand, when  $\mathcal{R}$  is a (*partial*) *order* on  $X$ , some appropriate extensions of Matkowski fixed point principle we just quoted were obtained in the 1986 papers by Turinici [60, 61]; two decades later, these results have been re-discovered—at the level of (Banach) contractive maps—by

Ran and Reurings [47]; see also Nieto and Rodriguez-Lopez [45]. Further, an extension—to the same framework—of Leader’s contribution was performed in Agarwal et al. [2] and Turinici [65]. Finally, when  $\mathcal{R}$  is a (general) relation, some results in this direction were obtained by Jachymski [34], within a graph setting, and by Samet and Turinici [51] under a general perspective.

- (II) The initial metric  $d$  is substituted by a separated family  $D = (d_\lambda; \lambda \in \Lambda)$  of semimetrics on  $X$ ; and the contractive condition is to written as

$$(G\text{-contr}) \quad d_\lambda(Tx, Ty) \leq k_\lambda d_{\varphi(\lambda)}(x, y), \text{ for all } x, y \in X, \text{ and all } \lambda \in \Lambda;$$

where  $(k_\lambda; \lambda \in \Lambda)$  is a family of positive numbers and  $\varphi : \Lambda \rightarrow \Lambda$  is a mapping. The first fixed point principle in this direction was obtained in 1967 by Gheorghiu [25]. Further refinement of it were obtained in Gheorghiu and Rotaru [27]; see also Gheorghiu [26]. Note that, by the very intricate form of contractive condition, genuine functional extensions of this principle were not yet obtained.

Recently, an interesting application of Gheorghiu contraction principle to the homotopic fixed point theory was obtained in Frigon [23]. It is our aim in the following to state a functional refinement of this result. We use here a metrical maximality principle as well as a metric interval, in contrast with—respectively—the Zorn-Bourbaki maximal one [9] and the standard (real) interval considered in that paper. Further aspects will be delineated elsewhere.

## 2 Conv-Cauchy Structures

Throughout this exposition, the axiomatic system in use is Zermelo-Fraenkel’s (abbreviated: ZF); cf. Cohen [16, Ch 2]. The notations and basic facts to be considered are standard. Some important ones are described below.

- (A) Let  $X$  be a nonempty set. By a *relation* over  $X$ , we mean any (nonempty) part  $\mathcal{R} \subseteq X \times X$ ; then,  $(X, \mathcal{R})$  will be referred to as a *relational structure*. For simplicity, we sometimes write  $(x, y) \in \mathcal{R}$  as  $x\mathcal{R}y$ . Note that  $\mathcal{R}$  may be regarded as a mapping between  $X$  and  $\exp[X]$  (=the class of all subsets in  $X$ ). In fact, denote for  $x \in X$ :

$$X(x, \mathcal{R}) = \{y \in X; x\mathcal{R}y\} \text{ (the } \textit{section} \text{ of } \mathcal{R} \text{ through } x);$$

then, the desired mapping representation is  $[\mathcal{R}(x) = X(x, \mathcal{R}), x \in X]$ .

A basic example of relational structure is to be constructed as below. Let  $N = \{0, 1, \dots\}$  be the set of *natural* numbers, endowed with the usual addition and (partial) order; note that

$(N, \leq)$  is well ordered: any (nonempty) subset of  $N$  has a first element.

Further, denote for  $p, q \in N, p \leq q$ ,

$$N[p, q] = \{n \in N; p \leq n \leq q\}, N]p, q[ = \{n \in N; p < n < q\},$$

$$N[p, q[ = \{n \in N; p \leq n < q\}, N]p, q] = \{n \in N; p < n \leq q\};$$

as well as, for  $r \in N$ ,

$$N[r, \infty[ = \{n \in N; r \leq n\}, N]r, \infty[ = \{n \in N; r < n\}.$$

By definition,  $N[0, r[ = N(r, >)$  is referred to as the *initial interval* (in  $N$ ) induced by  $r$ . Any set  $P$  with  $P \sim N$  (in the sense: there exists a bijection from  $P$  to  $N$ ) will be referred to as *effectively denumerable*. In addition, given some natural number  $n \geq 1$ , any set  $Q$  with  $Q \sim N(n, >)$  will be said to be *n-finite*; when  $n$  is generic here, we say that  $Q$  is *finite*. Finally, the (nonempty) set  $Y$  is called (at most) *denumerable* iff it is either effectively denumerable or finite.

Let  $X$  be a nonempty set. By a *sequence* in  $X$ , we mean any mapping  $x : N \rightarrow X$ , where  $N = \{0, 1, \dots\}$  is the set of *natural numbers*. For simplicity reasons, it will be useful to denote it as  $(x(n); n \geq 0)$ , or  $(x_n; n \geq 0)$ ; moreover, when no confusion can arise, we further simplify this notation as  $(x(n))$  or  $(x_n)$ , respectively. Also, any sequence  $(y_n := x_{i(n)}; n \geq 0)$  with

$$(i(n); n \geq 0) \text{ is } \textit{divergent} \text{ (} i(n) \rightarrow \infty \text{ as } n \rightarrow \infty \text{)}$$

will be referred to as a *subsequence* of  $(x_n; n \geq 0)$ . Note that, under such a convention, the relation “subsequence of” is transitive; i.e.:

$$(z_n) = \text{subsequence of } (y_n) \text{ and } (y_n) = \text{subsequence of } (x_n)$$

$$\text{imply } (z_n) = \text{subsequence of } (x_n).$$

(B) Remember that, an outstanding part of (ZF) is the *Axiom of Choice* (abbreviated: AC); which, in a convenient manner, may be written as

$$(AC) \text{ For each couple } (J, X) \text{ of nonempty sets and each } F : J \rightarrow \text{exp}(X),$$

$$\text{there exists a (selective) function } f : J \rightarrow X [f(v) \in F(v), \forall v \in J].$$

Here,  $\text{exp}(X)$  stands for the class of all nonempty elements in  $\text{exp}[X]$ . Sometimes, when the index set  $J$  is denumerable, the existence of such a selective function may be determined by using a weaker form of (AC), called: *Dependent Choice* principle (in short: DC). Call the relation  $\mathcal{R}$  over  $X$ , *proper* when

$$(X(x, \mathcal{R}) =) \mathcal{R}(x) \text{ is nonempty, for each } x \in X.$$

Then,  $\mathcal{R}$  is to be viewed as a mapping between  $X$  and  $\exp(X)$ ; and the couple  $(X, \mathcal{R})$  will be referred to as a *proper relational structure*. Further, given  $a \in X$ , let us say that the sequence  $(x_n; n \geq 0)$  in  $X$  is  $(a; \mathcal{R})$ -iterative, provided

$$x_0 = a, \text{ and } (x_n; n \geq 0) \text{ is } \mathcal{R}\text{-increasing } [x_n \mathcal{R} x_{n+1} \text{ (i.e.: } x_{n+1} \in \mathcal{R}(x_n)), \forall n].$$

**Proposition 1** *Let the relational structure  $(X, \mathcal{R})$  be proper. Then, for each  $a \in X$  there is at least an  $(a, \mathcal{R})$ -iterative sequence in  $X$ .*

This principle—proposed, independently, by Bernays [8] and Tarski [53]—is deductible from (AC), but not conversely; cf. Wolk [69]. Moreover, by the developments in Moskhovakis [44, Ch 8], and Schechter [52, Ch 6], the reduced system (ZF-AC+DC) is comprehensive enough so as to cover the “usual” mathematics; see also Moore [43, Appendix 2].

A basic consequence of (DC) is the so-called *Denumerable Axiom of Choice* [in short: AC(N)].

**Proposition 2** *Let  $F : N \rightarrow \exp(X)$  be a function. Then, for each  $a \in F(0)$  there exists a function  $f : N \rightarrow X$  with  $f(0) = a$  and  $(f(n) \in F(n), \forall n)$ .*

**Proof** Denote  $Q = N \times X$ ; and let us introduce the (proper) relation  $\mathcal{R}$  over it, according to:

$$\mathcal{R}(n, x) = \{n + 1\} \times F(n + 1), \quad n \geq 0, x \in X.$$

By an application of (DC) to the proper relational structure  $(Q, \mathcal{R})$  the conclusion follows; we do not give details.

As a consequence of the above facts,

$$(DC) \implies (AC(N)) \text{ in } (ZF-AC); \text{ or, equivalently:}$$

$$(AC(N)) \text{ is deductible in the system } (ZF-AC+DC).$$

The reciprocal of the written inclusion is not true; see Moskhovakis [44, Ch 8, Sect 8.25] for details.

(C) In the following, the concept of conv-Cauchy structure over a metric space is introduced.

Let  $X$  be a nonempty set. Further, let  $d : X \times X \rightarrow R_+$  be a mapping with

(m-1)  $d$  is *triangular*:  $d(x, z) \leq d(x, y) + d(y, z), \forall x, y, z \in X$

(m-2)  $d$  is *reflexive-sufficient*:  $d(x, y) = 0$  iff  $x = y$

(m-3)  $d$  is *symmetric*:  $d(x, y) = d(y, x)$ , for all  $x, y \in X$ .

We then say that  $d(., .)$  is a *metric* on  $X$ ; and the couple  $(X, d)$  will be then referred to as a *metric space*.

Given the sequence  $(x_n)$  in  $X$  and the point  $x \in X$ , we say that  $(x_n)$ ,  $d$ -converges to  $x$  (written as:  $x_n \xrightarrow{d} x$ ), provided  $d(x_n, x) \rightarrow 0$  as  $n \rightarrow \infty$ ; i.e.,

$$\forall \varepsilon > 0, \exists i = i(\varepsilon) : i \leq n \implies d(x_n, x) < \varepsilon; \text{ or, equivalently:}$$

$$\forall \varepsilon > 0, \exists i = i(\varepsilon) : i \leq n \implies d(x_n, x) \leq \varepsilon.$$

The set of all such points  $x$  will be denoted  $\lim_n(x_n)$ ; when it is nonempty, then  $(x_n)$  is called  $d$ -convergent. By this very definition, we have the properties:

(conv-1) ( $\xrightarrow{d}$  is hereditary)

$$x_n \xrightarrow{d} x \text{ implies } y_n \xrightarrow{d} x, \text{ for each subsequence } (y_n) \text{ of } (x_n)$$

(conv-2) ( $\xrightarrow{d}$  is reflexive)

$$(\forall u \in X): \text{ the constant sequence } (x_n = u; n \geq 0) \text{ fulfills } x_n \xrightarrow{d} u.$$

As a consequence,  $(\xrightarrow{d})$  has all properties required in Kasahara [36]; in addition—as  $d$  is triangular symmetric—the following extra property is holding here

(conv-3) ( $\xrightarrow{d}$  is separated (referred to as  $d$  is separated):

$$\lim_n(x_n) \text{ is an asingleton, for each sequence } (x_n) \text{ in } X.$$

The introduced concepts allow us to give a useful property.

**Proposition 3** *The mapping  $(x, y) \mapsto d(x, y)$  is  $d$ -Lipschitz, in the sense*

$$(23-1) |d(x, y) - d(u, v)| \leq d(x, u) + d(y, v), \forall (x, y), (u, v) \in X \times X.$$

*As a consequence, this map is  $d$ -continuous; i.e.,*

$$(23-2) x_n \xrightarrow{d} x, y_n \xrightarrow{d} y \text{ imply } d(x_n, y_n) \rightarrow d(x, y).$$

Further, call  $(x_n)$ ,  $d$ -Cauchy when  $d(x_m, x_n) \rightarrow 0$  as  $m, n \rightarrow \infty, m < n$ ; that is,

$$\forall \varepsilon > 0, \exists j = j(\varepsilon) : j \leq m < n \implies d(x_m, x_n) < \varepsilon; \text{ or, equivalently:}$$

$$\forall \varepsilon > 0, \exists j = j(\varepsilon) : j \leq m < n \implies d(x_m, x_n) \leq \varepsilon.$$

The class of all such sequences will be denoted as  $\text{Cauchy}(d)$ . As before, from this very definition one has the properties

(Cauchy-1) ( $\text{Cauchy}(d)$  is hereditary)

$$(x_n) \text{ is } d\text{-Cauchy implies } (y_n) \text{ is } d\text{-Cauchy, for each subsequence } (y_n) \text{ of } (x_n)$$

(Cauchy-2) ( $\text{Cauchy}(d)$  is reflexive)

$$(\forall u \in X): \text{ the constant sequence } (x_n = u; n \geq 0) \text{ is } d\text{-Cauchy.}$$

Hence,  $\text{Cauchy}(d)$  is a Cauchy structure, under the lines in Turinici [64].

Now—according to the quoted work—term the couple  $((\xrightarrow{d}), Cauchy(d))$ , a *conv-Cauchy structure* induced by  $d$ . The following regularity conditions about this structure are to be (optionally) considered

- (CC-1)  $d$  is *regular*: each  $d$ -convergent sequence in  $X$  is  $d$ -Cauchy
- (CC-2)  $d$  is *complete*: each  $d$ -Cauchy sequence in  $X$  is  $d$ -convergent.

Clearly, the former of these is always obtainable, via  $d$ -triangular symmetric; but the latter one is not in general valid.

(D) In the following, some  $d$ -Cauchy criteria will be stated.

Let us say that the sequence  $(x_n; n \geq 0)$  is  *$d$ -asymptotic*, provided

$$r_n := d(x_n, x_{n+1}) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Clearly, each  $d$ -Cauchy sequence is  $d$ -asymptotic too; the reciprocal of this is not in general true. This tells us that the  $d$ -Cauchy criteria we are looking for are to be sought in the class of  $d$ -asymptotic sequences. To get concrete examples of such properties, we need some conventions and auxiliary facts.

Given the  $d$ -asymptotic sequence  $(x_n; n \geq 0)$  and the number  $\varepsilon > 0$ , let us say that  $i \in N$  is  *$\varepsilon$ -regular*, provided

$$i \leq n \text{ implies } d(x_n, x_{n+1}) < \varepsilon.$$

The class  $\mathcal{Z}(\varepsilon)$  of all these ranks is nonempty; so that

- $(\forall \varepsilon > 0) : Z(\varepsilon) = \min \mathcal{Z}(\varepsilon)$  is well defined, as an element of  $N$ ;
- with, in addition:  $d(x_n, x_{n+1}) < \varepsilon$ , for all  $n \geq Z(\varepsilon)$ .

Define the subsets of  $N \times N$

$$(\leq; N) = \{(m, n) \in N \times N; m \leq n\}, \quad (<; N) = \{(m, n) \in N \times N; m < n\};$$

these are just the graph over  $N$  of the relations  $(\leq)$  and  $(<)$ , respectively.

Further, let us say that the subset  $\Theta$  of  $R_+^0 := ]0, \infty[$  is  *$(>)$ -cofinal* in  $R_+^0$ , when

$$\text{for each } \varepsilon > 0, \text{ there exists } \theta \in \Theta \text{ with } \varepsilon > \theta.$$

Given this subset, define (as before)

$$(\geq; \Theta) = \{(\beta, \gamma) \in \Theta \times \Theta; \beta \geq \gamma\}, \quad (>; \Theta) = \{(\beta, \gamma) \in \Theta \times \Theta; \beta > \gamma\};$$

these are just the graph over  $\Theta$  of the relations  $(\geq)$  and  $(>)$ , respectively. In addition, for each  $(\beta, \gamma) \in (>; \Theta)$ , denote

$$B(\beta, \gamma) = \{n \in N(2, \leq); 3^{-n} < \min\{\beta - \gamma, \gamma\}\};$$



clearly,  $B(\beta, \gamma)$  is nonempty and hereditary:

$$i \in B(\beta, \gamma) \text{ and } i \leq j \text{ imply } j \in B(\beta, \gamma).$$

In this case, it makes sense denoting

$$b(\beta, \gamma) = \min B(\beta, \gamma), (\beta, \gamma) \in (>; \Theta).$$

If no confusion arise, we write  $b(\beta, \gamma)$  as  $b$ , for simplicity; hence, by definition

$$(b \geq 2 \text{ and } k \in N(b, \leq)) \text{ implies } 3^{-k} < \min\{\beta - \gamma, \gamma\}.$$

From the perspective of our initial problem, the points  $\beta \in \Theta$  appearing here must be chosen according to some admissible properties. Some conventions are in order. Let  $(x_n)$  be a sequence in  $X$ . We say that  $\beta \in \Theta$  is  $(x_n)$ -admissible, when

$$\text{(rela-be)} \quad E(j; \beta) := \{(m, n) \in (<; N); j \leq m < n, d(x_m, x_n) > \beta\} \neq \emptyset, \forall j.$$

The class of all these points will be denoted as  $\text{adm}(\Theta; (x_n))$ .

**Proposition 4** *Let  $(x_n)$  be a sequence in  $X$  that is not  $d$ -Cauchy; and  $\Theta$  be a  $(>)$ -cofinal part of  $R_+^0$ . Then, necessarily,*

(24-1)  $\text{adm}(\Theta; (x_n))$  is (nonempty and)  $(>)$ -cofinal in  $R_+^0$

(24-2) for each  $\beta \in \text{adm}(\Theta; (x_n))$  and each  $\gamma \in \Theta$ , we have  $(\beta, \gamma) \in (>; \Theta)$  iff  $(\beta, \gamma) \in (>; \text{adm}(\Theta; (x_n)))$ .

**Proof**

(i): By definition, the  $d$ -Cauchy property of our sequence writes:

$$\forall \varepsilon \in R_+^0, \exists a \in N, \forall (m, n) \in (<; N) : a \leq m < n \implies d(x_m, x_n) \leq \varepsilon.$$

As  $\Theta$  is a  $(>)$ -cofinal part in  $R_+^0$ , this property may be also written as

$$\forall \theta \in \Theta, \exists \alpha \in N, \forall (m, n) \in (<; N) : \alpha \leq m < n \implies d(x_m, x_n) \leq \theta.$$

The negation of this property means: there exists  $\beta \in \Theta$  such that

$$E(j; \beta) := \{(m, n) \in (<; N); j \leq m < n, d(x_m, x_n) > \beta\} \neq \emptyset, \forall j;$$

which shows that  $\beta \in \text{adm}(\Theta; (x_n))$ . Moreover, as

$$(\forall j) : \beta_1 \geq \beta_2 \text{ implies } E(j; \beta_1) \subseteq E(j; \beta_2)$$

it is clear that (under our standard notations)

$$\beta \in \text{adm}(\Theta; (x_n)) \text{ implies } \Theta(\beta, >) \subseteq \Theta(\beta, \geq) \subseteq \text{adm}(\Theta; (x_n));$$

and proves (by means of (>)-cofinal property of  $\Theta$ ) the first conclusion.

(ii): Evident, in view of the preceding stage.

Finally, the following conventions are needed. For each sequence of ranks  $(\lambda(k); k \geq 0)$  in  $N(1, \leq)$ , define the property

$$(\lambda(k); k \geq 0) \text{ is } \textit{strictly ascending} : i < j \text{ implies } \lambda(i) < \lambda(j);$$

$$\text{hence, } (\lambda(k)) \text{ is divergent } (\lim_k \lambda(k) = \infty).$$

On the other hand, for each sequence  $(r_n)$  in  $R$  and each point  $r \in R$ , let us write

$$r_n \rightarrow r + \text{ (also written as: } \lim_n(r_n) = r + \text{) if } r_n \rightarrow r \text{ and } (r_n > r, \forall n)$$

$$r_n \rightarrow \rightarrow r + \text{ (also written as: } \lim_n(r_n) = r + + \text{) if } r_n \rightarrow r \text{ and } (r_n > r, \forall \forall n).$$

Here, given a property  $\pi(k)$  depending on  $k \in N$ , let us say that it holds for *nearly all*  $k$  [written:  $(\pi(k), \forall \forall k)$ ] provided

$$\text{there exists } c = c(\pi) \in N \text{ such that } (\pi(k) \text{ is true, for all } k \geq c).$$

In particular, given the sequence  $(w_k; k \geq 0)$  in  $X$ , the subset  $Y$  of  $X$ , and the property  $[\pi(k)$  holds iff  $w_k \in Y]$ , we introduce the convention

$$(\pi(k) \text{ holds for nearly all } k) \text{ is referred to as } ((w_k) \text{ is } \textit{nearly in } Y).$$

The following result, referred to as *Boyd-Wong non-Cauchy Criterion* (in short: (BW-n-CC)) is now available.

**Theorem 2** *Let the sequence  $(x_n; n \geq 0)$  in  $X$  be such that*

$$(21-i) \text{ } (x_n; n \geq 0) \text{ is } d\text{-asymptotic } (r_n := d(x_n, x_{n+1}) \rightarrow 0 \text{ as } n \rightarrow \infty)$$

$$(21-ii) \text{ } (x_n; n \geq 0) \text{ is not } d\text{-Cauchy.}$$

*Further, let the subset  $\Theta$  of  $R_+^0$  be (>)-cofinal in  $R_+^0$ ; hence (see above)*

$$\text{admc}(\Theta; (x_n)) \text{ (=the class of } (x_n)\text{-admissible couples) is nonempty.}$$

*Then, for each couple  $(\beta, \gamma) \in (>; \text{adm}(\Theta; (x_n)))$  (with the associated rank  $b = b(\beta, \gamma)$ ), and each strictly ascending rank sequence  $(\lambda(k); k \geq 0)$  in  $N(1, \leq)$ , there exists a rank sequence  $(J(k); k \geq 0)$  in  $N(1, \leq)$  and a couple of rank sequences  $(m(k); k \geq 0)$  and  $(n(k); k \geq 0)$  in  $N(1, \leq)$ , so that*

$$(21-a) \text{ } k + 1 \leq J(k) \leq m(k) < m(k) + 3\lambda(k) < n(k), \forall k$$

- (21-b)  $J(k) \leq m(k) < m(k) + 2\lambda(k) < n(k) - 1 < n(k)$ , and  
 $d(x_{m(k)}, x_{n(k)}) > \gamma$ ,  $d(x_{m(k)}, x_{n(k)-1}) \leq 3^{-b-k} + \gamma$ ,  $\forall k$ ,
- (21-c)  $U_k := d(x_{m(k)}, x_{n(k)}) \rightarrow \gamma + as$   $k \rightarrow \infty$
- (21-d) for each couple of rank sequences  $(\mu(k))$  and  $(\nu(k))$  in  $N$  with  
 $(\mu(k), \nu(k) \leq 3\lambda(k), \forall k)$ ,  $V_k := d(x_{m(k)+\mu(k)}, x_{n(k)+\nu(k)}) \rightarrow \gamma + as$   
 $k \rightarrow \infty$
- (21-e) for each couple of ranks  $(i, j)$  in  $N \times N$  with  
 $(i, j \leq 3\lambda(0))$ ,  $S_k := d(x_{m(k)+i}, x_{n(k)+j}) \rightarrow \gamma + as$   $k \rightarrow \infty$
- (21-f) for each couple of ranks  $(i, j)$  in  $N \times N$ , one has  
 $T_k := d(x_{m(k)+i}, x_{n(k)+j}); k \geq 0) \rightarrow \gamma + as$   $k \rightarrow \infty$ .

**Proof** Let the couple  $(\beta, \gamma) \in (>; \text{adm}(\Theta; (x_n)))$  be fixed in the sequel; so,  $[\beta \in \text{adm}(\Theta; (x_n))$  and  $\beta, \gamma) \in (>; \Theta)]$ ; remember that the first property means

$$(\text{rela-1}) \ E(j) := \{(m, n) \in (<; N); j \leq m < n, d(x_m, x_n) > \beta\} \neq \emptyset, \forall j.$$

Further, take a strictly ascending rank sequence  $(\lambda(k); k \geq 0)$  in  $N(1, \leq)$ . With the aid of these data, define the rank sequence  $(J(k); k \geq 0)$  in  $N(1, \leq)$ , according to

$$(J(k) = 1 + k + Z(3^{-2b-k}/\lambda(k)); k \geq 0);$$

where  $b = b(\beta, \gamma)$  (see above) and the mapping  $\varepsilon \mapsto Z(\varepsilon)$  is introduced by the  $d$ -asymptotic property of  $(x_n)$ . Then, denote

$$(A(k) = E(J(k)); k \geq 0); \text{ hence, by definition,}$$

$$A(k) := \{(m, n) \in (<; N); J(k) \leq m < n, d(x_m, x_n) > \beta\}, k \geq 0;$$

with, in addition:  $A(k)$  is a nonempty relation over  $N$ , for each  $k \geq 0$ .

By the triangle inequality (and the choice of  $b$ )

$$(\forall k) : (m, n) \in A(k) \text{ implies}$$

$$d(x_{m+s}, x_{n+t}) \geq d(x_m, x_n) - d(x_m, x_{m+s}) - d(x_n, x_{n+t}) > \\ \beta - 3^{-2b-k+1} - 3^{-2b-k+1} > \beta - 3^{-b-k} > \gamma, \forall s, t \in N[0, 3\lambda(k)];$$

which tells us that

$$(\forall k) : B(k) := \{(m, n) \in (<; N); J(k) \leq m < n,$$

$$d(x_{m+s}, x_{n+t}) > \gamma, \forall s, t \in N[0, 3\lambda(k)]\} \text{ is a nonempty relation over } N.$$

Having this precise, denote for each  $k \geq 0$

$$m(k) = \min \text{Dom}(B(k)), \ n(k) = \min B(k)(m(k)).$$

By this very definition, we get

(pro-1)  $(\forall k): k + 1 \leq J(k) \leq m(k) < n(k)$ ,

(pro-2)  $(\forall k): d(x_{m(k)+s}, x_{n(k)+t}) > \gamma, \forall s, t \in N[0, 3\lambda(k)]$ .

We claim that the couple  $(\beta, \gamma) \in (>; \text{adm}(\Theta; (x_n)))$ , the rank sequence  $(J(k))$  in  $N(1, \leq)$ , and the couple of rank-sequences  $[(m(k)), (n(k))]$  fulfill all desired conclusions.

(i): By (pro-1), it is clear that the first, second, and third relation of (21-a) holds.

(ii): Suppose by contradiction that

$$(m(k) <)n(k) \leq m(k) + 3\lambda(k), \text{ for some } k \geq 0.$$

By  $m(k) \geq J(k) \geq Z(3^{-2b-k}/\lambda(k))$ , the triangle inequality, and the choice of  $b$ ,

$$d(x_{m(k)}, x_{n(k)}) \leq 3^{-2b-k+1} < 3^{-b-k} < \gamma;$$

in contradiction with (pro-2); whence, the fourth relation in (21-a) holds too.

(iii): The first and second part of (21-b) are directly obtainable from the preceding stage and (pro-2), respectively. Concerning the third part of (21-b), let  $k \geq 0$  be arbitrary fixed. By definition,  $n(k)$  is the minimum of all ranks  $p \in N$  with

$(m(k), p) \in B(k)$ ; that is:

$$J(k) \leq m(k) < p \text{ and } d(x_{m(k)+s}, x_{p+t}) > \gamma, \forall s, t \in N[0, 3\lambda(k)].$$

As  $m(k) < m(k) + 2\lambda(k) < n(k) - 1$ , we must have (by this minimal property)

(pro-3)  $d(x_{m(k)+s}, x_{n(k)-1+t}) \leq \gamma$ , for some  $s, t \in N[0, 3\lambda(k)]$ .

But, in view of (pro-2) once again,

(pro-4)  $d(x_{m(k)+u}, x_{n(k)-1+v}) > \gamma, \forall u \in N[0, 3\lambda(k)], \forall v \in N[1, 3\lambda(k)]$ .

This, combined with (pro-3), tells us that, necessarily,

(pro-5)  $d(x_{m(k)+s}, x_{n(k)-1}) \leq \gamma$ , for some  $s \in N[0, 3\lambda(k)]$ .

By  $m(k) \geq Z(3^{-2b-k}/\lambda(k))$ , the triangle inequality, and  $b \geq 2$ , we then have

$$d(x_{m(k)}, x_{n(k)-1}) \leq d(x_{m(k)}, x_{m(k)+s}) + d(x_{m(k)+s}, x_{n(k)-1}) \leq 3^{-2b-k+1} + \gamma \leq 3^{-b-k} + \gamma$$

and the last conclusion of (21-b) follows.

(iv): From these facts and triangular inequality,

$$\gamma < d(x_{m(k)}, x_{n(k)}) \leq d(x_{m(k)}, x_{n(k)-1}) + r_{n(k)-1} \leq 3^{-b-k} + \gamma + r_{n(k)-1}, \forall k.$$

Passing to limit in this double inequality gives (21-c).

(v): By the very definition of  $[(\mu(k)), (\nu(k))]$ , and (pro-2),

$$V_k > \gamma, \text{ for all } k \geq 0; \text{ so, the first half of (21-d) holds.}$$

Moreover, from a metrical property of  $d$ , the very definition of  $(m(k); k \geq 0)$  and  $(n(k)); k \geq 0)$ , and the choice of  $b$

$$\begin{aligned} & |d(x_{m(k)}, x_{n(k)}) - d(x_{m(k)+\mu(k)}, x_{n(k)+\nu(k)})| \leq \\ & d(x_{m(k)}, x_{m(k)+\mu(k)}) + d(x_{n(k)}, x_{n(k)+\nu(k)}) \leq \\ & 3^{-2b-k+1} + 3^{-2b-k+1} < 3^{-2b-k+2} \leq 3^{-b-k}, \text{ for all } k \geq 0. \end{aligned}$$

Passing to limit in the relation between the first and the last member of this relation gives the second half of (21-d).

- (vi): By the strict ascending property of  $(\lambda(k); k \geq 0)$ , the sequences  $(\mu(k) = i; k \geq 0)$  and  $(\nu(k) = j; k \geq 0)$ , fulfill  $(\mu(k), \nu(k) \leq 3\lambda(k), \forall k)$ ; and this, along with the preceding stage, yields the desired conclusion.
- (vii): Let  $(i, j) \in N \times N$  be given. By the strict ascending property of  $(\lambda(k); k \geq 0)$ , there exists an index  $L = L(i, j) \in N$  with

$$i, j \leq 3\lambda(k), \text{ for all } k \geq L.$$

Then, define the couple of sequences  $(\mu(k); k \geq 0)$  and  $(\nu(k); k \geq 0)$  as

$$(\mu(k) = \nu(k) = 0; k \leq L); (\mu(k) = i, \nu(k) = j; k > L).$$

Clearly,  $(\mu(k), \nu(k) \leq 3\lambda(k), \forall k)$ ; and this, along with the preceding stage, yields

$$T_k^* = d(x_{m(k)+\mu(k)}, x_{n(k)+\nu(k)}) \rightarrow \gamma + \text{ as } n \rightarrow \infty.$$

It will now suffice observing that

$$T_k = T_k^*, \text{ for all } k > L \text{ (hence, for all } k \geq L + 1)$$

to get the written conclusion. The proof is complete.

In particular, when  $\Theta = R_+^0$  and  $(\lambda(k) = 1 + k; k \geq 0)$ , the obtained statement covers the 1969 one in Boyd and Wong [11]; so, it is natural that this result be referred to in the proposed way. Further aspects may be found in Reich [48]; see also Khan et al. [37].

### 3 Admissible Real Functions

In the following, some classes of admissible real functions are introduced. Their usefulness will become clear from our next developments.

(A) Denote for simplicity

- $\mathcal{F}_0(R_+)$  = the subclass of all  $\varphi \in \mathcal{F}(R_+)$ , with  $\varphi(0) = 0$
- $\mathcal{F}_0(re)(R_+)$  = the subclass of all  $\varphi \in \mathcal{F}_0(R_+)$ , with the *regressive* property:  $\varphi(t) < t, \forall t \in R_+^0$
- $\mathcal{F}_0(in)(R_+)$  = the subclass of all increasing functions  $\varphi \in \mathcal{F}_0(R_+)$
- $\mathcal{F}_0(re, in)(R_+) = \mathcal{F}_0(re)(R_+) \cap \mathcal{F}_0(in)(R_+)$ .

For each  $\varphi \in \mathcal{F}_0(re)(R_+)$ , let us introduce the sequential properties

- (M-a)  $\varphi$  is *Matkowski admissible*:  
for each  $(t_n)$  in  $R_+^0$  with  $(t_{n+1} \leq \varphi(t_n), \forall n)$  we have  $\lim_n t_n = 0$
- (n-d-a)  $\varphi$  is *non-diagonally admissible*:  
there are no sequences  $(t_n; n \geq 0)$  in  $R_+^0$   
and no elements  $\varepsilon \in R_+^0$  with  $t_n \rightarrow \varepsilon +, \varphi(t_n) \rightarrow \varepsilon +$
- (MK-a)  $\varphi$  is *Meir-Keeler admissible*:  
 $\forall \varepsilon > 0, \exists \delta > 0$ , such that  $\varepsilon < s < \varepsilon + \delta \implies \varphi(s) \leq \varepsilon$ .

The relationships between these properties are discussed in the statement below

**Theorem 3** For each  $\varphi \in \mathcal{F}_0(re)(R_+)$ , we have in (ZF-AC+DC)

$$(M-a) \implies (n-d-a) \implies (MK-a) \implies (M-a).$$

Hence, for each  $\varphi \in \mathcal{F}_0(re)(R_+)$ , the properties (M-a), (n-d-a), and (MK-a) are equivalent to each other.

**Proof** There are three cases to be discussed.

- (i) Suppose that  $\varphi$  is Matkowski admissible; we assert that  $\varphi$  is non-diagonally admissible. For, if  $\varphi$  is not endowed with such a property, there must be a sequence  $(t_n; n \geq 0)$  in  $R_+^0$  and a number  $\varepsilon > 0$ , such that

$$t_n \rightarrow \varepsilon + \text{ and } \varphi(t_n) \rightarrow \varepsilon +, \text{ as } n \rightarrow \infty.$$

Put  $i(0) = 0$ . As  $\varepsilon < \varphi(t_{i(0)})$  and  $t_n \rightarrow \varepsilon +$ , we have that

$$A(i(0)) := \{n > i(0); t_n < \varphi(t_{i(0)})\} \text{ is not empty;}$$

$$\text{hence, } i(1) := \min(A(i(0))) \text{ is an element of it, and } t_{i(1)} < \varphi(t_{i(0)}).$$

Likewise, as  $\varepsilon < \varphi(t_{i(1)})$  and  $t_n \rightarrow \varepsilon +$ , we have that

$A(i(1)) := \{n > i(1); t_n < \varphi(t_{i(1)})\}$  is not empty;

hence,  $i(2) := \min(A(i(1)))$  is an element of it, and  $t_{i(2)} < \varphi(t_{i(1)})$ .

This procedure may continue indefinitely; and yields (without any choice technique) a strictly ascending rank sequence  $(i(n); n \geq 0)$  (hence,  $i(n) \rightarrow \infty$  as  $n \rightarrow \infty$ ) for which the attached subsequence  $(s_n := t_{i(n)}; n \geq 0)$  of  $(t_n)$  fulfills

$$s_{n+1} < \varphi(s_n) (< s_n), \text{ for all } n.$$

On the other hand, by this very subsequence property,

$$(s_n > \varepsilon, \forall n) \text{ and } \lim_n s_n = \lim_n t_n = \varepsilon.$$

The obtained relations are in contradiction with the Matkowski property of  $\varphi$ ; hence, the working condition cannot be true; and we are done.

- (ii) Suppose that  $\varphi$  is non-diagonally admissible; we show that, necessarily,  $\varphi$  is Meir-Keeler admissible. For, if  $\varphi$  is not endowed with such a property, we must have (for some  $\varepsilon > 0$ )

$$H(\delta) := \{t \in R_+^0; \varepsilon < t < \varepsilon + \delta, \varphi(t) > \varepsilon\} \text{ is not empty, for each } \delta > 0.$$

Taking a strictly descending sequence  $(\delta_n; n \geq 0)$  in  $R_+^0$  with  $\delta_n \rightarrow 0$ , we get by the Denumerable Axiom of Choice (AC(N)) [deductible, as precise, in (ZF-AC+DC)], a sequence  $(t_n; n \geq 0)$  in  $R_+^0$ , so as

$$(\forall n) : t_n \text{ is an element of } H(\delta_n);$$

or, equivalently (by the very definition above and  $\varphi$ =regressive)

$$(\forall n) : \varepsilon < \varphi(t_n) < t_n < \varepsilon + \delta_n;$$

hence, in particular:  $\varphi(t_n) \rightarrow \varepsilon +$  and  $t_n \rightarrow \varepsilon +$ .

But, these relations are in contradiction with the non-diagonal admissible property of our function; hence, the assertion follows.

- (iii) Suppose that  $\varphi$  is Meir-Keeler admissible; we have to establish that  $\varphi$  is Matkowski admissible. Let  $(s_n; n \geq 0)$  be a sequence in  $R_+^0$  with the property  $(s_{n+1} \leq \varphi(s_n); n \geq 0)$ . Clearly,  $(s_n)$  is strictly descending in  $R_+^0$ ; hence,  $\sigma := \lim_n s_n$  exists in  $R_+$ . Suppose by contradiction that  $\sigma > 0$ ; and let  $\rho > 0$  be given by the Meir-Keeler admissible property of  $\varphi$ . By the above convergence relations, there exists some rank  $n(\rho)$ , such that

$$n \geq n(\rho) \text{ implies } \sigma < s_n < \sigma + \rho.$$

But then, under the notation  $(t_n := \varphi(s_n); n \geq 0)$ , we get (for the same ranks)

$$\sigma < s_{n+1} \leq t_n < s_n < \sigma + \rho;$$

in contradiction with the Meir-Keeler admissible property. Hence, necessarily,  $\sigma = 0$ ; and conclusion follows. The proof is complete.

In the following, some important examples of such objects will be given. For any  $\varphi \in \mathcal{F}_0(re)(R_+)$  and any  $s \in R_+^0$ , put

$$\Lambda^+ \varphi(s) = \inf_{\varepsilon > 0} \Phi(s+)(\varepsilon); \text{ where } \Phi(s+)(\varepsilon) = \sup \varphi(]s, s + \varepsilon[), \varepsilon > 0.$$

From the regressive property of  $\varphi$ , these quantities are finite; precisely,

$$0 \leq \Lambda^+ \varphi(s) \leq s, \forall s \in R_+^0.$$

The following completion of this will be useful.

**Proposition 5** *Let  $\varphi \in \mathcal{F}_0(re)(R_+)$  and  $s \in R_+^0$  be arbitrary fixed. Then,*

(31-1)  $\limsup_n(\varphi(t_n)) \leq \Lambda^+ \varphi(s)$ , for each sequence  $(t_n)$  in  $R_+^0$  with  $t_n \rightarrow s+$

(31-2) there exists a sequence  $(r_n)$  in  $R_+^0$  with  $r_n \rightarrow s+$  and  $\varphi(r_n) \rightarrow \Lambda^+ \varphi(s)$ .

**Proof** Denote, for simplicity,

$$\alpha = \Lambda^+ \varphi(s); \text{ hence, } \alpha = \inf_{\varepsilon > 0} \Phi(s+)(\varepsilon), \text{ and } 0 \leq \alpha \leq s.$$

(i): Given  $\varepsilon > 0$ , there exists a rank  $p(\varepsilon) \geq 0$  such that  $s < t_n < s + \varepsilon$ , for all  $n \geq p(\varepsilon)$ ; hence

$$\limsup_n(\varphi(t_n)) \leq \sup\{\varphi(t_n); n \geq p(\varepsilon)\} \leq \Phi(s+)(\varepsilon).$$

Passing to infimum over  $\varepsilon > 0$ , yields (see above)

$$\limsup_n(\varphi(t_n)) \leq \inf_{\varepsilon > 0} \Phi(s+)(\varepsilon) = \alpha; \text{ and the claim follows.}$$

(ii): Define  $(\beta_n := \Phi(s+)(2^{-n}); n \geq 0)$ ; this is a descending sequence in  $R_+$ , with

$$(\beta_n \geq \alpha, \forall n) \text{ and } \inf_n \beta_n = \alpha; \text{ hence } \lim_n \beta_n = \alpha.$$

By these properties, there may be constructed a sequence  $(\gamma_n; n \geq 0)$  in  $R$ , with

$$\gamma_n < \beta_n, \forall n; \lim_n \gamma_n = \lim_n \beta_n = \alpha.$$



(For example, we may take  $(\gamma_n = \beta_n - 3^{-n}; n \geq 0)$ ; we do not give details). Let  $n \geq 0$  be arbitrary fixed. By the supremum definition, there exists  $r_n \in ]s, s + 2^{-n}[$  such that  $\varphi(r_n) > \gamma_n, \forall n$ ; moreover (again by definition),  $\varphi(r_n) \leq \beta_n$ . The obtained sequence  $(r_n; n \geq 0)$  fulfills  $r_n \rightarrow s+$  and  $\varphi(r_n) \rightarrow \alpha$ ; wherefrom, all is clear.

We may now pass to the announced examples of such functions.

*Example 1* Call  $\varphi \in \mathcal{F}_0(re)(R_+)$ , *Boyd-Wong admissible* [11], if

$$\Lambda^+\varphi(s) < s, \text{ for all } s > 0.$$

In particular, this holds provided

$$\varphi \text{ is upper semicontinuous at the right on } R_+^0 : \Lambda^+\varphi(s) \leq \varphi(s), \forall s \in R_+^0.$$

Moreover, the written property is fulfilled when

$$\varphi \text{ is continuous at the right on } R_+^0;$$

for, in such a case,  $[\Lambda^+\varphi(s) = \varphi(s), \forall s \in R_+^0]$ .

Under these conditions, we have (cf. Meir and Keeler [42]):

$(\forall \varphi \in \mathcal{F}_0(re)(R_+)) : \varphi \text{ is Boyd-Wong admissible implies}$   
 $\varphi \text{ is Meir-Keeler admissible [or, equivalently: Matkowski admissible].}$

In fact, suppose that  $\varphi \in \mathcal{F}_0(re)(R_+)$  is Boyd-Wong admissible; and fix some  $\gamma > 0$ . As  $\Lambda^+\varphi(\gamma) < \gamma$ , there exists  $\beta = \beta(\gamma) > 0$  such that

$$\Phi(\gamma+)(\beta) < \gamma; \text{ wherefrom, } \gamma < t < \gamma + \beta \text{ implies } \varphi(t) < \gamma;$$

and this gives the desired property.

Concerning the reverse inclusion, let us consider the function  $\varphi \in \mathcal{F}_0(re)(R_+)$ , according to (for some  $r > 0$ ):

$$(\varphi(t) = t/2, \text{ if } 0 \leq t \leq r), (\varphi(t) = r, \text{ if } t > r).$$

Clearly,  $\varphi$  is Matkowski admissible, as it can be directly seen. On the other hand,

$$\Lambda^+\varphi(r) = r; \text{ whence, } \varphi \text{ is not Boyd-Wong admissible;}$$

proving that the reverse inclusion is not in general valid. For an extended example of this type, see Turinici [55] and the references therein.

*Example 2* According to its definition, the Matkowski admissible property of some  $\varphi \in \mathcal{F}_0(re, in)(R_+)$  writes (cf. Matkowski [41])

$$\varphi^n(t) \rightarrow 0 \text{ as } n \rightarrow \infty, \text{ for all } t > 0.$$

[Here, for each  $n \in N$ ,  $\varphi^n$  stands for the  $n$ -th iterate of  $\varphi$ ]. This, by a previous equivalence result, tells us that  $\varphi$  is Meir-Keeler admissible. A different way of proving this may be described as follows (cf. Jachymski [33]):

Assume that  $\varphi \in \mathcal{F}_0(re, in)(R_+)$  is Matkowski admissible. If the underlying property fails, then (for some  $\gamma > 0$ ):

$$\forall \beta > 0, \exists t \in ]\gamma, \gamma + \beta[, \text{ such that } \varphi(t) > \gamma.$$

As  $\varphi$  is increasing, this yields

$$\varphi(t) > \gamma, \forall t > \gamma; \text{ hence, by induction: } \varphi^n(t) > \gamma, \forall n \in N, \forall t > \gamma.$$

Taking some  $t > \gamma$  and passing to limit as  $n \rightarrow \infty$ , one gets  $0 \geq \gamma$ ; contradiction.

A sufficient condition for this property is to be obtained from the result above, by simply noting that

$$(\forall \varphi \in \mathcal{F}_0(re, in)(R_+)) : \Lambda^+ \varphi(s) = \varphi(s + 0), \forall s \in R_+^0.$$

Precisely, we have the practical characterization

$$(\forall \varphi \in \mathcal{F}_0(re, in)(R_+)) : \varphi \text{ is Boyd-Wong admissible iff } \varphi \text{ is strongly regressive } [\varphi(s + 0) < s, \text{ for all } s > 0].$$

*Example 3* Call  $\varphi \in \mathcal{F}_0(re)(R_+)$ , Geraghty admissible [24] provided each sequence  $(t_n)$  in  $R_+^0$  with  $\varphi(t_n)/t_n \rightarrow 1$  fulfills  $t_n \rightarrow 0$ .

The connection between this property and the preceding ones is described as:

$$\begin{aligned} & \text{(for each } \varphi \in \mathcal{F}_0(re)(R_+) \text{):} \\ & \varphi \text{ is Geraghty admissible implies } \varphi \text{ is Boyd-Wong admissible.} \end{aligned}$$

In fact, suppose that  $\varphi \in \mathcal{F}_0(re)(R_+)$  is not Boyd-Wong admissible. From a previous relation, there exists some  $s \in R_+^0$  with  $\Lambda^+ \varphi(s) = s$ . Combining with a preceding auxiliary fact, there exists a sequence  $(r_n; n \geq 0)$  in  $R_+^0$  with

$$r_n \rightarrow s + \text{ and } \varphi(r_n) \rightarrow s; \text{ whence } \varphi(r_n)/r_n \rightarrow 1;$$

i.e.:  $\varphi$  is not Geraghty admissible. The obtained contradiction proves our claim.

Concerning the reverse inclusion, note that, for the (continuous) Boyd-Wong admissible function  $[\varphi(t) = t(1 - e^{-t}), t \geq 0]$  in  $\mathcal{F}_0(re, in)(R_+)$ , and the sequence  $(t_n = n + 1; n \geq 0)$  in  $R_+^0$ , we have

$$\varphi(t_n)/t_n \rightarrow 1; \text{ but, evidently, } t_n \rightarrow \infty.$$

Hence,  $\varphi$  is not Geraghty admissible; so that the reciprocal is not in general true.

(B) A basic application involving these concepts may be described as below. Fix some  $\omega \in \mathcal{F}_0(re, in)(R_+)$  and let  $\omega^* \in \mathcal{F}(R_+)$  be defined as

$$(\omega^*(t) = t - \omega(t); t \in R_+) \text{ [in short: } \omega^* = I - \omega\text{]; called: the complement of } \omega.$$

We say that  $\omega$  is *complementary coercive*, in case

$$\begin{aligned} &\omega^* \text{ is coercive: } \omega^*(t) \rightarrow \infty \text{ as } t \rightarrow \infty; \text{ or, equivalently:} \\ &\{t \in R_+; \omega^*(t) \leq a\} \text{ is bounded, for each } a \geq 0. \end{aligned}$$

Note that, under such a condition, the function  $\gamma \in \mathcal{F}(R_+)$  given as

$$\gamma(s) = \sup A(s) \text{ where } A(s) = \{t \in R_+; t \leq \omega(s + t)\}, s \in R_+$$

[referred to as the *right complementary inverse* (in short: *rc-inverse*) associated to  $\omega$ ] is well defined. In fact, it will suffice noting that, for each  $s \in R_+$ ,

$$A(s) = \{t \in R_+; \omega^*(s + t) \leq s\}; \text{ whence, } A(s) \text{ is bounded in } R_+.$$

Clearly,  $\gamma(0) = 0$ . In addition, by the properties of  $\omega$ , we have

- (p-1) (Increasing):  $s_1 \leq s_2$  implies  $A(s_1) \subseteq A(s_2)$ ; whence  $\gamma$  is increasing
- (p-2) (Representation formula):  $\gamma(s) \leq \omega(s + \gamma(s))$ , for each  $s \in R_+$ .

In fact, let  $t \in A(s)$  be arbitrary fixed; hence,  $t \leq \omega(s + t)$ . By the definition of supremum,  $(t \leq \omega(s + \gamma(s)), \forall t \in A(s))$ ; wherefrom (again by the underlying definition) the conclusion follows.

Further properties of  $\gamma(\cdot)$  are available under extra properties of  $\omega(\cdot)$ .

**Proposition 6** *Suppose that, in addition,*

$$\omega \text{ is strongly regressive: } \omega(t + 0) < t, \forall t > 0.$$

*Then (in addition to the above)*

$$(\lim_{t \rightarrow 0+} \gamma(t) =) \gamma(0 + 0) = 0; \text{ that is: } s_n \rightarrow 0 \text{ implies } \gamma(s_n) \rightarrow 0.$$

**Proof** As  $\gamma$  is increasing,  $\delta := \gamma(0+0)$  exists. Assume by contradiction that  $\delta > 0$ . Let  $(s_n)$  be a strictly descending sequence in  $R_+^0$  with  $s_n \rightarrow 0$ ; whence,  $\gamma(s_n) \rightarrow \delta$ . From the preceding representation formula

$$\gamma(s_n) \leq \omega(s_n + \gamma(s_n)), \text{ for all } n.$$

As  $n \rightarrow \infty$ , we have  $s_n + \gamma(s_n) \rightarrow \delta+$ ; so that, passing to limit in this relation yields  $\delta \leq \omega(\delta + 0) < \delta$ ; contradiction. Hence,  $\gamma(0 + 0) = 0$ , as claimed.

Sometimes, we may ask for the associated rc-inverse  $\gamma(\cdot)$  a property like

$$\gamma \text{ is strongly regressive: } \gamma(t + 0) < t, \text{ for all } t \in R_+^0;$$

the initial function  $\omega$  will be referred to as *rc-regressive* in such a case. A useful answer in this direction is obtained by imposing regularity conditions upon

$$(\eta(t) = \omega(2t); t \geq 0) \text{ [the double function attached to } \omega\text{].}$$

Note that, by the increasing property of  $\omega$  (and  $\eta$ )

$$(\forall t > 0) : \eta(t + 0) = \inf_{s>0} \eta(t + s) = \inf_{s>0} \omega(2t + 2s) = \omega(2t + 0);$$

this will be useful in the sequel.

**Theorem 4** *Let the function  $\omega \in \mathcal{F}_0(re, in)(R_+)$  be such that*

- (32-i) *the double function  $\eta$  is strongly regressive ( $\eta(t+0) < t, \forall t > 0$ ); expressed as:  $\omega$  is double strongly regressive*
- (32-ii) *the double function  $\eta$  is complementary coercive ( $\eta^* := I - \eta$  is coercive); expressed as:  $\omega$  is double complementary coercive.*

Then,

- (32-a) *its associated complement  $\omega^* := I - \omega$  is coercive; so, the right complementary inverse*

$$\gamma(s) = \sup A(s) \text{ where } A(s) = \{t \in R_+; t \leq \omega(s + t)\}, s \in R_+$$

*is well defined as an element of  $\mathcal{F}_0(R_+)$*

- (32-b) *The right complementary inverse  $\gamma(\cdot)$  fulfills*

- (p-1)  *$\gamma$  is increasing and  $\gamma(s) \leq \omega(s + \gamma(s))$ , for all  $s \in R_+$*
- (p-2)  *$\gamma$  is zero-continuous:  $\gamma(t) \rightarrow 0 = \gamma(0 + 0)$  as  $t \rightarrow 0$*

- (32-c) *Finally, we have the properties*

- (p-3)  *$\gamma$  is strongly regressive; hence, in particular, regressive*
- (p-4)  *$\gamma$  is complementary coercive:  $\gamma^* := I - \gamma$  is coercive.*

**Proof** By a preceding relation,

$$(\forall t > 0) : \omega(t + 0) < t/2; \text{ whence, } t - \omega(t) > t - t/2 = t/2;$$

proving that  $\omega$  is complementary coercive:  $\omega^* := I - \omega$  is coercive.

- (i): Clearly, the right complementary inverse  $\gamma : R_+ \rightarrow R_+$  is well defined.
- (ii): The properties (p-1) and (p-2) of  $\gamma$  are directly obtainable by the preceding facts, in view of

$\omega$  is double strongly regressive implies  $\omega$  is strongly regressive.

(iii): Fix  $s > 0$ ; and let  $(s_n)$  be a strictly descending sequence with

$$s_n \rightarrow s + \text{ as } n \rightarrow \infty; \text{ so that, } \gamma(s_n) \rightarrow r := \gamma(s + 0) \text{ as } n \rightarrow \infty.$$

From the above representation formula,

$$\text{(rep-seq) } \gamma(s_n) \leq \omega(s_n + \gamma(s_n)), \forall n.$$

Denote for simplicity  $u = s + r$ . By the posed hypotheses,

$$s_n + \gamma(s_n) \rightarrow u+; \text{ whence } \omega(s_n + \gamma(s_n)) \rightarrow \omega(u + 0)$$

Passing to limit in (rep-seq) gives (as  $\omega$  is double strongly regressive)

$$r \leq \omega(u + 0) < (1/2)u = (1/2)(s + r); \text{ wherefrom, } r < s;$$

and the conclusion follows.

(iv): Let  $a \geq 0$  be arbitrary fixed. We have to establish that

$$A := \{t \in R_+; t \leq \gamma(t) + a\} \text{ is bounded.}$$

But, according to definition (and the properties of  $\omega$  and  $\gamma$ )

$$\begin{aligned} A \subseteq \{t \in R_+; t \leq \omega(t + \gamma(t)) + a\} &\subseteq \{t \in R_+; t \leq \omega(2t) + a\} = \\ \{t \in R_+; t \leq \eta(t) + a\} &= \{t \in R_+; \eta^*(t) \leq a\}; \end{aligned}$$

and this, combined with the last subset being bounded, ends our argument.

The obtained properties of the right complementary inverse  $\gamma \in \mathcal{F}(R_+)$  associated to  $\omega$  are basic tools for getting a lot of relative type statements involving real sequences, to be needed further. For simplicity reasons, we will express these as generic results (without explicitly mentioning the conditions under which these properties were derived).

**Theorem 5** *Let the sequences  $(a_n; n \geq 0)$  and  $(b_n; n \geq 0)$  in  $R_+$ , and the function  $\varphi \in \mathcal{F}_0(re, in)(R_+)$  be such that*

$$(33-i) \quad a_n \leq b_n + \varphi(a_n), \text{ for all } n.$$

*In addition, suppose that*

$$(33-ii) \quad \varphi \text{ is strongly regressive } (\varphi(t + 0) < t, \forall t > 0)$$

$$(33-iii) \quad \varphi \text{ is complementary coercive } (\varphi^* := I - \varphi \text{ is coercive}).$$

*Then, necessarily,*

$$(33-a) \quad (b_n) \text{ is bounded implies } (a_n) \text{ is bounded}$$

$$(33-b) \quad (b_n) \text{ is zero-convergent implies } (a_n) \text{ is zero-convergent.}$$

**Proof**

(i): By definition,  $\beta := \sup\{b_n; n \geq 0\}$  exists in  $R_+$ . This, by the imposed hypothesis, yields

$$(\forall n) : a_n \leq \beta + \varphi(a_n); \text{ that is } \varphi^*(a_n) \leq \beta.$$

Combining with the coerciveness of  $\varphi^*$ , yields the desired fact.

(ii): Suppose that  $(b_n)$  is zero-convergent. In particular,  $(b_n)$  is bounded; hence, by the preceding stage,  $(a_n)$  is bounded too. By a standard result,  $(a_n)$  has convergent subsequences. Let  $(a_{i(n)}; n \geq 0)$  be one of these; hence,  $a_{i(n)} \rightarrow a$  as  $n \rightarrow \infty$ , for some  $a \in R_+$ . Suppose by contradiction that  $a > 0$ . Passing to limit as  $n \rightarrow \infty$  in

$$a_{i(n)} \leq b_{i(n)} + \varphi(a_{i(n)}), n \geq 0$$

one derives (by the strongly regressive property of  $\varphi$ )

$$a \leq \varphi(a + 0) < a; \text{ absurd; whence, } a = 0.$$

In other words, all convergent subsequences  $(a_{i(n)}; n \geq 0)$  of  $(a_n)$  fulfill  $a_{i(n)} \rightarrow 0$  as  $n \rightarrow \infty$ . This necessarily gives  $a_n \rightarrow 0$  as  $n \rightarrow \infty$ ; and the conclusion follows.

A bi-dimensional version of this result may be stated as follows. By a *pseudometric* over  $N$ , we mean any mapping  $a : N \times N \rightarrow R_+$ . If, in addition,

$$a(., .) \text{ is reflexive: } a(n, n) = 0, \forall n$$

we say that  $a(., .)$  is a *r-pseudometric*. Given the r-pseudometric  $a(., .)$ , call it *Cauchy*, in case of  $[a(n, m) \rightarrow 0 \text{ as } n, m \rightarrow \infty, n \leq m]$ ; and *asymptotic*, provided  $[a(n, n + 1) \rightarrow 0 \text{ as } n \rightarrow \infty]$ . Clearly,

$$a(., .) \text{ is Cauchy implies } a(., .) \text{ is asymptotic;}$$

the reciprocal is not in general true.

**Theorem 6** *Let the r-pseudometrics  $a(., .)$  and  $b(., .)$  over  $N$  and the function  $\varphi \in \mathcal{F}_0(re, in)(R_+)$  be such that*

$$(34-i) \ a(n, m) \leq b(n, m) + \varphi(a(n, m)), \text{ for all } n, m \in N, n \leq m.$$

*In addition, suppose that*

$$(34-ii) \ \varphi \text{ is strongly regressive } (\varphi(t + 0) < t, \forall t > 0)$$

$$(34-iii) \ \varphi \text{ is complementary coercive } (\varphi^* := I - \varphi \text{ is coercive}).$$

Then, necessarily,

(34-a)  $b(., .)$  is asymptotic implies  $a(., .)$  is asymptotic

(34-b)  $b(., .)$  is Cauchy implies  $a(., .)$  is Cauchy.

**Proof** There are two steps to be passed.

**Step 1.** Suppose that

$b(., .)$  is asymptotic [ $(\beta_n := b(n, n + 1); n \geq 0)$  fulfills  $\beta_n \rightarrow 0$ ].

We have to establish that

$a(., .)$  is asymptotic [ $(\alpha_n := a(n, n + 1); n \geq 0)$  fulfills  $\alpha_n \rightarrow 0$ ].

This, however, by the sequential condition

(seq)  $\alpha_n \leq \beta_n + \varphi(\alpha_n)$ , for all  $n \in N$ ,

is deductible from the preceding statement; so that, our claim follows.

**Step 2.** Now, let us show that

$b(., .)$  is Cauchy implies  $a(., .)$  is Cauchy.

The last property means (by definition)

$\forall \varepsilon > 0$ , there exists  $j(\varepsilon)$ , such that  $j(\varepsilon) \leq n \leq m$  implies  $a(n, m) \leq \varepsilon$ ;

or, equivalently [passing to the successor  $J(\varepsilon) := j(\varepsilon) + 1$  and remembering that  $a(., .)$  is r-pseudometric]

$\forall \varepsilon > 0$ , there exists  $j(\varepsilon)$ , such that  $j(\varepsilon) < n < m$  implies  $a(n, m) \leq \varepsilon$ .

The negation of this means: there exists  $\varepsilon > 0$ , such that

$$C(j) = \{(n, m) \in N \times N; j < n < m, a(n, m) > \varepsilon\} \neq \emptyset, \forall j.$$

Denote, for each  $j$

$$n(j) = \min \text{Dom}(C(j)), m(j) = \min C(j)(n(j)).$$

Fix  $j(0) \in N$ . By this construction, there exists a couple of ranks ( $j(1) = n(j(0))$ ,  $j(2) = m(j(0))$ ) with  $j(0) < j(1) < j(2)$ ,  $a(j(1), j(2)) > \varepsilon$ . Further, given this index  $j(2)$  there exists a couple of ranks ( $j(3) = n(j(2))$ ,  $j(4) = m(j(2))$ ) with  $j(2) < j(3) < j(4)$ ,  $a(j(3), j(4)) > \varepsilon$ . The procedure may continue indefinitely (without any choice techniques); and yields a strictly ascending sequence of ranks ( $j(n); n \geq 0$ ) [hence,  $j(n) \rightarrow \infty$  as  $n \rightarrow \infty$ ], with

$$\alpha_n := a(j(2n + 1), j(2n + 2)) > \varepsilon, \forall n.$$

On the other hand, by the Cauchy property

$$\beta_n := b(j(2n + 1), j(2n + 2)) \rightarrow 0, \text{ as } n \rightarrow \infty.$$

This, in combination with the sequential condition

$$(\text{seqq}) \alpha_n \leq \beta_n + \varphi(\alpha_n), \text{ for all } n \in N$$

gives a contradiction with respect to the preceding statement. Hence, necessarily,  $a(., .)$  is Cauchy; as claimed.

An interesting question is that of such conclusions being retainable when one of the basic hypotheses about  $\varphi$  is to be dropped. The answer to this is negative, in general; we do not give details.

### 4 Topological Preliminaries

In the following, some basic concepts and results involving topological structures are given, with a special emphasis on gauge spaces.

Let  $N = \{0, 1, \dots\}$  denote the set of natural numbers. Given  $p \geq 1$ , each set  $M$  equivalent with  $N(p, >) := \{0, \dots, p - 1\}$  (in the sense: there exists a bijection between  $M$  and  $N(p, >)$ ) is called *p-finite*. For completeness, we accept that the empty set  $\emptyset$  is 0-finite. Finally, let us say that the set  $M$  is *finite*, provided it is *q-finite*, for some  $q \in N$ .

Let  $E$  be a nonempty set. Remember that  $\text{exp}[E]$  denotes the class of all subsets in  $E$ . Given the subset  $\mathcal{D} \subseteq \text{exp}[E]$  of (i.e.: a family of subsets in  $E$ ), remember that its union  $\cup \mathcal{D}$  is introduced, axiomatically, as

$$z \in \cup \mathcal{D} \text{ iff } z \in H, \text{ for some } H \in \mathcal{D}; \text{ hence, in particular, } \cup \emptyset = \emptyset.$$

On the other hand, whenever  $\mathcal{D}$  is nonempty, then  $\cap \mathcal{D}$  is introduced as

$$z \in \cap \mathcal{D} \text{ iff } z \in H, \text{ for each } H \in \mathcal{D}.$$

Finally, when  $\mathcal{D} = \emptyset$ , we put  $\cap \mathcal{D} = E$ . We must stress that this last definition is entirely *locally*; i.e.: it is restricted to the ambient set  $E$  and the class  $\text{exp}[E]$ .

**(A) [General aspects]**

Let  $X$  be a nonempty set. We say that the family  $\mathcal{G} \subseteq \text{exp}[X]$  is

$$\textit{semi-normal}, \text{ if } \emptyset \in \mathcal{G}; \textit{ normal}, \text{ if } \{\emptyset, X\} \subseteq \mathcal{G}.$$

By a *topology* on  $X$ , we mean any normal family  $\mathcal{T} \subseteq \text{exp}[X]$ , with the properties (cf. Bourbaki [10, Ch I, Sect 1])



- (top-1) the union of any subset in  $\mathcal{T}$  belongs to  $\mathcal{T}$
- (top-2) the intersection of any finite subset of  $\mathcal{T}$  is in  $\mathcal{T}$ .

In this case, the couple  $(X, \mathcal{T})$  will be called a *topological space*. Any  $D \in \mathcal{T}$  will be called *open*; and its complement  $X \setminus D$  is termed *closed* (with respect to  $\mathcal{T}$ ).

An equivalent way of introducing such a concept is the *complementary* one. Let us say that the normal family  $\mathcal{S} \subseteq \exp[X]$  is a *cotopology* on  $X$ , when

- (cotop-1) the intersection of any subset in  $\mathcal{S}$  belongs to  $\mathcal{S}$
- (cotop-2) the union of any finite subset of  $\mathcal{S}$  is in  $\mathcal{S}$ .

In this case, the couple  $(X, \mathcal{S})$  will be called a *cotopological space*.

Given the topology  $\mathcal{T}$  on  $X$ , the (normal) family  $\mathcal{S} \subseteq \exp[X]$  introduced as  $[E \in \mathcal{S} \text{ iff } X \setminus E \in \mathcal{T}]$  is a cotopology on  $X$ . This will be referred to as the *cotopology induced* by the topology  $\mathcal{T}$ ; denoted as:  $\mathcal{S} = X \setminus \mathcal{T}$ . Conversely, given the cotopology  $\mathcal{S}$  on  $X$ , the (normal) family  $\mathcal{T} \subseteq \exp[X]$  introduced as  $[D \in \mathcal{T} \text{ iff } X \setminus D \in \mathcal{S}]$  is a topology on  $X$ . This will be referred to as the *topology induced* by the cotopology  $\mathcal{S}$ ; denoted as:  $\mathcal{T} = X \setminus \mathcal{S}$ .

**(B) [Topological constructions]**

Let  $X$  be a nonempty set. There are several ways of constructing a topology on  $X$ ; the basic ones are described below.

- (I) We say that the selfmap  $A \mapsto \text{Klo}(A)$  of  $\exp[X]$  is a (*Kuratowski*) *closure* over  $X$ , provided (cf. Kuratowski [38, Ch I, Sect 4])

- (Klo-1)  $\emptyset = \text{Klo}(\emptyset)$ ,  $X = \text{Klo}(X)$
- (Klo-2)  $A \subseteq \text{Klo}(A)$ , for each  $A \in \exp[X]$
- (Klo-3)  $\text{Klo}(A \cup B) = \text{Klo}(A) \cup \text{Klo}(B)$ ,  $\forall A, B \in \exp[X]$
- (Klo-4)  $\text{Klo}(\text{Klo}(A)) = \text{Klo}(A)$ , for each  $A \in \exp[X]$ .

Note that, as a direct consequence of (Klo-3),

- (Klo-incr)  $A \subseteq B$  implies  $\text{Klo}(A) \subseteq \text{Klo}(B)$ .

In fact, let  $A, B \in \exp[X]$  be such that  $A \subseteq B$ . As  $B = A \cup (B \setminus A)$ , we must have  $[\text{Klo}(B) = \text{Klo}(A) \cup \text{Klo}(B \setminus A) \supseteq \text{Klo}(A)]$ ; wherefrom, the assertion follows.

Let  $\mathcal{T}$  be a topology over  $X$ . Define a selfmap  $A \mapsto \text{cl}(A)$  of  $\exp[X]$ , as:

- $(\forall A \in \exp[X]) : \text{cl}(A) = \text{the intersection of all } G \in \exp[X] \text{ with } A \subseteq G = \text{closed.}$

It is not hard to see that all properties (Klo-1)-(Klo-4) hold; whence,  $A \mapsto \text{cl}(A)$  is a Kuratowski closure over  $X$ ; referred to as the *closure operator* induced by  $\mathcal{T}$ . Conversely, suppose that  $A \mapsto \text{Klo}(A)$  is a Kuratowski closure. Then, the (normal) class  $\mathcal{S} = \{E \in \exp[X]; E = \text{Klo}(E)\}$  is a cotopology over  $X$  (see above); so that,  $\mathcal{T} = X \setminus \mathcal{S}$  is a topology over  $X$  in the above discussed sense. Moreover, with respect to this topology, we have

$$\text{cl}(A) = \text{Klo}(A), \text{ for each } A \in \exp[X].$$

(II) Suppose that, to each  $x \in X$  we attached a class  $\mathscr{W}(x) \subseteq \exp[X]$  (referred to as: *neighborhoods system* of  $x$ ) such that the family  $(\mathscr{W}(x); x \in X)$  fulfills

(ns-1) for each  $x \in X$  and each  $W \in \mathscr{W}(x)$ , we have  $x \in W$

(ns-2)  $(\forall x \in X): A \in \mathscr{W}(x)$  and  $A \subseteq B$  imply  $B \in \mathscr{W}(x)$ ;  
hence, in particular,  $X \in \mathscr{W}(x)$

(ns-3)  $(\forall x \in X): W_1, W_2 \in \mathscr{W}(x)$  imply  $W_1 \cap W_2 \in \mathscr{W}(x)$

(ns-4) for each  $x \in X$  and each  $A \in \mathscr{W}(x)$  there exists  $D \in \mathscr{W}(x)$  such that  $A \in \mathscr{W}(y)$ , for each  $y \in D$  (hence,  $A \supseteq D$ ).

We then say that  $(\mathscr{W}(x); x \in X)$  is a *neighborhoods system* over  $X$ . A concrete example is to be given as follows. Let  $\mathcal{T}$  be a topology over  $X$ . Given  $x \in X$ , let us say that  $Y \in \exp[X]$  is a  $\mathcal{T}$ -*neighborhood* of  $x$ , provided

$$x \in D \subseteq Y, \text{ for some } D \in \mathcal{T};$$

the class of all sets will be denoted as  $\mathscr{V}(x)$ . It is not hard to see that all properties (ns-1)-(ns-4) hold; whence,  $(\mathscr{V}(x); x \in X)$  is a neighborhoods system over  $X$ ; referred to as the neighborhoods system generated by  $\mathcal{T}$ . Concerning the reciprocal question, the following statement holds (cf. Costinescu [17, Ch II, Sect 1]):

**Proposition 7** *Suppose that  $(\mathscr{W}(x); x \in X)$  is a neighborhoods system over  $X$ ; and put  $\mathcal{T} =$  the set of all  $D \in \exp[X]$  with  $(D \in \mathscr{W}(x), \forall x \in D)$ . Then,*

(41-1)  $\mathcal{T}$  is (a normal family and) a topology over  $X$

(41-2) with respect to this topology,  $(\mathscr{V}(x) = \mathscr{W}(x), \text{ for all } x \in X)$ .

A variant of this construction is to be described as follows. Suppose that, to each  $x \in X$  we attached a class  $\mathscr{W}^*(x) \subseteq \exp[X]$  (referred to as: neighborhoods subsystem of  $x$ ) such that the family  $(\mathscr{W}^*(x); x \in X)$  fulfills the properties

(nss-1) for each  $x \in X$  and each  $W^* \in \mathscr{W}^*(x)$ , we have  $x \in W^*$

(nss-2)  $(\forall x \in X):$  for each  $W_1^*, W_2^* \in \mathscr{W}^*(x)$  there exists  $W_3^* \in \mathscr{W}^*(x)$ ,  
such that  $W_1^* \cap W_2^* \supseteq W_3^*$

(nss-3) for each  $x \in X$  and each  $A^* \in \mathscr{W}^*(x)$  there exists  $D^* \in \mathscr{W}^*(x)$ ,  
such that for each  $y \in D^*$ , there exists  $B_y^* \in \mathscr{W}^*(y)$  with  $A^* \supseteq B_y^*$ .

We then say that  $(\mathscr{W}^*(x); x \in X)$  is a *neighborhoods subsystem* over  $X$ . A concrete example is to be given as follows. Let  $\mathcal{T}$  be a topology over  $X$ ; and, for each  $x \in X$ , let  $\mathscr{V}(x)$  be the neighborhoods system over  $X$  generated by  $\mathcal{T}$ . Then, for each  $x \in X$ , let the subset  $\mathscr{V}^*(x)$  of  $\mathscr{V}(x)$  be taken as

$$\text{for each } V \in \mathscr{V}(x) \text{ where exists } V^* \in \mathscr{V}^*(x) \text{ with } V \supseteq V^*.$$

Then,  $(\mathscr{V}^*(x); x \in X)$  is a neighborhoods subsystem over  $X$ ; referred to as: the neighborhoods subsystem over  $X$  generated by  $\mathcal{T}$ . For example,  $(\mathscr{V}^*(x); x \in X)$  may be taken as the *open neighborhoods subsystem* over  $X$  generated by  $\mathcal{T}$ :

$$\mathscr{V}^*(x) = \{D \in \mathcal{T}; x \in D\}, x \in X.$$

Conversely, let  $(\mathcal{W}^*(x); x \in X)$  be a neighborhoods subsystem over  $X$ . Then, the class  $(\mathcal{W}(x); x \in X)$  introduced as

$$(\forall x \in X) : W \in \mathcal{W}(x) \text{ iff } W \supseteq W^*, \text{ for some } W^* \in \mathcal{W}^*(x)$$

is a neighborhoods system over  $X$ , as it can be directly seen. By a preceding statement, there exists a (uniquely determined) topology  $\mathcal{T}$  over  $X$  such that  $\mathcal{W}(x) = \mathcal{V}(x)$  (the neighborhood system on  $X$  generated by  $\mathcal{T}$ ). But then,  $(\mathcal{W}^*(x); x \in X)$  is a neighborhoods subsystem on  $X$  generated by  $\mathcal{T}$ .

(III) Let us say that the subset  $\mathcal{B} \subseteq \text{exp}[X]$  is a *basis* for  $\mathcal{T}$ , provided

$$\text{each } D \in \mathcal{T} \text{ is the union of a subset in } \mathcal{B}.$$

In particular, this means that  $\mathcal{B} \subseteq \mathcal{T}$ ; but, in general, we cannot have  $\mathcal{B} = \mathcal{T}$ .

A useful characterization of this concept is to be obtained by means of neighborhood subsystems generated by  $\mathcal{T}$ . Let  $\mathcal{B} \subseteq \mathcal{T}$  be a family of open sets. Denote

$$\mathcal{B}(x) = \{B \in \mathcal{B}; x \in B\}, \quad x \in X.$$

We have (see Costinescu [17, Ch II, Sect 2] for details)

**Proposition 8** *Let  $\mathcal{B} \subseteq \mathcal{T}$  be a family of open sets. The following are equivalent:*

(42-1)  $\mathcal{B}$  is a basis for the topology  $\mathcal{T}$

(42-2)  $(\mathcal{B}(x); x \in X)$  is a neighborhoods subsystem over  $X$  generated by  $\mathcal{T}$ :  
for each  $x \in X$  and each  $V \in \mathcal{V}(x)$  there exists  $B \in \mathcal{B}(x)$  with  $V \supseteq B$ .

Let  $\mathcal{B} \subseteq \text{exp}[X]$  be a class of sets. We may ask of to what extent there exists a topology  $\mathcal{T}$  on  $X$  such that  $\mathcal{B}$  is the basis for  $\mathcal{T}$ . To this end, note that if  $\mathcal{B} \subseteq \text{exp}[X]$  is a basis for the topology  $\mathcal{T}$ , then (by the properties of  $(\mathcal{B}(x); x \in X)$ )

(s-dir)  $\mathcal{B}$  is *strongly directed*: for each  $B_1, B_2 \in \mathcal{B}$  and each

$$x \in B_1 \cap B_2 \text{ there exists } B_3 \in \mathcal{B} \text{ such that } x \in B_3 \subseteq B_1 \cap B_2$$

(B-total)  $\mathcal{B}$  is *total*:  $X = \cup \mathcal{B}$  ( $\forall x \in X, \exists B \in \mathcal{B}$  such that  $x \in B$ ).

**Theorem 7** *Let the subset  $\mathcal{B} \subseteq \text{exp}[X]$  be strongly directed, total; and put*

$$\mathcal{T} = \text{the class of all } D \in \text{exp}[X] \text{ with } D = \text{union of a subset in } \mathcal{B}.$$

*Then,*

(41-a)  $\mathcal{T}$  is (a normal family and) a topology on  $X$

(41-b)  $\mathcal{B}$  is a basis of the topology  $\mathcal{T}$ .

**Proof** (cf. Engelking [22, Ch 1, Sect 1.2]). There are two steps to be passed.

**Step 1.** Clearly,  $\mathcal{T}$  is a normal family, by the total property. In addition, each union of elements in  $\mathcal{T}$  belongs to  $\mathcal{T}$ . It remains to establish that the intersection

of any finite subset in  $\mathcal{T}$  belongs to  $\mathcal{T}$ ; this, clearly, amounts to

$$D_1, D_2 \in \mathcal{T} \text{ imply } D_1 \cap D_2 \in \mathcal{T}.$$

Let  $D_1, D_2 \in \mathcal{T}$  be given. For the arbitrary fixed  $x \in D_1 \cap D_2$ , there exist  $B_1, B_2 \in \mathcal{B}$  such that  $x \in B_1 \subseteq D_1, x \in B_2 \subseteq D_2$ ; hence,  $x \in B_1 \cap B_2 \subseteq D_1 \cap D_2$ . From the strongly directed property, there exists  $B_3 \in \mathcal{B}$  such that  $x \in B_3 \subseteq B_1 \cap B_2 \subseteq D_1 \cap D_2$ . This, along with the arbitrariness of  $x$  in  $D_1 \cap D_2$ , tells us that

$$D_1 \cap D_2 = \text{union of a subset in } \mathcal{B}; \text{ whence, } D_1 \cap D_2 \in \mathcal{T};$$

so that,  $\mathcal{T}$  is a topology on  $X$ .

**Step 2.** By the preceding step,  $(B = \cup\{B\} \in \mathcal{T}, \forall B \in \mathcal{B})$ ; and this, along with definition of  $\mathcal{T}$ , gives the desired conclusion.

(IV) Given the topology  $\mathcal{T}$  over  $X$ , let us say that the family  $\mathcal{A} \subseteq \exp[X]$  is a *subbasis* of it, when

the class  $\mathcal{B}$  of intersections of finite subsets in  $\mathcal{A}$  is a basis for  $\mathcal{T}$ .

Clearly, this in particular means that  $\mathcal{A} \subseteq \mathcal{B}$ . Moreover, as  $\cup\mathcal{A} = \cup\mathcal{B}$  and  $\cup\mathcal{B} = X$ , we have  $\cup\mathcal{A} = X$ ; i.e.:  $\mathcal{A}$  is *total*.

Conversely, let  $\mathcal{A} \subseteq \exp[X]$  be a class of sets. As before, we may ask whether there exists a topology  $\mathcal{T}$  on  $X$  such that  $\mathcal{A}$  is a subbasis for  $\mathcal{T}$ . By the above developments, this happens when

$$\mathcal{B} = \text{the class of intersections of finite subsets in } \mathcal{A}$$

is endowed with the properties

(s-dir)  $\mathcal{B}$  is *strongly directed*: for each  $B_1, B_2 \in \mathcal{B}$  and each  $x \in B_1 \cap B_2$  there exists  $B_3 \in \mathcal{B}$  such that  $x \in B_3 \subseteq B_1 \cap B_2$

(B-total)  $\mathcal{B}$  is *total*:  $X = \cup\mathcal{B} (\forall x \in X, \exists B \in \mathcal{B} \text{ such that } x \in B)$ .

The former of these is evident (in our case), via

$$B_1, B_2 \in \mathcal{B} \text{ implies } B_1 \cap B_2 \in \mathcal{B}.$$

So, it remains to verify the latter; which, under  $\cup\mathcal{A} = \cup\mathcal{B}$ , means

(A-total)  $\mathcal{A}$  is *total*:  $X = \cup\mathcal{A} (\forall x \in X, \exists A \in \mathcal{A} \text{ such that } x \in A)$ .

Summing up, the following answer to the posed question is available.

**Theorem 8** *Let the family  $\mathcal{A} \subseteq \exp[X]$  be total; and put*

$$\mathcal{T} = \{D \in \exp[X]; D = \text{union of intersections of finite subsets in } \mathcal{A}\}.$$

Then,

(42-a)  $\mathcal{T}$  is (a normal family and) a topology on  $X$

(42-b)  $\mathcal{A}$  is a subbasis of the topology  $\mathcal{T}$ .

**(C) [Comparison criteria]**

In the following, some elementary facts involving comparison of topologies will be discussed.

Given the couple of topological spaces  $(X, \mathcal{T})$  and  $(Y, \mathcal{S})$ , let us say that  $f : X \rightarrow Y$  is  $(\mathcal{T}, \mathcal{S})$ -continuous, when

$$\text{for each } G \in \mathcal{S} \text{ we have } f^{-1}(G) \in \mathcal{T}.$$

Note that, an equivalent characterization of this property is by means of cotopologies; precisely,  $f : X \rightarrow Y$  is  $(\mathcal{T}, \mathcal{S})$ -continuous, iff

$$\text{for each } H \in Y \setminus \mathcal{S} \text{ we have } f^{-1}(H) \in X \setminus \mathcal{T}.$$

This concept may be used towards the comparison of topologies. Some preliminaries are needed. Let  $X$  be a nonempty set. Given the topologies  $\mathcal{T}, \mathcal{S}$  over  $X$ , let us introduce the relation

$$\mathcal{T} \subseteq \mathcal{S} : \text{ each } \mathcal{T}\text{-open set is } \mathcal{S}\text{-open.}$$

This will be referred to as:  $\mathcal{T}$  is *coarser* than  $\mathcal{S}$ ; or:  $\mathcal{S}$  is *finer* than  $\mathcal{T}$ .

Let  $i_X : X \rightarrow X$  stand for the *identity selfmap* ( $i_X(x) = x; x \in X$ ). The following result is now immediate; so, we do not give details.

**Proposition 9** *For the topologies  $\mathcal{T}, \mathcal{S}$  over  $X$ , the following are equivalent:*

(43-1)  $\mathcal{T}$  is coarser than  $\mathcal{S}$  (i.e.:  $\mathcal{S}$  is finer than  $\mathcal{T}$ )

(43-2) the identical application  $i_X : X \rightarrow X$  is  $(\mathcal{S}, \mathcal{T})$ -continuous.

This characterization is a handy tool for constructing topologies on a (nonempty) set  $X$ . Given a nonempty index set  $\Lambda$ , let  $((Y_\lambda, \mathcal{S}_\lambda); \lambda \in \Lambda)$  be a family of topological spaces, and  $(f_\lambda : X \rightarrow Y_\lambda; \lambda \in \Lambda)$  be a family of maps. We are interested to determine the (minimal) topology  $\mathcal{T}$  over  $X$  with respect to which

$$f_\lambda \text{ is } (\mathcal{T}, \mathcal{S}_\lambda)\text{-continuous, for each } \lambda \in \Lambda.$$

As we shall see, this topology is to be introduced by the subbase

$$\mathcal{A} = \cup\{f_\lambda^{-1}(\mathcal{S}_\lambda); \lambda \in \Lambda\}; \text{ where, by definition,}$$

$$f_\lambda^{-1}(\mathcal{S}_\lambda) = \{f_\lambda^{-1}(E_\lambda); E_\lambda \in \mathcal{S}_\lambda\}, \lambda \in \Lambda.$$

In fact,  $\mathcal{A}$  is total ( $X = \cup \mathcal{A}$ ); because  $[\emptyset, X \in f_\lambda^{-1}(\mathcal{S}_\lambda)$ , for each  $\lambda \in \Lambda$ ]. By the preceding statement,  $\mathcal{A}$  is a subbase of a topology  $\mathcal{T}$  over  $X$ , represented as

$$\mathcal{T} = \{D \in \exp[X]; D = \text{union of intersections of finite subsets in } \mathcal{A}\}.$$

The obtained topology is the desired one. Precisely, we have

**Theorem 9** *Under the above conventions,*

- (43-a)  $f_\lambda$  is  $(\mathcal{T}, \mathcal{S}_\lambda)$ -continuous, for each  $\lambda \in \Lambda$
- (43-b) *If the topology  $\mathcal{L}$  over  $X$  fulfills [ $f_\lambda$  is  $(\mathcal{L}, \mathcal{S}_\lambda)$ -continuous,  $\forall \lambda \in \Lambda$ ] then, necessarily,  $\mathcal{T} \subseteq \mathcal{L}$  ( $\mathcal{T}$  is coarser than  $\mathcal{L}$ ).*

*Remark 1* The following alternate representation of this topology is useful in many concrete cases. For each  $\lambda \in \Lambda$ , let  $\mathcal{S}_\lambda^*$  be a subbasis of  $\mathcal{S}_\lambda$ ; note that, by definition,  $\mathcal{S}_\lambda^*$  is total. Then, let us denote

$$\mathcal{A}^* = \cup \{f_\lambda^{-1}(\mathcal{S}_\lambda^*); \lambda \in \Lambda\}; \text{ where, by definition,}$$

$$f_\lambda^{-1}(\mathcal{S}_\lambda^*) = \{f_\lambda^{-1}(E_\lambda^*); E_\lambda^* \in \mathcal{S}_\lambda^*, \lambda \in \Lambda.$$

From the above observation,  $\mathcal{A}^*$  is total too; so, the formula

$$\mathcal{T}^* = \{D \in \exp[X]; D = \text{union of intersections of finite subsets in } \mathcal{A}^*\}$$

defines a topology over  $X$ . We claim that, necessarily,

$$\mathcal{T} = \mathcal{T}^*; \text{ i.e.: } \mathcal{T} \subseteq \mathcal{T}^* \text{ and } \mathcal{T}^* \subseteq \mathcal{T}.$$

The latter inclusion is clear, in view of  $\mathcal{A}^* \subseteq \mathcal{A}$ ; so, it remains to establish the former inclusion. This, in turn, reduces to

(f-int) each intersection of a finite subset in  $\mathcal{A}$  belongs to  $\mathcal{T}^*$ ;

for, in such a case,

$$D \in \mathcal{T} \text{ implies } D = \text{union of elements in } \mathcal{T}^*; \text{ whence } D \in \mathcal{T}^*.$$

Now, (f-int) means, ultimately:  $A, B \in \mathcal{A}$  implies  $A \cap B \in \mathcal{T}^*$ . To verify this, let  $A, B \in \mathcal{A}$ ; hence,

$$A = f_\lambda^{-1}(G_\lambda), B = f_\mu^{-1}(H_\mu), \text{ where } \lambda, \mu \in \Lambda, G_\lambda \in \mathcal{S}_\lambda, H_\mu \in \mathcal{S}_\mu.$$

Take some  $x \in A \cap B$ . As  $f_\lambda(x) \in G_\lambda$ , there exists  $G_\lambda^*$ =intersection of a finite subset in  $\mathcal{S}_\lambda^*$ , such that  $f_\lambda(x) \in G_\lambda^* \subseteq G_\lambda$ . Likewise, as  $f_\mu(x) \in H_\mu$ , there exists  $H_\mu^*$ =intersection of a finite subset in  $\mathcal{S}_\mu^*$ , such that  $f_\mu(x) \in H_\mu^* \subseteq H_\mu$ . Combining these, yields

$$x \in f_\lambda^{-1}(G_\lambda^*) \cap f_\mu^{-1}(H_\mu^*) \subseteq A \cap B; \text{ with}$$

$$f_\lambda^{-1}(G_\lambda^*) \cap f_\mu^{-1}(H_\mu^*) = \text{intersection of a finite subset in } \mathcal{A}^*;$$

and this, by the arbitrariness of  $x \in A \cap B$ , gives

$$A \cap B = \text{union of intersections of finite subsets in } \mathcal{A}^*; \text{ so, } A \cap B \in \mathcal{T}^*.$$

Further aspects may be found in Dugundji [20, Ch IV, Sect 1].

A basic particular case of these developments corresponds to the choice

$$Y_\lambda = X, f_\lambda = i_X, \text{ for all } \lambda \in \Lambda.$$

Precisely, let  $\Lambda$  be a nonempty index set; and  $(\mathcal{S}_\lambda; \lambda \in \Lambda)$  be a family of topologies over  $X$ . We have

**Theorem 10** *Under the described setting, there exists a unique topology  $\mathcal{T}$  over  $X$  having as subbase the class of sets  $\mathcal{A} = \cup\{\mathcal{S}_\lambda; \lambda \in \Lambda\}$ , and representable as*

$$\mathcal{T} = \{D \in \exp[X]; D = \text{union of intersections of finite subsets in } \mathcal{A}\}.$$

Moreover, the following properties hold:

(44-a)  $\mathcal{S}_\lambda \subseteq \mathcal{T}$ , for each  $\lambda \in \Lambda$

(44-b) If the topology  $\mathcal{L}$  over  $X$  fulfills  $(\mathcal{S}_\lambda \subseteq \mathcal{L}, \forall \lambda \in \Lambda)$  then,  $\mathcal{T} \subseteq \mathcal{L}$ .

In other words:

(44-c)  $\mathcal{T} = \sup\{\mathcal{S}_\lambda; \lambda \in \Lambda\}$  (in the inclusion sense).

*Remark 2* The following alternate representation of this topology is useful in many concrete cases. For each  $\lambda \in \Lambda$ , let  $\mathcal{S}_\lambda^*$  be a subbasis of  $\mathcal{S}_\lambda$ ; note that, by definition,  $\mathcal{S}_\lambda^*$  is total. Then, denote  $\mathcal{A}^* = \cup\{\mathcal{S}_\lambda^*; \lambda \in \Lambda\}$ . From the above observation,  $\mathcal{A}^*$  is also total; so, the formula

$$\mathcal{T}^* = \{D \in \exp[X]; D = \text{union of intersections of finite subsets in } \mathcal{A}^*\}$$

defines a topology over  $X$ . We claim that, necessarily,

$$\mathcal{T} = \mathcal{T}^*; \text{ i.e.: } \mathcal{T} \subseteq \mathcal{T}^* \text{ and } \mathcal{T}^* \subseteq \mathcal{T}.$$

The latter inclusion is clear, in view of  $\mathcal{A}^* \subseteq \mathcal{A}$ ; so, it remains to establish the former inclusion. This, in turn, reduces to

(f-int-id) each finite intersection of elements in  $\mathcal{A}$  belongs to  $\mathcal{T}^*$ ;

for, in such a case,

$$D \in \mathcal{T} \text{ implies } D = \text{union of elements in } \mathcal{T}^*; \text{ whence } D \in \mathcal{T}^*.$$

Now, (f-int-id) means, ultimately:  $A, B \in \mathcal{A}$  implies  $A \cap B \in \mathcal{T}^*$ . To verify this, let  $A, B \in \mathcal{A}$ ; hence,

$$A = G_\lambda, B = H_\mu, \text{ where } \lambda, \mu \in \Lambda, G_\lambda \in \mathcal{S}_\lambda, H_\mu \in \mathcal{S}_\mu.$$

Take some  $x \in A \cap B$ . As  $x \in G_\lambda$ , there exists  $G_\lambda^*$ =intersection of a finite subset in  $\mathcal{S}_\lambda^*$ , such that  $x \in G_\lambda^* \subseteq G_\lambda$ . Likewise, as  $x \in H_\mu$ , there exists  $H_\mu^*$ =intersection of a finite subset in  $\mathcal{S}_\mu^*$ , such that  $x \in H_\mu^* \subseteq H_\mu$ . Combining these, yields

$$x \in G_\lambda^* \cap H_\mu^* \subseteq A \cap B; \text{ with } G_\lambda^* \cap H_\mu^* \text{=intersection of a finite subset in } \mathcal{A}^*;$$

and this, by the arbitrariness of  $x \in A \cap B$ , gives

$$A \cap B \text{=union of intersections of finite subsets in } \mathcal{A}^*; \text{ so, } A \cap B \in \mathcal{T}^*.$$

**(D) [Basic concepts]**

In the following, some completions of the above facts will be provided. Let  $(X, \mathcal{T})$  be a topological space.

Given  $A \in \text{exp}[X]$ , let us say that  $x \in X$  is *interior* to  $A$ , when  $A \in \mathcal{V}(x)$ ; clearly,  $x \in A$ . The class of all these points will be denoted as  $\text{int}(A)$  (the *interior* of  $A$ ). A global characterization of this operator is given as

$$\text{int}(A) \text{=the union of all } D \in \mathcal{T} \text{ with } D \subseteq A.$$

Given  $A \in \text{exp}[X]$ , let us say that  $x \in X$  is *adherent* to  $A$ , when  $V \cap A \neq \emptyset$ , for all  $V \in \mathcal{V}(x)$ ; clearly,  $x \in A$  in not in general true. The class of all these points is just  $\text{cl}(A)$  (the *closure* (or, *adherence*) of  $A$ ); clearly,  $\text{int}(A) \subseteq \text{cl}(A)$ .

The introduced operators  $A \mapsto \text{int}(A)$  and  $A \mapsto \text{cl}(A)$  (from  $\text{exp}[X]$  to itself) are *duals*, in the sense

$$X \setminus \text{int}(A) = \text{cl}(X \setminus A), X \setminus \text{cl}(A) = \text{int}(X \setminus A), A \in \text{exp}[X].$$

This means that each property involving  $\text{int}(\cdot)$  has a dual property involving  $\text{cl}(\cdot)$ , and vice versa; we do not give details.

Given  $A \in \text{exp}[X]$ , we say that  $x \in X$  is *boundary* to  $A$ , when  $x \in \text{cl}(A) \cap \text{cl}(X \setminus A)$ ; as before, such a point need not be in  $A$ . The class of all these writes

$$\text{bd}(A) = \text{cl}(A) \cap \text{cl}(X \setminus A) = \text{cl}(A) \setminus \text{int}(A) \text{ (the } \textit{boundary} \text{ of } A);$$

clearly,  $\text{bd}(A) \subseteq \text{cl}(A)$ . As a consequence of this,

$$\text{cl}(A) = (\text{cl}(A) \setminus \text{int}(A)) \cup \text{int}(A) = \text{bd}(A) \cup \text{int}(A), A \in \text{exp}[X].$$



Let again  $(X, \mathcal{T})$  be a topological space. Define the properties

- (H-sep)  $(X, \mathcal{T})$  is (*Hausdorff*) *separated*:  $\forall x, y \in X$  with  $x \neq y$ , there exist  $V \in \mathcal{V}(x)$ ,  $W \in \mathcal{V}(y)$  with  $V \cap W = \emptyset$
- (q-comp)  $(X, \mathcal{T})$  is *quasi-compact*:  
 $X = \cup \mathcal{G}$  for  $\mathcal{G} \subseteq \mathcal{T}$  implies  $X = \cup \mathcal{H}$ , where  $\mathcal{H}$  =finite part of  $\mathcal{G}$ .

When both these properties hold, we say that  $(X, \mathcal{T})$  is *compact*. The basic properties of these concepts are well known; so, further details are not provided.

(E) [**Metrical structures**]

Some basic applications of these facts are to be given in a metrical context.

- (I) Let  $M$  be a nonempty set. Given the map (=pseudometric)  $d : M \times M \rightarrow R_+$ , consider the properties

- (m-1)  $d$  is *reflexive*:  $d(x, x) = 0$ , for all  $x \in X$
- (m-2)  $d$  is *triangular*:  $d(x, y) \leq d(x, z) + d(z, y)$ ,  $\forall x, y, z \in X$
- (m-3)  $d$  is *symmetric*:  $d(x, y) = d(y, x)$ , for all  $x, y \in X$
- (m-4)  $d$  is *sufficient*:  $d(x, y) = 0$  implies  $x = y$ .

When (m-1)-(m-3) hold, we say that  $d(., .)$  is a *semimetric* on  $M$ ; and  $(M, d)$  is called a *semimetric space*. On the other hand, when (m-1)-(m-4) hold, we say that  $d(., .)$  is a *metric* on  $M$ ; and  $(M, d)$  is called a *metric space*.

Let  $(M, d)$  be a semimetric space. The topology to be considered here is generated by the family of open spheres in  $M$ . Precisely, given  $a \in M$ ,  $\rho > 0$ , denote

$$M(a, \rho)(d) = \{x \in M; d(a, x) < \rho\}, \quad M[a, \rho](d) = \{x \in M; d(a, x) \leq \rho\};$$

these will be referred to as the *open* (respectively, *closed*) sphere with center  $a \in M$  and radius  $\rho > 0$ . (Clearly, both these spheres are nonempty; because  $a \in M$  is an element of them).

Define the family of (nonempty) sets

$$\mathcal{W}^*(a)(d) = \{M(a, \rho)(d); \rho > 0\}, \quad a \in M.$$

It is not hard to see that  $(\mathcal{W}^*(a)(d); a \in M)$  is a neighborhoods subsystem over  $M$ . Its attached family of (nonempty) sets  $(\mathcal{W}(a)(d); a \in M)$  defined as

$$(\forall a \in M) : W \in \mathcal{W}(a)(d) \text{ iff } W \supseteq W^*, \text{ for some } W^* \in \mathcal{W}^*(a)(d)$$

is a neighborhoods system over  $M$ . By a preceding statement, there exists a unique topology  $\mathcal{T}(d)$  over  $M$  such that

$$\mathcal{W}(a)(d) \text{ is just the } \mathcal{T}(d)\text{-neighborhoods system of } a, \text{ for each } a \in M.$$

According to its definition, we have, for each nonempty  $D \in \mathcal{T}(d)$  and each  $a \in D$ ,

$$D \in \mathcal{W}(a)(d); \text{ so, there exists } \rho > 0 \text{ such that } M(a, \rho)(d) \subseteq D.$$

Moreover,  $\mathcal{B}^*(d) = \{\mathcal{W}^*(a)(d); a \in M\}$  is a subbasis of  $\mathcal{T}(d)$ ; so,

$$(\text{for each } D \in \mathcal{T}(d)) : D = \text{union of intersections of finite subsets in } \mathcal{B}^*(d).$$

In particular,  $\mathcal{B}^*(d) \subseteq \mathcal{T}(d)$ ; whence,

$$M(a, \rho)(d) \text{ is open, for each } a \in M, \rho > 0.$$

On the other hand, by the semimetrical properties of  $d$ ,

$$(\forall a \in M, \forall \rho > 0) : M^e(a, \rho)(d) := \{x \in M; d(a, x) > \rho\} \text{ is open;}$$

$$\text{whence, } M[a, \rho](d) = M \setminus M^e(a, \rho)(d) \text{ is closed.}$$

Finally, when no confusion can arise, it would be convenient to drop the index  $(d)$  from all notations above.

(II) Let  $(M, d)$  be a metric space. The (topological) compactness property of this space was already introduced. A related notion is the following. Let us say that  $(M, d)$  is *sequentially compact*, provided

(s-comp) each sequence in  $M$  has a convergent subsequence.

**Proposition 10** *The following is valid, in (ZF-AC+DC)*

$$(M, d) \text{ is compact iff } (M, d) \text{ is sequentially compact.}$$

The verification is very similar with the one in Costinescu [17, Ch VII, Sect 3].

(III) Let  $(P, d)$  and  $(Q, e)$  be a couple of metric spaces; and  $f : P \rightarrow Q$  be a mapping. Let us say that  $f$  is *sequentially continuous*, when

$$\forall \text{ sequence } (x_n) \text{ in } P \text{ and } \forall \text{ element } x \in P : x_n \xrightarrow{d} x \text{ implies } f(x_n) \xrightarrow{e} f(x).$$

**Proposition 11** *We have, in (ZF-AC+DC),*

$$f \text{ is continuous iff } f \text{ is sequentially continuous.}$$

**Proof**

(i): Suppose that  $f$  is continuous (after the general definition); and let the sequence  $(x_n)$  in  $P$  and the point  $x \in P$  be such that  $x_n \xrightarrow{d} x$ . Let  $\varepsilon > 0$  be arbitrary fixed. By the topological definition of continuity, there exists  $\delta > 0$  such

that  $f(P(x, \delta)) \subseteq Q(f(x), \varepsilon)$ . Further, given this  $\delta > 0$ , there exists (by hypothesis) an index  $n(\delta) \in \mathbb{N}$  such that  $n \geq n(\delta)$  implies  $x_n \in P(x, \delta)$ . Combining these gives

$$n \geq n(\delta) \text{ implies } f(x_n) \in Q(f(x), \varepsilon);$$

and this, by the arbitrariness of  $\varepsilon > 0$ , tells us that  $f(x_n) \xrightarrow{e} f(x)$ .

(ii): Suppose that  $f$  is sequentially continuous; but  $f$  is not continuous (after the topological definition): there exists  $\varepsilon > 0$  such that

$$(\forall \delta > 0) : H(\delta) := f(P(x, \delta)) \setminus Q(f(x), \varepsilon) \text{ is nonempty.}$$

Let  $(\delta_n; n \geq 0)$  be a strictly descending sequence in  $R_+^0$  with  $\delta_n \rightarrow 0$ . (For, example, one may take  $(\delta_n = 2^{-n}; n \geq 0)$ ; but this is not the only possible choice). By the Denumerable Axiom of Choice (AC(N)) (deductible in (ZF-AC+DC)), there may be determined a sequence  $(x_n)$  in  $P$  with

$$(\forall n) : x_n \in H(\delta_n) \text{ (that is; } x_n \in P(x, \delta_n) \text{ and } f(x_n) \notin Q(f(x), \varepsilon)).$$

By the first half of this relation,  $x_n \xrightarrow{d} x$ ; so that, by the sequential continuity,

$$f(x_n) \xrightarrow{e} f(x); \text{ wherefrom:}$$

$$\text{there exists some index } n(\varepsilon) \geq 0, \text{ with } f(x_n) \in Q(f(x), \varepsilon), \forall n \geq n(\varepsilon).$$

This, however, is in contradiction with the second half of underlying relation; and then, the conclusion follows.

(IV) Let again  $(P, d)$  and  $(Q, e)$  be a couple of metric spaces; and  $f : P \rightarrow Q$  be a mapping. Let us say that  $f$  is *uniformly continuous*, when

$$\forall \varepsilon > 0, \exists \delta > 0, \text{ such that } d(x, y) \leq \delta \text{ implies } e(f(x), f(y)) \leq \varepsilon.$$

In this case, we may construct a mapping  $\psi : R_+ \rightarrow R_+ \cup \{\infty\}$ , according to

$$\psi(t) = \sup\{e(f(x), f(y)); d(x, y) \leq t\}, t \in R_+.$$

This will be referred to as the *uniform continuity modulus* of  $f$ ; the class of all these will be denoted as  $\mathcal{F}_u(R_+, R_+ \cup \{\infty\})$ . Clearly,

(uc-1)  $\psi$  is increasing and zero-continuous [ $\psi(0) = 0 = \psi(0 + 0)$ ],

(uc-2)  $e(f(x), f(y)) \leq \psi(d(x, y)), \forall x, y \in P$ ;

since the verification is immediate, we do not give details.

A basic example of this type is to be obtained in the compactness context.

**Theorem 11** *Suppose that  $(P, d)$  is sequentially compact. Then, the generic inclusion is valid, in  $(ZF-AC+DC)$ :*

$$(\forall f \in \mathcal{F}(P, Q)) : f \text{ is continuous implies } f \text{ is uniformly continuous.}$$

**Proof** Suppose that  $f$  is continuous, but not uniformly continuous. There exists then  $\varepsilon > 0$ , such that

$$C(\delta) := \{(x, y) \in P \times P; d(x, y) \leq \delta, e(f(x), f(y)) > \varepsilon\} \neq \emptyset, \text{ for each } \delta > 0.$$

Let  $(\delta_n; n \geq 0)$  be a sequence in  $R_+^0$  with  $\delta_n \rightarrow 0$ . (For example, we may take  $(\delta_n = 2^{-n}; n \geq 0)$ ; hence, this is not depending on (AC)). From the Denumerable Axiom of Choice (AC(N)) (deductible, as precise, in  $(ZF-AC+DC)$ ), we get a sequence  $((x_n, y_n); n \geq 0)$  in  $P \times P$ , such that

$$(x_n, y_n) \in C(\delta_n), \forall n; \text{ that is: } d(x_n, y_n) \leq \delta_n, e(f(x_n), f(y_n)) > \varepsilon, \forall n.$$

From the sequential compactness of  $(P, d)$ , there exists subsequences  $(x_n^*), (y_n^*)$  of  $(x_n)$  and  $(y_n)$  respectively, and elements  $x^*, y^* \in P$ , such that (from the first part of this relation)

$$d(x_n^*, y_n^*) \leq \delta_n, \forall n, x_n \xrightarrow{d} x^*, y_n^* \xrightarrow{d} y^*; \text{ whence, } x^* = y^*.$$

Combining with the second part of the same relation, we get (by a limit process)

$$0 = e(f(x^*), f(y^*)) \geq \varepsilon; \text{ contradiction.}$$

Hence, our working assumption is not acceptable; and the conclusion follows.

**(F) [Gauge spaces]**

Let  $X$  be a nonempty set; and  $\Lambda$  be some (nonempty) index set. For each  $\lambda \in \Lambda$ , let  $d_\lambda : X \times X \rightarrow R_+$  be a semimetric over  $X$ ; and  $\mathcal{T}(d_\lambda)$  stand for the associated topology (see above).

Technically speaking, these data generate two basic structures on  $X$ .

(I) The former of these is represented by the supremum topology of the family  $\{\mathcal{T}(d_\lambda); \lambda \in \Lambda\}$ . Precisely, note that for each  $\lambda \in \Lambda$ ,

$$\mathcal{A}_\lambda = \{X(x, \rho)(d_\lambda); x \in X, \rho > 0\} \text{ is a subbase of } \mathcal{T}(d_\lambda);$$

hence, in particular,  $\mathcal{A}_\lambda$  is total. Denote further

$$\mathcal{A} = \cup\{\mathcal{A}_\lambda; \lambda \in \Lambda\}; \text{ clearly, } \mathcal{A} \text{ is total too.}$$

By a preceding result, the formula

$$\mathcal{T} = \{D \in \exp[X]; D = \text{union of intersections of finite subsets in } \mathcal{A}\}$$

defines a topology over  $X$ . This, as precise, is just the supremum topology of the family  $\{\mathcal{T}(d_\lambda); \lambda \in \Lambda\}$ . Let  $\text{cl}(\cdot)$  stand for the associated closure operator. Note that, according to its definition, we have for each (nonempty) subset  $Y$  of  $X$

- (g-int)  $x \in \text{int}(Y)$  iff there exists an intersection  $W$  of a finite subset in  $\mathcal{A}$  with  $x \in W \subseteq Y$
- (g-clo)  $x \in \text{cl}(Y)$  iff each intersection  $W$  of a finite subset in  $\mathcal{A}$  with  $x \in W$  fulfills  $W \cap Y \neq \emptyset$ .

(II) The latter of these is represented by the (sequential) conv-Cauchy structure induced by the family  $D = (d_\lambda; \lambda \in \Lambda)$  of these semimetrics.

Take an arbitrary sequence  $(x_n; n \geq 0)$  in  $X$ . Given  $\lambda \in \Lambda$ , the  $d_\lambda$ -convergence of this sequence towards an  $x \in X$  [depicted as:  $x_n \xrightarrow{d_\lambda} x$ ], means:

$$d_\lambda(x_p, x) \rightarrow 0 \text{ as } p \rightarrow \infty$$

(i.e.:  $\forall \varepsilon > 0, \exists i = i(\lambda, \varepsilon)$ , such that  $i \leq p \implies d_\lambda(x_p, x) \leq \varepsilon$ ).

If this holds for all  $\lambda \in \Lambda$ , then  $(x_n; n \geq 0)$  is said to  $D$ -converge towards  $x$  [written as:  $x_n \xrightarrow{D} x$ ]. The set of all such points  $x$  will be denoted as  $D - \lim_n(x_n)$ ; when it is nonempty, then  $(x_n; n \geq 0)$  is called  $D$ -convergent. On the other hand, given  $\lambda \in \Lambda$ , the  $d_\lambda$ -Cauchy property of  $(x_n; n \geq 0)$  means:

$$d_\lambda(x_p, x_q) \rightarrow 0 \text{ as } p, q \rightarrow \infty, p \leq q$$

(i.e.:  $\forall \varepsilon > 0, \exists j := j(\lambda, \varepsilon)$ , such that  $j \leq p \leq q \implies d_\lambda(x_p, x_q) \leq \varepsilon$ ).

If this holds for each  $\lambda \in \Lambda$ , we say that  $(x_n)$  is  $D$ -Cauchy; the class of all such sequences will be denoted as  $\text{Cauchy}(D)$ . By definition, the triple  $(X, \Lambda; D)$  endowed with the conv-Cauchy structure  $((\xrightarrow{D}), \text{Cauchy}(D))$  and the regularity condition

$$D \text{ is sufficient : } x, y \in X \text{ and } (d_\lambda(x, y) = 0, \forall \lambda \in \Lambda) \text{ imply } x = y$$

is called a *gauge space*. Note that, in this setting, any  $D$ -convergent sequence is  $D$ -Cauchy too; the reciprocal is not in general valid.

Using the previous conventions, we may introduce a  $D$ -closure operator  $A \mapsto \text{Dcl}(A)$  from  $\exp[X]$  to itself, as: for each  $A \in \exp[X]$ ,

$$y \in \text{Dcl}(A) \text{ iff } y = D - \lim_n(y_n), \text{ for some sequence } (y_n) \text{ in } A.$$

It is not hard to see that  $A \mapsto \text{Dcl}(A)$  is a *semi-closure* over  $X$ , in the sense

- (Dcl-1) (identity)  $\emptyset = \text{Dcl}(\emptyset)$ ,  $X = \text{Dcl}(X)$ ,
- (Dcl-2) (progressiveness)  $Y \subseteq \text{Dcl}(Y)$ ,  $\forall Y \in \exp[X]$ ,
- (Dcl-3) (additivity)  $\text{Dcl}(U \cup V) = \text{Dcl}(U) \cup \text{Dcl}(V)$ ,  $\forall U, V \in \exp[X]$ .

Note that, as a direct consequence of these, we also have

- (Dcl-4) (monotonicity)  $Y_1 \subseteq Y_2$  implies  $\text{Dcl}(Y_1) \subseteq \text{Dcl}(Y_2)$ .

Further, call  $A \in \exp[X]$ , *D-closed* provided  $A = \text{Dcl}(A)$ ; note that in this case,

the *D*-limit of each sequence in  $A$  belongs to  $A$ .

Denote, for simplicity

$$\mathcal{K}[D] = \text{the class of all } D\text{-closed subsets in } X.$$

Clearly,  $\mathcal{K}[D] \subseteq \exp[X]$  is a normal family, with the finite union property

- (K-D-1) the union of any finite subset of  $\mathcal{K}[D]$  is in  $\mathcal{K}[D]$ .

However, the arbitrary intersection property

- (K-D-2) the intersection of any subset in  $\mathcal{K}[D]$  belongs to  $\mathcal{K}[D]$

is not in general true; so that,  $\mathcal{K}[D]$  is not a cotopology on  $X$ ; we then say that  $\mathcal{K}[D]$  is an *almost cotopology* on  $X$ . Equivalently, this tells us that

$$\mathcal{T}[D] = X \setminus \mathcal{K}[D] \text{ (in the sense: } Z \in \mathcal{T}[D] \text{ iff } X \setminus Z \in \mathcal{K}[D])$$

is a normal family with the finite intersection property

- (T-D-1) the intersection of any finite subset of  $\mathcal{T}[D]$  is in  $\mathcal{T}[D]$ .

However, the arbitrary union property

- (T-D-2) the union of any subset in  $\mathcal{T}[D]$  belongs to  $\mathcal{T}[D]$

is not in general true; so that,  $\mathcal{T}[D]$  is not topology on  $X$ ; we then say that  $\mathcal{T}[D]$  is an *almost topology* on  $X$  [referred to as the *gauge almost topology*]. An explanation of this bad property is due to the closure operator  $A \mapsto \text{Dcl}(A)$  being *not involutive*; i.e., a property like

- (Dcl-inv)  $\text{Dcl}(\text{Dcl}(A)) = \text{Dcl}(A)$ , for each  $A \in \exp[X]$

is not in general true; so that,  $A \mapsto \text{Dcl}(A)$  is not a closure over  $X$  according to Kuratowski [38, Ch I, Sect 4]. Further aspects of this problem may be found in Engelking [22, Ch 1, Sect 1.2].

Finally, let us compare the supremum topology  $\mathcal{T} = \sup\{\mathcal{T}(d_\lambda); \lambda \in \Lambda\}$  with the gauge almost topology  $\mathcal{T}[D]$ . In this direction, we have the properties

(p-1)  $(\forall Y \in \exp[X]): \text{Dcl}(Y) \subseteq \text{cl}(Y)$

(p-2)  $\mathcal{T} \subseteq \mathcal{T}[D];$  i.e.:  $Y=\text{closed}$  (open) implies  $Y=D\text{-closed}$  ( $D\text{-open}$ ).

In fact, letting  $x \in \text{Dcl}(Y)$ , there exists a sequence  $(x_n)$  in  $Y$  with

$$x_n \xrightarrow{D} x; \text{ hence, } x_n \xrightarrow{d_\lambda} x, \text{ for each } \lambda \in \Lambda.$$

Let  $W = \cap \{X(x, \rho_i)(d_{\lambda(i)}); i \in I\}$  be a finite intersection (of open spheres) that includes  $x$ ; where,  $I$  is a (nonempty) finite index set. By the above convergence property, there must be some index  $n(x)$ , such that

$$\{x_n; n \geq n(x)\} \subseteq W; \text{ whence, } Y \cap W \neq \emptyset.$$

This yields  $x \in \text{cl}(Y)$ ; and proves the first property. As a consequence,

$$Y=\text{closed} \text{ implies } \text{cl}(Y) = Y \subseteq \text{Dcl}(Y) \subseteq \text{cl}(Y); \text{ hence, } Y=D\text{-closed};$$

wherefrom, the second property holds too.

### 5 Gheorghiu Functional Contractions

Let  $X$  be a nonempty set; and  $\Lambda$  be some nonempty (index) set. For each  $\lambda \in \Lambda$ , let  $d_\lambda : X \times X \rightarrow R_+$  be a semimetric over  $X_\lambda$ ; and  $\mathcal{T}(d_\lambda)$  stand for the associated to  $d_\lambda$  topology over  $X$ .

Technically speaking, these data generate two basic structures on  $X$ .

- (I) The former of these is represented by the supremum topology of the family  $\{\mathcal{T}(d_\lambda); \lambda \in \Lambda\}$ . For the moment, this topology is not essential for us; however, in the homotopical fixed point theory, a limited use of it will be made in the reasonings to be developed.
- (II) The latter of these is represented by the (sequential) conv-Cauchy structure induced by the family  $D = (d_\lambda; \lambda \in \Lambda)$  of these semimetrics. This, essentially, may be described in the following way.

Take an arbitrary sequence  $(x_n; n \geq 0)$  in  $X$ . Given  $\lambda \in \Lambda$ , the  $d_\lambda$ -convergence of this sequence towards an  $x \in X$  [depicted as:  $x_n \xrightarrow{d_\lambda} x$ ], means:

$$d_\lambda(x_p, x) \rightarrow 0 \text{ as } p \rightarrow \infty$$

$$(\text{i.e.: } \forall \varepsilon > 0, \exists i = i(\lambda, \varepsilon), \text{ such that } i \leq p \implies d_\lambda(x_p, x) \leq \varepsilon).$$

If this holds for all  $\lambda \in \Lambda$ , then  $(x_n; n \geq 0)$  is said to  $D$ -converge towards  $x$  [written as:  $x_n \xrightarrow{D} x$ ]. The set of all such points  $x$  will be denoted as  $D - \lim_n(x_n)$ ; when it is nonempty, then  $(x_n; n \geq 0)$  is called  $D$ -convergent. On the other hand, given  $\lambda \in \Lambda$ , the  $d_\lambda$ -Cauchy property of  $(x_n; n \geq 0)$  means:

$$d_\lambda(x_p, x_q) \rightarrow 0 \text{ as } p, q \rightarrow \infty, p \leq q$$

(i.e.:  $\forall \varepsilon > 0, \exists j := j(\lambda, \varepsilon)$ , such that  $j \leq p \leq q \implies d_\lambda(x_p, x_q) \leq \varepsilon$ ).

If this holds for each  $\lambda \in \Lambda$ , we say that  $(x_n; n \geq 0)$  is *D-Cauchy*. Note that any *D-convergent* sequence is *D-Cauchy* too; the reciprocal is not in general valid. By definition, the triple  $(X, \Lambda; D)$  endowed with this conv-Cauchy structure and with the regularity condition

$$D \text{ is sufficient : } x, y \in X \text{ and } (d_\lambda(x, y) = 0, \forall \lambda \in \Lambda) \text{ imply } x = y$$

is called a *gauge space*.

Let  $(X, \Lambda; D)$  be a gauge space. Remember that the subset  $Y$  of  $X$ , is called *asingleton*, if  $[y_1, y_2 \in Y \text{ imply } y_1 = y_2]$ ; and *singleton* if, in addition,  $Y$  is nonempty; note that, in this case  $Y = \{y\}$ , for some  $y \in X$ .

Further, let  $T : X \rightarrow X$  be a selfmap of  $X$ . As usual,  $\text{Fix}(T) := \{z \in X; z = Tz\}$  stands for the set of all fixed points of  $T$  in  $X$ . These are to be determined in the context below, comparable with the one in Rus [50, Ch 2, Sect 2.2]:

- (gpic-1) We say that  $T$  is *fix-asingleton*, if  $\text{Fix}(T)$  is an asingleton; and *fix-singleton*, if  $\text{Fix}(T)$  is a singleton
- (gpic-2) We say that  $x \in X$  is a *Picard point* (modulo  $(D, T)$ ), when the iterative sequence  $(T^n x; n \geq 0)$  is *D-Cauchy*. If this property holds for all  $x \in X$ , we say that  $T$  is a *Picard operator* (modulo  $D$ )
- (gpic-3) We say that  $x \in X$  is a *strongly Picard point* (modulo  $(D, T)$ ), when  $(T^n x; n \geq 0)$  is *D-convergent* with  $\lim_n(T^n x) \in \text{Fix}(T)$ . If this property holds for all  $x \in X$ , we say that  $T$  is a *strongly Picard operator* (modulo  $D$ ).

The sufficient (regularity) conditions for such properties are being founded on *orbital* concepts (in short: o-concepts). Given  $x \in X$ , let us say that the sequence  $(z_n; n \geq 0)$  in  $X$  is *T-orbital* with respect to  $x$ , when  $(z_n = T^n x; n \geq 0)$ ; if  $x \in X$  is generic, we then say that  $(z_n; n \geq 0)$  is a *T-orbital* (or, equivalently: an *o-sequence*).

- (greg-1) Call  $X, (o;D)$ -complete, provided (for each o-sequence) *D-Cauchy*  $\implies$  *D-convergent*
- (greg-2) We say that  $T$  is *(o;D)-continuous*, if  $[(z_n)=\text{o-sequence and } z_n \xrightarrow{D} z]$  imply  $Tz_n \xrightarrow{D} Tz$ .

As a completion of these, we must now formulate the metrical contractive conditions upon our data. Let  $\Omega = (\omega_\lambda; \lambda \in \Lambda)$  be a family over  $\mathcal{F}(R_+)$ . Then, let us introduce the mappings: for each  $x, y \in X$ , and each  $\lambda \in \Lambda$

$$M_1(x, y; T; d_\lambda) = d_\lambda(x, y),$$

$$M_2(x, y; T; d_\lambda) = \max\{d_\lambda(x, Tx), d_\lambda(y, Ty)\},$$



$$M_3(x, y; T; d_\lambda) = (1/2)[d_\lambda(x, Ty) + d_\lambda(y, Tx)],$$

$$M(x, y; T; d_\lambda) = \max\{M_1(x, y; T; d_\lambda), M_2(x, y; T; d_\lambda), M_3(x, y; T; d_\lambda)\}.$$

We say that  $T$  is  $(\Omega, M)$ -contractive, if

$$d_\lambda(Tx, Ty) \leq \omega_\lambda(M(x, y; T; d_\lambda)), \forall x, y \in X, \forall \lambda \in \Lambda.$$

The deep part of the main result to be stated is an auxiliary fact involving almost Picard properties over semimetric structures. Some preliminaries are needed.

Let  $X$  be a nonempty set; and  $d : X \times X \rightarrow R_+$  be a semimetric over  $X$ ; the couple  $(X, d)$  will be referred to as a *semimetric space*. Further, take some  $T \in \mathcal{F}(X)$ . In the following, sufficient conditions are given so that

- (AP)  $T$  is almost Picard (modulo  $d$ ); in the sense
- (ap-1)  $(\forall x \in X)$ : the iterative sequence  $(T^n x; n \geq 0)$  is  $d$ -Cauchy
- (ap-2)  $(\forall x \in X, \forall z \in X)$ :  $T^n x \xrightarrow{d} z$  implies  $d(z, Tz) = 0$ .

The conditions in question are to be written as follows. Let the functional classes

$$\mathcal{F}_0(R_+), \mathcal{F}_0(re)(R_+), \mathcal{F}_0(in)(R_+), \mathcal{F}_0(re, in)(R_+)$$

be introduced as before. Then, let us define the mappings: for each  $x, y \in X$ ,

$$M_1(x, y; T; d) = d(x, y),$$

$$M_2(x, y; T; d) = \max\{d(x, Tx), d(y, Ty)\},$$

$$M_3(x, y; T; d) = (1/2)[d(x, Ty) + d(y, Tx)],$$

$$M(x, y; T; d) = \max\{M_1(x, y; T; d), M_2(x, y; T; d), M_3(x, y; T; d)\}.$$

Given  $\omega \in \mathcal{F}(R_+)$ , we say that  $T$  is  $(\omega, M)$ -contractive, provided

$$d(Tx, Ty) \leq \omega(M(x, y; T; d)), \forall x, y \in X.$$

The following almost Picard type statement involving these data is available.

**Theorem 12** *Suppose that the selfmap  $T$  is  $(\omega, M)$ -contractive, where the function  $\omega \in \mathcal{F}_0(re, in)(R_+)$  fulfills*

$$\omega \text{ is Meir-Keeler admissible; or, equivalently: Matkowski admissible.}$$

*Then,  $T$  is an almost Picard operator (see above).*

A direct verification was provided in Leader [39]. For completeness reasons, we will however develop the argument, with certain modifications.

**Proof** Fix some point  $x_0 \in X$ , and put  $(x_n = T^n x_0; n \geq 0)$ .

**Step 1.** From the contraction hypothesis, we have

$$d(x_{n+1}, x_{n+2}) \leq \omega(\max\{d(x_n, x_{n+1}), d(x_{n+1}, x_{n+2}), (1/2)d(x_n, x_{n+2})\}), \forall n.$$

In view of

$$(1/2)d(x_n, x_{n+2}) \leq \max\{d(x_n, x_{n+1}), d(x_{n+1}, x_{n+2})\}, \forall n,$$

the underlying relation becomes

$$d(x_{n+1}, x_{n+2}) \leq \omega(\max\{d(x_n, x_{n+1}), d(x_{n+1}, x_{n+2})\}), \forall n.$$

If, for some  $n$ , one has

$$d(x_{n+1}, x_{n+2}) > d(x_n, x_{n+1})(\geq 0),$$

then, by the above relation (and regressiveness of  $\omega$ )

$$d(x_{n+1}, x_{n+2}) \leq \omega(d(x_{n+1}, x_{n+2})) < d(x_{n+1}, x_{n+2}); \text{ a contradiction.}$$

Hence, necessarily,

$$d(x_{n+1}, x_{n+2}) \leq d(x_n, x_{n+1}); \text{ whence, } d(x_{n+1}, x_{n+2}) \leq \omega(d(x_n, x_{n+1})), \forall n;$$

and this (along with  $\omega$  being Matkowski admissible) yields

$$(x_n) \text{ is } d\text{-asymptotic: } r_n := d(x_n, x_{n+1}) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

**Step 2.** Let  $\varepsilon > 0$  be arbitrary fixed; and  $\delta > 0$  be attached to it, by the Meir-Keeler admissible property of  $\omega$ :

$$\begin{aligned} \varepsilon < t < \varepsilon + \delta &\text{ implies } \omega(t) \leq \varepsilon; \text{ or, equivalently} \\ t < \varepsilon + \delta &\text{ implies } \omega(t) \leq \varepsilon \text{ (as } \omega\text{-regressive, } \omega(0) = 0); \end{aligned}$$

clearly, without loss, one may assume that  $\delta < \varepsilon$ . By the  $d$ -asymptotic property we just established, there exists a rank  $n(\delta)$ , such that

$$(d\text{-asy}) (\forall n \geq n(\delta)): d(x_n, x_{n+1}) < \delta/4; \text{ hence, } d(x_n, x_{n+2}) < \delta/2.$$

We prove by induction that, for each  $j \geq 1$ , the following relation holds

$$(d\text{-C};j) \quad d(x_n, x_{n+j}) < \varepsilon + \delta/2, \forall n \geq n(\delta);$$

wherefrom,  $(x_n; n \geq 0)$  is  $d$ -Cauchy. The case  $j \in \{1, 2\}$  is evident, by the evaluation (d-asy). Assume that (d-C;j) holds for all  $j \in \{1, \dots, p\}$ , where

$p \geq 2$ ; we must establish that (d-C;p+1) holds too. Let  $n \geq n(\delta)$  be arbitrary fixed. From the inductive hypothesis (and our asymptotic relation)

$$\begin{aligned} d(x_n, x_{n+p}), d(x_{n+1}, x_{n+p}) &< \varepsilon + \delta/2 \\ d(x_n, x_{n+1}), d(x_{n+p}, x_{n+p+1}) &< \delta/4 < \varepsilon + \delta/2. \end{aligned}$$

This, along with the triangular inequality, gives us

$$\begin{aligned} d(x_n, x_{n+p+1}) &\leq d(x_n, x_{n+p}) + d(x_{n+p}, x_{n+p+1}) < \varepsilon + 3\delta/4; \\ \text{so that, by definition, } M(x_n, x_{n+p}; T; d) &< \varepsilon + \delta. \end{aligned}$$

Combining with contractive hypothesis and Meir-Keeler property of  $\omega$ , gives

$$d(x_{n+1}, x_{n+p+1}) = d(Tx_n, Tx_{n+p}) \leq \omega(M(x_n, x_{n+p}; T; d)) \leq \varepsilon;$$

wherefrom (again by the triangular inequality)

$$d(x_n, x_{n+p+1}) \leq d(x_n, x_{n+1}) + d(x_{n+1}, x_{n+p+1}) < \varepsilon + \delta/4 < \varepsilon + \delta/2;$$

and our claim follows.

**Step 3.** Suppose that  $z \in X$  is such that

$$x_n \xrightarrow{d} z \text{ (i.e.: } d(x_n, z) \rightarrow 0), \text{ as } n \rightarrow \infty.$$

We claim that, necessarily,  $d(z, Tz) = 0$ . Suppose not:  $b := d(z, Tz) > 0$ . From the contractive condition,

$$d(x_{n+1}, Tz) \leq \omega(M(x_n, z; T; d)), \forall n.$$

By the convergence and  $d$ -asymptotic properties, there exists a rank  $n(b)$ , with

$$\begin{aligned} (\forall n \geq n(b)) : d(x_n, z), d(x_n, Tx_n), d(Tx_n, z) &< b/2; \\ \text{whence, } M_1(x_n, z; T; d) &< b/2, M_2(x_n, z; T; d) = b. \end{aligned}$$

This, along with the triangle inequality, gives

$$\begin{aligned} (\forall n \geq n(b)) : d(x_n, Tz) &\leq d(x_n, z) + d(z, Tz) < 3b/2; \\ \text{wherefrom } M_3(x_n, z; T; d) &< b, M(x_n, z; T; d) = b. \end{aligned}$$

Replacing these into the contractive condition, we derive

$$d(x_{n+1}, Tz) \leq \omega(b), \forall n \geq n(b).$$

Passing to limit as  $n \rightarrow \infty$  gives (by a previous auxiliary fact)

$$b \leq \omega(b) < b; \text{ contradiction.}$$

Hence, our initial assumption about  $z$  cannot be true; so that,  $d(z, Tz) = 0$ . The proof is thereby complete.

*Remark 3* Remember that  $\omega \in \mathcal{F}_0(re)(R_+)$  is Matkowski admissible, when

$$\omega \text{ is Boyd-Wong admissible: } \Lambda^+ \omega(s) < s, \forall s \in R_+^0.$$

In particular, when  $\omega \in \mathcal{F}_0(re, in)(R_+)$ , this last condition means

$$\omega \text{ is strongly regressive } [\omega(s + 0) < s, \text{ for all } s > 0].$$

It is worth noting that, under the Boyd-Wong property, there is a specific way of completing the second part of statement above. For technical reasons, we will describe it in what follows.

Suppose that the  $d$ -asymptotic sequence  $(x_n)$  is not  $d$ -Cauchy; and fix some couple  $(\beta, \gamma) \in (> ; \text{adm}(R_+^0; (x_n)))$  as well as the rank sequence  $(\lambda(k) = 1; k \geq 0)$ . By a previous auxiliary statement, there exist  $b \in R_+^0$ , a rank sequence  $(J(k); k \geq 0)$  in  $N(1, \leq)$  and a couple of rank-sequences  $(m(k); k \geq 0)$ ,  $(n(k); k \geq 0)$ , with

$$(p-1) \quad k + 1 \leq J(k) \leq m(k) < m(k) + 3\lambda(k) < n(k), \forall k \geq 0$$

$$(p-2) \quad \forall p, q \in N[0, 3\lambda(0)], u_k(p, q) := d(x_{m(k)+p}, x_{n(k)+q}) \rightarrow b + \text{ as } k \rightarrow \infty.$$

From the contractive hypothesis,

$$u_k(1, 1) \leq \omega(v_k), \forall k; \text{ where, by definition}$$

$$v_k = \max\{u_k(0, 0), r_{m(k)}, r_{n(k)}, (1/2)[u_k(0, 1) + u_k(1, 0)]\}, k \geq 0.$$

By the above convergence properties,

$$v_k \rightarrow b + \text{ as } k \rightarrow \infty; \text{ because } u_k(0, 0), u_k(0, 1), u_k(1, 0) \rightarrow b + \text{ as } k \rightarrow \infty.$$

Passing to  $\limsup$  as  $k \rightarrow \infty$ , in the previous relation, gives (cf. a preceding result)

$$b \leq \Lambda^+ \omega(b) < b; \text{ contradiction.}$$

Hence, our working assumption is not acceptable; and the conclusion follows.

Under these preliminaries, we may now pass to our basic fixed point result in this exposition.

**Theorem 13** *Suppose that  $T$  is  $(\Omega, M)$ -contractive, where  $\Omega = (\omega_\lambda; \lambda \in \Lambda)$  is a family over  $\mathcal{F}_0(re, in)(R_+)$ , fulfilling*

$$(\forall \lambda \in \Lambda): \omega_\lambda \text{ is Meir-Keeler admissible}$$

*(or, equivalently: Matkowski admissible).*

In addition, let  $X$  be  $(o;D)$ -complete. Then,

(52-a)  $T$  is fix-asingleton

(52-b)  $T$  is strongly Picard (modulo  $D$ ).

**Proof** There are several steps to be passed.

**Part 1** We prove that  $T$  is fix-asingleton. Let  $z_1, z_2 \in \text{Fix}(T)$  be arbitrary fixed. By the contractive condition,

$$d_\lambda(z_1, z_2) \leq \omega_\lambda(d_\lambda(z_1, z_2)), \lambda \in \Lambda.$$

This, by the regressive property of the functions in  $\Omega$ , gives

$$d_\lambda(z_1, z_2) = 0, \forall \lambda \in \Lambda; \text{ wherefrom (as } D \text{ is sufficient) } z_1 = z_2.$$

**Part 2** Take some  $x_0 \in X$ , and put  $(x_n := T^n x_0; n \geq 0)$ ; clearly,  $(x_n)$  is orbital. By the preceding semimetric statement, one derives that

(52-c)  $(x_n)$  is  $D$ -Cauchy (i.e.:  $d_\lambda$ -Cauchy, for each  $\lambda \in \Lambda$ ).

**Part 3** Summing up,  $(x_n; n \geq 0)$  is an orbital  $D$ -Cauchy sequence. As  $X$  is  $(o;D)$ -complete, there must be some (uniquely determined)  $z \in X$ , with

$$x_n \xrightarrow{D} z \text{ (hence, } x_n \xrightarrow{d_\lambda} z, \text{ for each } \lambda \in \Lambda).$$

Again by the preceding semimetric statement, we get

(52-d)  $d_\lambda(z, Tz) = 0$ , for each  $\lambda \in \Lambda$ .

This, along with  $D$ =sufficient, gives  $z = Tz$ ; i.e.:  $z \in \text{Fix}(T)$ .

In particular, when  $\Lambda$  is a singleton, this main result covers the one in Leader [39]. Passing to the general case, note that our main result includes the trivial quasi-order version of the related 2017 one in Turinici [67]; and has large overlaps with the fixed point result in Agarwal et al. [1]. In fact, it is not difficult to get a quasi-order version of our main result so as to include both these statements; we do not give details. An interesting problem to be posed is that of this statement having Maia type extensions like in Gheorghiu [26]; further aspects will be discussed elsewhere. Some applications of such results to integral equations may be found in Gheorghiu and Turinici [28] or Cherichi and Samet [15]; see also the Angelov's monograph [4, Ch II, Sect 2.2].

## 6 Homotopic Fixed Points

Let  $X$  be a nonempty set; and  $\Lambda$  be some nonempty (index) set. For each  $\lambda \in \Lambda$ , let  $d_\lambda : X \times X \rightarrow R_+$  be a semimetric over  $X$ ; and  $\mathcal{T}(d_\lambda)$  stand for the associated topology. Remember that there are two basic structures on  $X$  to be used further.

- (I) The former of these is represented by the supremum topology of the family  $\{\mathcal{T}(d_\lambda); \lambda \in \Lambda\}$ . Precisely, note that for each  $\lambda \in \Lambda$ ,

$$\mathcal{A}_\lambda = \{X(x, \rho)(d_\lambda); x \in X, \rho > 0\} \text{ is a subbase of } \mathcal{T}(d_\lambda);$$

hence, in particular,  $\mathcal{A}_\lambda$  is total. Denote further

$$\mathcal{A} = \cup\{\mathcal{A}_\lambda; \lambda \in \Lambda\}; \text{ clearly, } \mathcal{A} \text{ is total too.}$$

By a preceding result, the formula

$$\mathcal{T} = \{D \in \text{exp}[X]; D = \text{union of intersections of finite subsets in } \mathcal{A}\}$$

defines a topology over  $X$ . This, as precise, is just the supremum topology of the family  $\{\mathcal{T}(d_\lambda); \lambda \in \Lambda\}$ . Note that, according to its definition, we have for each (nonempty) subset  $Y$  of  $X$  and each  $x \in X$ ,

- (g-int)  $x \in \text{int}(Y)$  iff there exists an intersection  $W$  of a finite subset in  $\mathcal{A}$  with  $x \in W \subseteq Y$
- (g-clo)  $x \in \text{cl}(Y)$  iff each intersection  $W$  of a finite subset in  $\mathcal{A}$  with  $x \in W$  fulfills  $W \cap Y \neq \emptyset$ .

- (II) The latter of these is represented by the (sequential) conv-Cauchy structure induced by the family  $D = (d_\lambda; \lambda \in \Lambda)$  of these semimetrics. This, essentially, may be described in the following way.

Take an arbitrary sequence  $(x_n; n \geq 0)$  in  $X$ . Given  $\lambda \in \Lambda$ , the  $d_\lambda$ -convergence of this sequence towards an  $x \in X$  [depicted as:  $x_n \xrightarrow{d_\lambda} x$ ], means:

$$d_\lambda(x_p, x) \rightarrow 0 \text{ as } p \rightarrow \infty$$

(i.e.:  $\forall \varepsilon > 0, \exists i = i(\lambda, \varepsilon)$ , such that  $i \leq p \implies d_\lambda(x_p, x) \leq \varepsilon$ ).

If this holds for all  $\lambda \in \Lambda$ , then  $(x_n; n \geq 0)$  is said to  $D$ -converge towards  $x$  [written as:  $x_n \xrightarrow{D} x$ ]. The set of all such points  $x$  will be denoted as  $D - \lim_n(x_n)$ ; when it is nonempty, then  $(x_n; n \geq 0)$  is called  $D$ -convergent. On the other hand, given  $\lambda \in \Lambda$ , the  $d_\lambda$ -Cauchy property of  $(x_n; n \geq 0)$  means:

$$d_\lambda(x_p, x_q) \rightarrow 0 \text{ as } p, q \rightarrow \infty, p \leq q$$

(i.e.:  $\forall \varepsilon > 0, \exists j := j(\lambda, \varepsilon)$ , such that  $j \leq p \leq q \implies d_\lambda(x_p, x_q) \leq \varepsilon$ ).

If this holds for each  $\lambda \in \Lambda$ , we say that  $(x_n; n \geq 0)$  is  $D$ -Cauchy. Note that any  $D$ -convergent sequence is  $D$ -Cauchy too; the reciprocal is not in general valid. By definition, the triple  $(X, \Lambda; D)$  endowed with this conv-Cauchy structure and the regularity condition

$D$  is *sufficient* :  $x, y \in X$  and  $(d_\lambda(x, y) = 0, \forall \lambda \in \Lambda)$  imply  $x = y$

is called a *gauge space*.

Suppose that we fixed such a structure. Define a  $D$ -closure operator  $A \mapsto \text{Dcl}(A)$  from  $\exp[X]$  to itself, as: for each  $A \in \exp[X]$ ,

$$y \in \text{Dcl}(A) \text{ iff } y = D - \lim_n(y_n), \text{ for some sequence } (y_n) \text{ in } A.$$

It is not hard to see that  $A \mapsto \text{Dcl}(A)$  is a *semi-closure* over  $X$ , in the sense

(Dcl-1) (identity)  $\emptyset = \text{Dcl}(\emptyset), X = \text{Dcl}(X)$ ,

(Dcl-2) (progressiveness)  $Y \subseteq \text{Dcl}(Y), \forall Y \in \exp[X]$ ,

(Dcl-3) (additivity)  $\text{Dcl}(U \cup V) = \text{Dcl}(U) \cup \text{Dcl}(V), \forall U, V \in \exp[X]$ .

Note that, as a direct consequence of these, we also have

(Dcl-4) (monotonicity)  $Y_1 \subseteq Y_2$  implies  $\text{Dcl}(Y_1) \subseteq \text{Dcl}(Y_2)$ .

Unfortunately,  $A \mapsto \text{Dcl}(A)$  is not involutive; i.e.,

(Dcl-inv)  $\text{Dcl}(\text{Dcl}(A)) = \text{Dcl}(A)$ , for each  $A \in \exp[X]$

is not in general true; so that,  $A \mapsto \text{Dcl}(A)$  is not a closure over  $X$  according to Kuratowski [38, Ch I, Sect 4].

Further, let us say that  $A \in \exp[X]$  is *D-closed*, provided  $A = \text{Dcl}(A)$ ; note that

the  $D$ -limit of each sequence in  $A$  belongs to  $A$ .

Having these precise, let  $(J, g)$  be a metric space; and  $(\leq)$  be a (partial) order over  $J$ ; the triple  $(J, g, \leq)$  will be then referred to as an *ordered metric space*. Let  $(<)$  stand for the *attached strict order*:

$$t < s \text{ iff } t \leq s \text{ and } t \neq s \text{ [clearly, } (<) \text{ is irreflexive and transitive].}$$

Call  $r \in J$ ,  $(\leq)$ -*maximal* provided

$$J(r, \leq) = \{r\}; \text{ or, equivalently, } J(r, <) = \emptyset;$$

the class of all these will be denoted as  $\max(J, \leq)$ . The negation of this property ( $J(r, <) \neq \emptyset$ ) will be referred to as  $r$  is  $(\leq)$ -*nonmaximal*; denote the class of all such points as  $\text{nmax}(J, \leq)$ . The following properties upon  $(J, g, \leq)$  will be admitted:

(Jgle-1)  $(J, g, \leq)$  is *conv-bd-regular*: each ascending sequence  $(t_n)$  in  $J$  is  $g$ -convergent (hence,  $g$ -Cauchy) in  $J$ , with  $t_n \leq \lim_n(t_n), \forall n$

(Jgle-2)  $(J, g, \leq)$  is *nonmaximal accessible*:

each  $r \in \text{nmax}(J, \leq)$  is the  $g$ -limit of some sequence  $(p_n)$  in  $J(r, <)$ .

Further, take a mapping  $(t, x) \mapsto H(t, x) = H(t)x$  from  $J \times X$  into  $X$ ; hence, (s-map)  $H(t)$  is a selfmap of  $X$ , for each  $t \in J$ .

Given the subset  $U \in \exp(X)$ , the family  $\Psi = (\psi_\lambda; \lambda \in \Lambda)$  over  $\mathcal{F}(R_+)$ , and the family  $\Omega = (\omega_\lambda; \lambda \in \Lambda)$  over  $\mathcal{F}(R_+)$ , let us say that  $H$  is a  $(\Psi, \Omega; M; U)$ -homotopic mapping, when the following conditions hold

- (hom-1)  $H$  is first variable  $\Psi$ -contractive:  
 $d_\lambda(H(t, x), H(s, x)) \leq \psi_\lambda(g(t, s))$ , for all  $t, s \in J$  and all  $x \in U$
- (hom-2)  $H$  is second variable  $\Omega$ -contractive:  
 $d_\lambda(H(t, x), H(t, y)) \leq \omega_\lambda(M(x, y; H(t); d_\lambda))$ ,  $\forall t \in J, \forall x, y \in U$ .

Here, letting  $d$  be any semimetric over  $X$  and  $T : X \rightarrow X$  be any selfmap of  $X$ , we introduced the mappings  $(M_1, M_2, M_3, M)$  as: for each  $x, y \in X$ ,

$$\begin{aligned}
 M_1(x, y; T; d) &= d(x, y), \\
 M_2(x, y; T; d) &= \max\{d(x, Tx), d(y, Ty)\}, \\
 M_3(x, y; T; d) &= (1/2)[d(x, Ty) + d(y, Tx)], \\
 M(x, y; T; d) &= \max\{M_1(x, y; T; d), M_2(x, y; T; d), M_3(x, y; T; d)\}.
 \end{aligned}$$

Finally, let the functional classes

$$\mathcal{F}_0(R_+), \mathcal{F}_0(re)(R_+), \mathcal{F}_0(in)(R_+), \mathcal{F}_0(re, in)(R_+)$$

be introduced as before. Given  $\omega \in \mathcal{F}_0(re, in)(R_+)$ , let  $(\eta(t) = \omega(2t); t \geq 0)$  stand for the *double function* attached to it; clearly,  $\eta \in \mathcal{F}_0(in)(R_+)$ . The conditions to be considered are essentially related to this last function; precisely

- (s-reg)  $\eta$  is strongly regressive ( $\eta(t + 0) < t, \forall t > 0$ );  
referred to as:  $\omega$  is *double strongly regressive*
- (c-co)  $\eta$  is complementary coercive [ $\eta^* = I - \eta$  is coercive];  
referred to as:  $\omega$  is *double complementary coercive*.

Denote for simplicity

$$J(H; U) = \{t \in J; \text{Fix}(H(t)) \cap U \neq \emptyset\}.$$

We say that the homotopic map  $H$  is *transversal*, provided

$$J(H; U) \neq \emptyset \text{ implies } J(H; U) \cap \max(J, \leq) \neq \emptyset;$$

or, in other words: the generic property (involving points of  $J$ )

$$P(t) : \text{Fix}(H(t)) \cap U \text{ is nonempty}$$

may jump from the points of  $J$  to the points of  $\max(J, \leq)$ .

The following homotopic fixed point theorem is available.



**Theorem 14** *Let the ordered metric space  $(J, g, \leq)$  with*

(61-i)  *$(J, g, \leq)$  is conv-bd-regular and nonmaximal accessible*

*and the subset  $U \in \exp(X)$  be such that the map  $(t, x) \mapsto H(t, x) = H(t)x$  in  $\mathcal{F}(J \times X, X)$  is  $(\Psi, \Omega; M; U)$ -homotopic, where the family  $\Psi = (\psi_\lambda; \lambda \in \Lambda)$  over  $\mathcal{F}(R_+)$ , and the family  $\Omega = (\omega_\lambda; \lambda \in \Lambda)$  over  $\mathcal{F}_0(re, in)(R_+)$ , fulfill*

(61-ii)  $(\forall \lambda \in \Lambda)$ :  $\psi_\lambda$  is zero-continuous ( $\psi_\lambda(t) \rightarrow 0 = \psi_\lambda(0)$  as  $t \rightarrow 0$ )

(61-iii)  $(\forall \lambda \in \Lambda)$ : the double function  $(\eta_\lambda(t) = \omega_\lambda(2t); t \geq 0)$  is strongly regressive and complementary coercive.

*In addition, suppose that*

(61-iv)  *$X$  is  $D$ -complete and  $U$  is  $D$ -closed*

(61-v)  $J(H; U \setminus \text{int}(U)) = \emptyset$ ; i.e.:  $\text{Fix}(H(t)) \cap (U \setminus \text{int}(U)) = \emptyset$ , for all  $t \in J$ .

*Then, the homotopic map  $H$  is transversal in  $(ZF\text{-}AC+DC)$ ; i.e.,*

(61-a)  $J(H; U) \neq \emptyset$  implies  $J(H; U) \cap \max(J, \leq) \neq \emptyset$ .

**Proof** There are several steps to be passed.

**Step 1.** By the Zorn-Bourbaki metric maximality principle (ZB-m-mp) (see below),  $\max(J, \leq)$  is nonempty. Moreover, by the posed hypothesis,  $J(H; U)$  is nonempty. We claim that

(61-b) the ordered metric subspace  $(J(H; U), g, \leq)$  is conv-bd-regular.

To this end, let  $(r_n)$  be an ascending sequence in  $J(H; U)$ ; hence

$$((r_n)=\text{ascending and}) \text{Fix}(H(r_n)) \cap U \text{ is nonempty, for each } n.$$

By the Denumerable Axiom of Choice (deductible, as above said, in  $(ZF\text{-}AC+DC)$ ), there exists a sequence  $(x_n)$  in  $U$ , with

$$(\forall n): x_n \in \text{Fix}(H(r_n)); \text{ that is: } x_n = H(r_n)x_n.$$

As  $(J, g, \leq)$  is conv-bd-regular, there exists  $r_\infty := \lim_n(r_n)$  in  $J$  with  $(r_n \leq r_\infty, \forall n)$ . We now claim that  $r_\infty \in J(H; U)$ ; and, from this, all is clear. Remember that, by a previous auxiliary fact, we have that, for each  $\lambda \in \Lambda$ ,

(p-1) right complementary inverse  $(\gamma_\lambda(s) = \sup\{t \in R_+; t \leq \omega_\lambda(s + t)\}, s \in R_+)$  belongs to  $\mathcal{F}_0(re, in)(R_+)$

(p-2) moreover,  $\gamma_\lambda$  is strongly regressive ( $\gamma_\lambda(s + 0) < s, \forall s > 0$ ) as well as complementary coercive ( $\gamma_\lambda^* := I - \gamma_\lambda$  is coercive).

For the arbitrary fixed  $\lambda \in \Lambda$  and each couple of ranks  $(n, m)$  with  $n \leq m$ , we have (by the contractive conditions)

$$(61-c) \begin{aligned} d_\lambda(x_n, x_m) &= d_\lambda(H(r_n)x_n, H(r_m)x_m) \leq \\ &d_\lambda(H(r_n)x_n, H(r_m)x_n) + d_\lambda(H(r_m)x_n, H(r_m)x_m) \leq \\ &\psi_\lambda(g(r_n, r_m)) + d_\lambda(H(r_m)x_n, H(r_m)x_m). \end{aligned}$$

Let us now evaluate the last expression. By definition (and preceding facts)

$$\begin{aligned}
 M_1(x_n, x_m; H(r_m); d_\lambda) &= d_\lambda(x_n, x_m), \\
 M_2(x_n, x_m; H(r_m); d_\lambda) &= \max\{d_\lambda(x_n, H(r_m)x_n), d(x_m, H(r_m)x_m)\} = \\
 d_\lambda(x_n, H(r_m)x_n) &= d_\lambda(H(r_n)x_n, H(r_m)x_n) \leq \psi_\lambda(g(r_n, r_m)), \\
 M_3(x_n, x_m; H(r_m); d_\lambda) &= (1/2)[d_\lambda(x_n, H(r_m)x_m) + d_\lambda(x_m, H(r_m)x_n)] = \\
 (1/2)[d_\lambda(x_n, x_m) + d_\lambda(H(r_m)x_m, H(r_m)x_n)].
 \end{aligned}$$

Denote for each couple  $(n, m)$ ,

$$\begin{aligned}
 a_\lambda(n, m) &= d_\lambda(x_n, x_m), \quad b_\lambda(n, m) = \psi_\lambda(g(r_n, r_m)), \\
 C_\lambda(n, m) &= d_\lambda(H(r_m)x_n, H(r_m)x_m)
 \end{aligned}$$

By the above evaluations,

$$M(x_n, x_m; H(r_m); d_\lambda) \leq a_\lambda(n, m) + b_\lambda(n, m) + C_\lambda(n, m)$$

wherefrom (as  $\omega_\lambda$  is increasing)

$$\begin{aligned}
 C_\lambda(n, m) &\leq \omega_\lambda(M(x_n, x_m; H(r_m); d_\lambda)) \leq \\
 \omega_\lambda(a_\lambda(n, m) + b_\lambda(n, m) + C_\lambda(n, m)), &\text{ for } n \leq m.
 \end{aligned}$$

Combining with the complementary coerciveness of  $\omega_\lambda$ , gives

$$C_\lambda(n, m) \leq \gamma_\lambda(a_\lambda(n, m) + b_\lambda(n, m));$$

and this, along with (61-c), implies

$$a_\lambda(n, m) \leq b_\lambda(n, m) + \gamma_\lambda(a_\lambda(n, m) + b_\lambda(n, m));$$

wherefrom (by a simple addition)

$$a_\lambda(n, m) + b_\lambda(n, m) \leq 2b_\lambda(n, m) + \gamma_\lambda(a_\lambda(n, m) + b_\lambda(n, m)).$$

As  $(r_n)$  is  $g$ -Cauchy and  $\psi_\lambda$  is zero-continuous, it is clear that  $b_\lambda(., .)$  (hence,  $2b_\lambda(., .)$  as well) appears as Cauchy. This, along with  $\gamma_\lambda$  being strongly regressive and complementary coercive tells us (by a preceding auxiliary fact) that  $a_\lambda(., .) + b_\lambda(., .)$  is Cauchy too; wherefrom

$$(\forall \lambda \in A): d_\lambda(x_n, x_m) = a_\lambda(n, m) \rightarrow 0 \text{ as } n, m \rightarrow \infty, n \leq m;$$

proving that  $(x_n)$  is  $D$ -Cauchy; so that (by the properties of  $(X, U)$ )

$$x_n \xrightarrow{D} x_\infty \text{ as } n \rightarrow \infty, \text{ for some } x_\infty \in U.$$

**Step 2.** For each  $\lambda \in \Lambda$  and each  $n \in N$ , we have (by the contractive conditions)

$$\begin{aligned}
 (61-d) \quad & d_\lambda(x_\infty, H(r_\infty)x_\infty) - d_\lambda(x_\infty, x_n) \leq \\
 & d_\lambda(x_n, H(r_\infty)x_\infty) = d_\lambda(H(r_n)x_n, H(r_\infty)x_\infty) \leq \\
 & d_\lambda(H(r_n)x_n, H(r_\infty)x_n) + d_\lambda(H(r_\infty)x_n, H(r_\infty)x_\infty) \leq \\
 & \psi_\lambda(g(r_n, r_\infty)) + d_\lambda(H(r_\infty)x_n, H(r_\infty)x_\infty).
 \end{aligned}$$

Let us now evaluate the last expression. By definition (and preceding facts)

$$\begin{aligned}
 M_1(x_n, x_\infty; H(r_\infty); d_\lambda) &= d_\lambda(x_n, x_\infty), \\
 M_2(x_n, x_\infty; H(r_\infty); d_\lambda) &= \max\{d_\lambda(x_n, H(r_\infty)x_n), d_\lambda(x_\infty, H(r_\infty)x_\infty)\} = \\
 &= \max\{d_\lambda(H(r_n)x_n, H(r_\infty)x_n), d_\lambda(x_\infty, H(r_\infty)x_\infty)\} \leq \\
 &= \max\{\psi_\lambda(g(r_n, r_\infty)), d_\lambda(x_\infty, H(r_\infty)x_\infty)\}, \\
 M_3(x_n, x_\infty; H(r_\infty); d_\lambda) &= (1/2)[d_\lambda(x_n, H(r_\infty)x_\infty) + d_\lambda(x_\infty, H(r_\infty)x_n)] \leq \\
 &= d_\lambda(x_n, x_\infty) + (1/2)[d_\lambda(x_\infty, H(r_\infty)x_\infty) + d_\lambda(H(r_n)x_n, H(r_\infty)x_n)] \leq \\
 &= d_\lambda(x_n, x_\infty) + (1/2)[d_\lambda(x_\infty, H(r_\infty)x_\infty) + \psi_\lambda(g(r_n, r_\infty))] \leq \\
 &= d_\lambda(x_n, x_\infty) + \max\{\psi_\lambda(g(r_n, r_\infty)), d_\lambda(x_\infty, H(r_\infty)x_\infty)\}.
 \end{aligned}$$

Denote, for  $n \in N$

$$\begin{aligned}
 E_\lambda &= d_\lambda(x_\infty, H(r_\infty)x_\infty), a_\lambda(n) = \psi_\lambda(g(r_n, r_\infty)), \\
 b_\lambda(n) &= d_\lambda(x_n, x_\infty), c_\lambda(n) = a_\lambda(n) + b_\lambda(n) + E_\lambda.
 \end{aligned}$$

By the above evaluations,

$$\begin{aligned}
 M(x_n, x_\infty; H(r_\infty); d_\lambda) &\leq b_\lambda(n) + \max\{a_\lambda(n), E_\lambda\}; \\
 \text{whence, } M(x_n, x_\infty; H(r_\infty); d_\lambda) &\leq c_\lambda(n);
 \end{aligned}$$

and this, along with  $\omega_\lambda$ =increasing, gives

$$d_\lambda(H(r_\infty)x_n, H(r_\infty)x_\infty) \leq \omega_\lambda(M(x_n, x_\infty; H(r_\infty); d_\lambda)) \leq \omega_\lambda(c_\lambda(n)).$$

Combining with (61-d), one gets

$$\begin{aligned}
 (61-e) \quad (\forall n): E_\lambda &\leq a_\lambda(n) + b_\lambda(n) + \omega_\lambda(c_\lambda(n)); \text{ wherefrom} \\
 c_\lambda(n) &\leq 2(a_\lambda(n) + b_\lambda(n)) + \omega_\lambda(c_\lambda(n)).
 \end{aligned}$$

As  $2(a_\lambda(n) + b_\lambda(n)) \rightarrow 0$  if  $n \rightarrow \infty$  and  $\omega_\lambda$  is strongly regressive and complementary coercive, one gets (by a previous auxiliary fact)

$$(\forall \lambda \in \Lambda) : c_\lambda(n) \rightarrow 0 \text{ as } n \rightarrow \infty; \text{ so that, } E_\lambda = 0.$$

This, by definition, yields

$$d_\lambda(x_\infty, H(r_\infty)x_\infty) = 0, \forall \lambda \in \Lambda; \text{ whence, } x_\infty = H(r_\infty)x_\infty;$$

telling us that  $r_\infty \in J(H; U)$ .

**Step 3.** Summing up, the Zorn-Bourbaki metric maximality principle (ZB-m-mp) applies to  $(J(H; U), g, \leq)$ . Hence, for the fixed  $q_0 \in J(H; U)$  (assured by hypothesis), there exists an ascending sequence  $(q_n)$  in  $J(H; U)$ , with

$$q^* := \lim_n(q_n) \text{ exists in } J(H; U) \text{ and } q^* \in \max(J(H; U), \leq);$$

$$\text{that is: } q^* \leq r \in J(H; U) \text{ implies } q^* = r.$$

Note that, as  $q^* \in J(H; U)$ , one has

$$y^* = H(q^*)y^*, \text{ for some } y^* \in U.$$

We show that, necessarily,  $q^* \in \max(J, \leq)$ . Assume not:  $J(q^*, <)$  is not empty. By the imposed hypotheses

$$\text{there exists a sequence } (p_m; m \geq 0) \text{ in } J(q^*, <) \text{ with } \lim_m(p_m) = q^*.$$

On the other hand (again by hypothesis)

$$y^* \in \text{Fix}(H(q^*)) \cap U \text{ implies } y^* \notin U \setminus \text{int}(U); \text{ so that, } y^* \in \text{int}(U);$$

where  $\text{int}(\cdot)$  means: the interior with respect to the supremum topology  $\mathcal{T}$  (see above). From the definition of underlying topology, there exists a finite part  $\Lambda^*$  of  $\Lambda$  and a number  $s = s(\Lambda^*)$  in  $\mathcal{R}_+^0$ , such that

$$V := \cap\{X[y^*, s](d_\lambda); \lambda \in \Lambda^*\} \subseteq U.$$

We now claim that there exists some index  $i = i(V^*; s)$ , such that

$$y_m := H(p_m)y \in V, \text{ for all } m \geq i, y \in V.$$

In fact, for each  $\lambda \in \Lambda^*, m \in N, y \in V$  we have (in view of  $\omega_\lambda$ =increasing)

$$\begin{aligned} (61\text{-f}) \quad d_\lambda(y_m, y^*) &= d_\lambda(H(p_m)y, H(q^*)y^*) \leq \\ &d_\lambda(H(p_m)y, H(q^*)y) + d_\lambda(H(q^*)y, H(q^*)y^*) \leq \\ &\psi_\lambda(g(p_m, q^*)) + d_\lambda(H(q^*)y, H(q^*)y^*). \end{aligned}$$

Let us now evaluate the last expression. By definition (and preceding facts)

$$\begin{aligned}
 M_1(y, y^*; H(q^*); d_\lambda) &= d_\lambda(y, y^*) \leq s \\
 M_2(y, y^*; H(q^*); d_\lambda) &= \max\{d_\lambda(y, H(q^*)y), d_\lambda(y^*, H(q^*)y^*)\} = \\
 &= d_\lambda(y, H(q^*)y) \leq d_\lambda(y, y^*) + d_\lambda(H(q^*)y, H(q^*)y^*) \leq \\
 &= s + d_\lambda(H(q^*)y, H(q^*)y^*), \\
 M_3(y, y^*; H(q^*); d_\lambda) &= (1/2)[d_\lambda(y, H(q^*)y^*) + d_\lambda(y^*, H(q^*)y)] = \\
 &= (1/2)[d_\lambda(y, y^*) + d_\lambda(H(q^*)y, H(q^*)y^*)] \leq (1/2)[s + d_\lambda(H(q^*)y, H(q^*)y^*)].
 \end{aligned}$$

This yields directly

$$M(y, y^*; H(q^*); d_\lambda) \leq s + d_\lambda(H(q^*)y, H(q^*)y^*);$$

wherefrom (as  $\omega_\lambda$  is increasing)

$$\begin{aligned}
 d_\lambda(H(q^*)y, H(q^*)y^*) &\leq \omega_\lambda(M(y, y^*; H(q^*); d_\lambda)) \leq \\
 &\leq \omega_\lambda(s + d_\lambda(H(q^*)y, H(q^*)y^*)).
 \end{aligned}$$

As  $\omega_\lambda$  is strongly regressive and complementary coercive, one gets (see above)

$$d_\lambda(H(q^*)y, H(q^*)y^*) \leq \gamma_\lambda(s);$$

and this, along with (61-f), yields

$$d_\lambda(y_m, y^*) \leq \psi_\lambda(g(p_m, q^*)) + \gamma_\lambda(s).$$

In view of

$$\gamma_\lambda(s) < s, \lim_m \psi_\lambda(g(p_m, q^*)) = 0, \forall \lambda \in \Lambda^*,$$

there must be some index  $i = i(V^*; s)$ , such that

$$\psi_\lambda(g(p_m, q^*)) + \gamma_\lambda(s) < s, \forall m \geq i, \forall \lambda \in \Lambda^*.$$

But then, from the preceding evaluation,

$$y_m := H(p_m)y \in V, \text{ for all } m \geq i, y \in V,$$

and our assertion follows.

By the obtained facts,  $H(p_i)$  is a  $\Omega$ -contractive selfmap of the  $D$ -closed subset  $V$  of  $U$ , with (cf. the hypotheses)

$(\forall \lambda \in \Lambda) : \omega_\lambda$  is strongly regressive; hence, Boyd-Wong admissible.

Combining with our basic fixed point result it follows that, for the starting point  $y^* \in V$ , there exists another point  $z^* \in V$  with

$$z^* = \lim_n H(p_i)^n y^*, \text{ and } z^* = H(p_i)z^* \text{ (whence, } p_i \in J(H; U)).$$

This, in view of  $q^* < p_i$ , contradicts the maximality of  $q^*$  in  $J(H; U)$ . Hence, necessarily,  $q^* \in \max(J, \leq)$ ; and the conclusion follows.

A basic particular case of these developments corresponds to

$$(\forall \lambda \in \Lambda) : \omega_\lambda(t) = k_\lambda t, t \in R_+, \text{ for some } k_\lambda \in R_+.$$

Precisely, let  $(J, g, \leq)$  be an ordered metric space; and take a mapping  $(t, x) \mapsto H(t, x) = H(t)x$  from  $J \times X$  into  $X$ . Given the (nonempty) subset  $U \in \exp(X)$ , the family  $\Psi = (\psi_\lambda; \lambda \in \Lambda)$  over  $\mathcal{F}(R_+)$ , and the family  $k = (k_\lambda; \lambda \in \Lambda)$  over  $R_+$  suppose that the following conditions hold

(homm-1)  $H$  is first variable  $\Psi$ -contractive:

$$d_\lambda(H(t, x), H(s, x)) \leq \psi_\lambda(g(t, s)), \text{ for all } t, s \in J \text{ and all } x \in U$$

(homm-2)  $H$  is second variable  $k$ -contractive:

$$d_\lambda(H(t, x), H(t, y)) \leq k_\lambda M(x, y; H(t); d_\lambda), \text{ for all } t \in J \text{ and all } x, y \in U.$$

We then say that  $H$  is a  $(\Psi, k; M; U)$ -homotopic mapping.

The following (linear) homotopic fixed point theorem is available.

**Theorem 15** *Let the ordered metric space  $(J, g, \leq)$  with*

(62-i)  *$(J, g, \leq)$  is conv-bd-regular and nonmaximal accessible*

*and the subset  $U \in \exp(X)$  be such that the map  $(t, x) \mapsto H(t, x) = H(t)x$  from  $J \times X$  into  $X$  be  $(\Psi, k; M; U)$ -homotopic, where the family  $\Psi = (\psi_\lambda; \lambda \in \Lambda)$  over  $\mathcal{F}(R_+)$  and the family  $k = (k_\lambda; \lambda \in \Lambda)$  over  $R_+$  are such that*

(62-ii)  $(\forall \lambda \in \Lambda)$ :  $\psi_\lambda$  is zero-continuous ( $\psi_\lambda(t) \rightarrow 0 = \psi(0)$  as  $t \rightarrow 0$ )

(62-iii)  $(\forall \lambda \in \Lambda)$ :  $k_\lambda$  is double subunitary ( $0 \leq k_\lambda < 1/2$ ).

*In addition, suppose that*

(62-iv)  *$X$  is  $D$ -complete and  $U$  is  $D$ -closed*

(62-v)  $J(H; U \setminus \text{int}(U)) = \emptyset$ ; i.e.:  $\text{Fix}(H(t)) \cap (U \setminus \text{int}(U)) = \emptyset$ , for all  $t \in J$ .

*Then, the homotopic map  $H$  is transversal in  $(ZF-AC+DC)$ ; i.e.,*

$$J(H; U) \neq \emptyset \text{ implies } J(H; U) \cap \max(J, \leq) \neq \emptyset.$$

In particular, when the data  $(X, U)$  and  $(J, g, \leq)$  are taken as

(p-1)  $X$  is  $D$ -complete and  $U$  is  $\mathcal{F}$ -closed

(p-2)  $J = [0, 1]$ ,  $(g, \leq)$ =the usual metric and ordering in  $R$

the obtained result is comparable with a similar one due to Ariza-Ruiz and Jimenez-Melado [5, 6]; which, in turn, extends a related one in Frigon [23]. Concerning this aspect, two remarks are in order.

(I) The first variable contraction upon  $H$  considered in the quoted papers is

$$(\forall \lambda \in \Lambda): d_\lambda(H(t, x), H(s, x)) \leq |\varphi_\lambda(t) - \varphi_\lambda(s)|,$$

for all  $t, s \in J$  and all  $x \in U$ , where  $\varphi_\lambda : J \rightarrow R$  is continuous.

But, in view of any continuous function on the compact  $J = [0, 1]$  being uniformly continuous, we must have,  $\forall \lambda \in \Lambda$ ,

$$|\varphi_\lambda(t) - \varphi_\lambda(s)| \leq \psi_\lambda(|t - s|), t, s \in J,$$

where  $\psi_\lambda \in \mathcal{F}(R_+)$  is zero-continuous;

and, from this, we arrive at the posed first variable condition above.

(II) On the other hand, the linear homotopic result above holds over the reduced system (ZF-AC+DC) by working with metrical type intervals; while, the quoted result is holding over the complete system (ZF) and having as essential tool a connectedness characterization of (nonempty) real intervals.

As a consequence of this, the homotopic results above include the precise one. Note that, the metrical and order character of the interval space  $(J, g, \leq)$  may be put in an extended context, so as to include a related statement in O'Regan and Precup [46]. Finally, one may ask whether such results are applicable to practical Ulam-Hyers-Rassias stability for nonlinear equations over Banach spaces, under the lines in Wang and Fečkan [68]; we conjecture that a positive answer to this is available.

## 7 Kang-Park Principles

Let  $X$  be a nonempty set. By a *pseudometric* over  $X$  we mean any map  $d : X \times X \rightarrow R_+$ . Fix such a map, endowed with

(ref)  $d$  is *reflexive*:  $d(x, x) = 0, \forall x \in X$ ;

we then say that it is a *r-pseudometric* on  $X$ . Further, let  $(\preceq)$  be a *quasi-order* (i.e.: reflexive and transitive relation) over  $X$ ; the triple  $(X, d, \preceq)$  will be referred to as a *quasi-ordered r-pseudometric space*. Given  $M \in \exp(X)$ , call  $z \in X$ ,

(max-1)  $(d, \preceq)$ -*maximal* over  $M$ , if  $(u, v \in M, z \preceq u \leq v) \implies d(u, v) = 0$

(max-2)  $(\preceq)$ -*maximal* over  $M$ , if  $(u \in M, z \preceq u) \implies d(z, u) = 0$ .

It is our aim in the following to give sufficient conditions for such properties. Then, an application is given to Ekeland variational principles.

(A) Call the sequence  $(x_n)$  in  $X$ ,  $d$ -Cauchy when  $d(x_m, x_n) \rightarrow 0$  as  $m, n \rightarrow \infty$ ,  $m \leq n$ ; that is

$$\forall \varepsilon > 0, \exists n(\varepsilon), \text{ such that } n(\varepsilon) \leq p \leq q \implies d(x_p, x_q) \leq \varepsilon;$$

or, equivalently (by the reflexive property)

$$\forall \varepsilon > 0, \exists n(\varepsilon), \text{ such that } n(\varepsilon) < p < q \implies d(x_p, x_q) \leq \varepsilon;$$

and  $d$ -asymptotic, if  $\lim_n d(x_n, x_{n+1}) = 0$ ; that is

$$\forall \varepsilon > 0, \exists n(\varepsilon), \text{ such that } n(\varepsilon) \leq p \implies d(x_p, x_{p+1}) \leq \varepsilon.$$

Then, let us consider the global conditions

(C-reg)  $(M, d, \preceq)$  is *Cauchy regular*:

each  $(\preceq)$ -ascending sequence in  $M$  is  $d$ -Cauchy

(A-reg)  $(M, d, \preceq)$  is *asymptotic regular*:

each  $(\preceq)$ -ascending sequence in  $M$  is  $d$ -asymptotic.

As each  $d$ -Cauchy sequence is  $d$ -asymptotic too, it follows that (C-reg)  $\implies$  (A-reg).

The reverse implication also holds, in the sense

**Proposition 12** *We have, in (ZF-AC),*

$$(A\text{-reg}) \implies (C\text{-reg}); \text{ whence, } (A\text{-reg}) \iff (C\text{-reg}).$$

**Proof** Suppose that (A-reg) holds; but some  $(\preceq)$ -ascending  $(x_n)$  is not entitled with the  $d$ -Cauchy property; i.e. (for some  $\varepsilon > 0$ )

$$C(n) = \{(p, q) \in N \times N; n < p < q, d(x_p, x_q) > \varepsilon\} \neq \emptyset, \forall n.$$

Denote, for simplicity

$$p(n) = \min \text{Dom}(C(n)), q(n) = \max(C(n)(p(n)), n \in N;$$

clearly, no choice techniques are used in this construction. Fix some rank  $i(0)$ . By this assumption, there exist  $i(1) = p(i(0)), i(2) = q(i(0))$  with  $i(0) < i(1) < i(2)$ ,  $d(x_{i(1)}, x_{i(2)}) > \varepsilon$ . Further, given the rank  $i(2)$ , there exist  $i(3) = p(i(2)), i(4) = q(i(2))$  with  $i(2) < i(3) < i(4)$ ,  $d(x_{i(3)}, x_{i(4)}) > \varepsilon$ , and so on. By induction, we get a  $(\preceq)$ -ascending subsequence  $(y_n = x_{i(n)})$  of  $(x_n)$  with  $d(y_{2n+1}, y_{2n+2}) > \varepsilon$ , for all  $n$ . This contradicts (A-reg); hence the claim.

By definition, either of these conditions will be referred to as:  $(M, d, \preceq)$  is *regular*. A basic consequence of this property is the following.



**Proposition 13** *We have, in (ZF-AC+DC)*

$(M, d, \preceq)$  *is regular implies*

$(M, d, \preceq)$  *is weakly regular:  $\forall x \in M, \forall \varepsilon > 0, \exists y = y(x, \varepsilon) \in M(x, \preceq)$ ,*

*such that:  $(u, v \in M, y \preceq u \preceq v) \implies d(u, v) \leq \varepsilon$ .*

**Proof** Assume this would be false; that is (for some  $x \in M, \varepsilon > 0$ )

for each  $y \in M(x, \preceq)$  there exist  $u, v \in M$  with  $y \preceq u \preceq v, d(u, v) > \varepsilon$ .

This, by definition, yields (for the same  $(x, \varepsilon)$ ):

$$\forall y \in M(x, \preceq), \exists (u, v) \in (\preceq) : y \preceq u, d(u, v) > \varepsilon;$$

where  $(\preceq) := \{(a, b) \in M \times M; a \preceq b\}$ . Put  $Q := \{(a, b) \in (\preceq); x \preceq a\}$ ; and fix a couple  $(y_0, y_1) \in Q$ ; for example,  $y_0 = y_1 = x$ . Define a relation  $\mathcal{R} = \mathcal{R}(\varepsilon)$  on  $Q$  as

$$(a_1, b_1) \mathcal{R} (a_2, b_2) \text{ if and only if } b_1 \preceq a_2, d(a_2, b_2) > \varepsilon.$$

From the imposed condition,  $Q((a, b), \mathcal{R}) \neq \emptyset, \forall (a, b) \in Q$ . So, by (DC), it follows that, for the starting point  $w_0 = (y_0, y_1)$  in  $Q$  there exists a sequence  $(w_n := (y_{2n}, y_{2n+1}); n \geq 0)$  in  $Q$  with

$$w_n \mathcal{R} w_{n+1}, \text{ for all } n; \text{ hence, by definition,}$$

$$y_{2n+1} \preceq y_{2n+2}, d(y_{2n+2}, y_{2n+3}) > \varepsilon, \text{ for all } n.$$

As a consequence,  $(y_n; n \geq 0)$  is  $(\preceq)$ -ascending and not  $d$ -asymptotic; in contradiction with the regularity of  $(M, d, \preceq)$ ; hence the claim.

Starting from the quasi-order  $(\preceq)$ , let us introduce a *quasi-order convergence*  $\mathcal{B} = \mathcal{B}(\preceq)$  over  $X$  as follows: given the sequence  $(x_n)$  in  $X$  and the point  $x \in X$ ,

$$x_n \xrightarrow{\mathcal{B}} x \text{ if } \exists m(x) \in \mathbb{N}, \text{ such that } n \geq m(x) \text{ implies } x_n \preceq x;$$

also referred to as:  $x$  is a  $\mathcal{B}$ -limit of  $(x_n)$ . The class of all these will be denoted as  $\mathcal{B} - \lim_n(x_n)$ ; when it is nonempty, we say that  $(x_n)$  is  $\mathcal{B}$ -convergent. Given  $Y \in \text{exp}[X]$ , let us say that  $w \in X$  is a  $(\mathcal{B}, \preceq)$ -adherence point of it when

$$w \in \mathcal{B} - \lim_n(z_n), \text{ for some } (\preceq)\text{-ascending sequence } (z_n) \text{ of } Y;$$

the class of all these will be denoted as  $\text{ocl}(Y)$ . It is not hard to see that  $Y \mapsto \text{ocl}(Y)$  is a semi-closure over  $X$ , in the sense

- (ocl-1)  $\text{ocl}(\emptyset) = \emptyset, \text{ocl}(X) = X,$
- (ocl-2)  $\text{ocl}(U \cup V) = \text{ocl}(U) \cup \text{ocl}(V), \forall U, V \in \exp[X]$
- (ocl-3)  $Y \subseteq \text{ocl}(Y), \forall Y \in \exp[X].$

Unfortunately,  $Y \mapsto \text{ocl}(Y)$  is not involutive; i.e.,

$$\text{(ocl-inv)} \quad \text{ocl}(\text{ocl}(Y)) = \text{ocl}(Y), \text{ for each } Y \in \exp[X]$$

is not in general true; so that,  $Y \mapsto \text{ocl}(Y)$  is not a closure over  $X$  according to Kuratowski [38, Ch I, Sect 4]. In this context, let us say that  $Y \in \exp[X]$  is *o-closed*, provided  $Y = \text{ocl}(Y)$ ; note that, in this case,

the  $\mathcal{B}$ -limit of each  $(\preceq)$ -ascending sequence in  $Y$  is included in  $Y$ .

Putting these together, the following maximal principle (referred to as: Kang-Park pseudometric maximal principle; in short: (KP-p-mp)) is available.

**Theorem 16** *Let the quasi-ordered r-pseudometric space  $(X, d, \preceq)$  and the subset  $M \in \exp(X)$  be such that*

- (71-i)  $(M, d, \preceq)$  is regular: each  $(\preceq)$ -ascending sequence in  $M$  is  $d$ -Cauchy
- (71-ii)  $(M, d, \preceq)$  is  $\mathcal{B}$ -complete:  
each ascending  $d$ -Cauchy sequence in  $M$  is  $\mathcal{B}$ -convergent (in  $X$ ).

Then, the conclusion below holds in (ZF-AC+DC)

$$\forall u \in M, \exists v \in \text{ocl}(M), \text{ with } u \preceq v \text{ and } v \text{ is } (d, \preceq)\text{-maximal over } M.$$

**Proof** By an auxiliary fact above,  $(M, d, \preceq)$  is weakly regular. Hence, given  $u \in M$ , one may construct a  $(\preceq)$ -ascending sequence  $(u_n)$  in  $M$  with

$$u \preceq u_0, \text{ and } [(\forall n), (\forall y, z \in M)]: u_n \preceq y \preceq z \implies d(y, z) \leq 2^{-n}.$$

In particular, this tells us that  $(u_n)$  is  $d$ -Cauchy (in  $M$ ). Combining with the completeness hypothesis, there exists  $v \in \text{ocl}(M)$  such that

$$u_n \xrightarrow{\mathcal{B}} v; \text{ whence, } u_n \preceq v, \forall n.$$

This, via inclusion above, tells us that  $v$  is  $(d, \preceq)$ -maximal over  $M$ .

Formally, this result is comparable with the 1990 one in Kang and Park [35]. However, its basic lines were already set up in the papers by Turinici [58, 62]. Note that, both these results extend the 1976 Brezis-Browder ordering principle [12]. Further aspects may be found in Altman [3].

Technically speaking, the completeness hypothesis above is not very appropriate for the local version of Kang-Park pseudometric maximal principle (KP-p-mp); because, for many subsets  $M \in \exp(X)$ , the points in  $\text{ocl}(M)$  may be pretty far from the ones in  $M$ . To remove this inconvenient, we need some preliminaries.

Define a  $d$ -convergence structure on  $X$  under the precise way. Given  $Y \in \exp[X]$ , let us say that  $w \in X$  is a  $(d, \preceq)$ -adherence point of it when

$$w \in d - \lim_n(z_n), \text{ for some } (\preceq)\text{-ascending sequence } (z_n) \text{ of } Y;$$

the class of all these will be denoted as  $\text{docl}(Y)$ . It is not hard to see that  $Y \mapsto \text{docl}(Y)$  is a semi-closure over  $X$ , in the sense

- (docl-1)  $\text{docl}(\emptyset) = \emptyset, \text{docl}(X) = X,$
- (docl-2)  $\text{docl}(U \cup V) = \text{docl}(U) \cup \text{docl}(V), U, V \in \exp[X]$
- (docl-3)  $Y \subseteq \text{docl}(Y), \forall Y \in \exp[X].$

Unfortunately,  $Y \mapsto \text{docl}(Y)$  is not involutive; i.e.,

$$(\text{docl-inv}) \text{docl}(\text{docl}(Y)) = \text{docl}(Y), \text{ for each } Y \in \exp[X]$$

is not in general true; so that,  $Y \mapsto \text{docl}(Y)$  is not a closure over  $X$  according to Kuratowski [38, Ch I, Sect 4]. Further, let us say that  $Y \in \exp[X]$  is *do-closed*, provided  $Y = \text{docl}(Y)$ ; note that, in this case,

the  $d$ -limit of each  $(\preceq)$ -ascending sequence in  $Y$  is included in  $Y$ .

The following version of the statement above (referred to as: Kang-Park strong pseudometric maximal principle; in short: (KP-sp-mp)) is now available.

**Theorem 17** *Let the quasi-ordered  $r$ -pseudometric space  $(X, d, \preceq)$  and the subset  $M \in \exp(X)$  be such that*

- (72-i)  $(M, d, \preceq)$  is regular: each  $(\preceq)$ -ascending sequence in  $M$  is  $d$ -Cauchy
- (72-ii)  $(M, d, \preceq)$  is  $(d, \mathcal{B})$ -complete: each ascending  $d$ -Cauchy sequence  $(x_n)$  in  $M$  is  $d$ -convergent and  $\mathcal{B}$ -convergent (in  $X$ ) with (in addition)  $(d - \lim_n(x_n)) \cap (\mathcal{B} - \lim_n(x_n)) \neq \emptyset.$

Then, the conclusion below holds in (ZF-AC+DC)

$$\forall u \in M, \exists v \in \text{docl}(M), \text{ with } u \preceq v \text{ and } v \text{ is } (d, \preceq)\text{-maximal over } M.$$

A basic particular case of these corresponds to  $d$ -metric on  $X$ .

(I) As a first consequence of this choice, we have

$$(\forall M \in \exp(X), \forall z \in M) : z \text{ is } (d, \preceq)\text{-maximal over } M \text{ iff } z \text{ is } (\preceq)\text{-maximal over } M [z \preceq w \in M \text{ implies } z = w].$$

In fact, the left to right inclusion follows via  $z \preceq z \preceq w$  and  $z, w \in M$  imply  $z = w$ . On the other hand, the right to left inclusion is immediate, by the very definition of  $(\preceq)$ -maximal element, combined with

$$z \preceq u \preceq v \implies (z \preceq u \text{ and } z \preceq v).$$

(II) A second consequence of the same choice is contained in

**Proposition 14** *Suppose that  $(X, d, \preceq)$  is regular. Then, necessarily,*

$(\preceq)$  *is antisymmetric; hence, an ordering.*

**Proof** (cf. Hamel [30, Ch 4, Sect 4.1]) Let  $u, v \in X$  be such that  $u \preceq v$  and  $v \preceq u$ . The sequence  $(y_{2n} = u, y_{2n+1} = v; n \geq 0)$  is ascending; hence,  $d$ -asymptotic by hypothesis. This yields  $d(u, v) = 0$ ; whence (as  $d$ -metric),  $u = v$ .

(III) For the last consequence of our choice let us introduce the usual  $d$ -closure operator  $Y \mapsto \text{dcl}(Y)$ , as:

$$z \in \text{dcl}(Y) \text{ iff } z = \lim_n(y_n), \text{ for some sequence } (y_n) \text{ in } Y.$$

Note that, an equivalent way of describing this operator is

$$z \in \text{dcl}(Y) \text{ iff } \forall \varepsilon > 0, \exists y = y(\varepsilon) \in Y: d(z, y) < \varepsilon.$$

Starting from this, it is not hard to see that  $Y \mapsto \text{dcl}(Y)$  has the properties

- (dcl-1)  $\text{dcl}(\emptyset) = \emptyset, \text{dcl}(X) = X,$
- (dcl-2)  $\text{dcl}(U \cup V) = \text{dcl}(U) \cup \text{dcl}(V), U, V \in \text{exp}[X]$
- (dcl-3)  $Y \subseteq \text{dcl}(Y), \forall Y \in \text{exp}[X]$
- (dcl-4)  $\text{dcl}(\text{dcl}(Y)) = \text{dcl}(Y),$  for each  $Y \in \text{exp}[X];$

so,  $Y \mapsto \text{dcl}(Y)$  is a closure over  $X$  according to Kuratowski [38, Ch I, Sect 4]. In this case,  $Y \in \text{exp}[X]$  is called  $d$ -closed, provided  $Y = \text{dcl}(Y)$ .

Putting these together, gives the following maximal statement (referred to as: Zorn-Bourbaki metric maximal principle; in short: (ZB-m-mp)).

**Theorem 18** *Let the ordered metric space  $(X, d, \preceq)$  and the subsets  $M, Y \in \text{exp}(X)$  be such that*

- (73-i)  $M \subseteq Y$  and  $Y$  is  $d$ -closed ( $Y = \text{dcl}(Y)$ )
- (73-ii)  $(Y, d, \preceq)$  is  $(\text{conv}, \mathcal{B})$ -regular: each  $(\preceq)$ -ascending sequence  $(x_n)$  in  $Y$  is  $d$ -convergent in  $X$  and  $x_n \preceq \lim_n(x_n), \forall n.$

*Then, the conclusion below holds in (ZF-AC+DC)*

*for each  $u \in M$  there exists some other point  $v \in Y$ , with  $u \preceq v$  and  $v$  is  $(\preceq)$ -maximal over  $M: v \preceq y \in M$  imply  $v = y.$*

**Proof** By the posed hypothesis,

(73-a)  $(Y, d, \preceq)$  is regular: each  $(\preceq)$ -ascending sequence in  $Y$  is  $d$ -Cauchy.

Moreover, under the same condition, it is clear that

(73-b)  $(Y, d, \preceq)$  is  $(d, \mathcal{B})$ -complete: each  $(\preceq)$ -ascending  $d$ -Cauchy sequence  $(x_n)$  in  $Y$  is  $d$ -convergent in  $X$  and  $\lim_n(x_n) \in \mathcal{B} - \lim_n(x_n).$

Putting these together, we have that the Kang-Park strong pseudometric maximal principle (KP-sp-mp) is applicable here. Hence, for the starting  $u \in M \subseteq Y$  there exists some other point  $v \in \text{docl}(Y)$  with

$$u \preceq v \text{ and } v \text{ is } (d, \preceq)\text{-maximal over } Y.$$

But, evidently,  $\text{docl}(Y) \subseteq Y$  (as  $Y=d$ -closed); and this yields

$$v \in Y; \text{ hence (see above) } v \text{ is } (\preceq)\text{-maximal over } Y.$$

As a direct consequence of definition, we get

$$v \preceq y \in M \implies v \preceq y \in Y; \text{ whence } v = y;$$

and the proof is complete.

This result may be viewed as a metrical variant of the Zorn-Bourbaki maximal principle [9]. An early version of it was formulated in Turinici [57]; note that it includes the one due to Dancs et al. [18]. Further aspects may be found in the related papers by Turinici [56, 59].

(B) A basic application of these facts is to (local) variational principles. Let  $X$  be a nonempty set; and  $d(., .)$  be a metric on  $X$ ; the couple  $(X, d)$  will be referred to as a *metric space*. Further, let the function  $\varphi : X \rightarrow R \cup \{\infty\}$  be such that the following admissible conditions hold:

- (adm-1)  $\varphi$  is inf-proper ( $\text{Dom}(\varphi) \neq \emptyset$  and  $\inf[\varphi(X)] > -\infty$ )
- (adm-2)  $\varphi$  is  $d$ -lsc:  $[\varphi \leq t] := \{x \in X; \varphi(x) \leq t\}$  is  $d$ -closed,  $\forall t \in R$ .

The following variational statement (referred to as: Ekeland variational principle on metric spaces; in short: (EVP-m)) is available.

**Theorem 19** *Let the metric space  $(X, d)$  and the function  $\varphi : X \rightarrow R \cup \{\infty\}$  be taken according to*

$$(X, d) \text{ is complete and } \varphi \text{ is admissible.}$$

*Further, take some  $u \in \text{Dom}(\varphi)$ . Then, in the reduced system (ZF-AC+DC), there exists  $v \in \text{Dom}(\varphi)$ , with*

- (74-a)  $d(u, v) \leq \varphi(u) - \varphi(v)$  (hence  $\varphi(u) \geq \varphi(v)$ )
- (74-b)  $d(v, x) > \varphi(v) - \varphi(x)$ , for each  $x \in X \setminus \{v\}$ .

**Proof** Denote  $X[u] = \{x \in X; \varphi(u) \geq \varphi(x)\}$ ; note that, by the posed conditions,  $X[u]$  is a nonempty closed subset of  $X$ , with  $X[u] \subseteq \text{Dom}(\varphi)$ . Let  $(\preceq)$  be the relation on  $X$  defined as

$$x \preceq y \text{ iff } d(x, y) + \varphi(y) \leq \varphi(x).$$

Clearly,  $(\preceq)$  appears as a quasi-order on  $X$ ; moreover,  $(\preceq)$  is antisymmetric (hence an order) on  $\text{Dom}(\varphi)$ ; hence, on  $X[u]$  as well. Let also  $\mathcal{B}$  stand for the order convergence structure (over  $X$ ). We claim that the Zorn-Bourbaki metric maximal principle (ZB-m-mp) is applicable over the ordered metric space  $(X[u], d, \preceq)$ , under the choice  $M = Y = X[u]$ ; and this will complete the argument. In fact, let  $(x_n)$  be a  $(\preceq)$ -ascending sequence in  $X[u]$ :

$$(74\text{-c}) \quad d(x_n, x_m) \leq \varphi(x_n) - \varphi(x_m), \text{ if } n \leq m.$$

The sequence  $(\varphi(x_n))$  is descending bounded from below; hence a Cauchy one. This, along with the preceding relation, shows that  $(x_n)$  is a  $(\preceq)$ -ascending  $d$ -Cauchy sequence in  $X[u]$ . By the completeness hypothesis (and closeness of  $X[u]$ ) it follows that  $x := \lim_n(x_n)$  exists in  $X[u]$ ; moreover,

$$\varphi(x_n) \geq \varphi(x), \text{ for all } n \text{ [as } \varphi \text{ is } d\text{-lsc]}.$$

Finally, take some rank  $n$ . From (74-c) (and the triangular property)

$$d(x_n, x) \leq d(x_n, x_m) + d(x_m, x) \leq \varphi(x_n) - \varphi(x) + d(x_m, x), \quad \forall m \geq n;$$

wherefrom (passing to limit as  $m \rightarrow \infty$ )

$$d(x_n, x) \leq \varphi(x_n) - \varphi(x) \text{ (i.e.: } x_n \preceq x), \quad \forall n;$$

and our claim follows. Summing up, the Zorn-Bourbaki metric maximal principle (ZB-m-mp) is indeed applicable to our data. According to the quoted result, we have that, for the starting  $u \in X[u]$ , there exists another point  $v \in X[u]$ , with

$$(74\text{-d}) \quad u \preceq v \text{ and } v \text{ is } (\preceq)\text{-maximal on } X[u]: v \preceq w \in X[u] \text{ implies } v = w.$$

The first half of this gives our first conclusion. Concerning our second conclusion, suppose by absurd that

$$d(v, y) \leq \varphi(v) - \varphi(y) \text{ (i.e.: } v \preceq y), \text{ for some } y \in X, y \neq v.$$

As  $u \preceq v$ , one gets  $u \preceq y$ ; hence,  $y \in X[u]$ . This, under the second half of (74-d), yields  $v \preceq y \in X[u]$ ; hence,  $v = y$ ; a contradiction; so that, we are done.

*Remark 4* A local version of this result is immediately obtainable from the Zorn-Bourbaki metric maximal principle (ZB-m-mp) by simply taking the couple  $(M, Y)$  appearing there as  $(M_u, Y_u)$ , where

$$M_u := M \cap X[u], Y_u := \text{dcl}(M) \cap X[u];$$

here  $M \in \text{exp}(X)$  and  $u \in M \cap \text{Dom}(\varphi)$  are fixed elements.

This principle, due to Ekeland [21], found some basic applications to control and optimization, generalized differential calculus, critical point theory and global

analysis; we refer to the quoted paper for a survey of these. So, it cannot be surprising that, soon after its formulation, many extensions of (EVP-m) were proposed. For example, the *dimensional* way of extension refers to the ambient space ( $R$ ) of  $\varphi(X)$  being substituted by a (topological or not) vector space. An account of the results in this area is to be found in Goepfert et al. [29, Ch 3]. The metrical extension of the same consists in conditions imposed upon our metric being relaxed. Some of these extensions were already stated; for the remaining ones, we refer to Hyers, Isac and Rassias [31, Ch 5]; see also Turinici [63].

By the developments above, we therefore have the implications:

$$(DC) \implies (KP\text{-}p\text{-}mp) \implies (KP\text{-}sp\text{-}mp) \implies (ZB\text{-}m\text{-}mp) \implies (EVP\text{-}m).$$

So, we may ask whether these may be reversed. Clearly, the natural setting for solving this problem is the strongly reduced system (ZF-AC).

Let  $X$  be a nonempty set; and  $(\leq)$  be a (partial) order on it. We say that  $(\leq)$  has the *inf-lattice* property, provided:  $x \wedge y := \inf(x, y)$  exists, for all  $x, y \in X$ . Remember that  $z \in X$  is a  $(\leq)$ -*maximal* element if  $X(z, \leq) = \{z\}$ ; the class of all these points will be denoted as  $\max(X, \leq)$ . Call  $(\leq)$ , a *Zorn order* when

$$\begin{aligned} \max(X, \leq) \text{ is nonempty and } & \textit{cofinal} \text{ in } X \\ \text{(for each } u \in X \text{ there exists a } & (\leq)\text{-maximal } v \in X \text{ with } u \leq v). \end{aligned}$$

Further aspects are to be described in a metric setting. Let  $d(., .)$  be a metric over  $X$ ; and  $\varphi : X \rightarrow R_+$  be some function. Then, the natural choice for  $(\leq)$  above is

$$x \leq_{(d, \varphi)} y \text{ iff } d(x, y) \leq \varphi(x) - \varphi(y);$$

referred to as the *Brøndsted order* [13] attached to  $(d, \varphi)$ . Denote

$$X(x, \rho) = \{u \in X; d(x, u) < \rho\}, x \in X, \rho > 0$$

[the *open* sphere with center  $x$  and radius  $\rho$ ]. Call  $(X, d)$ , *discrete* when

$$\text{for each } x \in X \text{ there exists } \rho = \rho(x) > 0 \text{ such that } X(x, \rho) = \{x\}.$$

Note that, under such an assumption, any function  $\psi : X \rightarrow R$  is continuous over  $X$ . However, this is not extendable to the *d-Lipschitz property*

$$|\psi(x) - \psi(y)| \leq Ld(x, y), x, y \in X, \text{ for some } L > 0;$$

hence, all the more, to the *d-nonexpansive property* ( $L = 1$ ).

Now, the statement below is a particular case of (EVP-m):

**Theorem 20** *Let the metric space  $(X, d)$  and the function  $\varphi : X \rightarrow R_+$  satisfy*

- (75-i)  $(X, d)$  is discrete bounded and complete
- (75-ii)  $(\leq_{(d,\varphi)})$  has the inf-lattice property
- (75-iii)  $\varphi$  is  $d$ -nonexpansive and  $\varphi(X)$  is countable.

*Then,  $(\leq_{(d,\varphi)})$  is a Zorn order.*

We shall refer to it as: the discrete Lipschitz countable version of (EVP-m) (in short: (EVP-m-dLc)). Clearly, (EVP-m)  $\implies$  (EVP-m-dLc). The remarkable fact is that this last principle yields (DC); and completes the circle between all these.

**Proposition 15** *We have the inclusion (EVP-m-dLc)  $\implies$  (DC) [in the strongly reduced system (ZF-AC)]. So (by the above),*

- (74-1) *the maximal/variational principles (KP-p-mp), (KP-sp-mp), (ZB-m-mp) and (EVP-m) are all equivalent with (DC); hence, mutually equivalent*
- (74-2) *each intermediary maximal/variational statement (VP) with (DC)  $\implies$  (VP)  $\implies$  (EVP-m) is equivalent with both (DC) and (EVP-m).*

For a complete proof, see Turinici [66]. In particular, when the inf-lattice, nonexpansive, and countable properties are ignored in (EVP-m-dLc), the last result above reduces to the one in Brunner [14]. Note that, in the same particular setting, a different proof of the underlying inclusion was provided in Dodu and Morillon [19]; see also Schechter [52, Ch 19, Sect 19.51]. For a number of finite dimensional stability criteria involving the maximal points in question, we refer to the monograph by Hyers et al. [32, Ch 11].

## References

1. R.P. Agarwal, D. O'Regan, N. Shahzad, Fixed point theory for generalized contractive maps of Meir-Keeler type. *Math. Nachr.* **276**, 3–22 (2004)
2. R.P. Agarwal, M.A. El-Gebeily, D. O'Regan, Generalized contractions in partially ordered metric spaces. *Appl. Anal.* **87**, 109–116 (2008)
3. M. Altman, A generalization of the Brezis-Browder principle on ordered sets. *Nonlin. Anal.* **6**, 157–165 (1982)
4. V.G. Angelov, *Fixed Points in Uniform Spaces and Applications* (Cluj University Press, Cluj-Napoca, 2009)
5. D. Ariza-Ruiz, A. Jimenez-Melado, A continuation method for weakly contractive mappings under the interior condition. *Fixed Point Th. Appl.* **2009**, Article ID 809315 (2009)
6. D. Ariza-Ruiz, A. Jimenez-Melado, A continuation method for weakly Kannan maps. *Fixed Point Th. Appl.* **2010**, Article ID 321594 (2010)
7. S. Banach, Sur les opérations dans les ensembles abstraits et leur application aux équations intégrales. *Fund. Math.* **3**, 133–181 (1922)
8. P. Bernays, A system of axiomatic set theory: Part III. Infinity and enumerability analysis. *J. Symbolic Logic* **7**, 65–89 (1942)
9. N. Bourbaki, Sur le théorème de Zorn. *Archiv Math.* **2**, 434–437 (1949/1950)
10. N. Bourbaki, *General Topology (Chs 1–4)* (Springer, Berlin, 1995)
11. D.W. Boyd, J.S.W. Wong, On nonlinear contractions. *Proc. Amer. Math. Soc.* **20**, 458–464 (1969)



12. H. Brezis, F.E. Browder, A general principle on ordered sets in nonlinear functional analysis. *Adv. Math.* **21**, 355–364 (1976)
13. A. Brøndsted, Fixed points and partial orders. *Proc. Amer. Math. Soc.* **60**, 365–366 (1976)
14. N. Brunner, Topologische Maximalprinzipien. *Zeitschr. Math. Logik Grundl. Math.* **33**, 135–139 (1987)
15. M. Cherichi, B. Samet, Fixed point theorems on ordered gauge spaces with applications to integral equations. *Fixed Point Th. Appl.* **2012**, 13 (2012)
16. P.J. Cohen, *Set Theory and the Continuum Hypothesis* (Benjamin, New York, 1966)
17. O. Costinescu, *Elements of General Topology* (Romanian), Ed. (Tehnică, București, 1970)
18. S. Dancs, M. Hegedus, P. Medvegyev, A general ordering and fixed-point principle in complete metric space. *Acta Sci. Math. (Szeged)* **46**, 381–388 (1983)
19. J. Dodu, M. Morillon, The Hahn-Banach property and the axiom of choice. *Math. Logic Q.* **45**, 299–314 (1999)
20. J. Dugundji, *Topology* (Allyn and Bacon, Boston, 1966)
21. I. Ekeland, Nonconvex minimization problems. *Bull. Amer. Math. Soc. (New Series)* **1**, 443–474 (1979)
22. R. Engelking, *General Topology* (Heldermann Verlag, Berlin, 1989)
23. M. Frigon, Fixed point results in Fréchet and gauge spaces and applications to differential and integral equations. *Lect. Notes. Nonlin. Anal.* **16**, 9–91 (2017)
24. M.A. Geraghty, On contractive mappings. *Proc. Amer. Math. Soc.* **40**, 604–608 (1973)
25. N. Gheorghiu, Contraction theorem in uniform spaces (Romanian). *Stud. Cerc. Mat.* **19**, 119–122 (1967)
26. N. Gheorghiu, Fixed point theorems in uniform spaces. *An. Șt. Univ. Al. I. Cuza Iași (Sect. I-a, Mat.)* **28**, 17–18 (1982)
27. N. Gheorghiu, E. Rotaru, A fixed point theorem in uniform spaces. *An. Șt. Univ. Al. I. Cuza Iasi (Sect. I-a, Mat)* **18**, 311–314 (1972)
28. N. Gheorghiu, M. Turinici, Equations intégrales dans les espaces localement convexes. *Rev. Roum. Math. Pures Appl.* **23**, 33–40 (1978)
29. A. Goepfert, H. Riahi, C. Tammer, C. Zălinescu, *Variational Methods in Partially Ordered Spaces*. *Canad. Math. Soc. Books Math.*, vol. 17 (Springer, New York, 2003)
30. A. Hamel, *Variational Principles on Metric and Uniform Spaces*. Habilitation Thesis. Martin-Luther University, Halle-Wittenberg, 2005
31. D.H. Hyers, G. Isac, Th. M. Rassias, *Topics in Nonlinear Analysis and Applications* (World Scientific Publishing, Singapore, 1997)
32. D.H. Hyers, G. Isac, Th.M. Rassias, *Stability of Functional Equations in Several Variables* (Birkhäuser, Boston, 1998)
33. J. Jachymski, Common fixed point theorems for some families of mappings. *Indian J. Pure Appl. Math.* **25**, 925–937 (1994)
34. J. Jachymski, The contraction principle for mappings on a metric space with a graph. *Proc. Amer. Math. Soc.* **136**, 1359–1373 (2008)
35. B.G. Kang, S. Park, On generalized ordering principles in nonlinear analysis. *Nonlin. Anal.* **14**, 159–165 (1990)
36. S. Kasahara, On some generalizations of the Banach contraction theorem. *Publ. Res. Inst. Math. Sci. Kyoto Univ.* **12**, 427–437 (1976)
37. M.S. Khan, M. Swaleh, S. Sessa, Fixed point theorems by altering distances between the points. *Bull. Austral. Math. Soc.* **30**, 1–9 (1984)
38. K. Kuratowski, *Topology* (Academic Press, New York, 1966)
39. S. Leader, Fixed points for general contractions in metric spaces. *Math. Japonica* **24**, 17–24 (1979)
40. J. Matkowski, *Integrable Solutions of Functional Equations*. *Dissertationes Math.*, vol. 127 (Polish Scientific Publishers, Warsaw, 1975)
41. J. Matkowski, Fixed point theorems for mappings with a contractive iterate at a point. *Proc. Amer. Math. Soc.* **62**, 344–348 (1977)

42. A. Meir, E. Keeler, *A theorem on contraction mappings*. J. Math. Anal. Appl. **28**, 326–329 (1969)
43. G.H. Moore, *Zermelo's Axiom of Choice: its Origin, Development and Influence* (Springer, New York, 1982)
44. Y. Moskhovakis, *Notes on Set Theory* (Springer, New York, 2006)
45. J.J. Nieto, R. Rodríguez-Lopez, Contractive mapping theorems in partially ordered sets and applications to ordinary differential equations. Order **22**, 223–239 (2005)
46. D. O'Regan, R. Precup, Continuation theory for contractions on spaces with two vector-valued metrics. Applicable Anal. **82**, 131–144 (2003)
47. A.C.M. Ran, M.C. Reurings, A fixed point theorem in partially ordered sets and some applications to matrix equations. Proc. Amer. Math. Soc. **132**, 1435–1443 (2004)
48. S. Reich, Fixed points of contractive functions. Boll. Un. Mat. Ital. **5**, 26–42 (1972)
49. B.E. Rhoades, A comparison of various definitions of contractive mappings. Trans. Amer. Math. Soc. **226**, 257–290 (1977)
50. I.A. Rus, *Generalized Contractions and Applications* (Cluj University Press, Cluj-Napoca, 2001)
51. B. Samet, M. Turinici, Fixed point theorems on a metric space endowed with an arbitrary binary relation and applications. Commun. Math. Anal. **13**, 82–97 (2012)
52. E. Schechter, *Handbook of Analysis and its Foundation* (Academic Press, New York, 1997)
53. A. Tarski, Axiomatic and algebraic aspects of two theorems on sums of cardinals. Fund. Math. **35**, 79–104 (1948)
54. M. Turinici, Fixed points of implicit contraction mappings. An. Șt. Univ. Al. I. Cuza Iași (Sect. I-a, Mat.) **22**, 177–180 (1976)
55. M. Turinici, Nonlinear contractions and applications to Volterra functional equations. An. Șt. Univ. Al. I. Cuza Iași (Sect. I-a, Mat) **23**, 43–50 (1977)
56. M. Turinici, Maximal elements in a class of order complete metric spaces. Math. Japonica **25**, 511–517 (1980)
57. M. Turinici, Maximality principles and mean value theorems. An. Acad. Brasil. Cienc. **53**, 653–655 (1981)
58. M. Turinici, A generalization of Brezis-Browder's ordering principle. An. Șt. Univ. Al. I. Cuza Iași (Sect. I-a, Mat) **28**, 11–16 (1982)
59. M. Turinici, A generalization of Altman's ordering principle. Proc. Amer. Math. Soc. **90**, 128–132 (1984)
60. M. Turinici, Fixed points for monotone iteratively local contractions. Dem. Math. **19**, 171–180 (1986)
61. M. Turinici, Abstract comparison principles and multivariable Gronwall-Bellman inequalities. J. Math. Anal. Appl. **117**, 100–127 (1986)
62. M. Turinici, Pseudometric extensions of the Brezis-Browder ordering principle. Math. Nachrichten **130**, 91–103 (1987)
63. M. Turinici, A monotone version of the variational Ekeland's principle. An. Șt. Univ. Al. I. Cuza Iași (Sect. I-a, Mat.) **36**, 329–352 (1990)
64. M. Turinici, Function pseudometric VP and applications. Bul. Inst. Polit. Iași (S. Mat., Mec. Teor., Fiz.) **53**(57), 393–411 (2007)
65. M. Turinici, Ran-Reurings theorems in ordered metric spaces. J. Indian Math. Soc. **78**, 207–214 (2011)
66. M. Turinici, Sequential maximality principles, in *Mathematics Without Boundaries*, ed. by T.M. Rassias, P.M. Pardalos (Springer, New York, 2014), pp. 515–548
67. M. Turinici, *Fixed Point Results in Ordered Gauge Spaces*. Selected Topics in Metrical Fixed Point Theory, Paper 3–4 (Pim Editorial House, Iași, 2017)
68. J.R. Wang, M. Fečkan, Practical Ulam-Hyers-Rassias stability for nonlinear equations. Math. Bohemica **142**, 47–56 (2017)
69. E.S. Wolk, On the principle of dependent choices and some forms of Zorn's lemma. Canad. Math. Bull. **26**, 365–367 (1983)

# Analytic Methods in Rhoades Contractions Theory



Mihai Turinici

**Abstract** A lot of implicit analytic methods is proposed for the study of Rhoades contractions over a class of relational metric spaces, via regulated functions. Technical connections with some particular statements in the area due to Vujaković et al. (Mathematics 767:8 (2020)) are also discussed.

## 1 Introduction

Let  $X$  be a nonempty set. Call the subset  $Y$  of  $X$ , *almost singleton* (in short: *asingleton*), provided  $[y_1, y_2 \in Y \text{ implies } y_1 = y_2]$ ; and *singleton* if, in addition,  $Y$  is nonempty; note that in this case  $Y = \{y\}$ , for some  $y \in X$ . Further, take a metric  $d : X \times X \rightarrow R_+ := [0, \infty[$  over  $X$ ; the couple  $(X, d)$  will be referred to as a *metric space*. Finally, let  $T \in \mathcal{F}(X)$  be a selfmap of  $X$ . [Here, for each couple  $A, B$  of nonempty sets,  $\mathcal{F}(A, B)$  denotes the class of all functions from  $A$  to  $B$ ; when  $A = B$ , we write  $\mathcal{F}(A)$  in place of  $\mathcal{F}(A, A)$ ]. Denote  $\text{Fix}(T) = \{x \in X; x = Tx\}$ ; each point of this set is referred to as *fixed* under  $T$ . The determination of such points is carried out in the context below, comparable with the one in Rus [32, Ch 2, Sect 2.2]:

**pic-1)** We say that  $T$  is a *Picard operator* (modulo  $d$ ) if, for each  $x \in X$ , the iterative sequence  $(T^n x; n \geq 0)$  is  $d$ -Cauchy

**pic-2)** We say that  $T$  is a *strong Picard operator* (modulo  $d$ ) if, for each  $x \in X$ ,  $(T^n x; n \geq 0)$  is  $d$ -convergent with  $\lim_n(T^n x) \in \text{Fix}(T)$ .

**pic-3)** We say that  $T$  is *fix-asingleton* if  $\text{Fix}(T)$  is asingleton; and *fix-singleton*, provided  $\text{Fix}(T)$  is singleton.

---

**AMS Subject Classification:** 47H10 (Primary), 54H25 (Secondary)

---

M. Turinici (✉)

A. Myller Mathematical Seminar, A. I. Cuza University, Iași, Romania

e-mail: [mturi@uaic.ro](mailto:mturi@uaic.ro)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,

Springer Optimization and Its Applications 173,

[https://doi.org/10.1007/978-3-030-72563-1\\_25](https://doi.org/10.1007/978-3-030-72563-1_25)

705

Concerning the existence and uniqueness of such points, a basic result (referred to as: Banach fixed point theorem; in short: (B-fpt)) may be stated as follows. Call the selfmap  $T$ ,  $(d; \alpha)$ -contractive (where  $\alpha \geq 0$ ), if

$$(con) \quad d(Tx, Ty) \leq \alpha d(x, y), \text{ for all } x, y \in X.$$

**Theorem 1** *Suppose that  $T$  is  $(d; \alpha)$ -contractive, for some  $\alpha \in [0, 1[$ . In addition, let  $X$  be  $d$ -complete. Then,  $T$  is a strong Picard operator and fix-asingleton (or, equivalently: fix-singleton).*

This result, established in 1922 by Banach [2], found some important applications to the operator equations theory. Consequently, a multitude of extensions for (B-fpt) were proposed. From the perspective of this exposition, the set implicit ones are of interest. Denote, for  $x, y \in X$

$$\begin{aligned} Q_1(x, y) &= d(x, Tx), \quad Q_2(x, y) = d(x, y), \\ Q_3(x, y) &= d(x, Ty), \quad Q_4(x, y) = d(Tx, y), \\ Q_5(x, y) &= d(Tx, Ty), \quad Q_6(x, y) = d(y, Ty), \\ \mathcal{Q}(x, y) &= (Q_1(x, y), Q_2(x, y), Q_3(x, y), Q_4(x, y), Q_5(x, y), Q_6(x, y)). \end{aligned}$$

Then, the underlying contractions may be written as

$$(i\text{-s-con}) \quad \mathcal{Q}(x, y) \in \Upsilon, \text{ for all } x, y \in X \text{ with } x\mathcal{R}y;$$

where  $\Upsilon \subseteq R_+^6$  is a (nonempty) subset and  $\mathcal{R}$  is a relation over  $X$ . In particular, when  $\Upsilon$  is the zero-section of a certain function  $F : R_+^6 \rightarrow R$ ; i.e.,

$$\Upsilon = \{(t_1, \dots, t_6) \in R_+^6; F(t_1, \dots, t_6) \leq 0\},$$

the implicit contractive condition above has the functional form:

$$(i\text{-f-con}) \quad F(\mathcal{Q}(x, y)) \leq 0, \text{ for all } x, y \in X \text{ with } x\mathcal{R}y.$$

The natural case discussed in the immense majority of papers is  $\mathcal{R} = X \times X$  (the *trivial* relation over  $X$ ). Concerning the “genuine” implicit case, some recent contributions in the area may be found in Akkouchi [1], Berinde and Vetro [3], Nashine et al. [22], or Popa and Mocanu [26]. We stress that in almost all papers based on implicit techniques—including the ones we just quoted—it is claimed that the starting point in the area is represented by the contribution due to Popa [25]. Unfortunately, this claim is not true: fixed point results based on implicit techniques were obtained more than two decades ago in two papers by Turinici [36, 37]. But, we must note that some partial aspects of the set-implicit theory have been discussed in the (classical by now) 1969 Meir-Keeler fixed point principle [18].

On the other hand, for the explicit case, some basic contributions were obtained in Boyd and Wong [5], Leader [16], Reich [29], or Matkowski [17]; see also the survey paper by Rhoades [30].

Having these precise, it is our aim in the following to give a lot of fixed point results—involving six-dimensional contractive conditions—for selfmaps acting upon relational metric spaces. Precisely, as a by-product of our developments, two classes of fixed point statements were build up for contractive conditions of the type discussed in the papers by Rhoades [31] and Dutta and Choudhury [12]; namely:

$$(RDC) \quad \psi(u \circ \mathcal{Q}(x, y)) \leq \psi(v \circ \mathcal{Q}(x, y)) - \varphi(w \circ \mathcal{Q}(x, y)), \quad x, y \in X, x\mathcal{R}y,$$

where  $(\psi, \varphi)$  is a functional couple over  $\mathcal{F}(R_+^0, R)$ ,  $(u, v, w)$  is a triple of functions over  $\mathcal{F}(R_+^6, R_+)$ , and  $\mathcal{R}$  is a relation over  $X$ . Note that, the fixed point results of the first class involve contractions like before based on  $\psi$  being a regulated function—hence, not in general increasing. Then, as another by-product of these, some fixed point results are given for a class of contractive maps introduced by Cosentino and Vetro [9] and refined by Vujaković et al. [44]. Finally, as a special application of the used asymptotic techniques, an almost regulated version is given for the standard result in the area due to Wardowski [46]. Some anticipative type variants of these results, including the 2015 contribution in the area due to Dung and Hang [11] will be discussed elsewhere.

## 2 Dependent Choice Principle

Throughout this exposition, the axiomatic system in use is Zermelo-Fraenkel’s (abbreviated: ZF), as described by Cohen [8, Ch 2]. The notations and basic facts to be considered in this system are more or less standard. Some important ones are described below.

(A) Let  $X$  be a nonempty set. By a *relation* over  $X$ , we mean any nonempty part  $\mathcal{R} \subseteq X \times X$ ; then,  $(X, \mathcal{R})$  will be referred to as a *relational structure*. For simplicity, we sometimes write  $(x, y) \in \mathcal{R}$  as  $x\mathcal{R}y$ . Note that  $\mathcal{R}$  may be regarded as a mapping between  $X$  and  $\exp[X]$  (=the class of all subsets in  $X$ ). In fact, denote for  $x \in X$ :

$$X(x, \mathcal{R}) = \{y \in X; x\mathcal{R}y\} \text{ (the section of } \mathcal{R} \text{ through } x\text{);}$$

then, the desired mapping representation is  $[\mathcal{R}(x) = X(x, \mathcal{R}), x \in X]$ .

A basic example of relational structure is to be constructed as below. Let  $N = \{0, 1, \dots\}$  be the set of *natural* numbers, endowed with the usual addition and (partial) order; note that

$(N, \leq)$  is well ordered: any (nonempty) subset of  $N$  has a first element.

Further, denote for  $p, q \in N, p \leq q$ ,

$$N[p, q] = \{n \in N; p \leq n \leq q\}, N]p, q[ = \{n \in N; p < n < q\},$$

$$N[p, q[ = \{n \in N; p \leq n < q\}, N]p, q] = \{n \in N; p < n \leq q\};$$

as well as, for  $r \in N$ ,

$$N[r, \infty[ = \{n \in N; r \leq n\}, N]r, \infty[ = \{n \in N; r < n\}.$$

By definition,  $N[0, r[ = N(r, >)$  is referred to as the *initial interval* (in  $N$ ) induced by  $r$ . Any set  $P$  with  $P \sim N$  (in the sense: there exists a bijection from  $P$  to  $N$ ) will be referred to as *effectively denumerable*. In addition, given some natural number  $n \geq 1$ , any set  $Q$  with  $Q \sim N(n, >)$  will be said to be *n-finite*; when  $n$  is generic here, we say that  $Q$  is *finite*. Finally, the (nonempty) set  $Y$  is called (at most) *denumerable* iff it is either effectively denumerable or finite.

Let  $X$  be a nonempty set. By a *sequence* in  $X$ , we mean any mapping  $x : N \rightarrow X$ , where  $N = \{0, 1, \dots\}$  is the set of *natural* numbers. For simplicity reasons, it will be useful to denote it as  $(x(n); n \geq 0)$ , or  $(x_n; n \geq 0)$ ; moreover, when no confusion can arise, we further simplify this notation as  $(x(n))$  or  $(x_n)$ , respectively. Also, any sequence  $(y_n := x_{i(n)}; n \geq 0)$  with

$$(i(n); n \geq 0) \text{ is divergent } (i(n) \rightarrow \infty \text{ as } n \rightarrow \infty)$$

will be referred to as a *subsequence* of  $(x_n; n \geq 0)$ . Note that, under such a convention, the relation “subsequence of” is transitive; i.e.:

$$(z_n) = \text{subsequence of } (y_n) \text{ and } (y_n) = \text{subsequence of } (x_n)$$

$$\text{imply } (z_n) = \text{subsequence of } (x_n).$$

**(B)** Remember that, an outstanding part of (ZF) is the *Axiom of Choice* (abbreviated: AC); which, in a convenient manner, may be written as

(AC) For each couple  $(J, X)$  of nonempty sets and each function

$F : J \rightarrow \exp(X)$ , there exists a (selective) function

$f : J \rightarrow X$ , with  $f(v) \in F(v)$ , for each  $v \in J$ .

Here,  $\exp(X)$  stands for the class of all nonempty elements in  $\exp[X]$ . Sometimes, when the index set  $J$  is denumerable, the existence of such a selective function may be determined by using a weaker form of (AC), called: *Dependent Choice* principle (in short: DC). Call the relation  $\mathcal{R}$  over  $X$ , *proper* when

$$(X(x, \mathcal{R}) =) \mathcal{R}(x) \text{ is nonempty, for each } x \in X.$$

Then,  $\mathcal{R}$  is to be viewed as a mapping between  $X$  and  $\exp(X)$ ; and the couple  $(X, \mathcal{R})$  will be referred to as a *proper relational structure*. Further, given  $a \in X$ , let us say that the sequence  $(x_n; n \geq 0)$  in  $X$  is  $(a, \mathcal{R})$ -*iterative*, provided

$x_0 = a$ , and  $(x_n; n \geq 0)$  is  $\mathcal{R}$ -increasing :  
 $x_n \mathcal{R} x_{n+1}$  (i.e.:  $x_{n+1} \in \mathcal{R}(x_n)$ ), for all  $n$ .

**Proposition 1** *Let the relational structure  $(X, \mathcal{R})$  be proper. Then, for each  $a \in X$  there is at least an  $(a, \mathcal{R})$ -iterative sequence in  $X$ .*

This principle—proposed, independently, by Bernays [4] and Tarski [35]—is deductible from (AC), but not conversely; cf. Wolk [47]. Moreover, by the developments in Moskhovakis [20, Ch 8], and Schechter [33, Ch 6], the reduced system (ZF-AC+DC) is comprehensive enough so as to cover the “usual” mathematics; see also Moore [19, Appendix 2].

A basic consequence of (DC) is the so-called *Denumerable Axiom of Choice* [in short: AC(N)].

**Proposition 2** *Let  $F : N \rightarrow \exp(X)$  be a function. Then, for each  $a \in F(0)$  there exists a function  $f : N \rightarrow X$  with  $f(0) = a$  and  $f(n) \in F(n)$ ,  $\forall n \geq 0$ .*

**Proof** Denote  $Q = N \times X$ ; and let us introduce the (proper) relation  $\mathcal{R}$  over it, according to:

$$\mathcal{R}(n, x) = \{n + 1\} \times F(n + 1), \quad n \geq 0, x \in X.$$

By an application of (DC) to the proper relational structure  $(Q, \mathcal{R})$  the conclusion follows; we do not give details.

As a consequence of the above facts,

$$\begin{aligned} \text{(DC)} &\implies \text{(AC(N)) in (ZF-AC); or, equivalently :} \\ &\text{(AC(N)) is deductible in the system (ZF-AC + DC).} \end{aligned}$$

The reciprocal of the written inclusion is not true; see Moskhovakis [20, Ch 8, Sect 8.25] for details.

### 3 Conv-Cauchy Structures

Let  $X$  be a nonempty set. Further, let  $d : X \times X \rightarrow R_+$  be a mapping with

- (m-1)  $d$  is *triangular*:  $d(x, z) \leq d(x, y) + d(y, z)$ ,  $\forall x, y, z \in X$
- (m-2)  $d$  is *reflexive-sufficient*:  $d(x, y) = 0$  iff  $x = y$
- (m-3)  $d$  is *symmetric*:  $d(x, y) = d(y, x)$ , for all  $x, y \in X$ .

We then say that  $d(., .)$  is a *metric* on  $X$ ; and the couple  $(X, d)$  will be referred to as a *metric space*.

(A) We introduce a  $d$ -convergence and  $d$ -Cauchy structure on  $X$  as follows. Given the sequence  $(x_n)$  in  $X$  and the point  $x \in X$ , we say that  $(x_n)$ ,  $d$ -converges to  $x$  (written as:  $x_n \xrightarrow{d} x$ ), provided  $d(x_n, x) \rightarrow 0$  as  $n \rightarrow \infty$ ; i.e.,

$$\forall \varepsilon > 0, \exists i = i(\varepsilon) : i \leq n \implies d(x_n, x) < \varepsilon;$$

or, equivalently:

$$\forall \varepsilon > 0, \exists i = i(\varepsilon) : i \leq n \implies d(x_n, x) \leq \varepsilon.$$

The set of all such points  $x$  will be denoted  $\lim_n(x_n)$ ; when it is nonempty, then  $(x_n)$  is called  $d$ -convergent. By this very definition, we have the properties:

(conv-1) ( $\xrightarrow{d}$ ) is *hereditary*)

$x_n \xrightarrow{d} x$  implies  $y_n \xrightarrow{d} x$ , for each subsequence  $(y_n)$  of  $(x_n)$

(conv-2) ( $\xrightarrow{d}$ ) is *reflexive*)

$(\forall u \in X) : \text{the constant sequence } (x_n = u; n \geq 0) \text{ fulfills } x_n \xrightarrow{d} u.$

As a consequence,  $(\xrightarrow{d})$  has all properties required in Kasahara [14]; in addition— as  $d$  is triangular symmetric—the following extra property is holding here

(conv-3) ( $\xrightarrow{d}$ ) is *separated* (referred to as  $d$  is *separated*):  
 $\lim_n(x_n)$  is an asingleton, for each sequence  $(x_n)$  in  $X$ .

The introduced concepts allow us to give a useful property.

**Proposition 3** *The mapping  $(x, y) \mapsto d(x, y)$  is  $d$ -Lipschitz, in the sense*

$$(31-1) \quad |d(x, y) - d(u, v)| \leq d(x, u) + d(y, v), \quad \forall (x, y), (u, v) \in X \times X.$$

As a consequence, this map is  $d$ -continuous; i.e.,

$$(31-2) \quad x_n \xrightarrow{d} x, y_n \xrightarrow{d} y \text{ imply } d(x_n, y_n) \rightarrow d(x, y).$$

Further, call the sequence  $(x_n)$ ,  $d$ -Cauchy when  $d(x_m, x_n) \rightarrow 0$  as  $m, n \rightarrow \infty$ ,  $m < n$ ; that is,

$$\forall \varepsilon > 0, \exists j = j(\varepsilon) : j \leq m < n \implies d(x_m, x_n) < \varepsilon;$$

or, equivalently:

$$\forall \varepsilon > 0, \exists j = j(\varepsilon) : j \leq m < n \implies d(x_m, x_n) \leq \varepsilon.$$



the class of all these will be denoted as  $Cauchy(d)$ . As before, from this very definition one has the properties

(Cauchy-1) ( $Cauchy(d)$  is hereditary)  
 $(x_n)$  is  $d$ -Cauchy implies  $(y_n)$  is  $d$ -Cauchy,  
 for each subsequence  $(y_n)$  of  $(x_n)$

(Cauchy-2) ( $Cauchy(d)$  is reflexive)  
 $(\forall u \in X)$ : the constant sequence  $(x_n = u; n \geq 0)$  is  $d$ -Cauchy.

Hence,  $Cauchy(d)$  is a Cauchy structure, under the lines in Turinici [38].

Now—according to the quoted work—term the couple  $((\xrightarrow{d}), Cauchy(d))$ , a *conv-Cauchy structure* induced by  $d$ . The following regularity conditions about this structure are to be (optionally) considered

(CC-1)  $d$  is regular: each  $d$ -convergent sequence in  $X$  is  $d$ -Cauchy  
 (CC-2)  $d$  is complete: each  $d$ -Cauchy sequence in  $X$  is  $d$ -convergent.

Clearly, the former of these is always obtainable, via  $d$ -triangular symmetric; but the latter one is not in general valid.

**(B)** In the following, some  $d$ -Cauchy criteria will be stated. Some preliminaries are needed.

Let us say that  $(x_n; n \geq 0)$  is  *$d$ -asymptotic*, provided

$$r_n := d(x_n, x_{n+1}) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Clearly, each  $d$ -Cauchy sequence is  $d$ -asymptotic too; the reciprocal of this is not in general true. This tells us that the  $d$ -Cauchy criteria we are looking for are to be sought in the class of  $d$ -asymptotic sequences. Precisely, suppose that the following inclusion is valid

(for each  $d$ -asymptotic sequence  $(x_n)$ ) :  
 $(x_n)$  is not  $d$ -Cauchy implies  $(x_n)$  has the property  $\pi$ .

Then, the  $d$ -Cauchy criterion in question writes

(for each  $d$ -asymptotic sequence  $(x_n)$ ) :  
 $(x_n)$  is not endowed with the property  $\pi$  implies  $(x_n)$  is  $d$ -Cauchy.

To get concrete examples of such properties, we need some other conventions. Given the  $d$ -asymptotic sequence  $(x_n; n \geq 0)$  and  $\varepsilon > 0$ , let us say that  $i \in N$  is  $\varepsilon$ -regular, provided

$$i \leq n \text{ implies } d(x_n, x_{n+1}) < \varepsilon.$$

The class  $\mathcal{Z}(\varepsilon)$  of all these numbers is nonempty; so that

$(\forall \varepsilon > 0) : Z(\varepsilon) = \min \mathcal{Z}(\varepsilon)$  is well defined, as an element of  $N$ ;

with, in addition:  $d(x_n, x_{n+1}) < \varepsilon$ , for all  $n \geq Z(\varepsilon)$ .

Further, note that by this very definition,

$(\forall h \geq 1) : \Gamma_n(h) := \max\{d(x_n, x_{n+i}); i \in N[0, h]\} < \varepsilon$ ,

for all  $n \geq Z(\varepsilon/h)$ ; whence,  $\Gamma_n(h) \rightarrow 0$  as  $n \rightarrow \infty$ .

Define the subsets of  $N \times N$

$(\leq; N) = \{(m, n) \in N \times N; m \leq n\}$ ,  $(<; N) = \{(m, n) \in N \times N; m < n\}$ ;

these are just the graphs of the relations  $(\leq)$  and  $(<)$  introduced over  $N$ .

Further, let us say that the subset  $\Theta$  of  $R_+^0 := ]0, \infty[$  is  $(>)$ -cofinal in  $R_+^0$ , when for each  $\varepsilon > 0$ , there exists  $\theta \in \Theta$  with  $\varepsilon > \theta$ .

Finally, given the sequence  $(r_n; n \geq 0)$  in  $R$  and the point  $r \in R$ , let us write

$r_n \rightarrow r +$  (also written as:  $\lim_n r_n = r +$ ),

if  $r_n \rightarrow r$  and  $r_n > r$ , for all  $n \geq 0$ .

The following result, referred to as *Boyd-Wong Criterion* is now available.

**Theorem 2** *Let the sequence  $(x_n; n \geq 0)$  in  $X$  be such that*

(31-i)  $(x_n; n \geq 0)$  is  $d$ -asymptotic ( $r_n := d(x_n, x_{n+1}) \rightarrow 0$  as  $n \rightarrow \infty$ )

(31-ii)  $(x_n; n \geq 0)$  is not  $d$ -Cauchy.

Further, let the subset  $\Theta$  of  $R_+^0$  be  $(>)$ -cofinal in  $R_+^0$ ; and  $h \geq 1$  be some rank. There exist then a number  $\gamma \in \Theta$ , a rank  $J := J(\gamma, h) \geq 1$ , and a couple of rank-sequences  $(m(k); k \geq 0)$ ,  $(n(k); k \geq 0)$ , with

(31-a)  $J + k \leq m(k) < m(k) + 3h < n(k)$ ,  $\forall k$

(31-b)  $J + k \leq m(k) < m(k) + 2h < n(k) - 1 < n(k)$ , and  $d(x_{m(k)}, x_{n(k)}) > \gamma$ ,

$d(x_{m(k)}, x_{n(k)-1}) \leq \Gamma_{m(k)}(3h) + \gamma$ ,  $\forall k$

(31-c)  $U_k := d(x_{m(k)}, x_{n(k)}) \rightarrow \gamma +$  as  $k \rightarrow \infty$

(31-d)  $\forall s, t \in N[0, 3h]: V_k(s, t) := d(x_{m(k)+s}, x_{n(k)+t}) \rightarrow \gamma +$  as  $k \rightarrow \infty$

(31-e)  $\forall p, q \in N: V_k(p, q) := d(x_{m(k)+p}, x_{n(k)+q}) \rightarrow \gamma$  as  $k \rightarrow \infty$ .

**Proof** By definition, the  $d$ -Cauchy property of our sequence writes:

$\forall \varepsilon \in R_+^0, \exists a \in N, \forall (m, n) \in (<; N) : a \leq m < n \implies d(x_m, x_n) \leq \varepsilon$ .

As  $\Theta$  is a  $(>)$ -cofinal part in  $R_+^0$ , this property may be also written as

$\forall \theta \in \Theta, \exists \alpha \in N, \forall (m, n) \in (<; N) : \alpha \leq m < n \implies d(x_m, x_n) \leq \theta$ .

The negation of this property means: there exists  $\beta \in \Theta$  such that

$$\text{(rela-1) } E(j) := \{(m, n) \in (<; N); j \leq m < n, d(x_m, x_n) > \beta\} \neq \emptyset, \forall j.$$

As  $(x_n; n \geq 0)$  is  $d$ -asymptotic, we must have (under  $d$ =metric)

$$d(x_n, x_{n+i}) \rightarrow 0 \text{ as } n \rightarrow \infty, \text{ for each } i \geq 0.$$

Let the number  $\gamma \in \Theta$  be given according to

$$\beta > 3\gamma \text{ (possible, since } \Theta \text{ is } (>)\text{-cofinal in } \mathbb{R}_+^0\text{)}.$$

Further; let us take

$$\text{(rela-2) } J(\gamma, h) = Z(\gamma/3h); \text{ hence, in particular}$$

$$\Gamma_n(3h) := \max\{d(x_n, x_{n+i}); i \in N[0, 3h]\} < \gamma, \forall n \geq J(\gamma, h).$$

Denote for simplicity  $J = J(\gamma, h)$  and

$$(A(k) = E(J + k); k \geq 0); \text{ hence, by definition,}$$

$$A(k) := \{(m, n) \in (<; N); J + k \leq m < n, d(x_m, x_n) > \beta\}, k \geq 0;$$

with, in addition:  $A(k)$  is nonempty, for each  $k \geq 0$ .

By the triangle inequality, we have for each  $(m, n) \in A(k)$

$$\begin{aligned} d(x_{m+s}, x_{n+t}) &\geq d(x_m, x_n) - d(x_m, x_{m+s}) - d(x_n, x_{n+t}) > \\ &\beta - \gamma - \gamma = \beta - 2\gamma > \gamma, \forall s, t \in N[0, 3h]; \end{aligned}$$

which tells us that

$$B(k) := \{(m, n) \in (<; N); J + k \leq m < n, d(x_{m+s}, x_{n+t}) > \gamma, \forall s, t \in N[0, 3h]\}$$

is nonempty, for all  $k \geq 0$ .

Having this precise, denote for each  $k \geq 0$

$$m(k) = \min \text{Dom}(B(k)), n(k) = \min B(k)(m(k)).$$

By this very definition, we get

$$\text{(pro-1) } J + k \leq m(k) < n(k), \forall k \geq 0$$

$$\text{(pro-2) } d(x_{m(k)+s}, x_{n(k)+t}) > \gamma, \forall s, t \in N[0, 3h], \forall k \geq 0.$$

We claim that the couple  $(\gamma, J)$  and the couple of rank-sequences  $(m(k); k \geq 0)$  and  $(n(k); k \geq 0)$  fulfill all conclusions in the statement.

i): By (pro-1), it is clear that the first and the second relation of (31-a) holds. Concerning the third relation of the same, suppose by contradiction that

$$(m(k) <)n(k) \leq m(k) + 3h, \text{ for some } k \geq 0.$$

Then by (rela-2) (and  $m(k) \geq J$ )

$$d(x_{m(k)}, x_{n(k)}) \leq \Gamma_{m(k)}(3h) < \gamma;$$

in contradiction with (pro-2); whence, the third relation in (31-a) holds too.

ii): The first and second relation of (31-b) are directly obtainable from the preceding stage and (pro-2), respectively. Concerning the third relation of (31-b), let  $k \geq 0$  be arbitrary fixed. By definition,  $n(k)$  is the minimum of all ranks  $p \in N$  with

$$(m(k), p) \in B(k); \text{ that is:}$$

$$J + k \leq m(k) < p \text{ and } d(x_{m(k)+s}, x_{p+t}) > \gamma, \forall s, t \in N[0, 3h].$$

As  $m(k) < m(k) + 2h < n(k) - 1$  we must have (by this minimal property)

$$\text{(pro-3) } d(x_{m(k)+s}, x_{n(k)-1+t}) \leq \gamma, \text{ for some } s, t \in N[0, 3h].$$

But, in view of (pro-2) once again,

$$\text{(pro-4) } d(x_{m(k)+u}, x_{n(k)-1+v}) > \gamma, \text{ for all } u \in N[0, 3h], v \in N[1, 3h].$$

This, combined with (pro-3), tells us that, necessarily,

$$\text{(pro-5) } d(x_{m(k)+s}, x_{n(k)-1}) \leq \gamma, \text{ for some } s \in N[0, 3h].$$

By the triangular inequality, we have (under the precise index  $s$ )

$$d(x_{m(k)}, x_{n(k)-1}) \leq d(x_{m(k)}, x_{m(k)+s}) + d(x_{m(k)+s}, x_{n(k)-1}) \leq \Gamma_{m(k)}(3h) + \gamma;$$

and the last conclusion of (31-b) follows.

iii): By the very definition of  $(B(k); k \geq 0)$ ,

$$U_k > \gamma, \text{ for all } k \geq 0.$$

Moreover, taking the triangular inequality into account,

$$\gamma < d(x_{m(k)}, x_{n(k)}) \leq d(x_{m(k)}, x_{n(k)-1}) + r_{n(k)-1} \leq \Gamma_{m(k)}(3h) + \gamma + r_{n(k)-1}, \forall k.$$

Passing to limit in this double inequality gives (31-c).

iv): Let  $s, t \in N[0, 3h]$  be arbitrary fixed. By the very definition of  $(B(k); k \geq 0)$ ,

$$V_k(s, t) > \gamma, \text{ for all } k \geq 0.$$

Moreover, from a metrical property of  $d$ ,

$$|d(x_{m(k)}, x_{n(k)}) - d(x_{m(k)+s}, x_{n(k)+t})| \leq d(x_{m(k)}, x_{m(k)+s}) + d(x_{n(k)}, x_{n(k)+t}) \leq \Gamma_{m(k)}(3h) + \Gamma_{n(k)}(3h) < 2\gamma, \text{ for all } k \geq 0.$$

Passing to limit in the relation between the first and the third member of this relation gives (31-d).

v): Fix  $p, q \in N$ . From the metrical property of  $d$  we just evoked,

$$|d(x_{m(k)}, x_{n(k)}) - d(x_{m(k)+p}, x_{n(k)+q})| \leq d(x_{m(k)}, x_{m(k)+p}) + d(x_{n(k)}, x_{n(k)+q}) \leq \Gamma_{m(k)}(p) + \Gamma_{n(k)}(q), \text{ for all } k \geq 0.$$

Passing to limit in the relation between the first and the third member of this relation and noting that

$$\Gamma_{m(k)}(p) \rightarrow 0 \text{ and } \Gamma_{n(k)}(q) \rightarrow 0 \text{ as } k \rightarrow \infty$$

gives the last conclusion (31-e). The proof is complete.

*Remark 1* A natural problem to be posed is that of the last conclusion above being retainable in terms of right convergence; that is,

$$\forall p, q \in N : V_k(p, q) := d(x_{m(k)+p}, x_{n(k)+q}) \rightarrow \gamma + \text{ as } k \rightarrow \infty.$$

This is not in general true; because, fixing some rank  $i$ , and putting  $p = n(i) - m(i)$ , we have

$$d(x_{m(i)+p}, x_{n(i)}) = 0 < \gamma;$$

so that, a relation like  $d(x_{m(k)+p}, x_{n(k)}) \rightarrow \gamma +$  as  $k \rightarrow \infty$  is not possible.

This contradicts an affirmation in Vujaković et al. [44, Lemma 1]; but, fortunately, it has no impact upon the remaining statements in that paper.

By definition, the quadruple  $[\gamma; J; (m(k); k \geq 0); (n(k); k \geq 0)]$  given by this result will be referred to as a *Boyd-Wong*  $(\Theta, h)$ -system attached to  $(x_n)$ . In this case, the result above may be expressed as below.

**Theorem 3** *Let the sequence  $(x_n; n \geq 0)$  in  $X$  be  $d$ -asymptotic. Then, the following conditions/properties are equivalent*

(32-a)  $(x_n; n \geq 0)$  is not  $d$ -Cauchy

(32-b) for each  $(>)$ -cofinal subset  $\Theta$  of  $R_+^0$  and each  $h \geq 1$ , there is at least one *Boyd-Wong*  $[\Theta, h]$ -system attached to  $(x_n)$ .

**Proof** By the preceding result, the first condition includes the second one. Conversely, if the second condition holds then, under the choice  $\Theta = R_+^0$  and  $h = 1$ , there must be one Boyd-Wong  $[\Theta, h]$ -system attached to  $(x_n)$ ; hence, in particular

$$\lim_k d(x_{m(k)}, x_{n(k)}) = \gamma+, \text{ where } \gamma > 0.$$

This necessarily gives us that  $(x_n; n \geq 0)$  is not  $d$ -Cauchy; for otherwise—under a  $d$ -Cauchy condition upon our starting sequence—it follows that

$$\lim_k d(x_{m(k)}, x_{n(k)}) = 0;$$

in contradiction with the limiting property above.

In particular, when  $\Theta = R_+^0$ , the obtained statement covers the 1969 one in Boyd and Wong [5]; so, it is natural that this result be referred to in the proposed way. Further aspects may be found in Reich [29]; see also Khan et al. [15].

### 4 Statement of the Problem

Let  $(X, d)$  be a metric space. For each relation  $\mathcal{A}$  on  $X$  and each (nonempty) subset  $Z$  of  $X$ , let us introduce the convention

$$(\mathcal{A}; Z) = \mathcal{A} \cap (Z \times Z) \text{ (the restriction of } \mathcal{A} \text{ to } Z).$$

In particular,  $\mathcal{A}$  is identical with  $(\mathcal{A}; X)$ .

Let  $T \in \mathcal{F}(X)$  be a selfmap of  $X$ . In the following, sufficient conditions are given for the existence and/or uniqueness of elements in  $\text{Fix}(T)$ .

**4-I)** The proposed problem will be developed in the setting of

(it-seq) the fixed points of  $T$  are ultimately chosen among the limit points (if any)  $T^\omega x_0 := \lim_n (T^n x_0)$ , where  $x_0 \in X$  is arbitrary fixed.

(Here,  $\omega$  is the first transfinite ordinal). To do this, a lot of technical facts about iterative processes is needed.

Let  $x_0$  be some point in  $X$ . By an *orbital* (or: *iterative*) sequence attached to  $x_0$  (and  $T$ ), we mean any sequence  $X_0 := (x_n; n \geq 0)$  (or, simply,  $X_0 = (x_n)$ ) defined as  $(x_n = T^n x_0; n \geq 0)$ . When  $x_0$  is generic here, the resulting object will be referred to as an *orbital* sequence (in short: *o-sequence*) on  $X$ .

Fix in the following such an object  $X_0 = (x_n)$ . Then, denote

$$[X_0] = \{x_n; n \geq 0\} \text{ (the } \textit{trajectory} \text{ attached to } X_0 = (x_n))$$

$$[[X_0]] := \text{cl}([X_0]) \text{ (the } \textit{complete trajectory} \text{ attached to } X_0 = (x_n)).$$

**Proposition 4** *Under these conventions,*

$$[[X_0]] = [X_0] \cup \{z\}, \text{ whenever } z := \lim_n(x_n) \text{ exists.}$$

**Proof** Denote for simplicity  $U_0 = [X_0]$ ,  $V_0 = [[X_0]]$ ; hence,  $V_0 = \text{cl}(U_0)$ . Clearly,  $V_0 \supseteq U_0 \cup \{z\}$ . Suppose that there exists  $v \in V_0$  that is outside  $U_0 \cup \{z\}$ . By the limit definition, there exists  $\sigma > 0$  such that

$$X(v, \sigma) := \{x \in X; d(v, x) < \sigma\} \text{ is disjoint from } U_0 \cup \{z\}.$$

In particular, this tells us that  $v$  cannot belong to  $\text{cl}(U_0) = V_0$ ; contradiction. Consequently,  $V_0 = U_0 \cup \{z\}$ ; and we are done.

**4-II)** Passing to the basic part of our setting, denote for  $x, y \in X$

$$\begin{aligned} Q_1(x, y) &= d(x, Tx), \quad Q_2(x, y) = d(x, y), \\ Q_3(x, y) &= d(x, Ty), \quad Q_4(x, y) = d(Tx, y), \\ Q_5(x, y) &= d(Tx, Ty), \quad Q_6(x, y) = d(y, Ty), \\ \mathcal{Q}(x, y) &= (Q_1(x, y), Q_2(x, y), Q_3(x, y), Q_4(x, y), Q_5(x, y), Q_6(x, y)). \end{aligned}$$

Further, let us construct the family of functions [for  $x, y \in X$ ]

$$\begin{aligned} B_0(x, y) &= \min\{Q_2, Q_5\}(x, y), \\ B_1(x, y) &= \min\{Q_1, Q_2, Q_5, Q_6\}(x, y) \\ B_2(x, y) &= \min\{Q_2, Q_3, Q_4, Q_5\}(x, y), \\ B_3(x, y) &= \min\{Q_1, Q_2, Q_3, Q_4, Q_5, Q_6\}(x, y), \end{aligned}$$

and introduce the relations (over  $X$ )

$$(\text{rela-B}) \quad (B_i > 0) = \{(x, y) \in X \times X; B_i(x, y) > 0\}, i \in \{0, 1, 2, 3\}.$$

Finally, let  $\mathcal{Y}$  be a nonempty subset of  $R_+^6$ ; and let  $\mathcal{R} = \mathcal{R}(\mathcal{Y})$  stand for the relation

$$\mathcal{R} = \mathcal{Q}^{-1}(\mathcal{Y}); \text{ hence, } \mathcal{Q}(\mathcal{R}) \subseteq \mathcal{Y}.$$

Technically, the couple  $(\mathcal{Y}, \mathcal{R})$  will suffice for setting up a class of (implicit) fixed point statements over the ambient space based on the contractive property

(R-contr)  $T$  is  $(d, \mathcal{R}; \mathcal{Y})$ -contractive, provided

$$\mathcal{Q}(x, y) \in \mathcal{Y}, \forall (x, y) \in \mathcal{R}; \text{ that is: } \mathcal{Q}(\mathcal{R}) \subseteq \mathcal{Y}.$$

The starting point of it is the construction of an orbital sequence  $X_0 = (x_n)$  with

(R-asc)  $X_0 = (x_n)$  is  $\mathcal{R}$ -ascending:  $x_n \mathcal{R} x_{n+1}$ , for all  $n$ .

The general aspects of this method will be discussed elsewhere. Here, we will consider a particular case of it—that, in fact, includes a large number of such contractions—based on the restrictive condition

(B1-adm)  $\mathcal{R}$  is  $B_1$ -admissible:  $(B_1 > 0) \subseteq \mathcal{R}$ .

This tells us that, in place of our general contraction, we study its restrictive part

(B1-contr)  $T$  is  $(d, (B_1 > 0); \mathcal{Y})$ -contractive, provided  
 $\mathcal{Q}(x, y) \in \mathcal{Y}, \forall (x, y) \in (B_1 > 0)$ .

Another condition to be used, in certain moments of our exposition, is

(B0-adm)  $\mathcal{R}$  is  $B_0$ -admissible:  $(B_0 > 0) \subseteq \mathcal{R}$ ;

with the associated (stronger) contractive condition

(B0-contr)  $T$  is  $(d, (B_0 > 0); \mathcal{Y})$ -contractive, provided  
 $\mathcal{Q}(x, y) \in \mathcal{Y}, \forall (x, y) \in (B_0 > 0)$ .

This, as we will see, has an impact upon the uniqueness property.

Returning to the  $B_1$ -admissible setting we stress that, in this case, our initial objective of constructing a  $\mathcal{R}$ -ascending orbital sequence is directly attainable, in terms of the coarser relation  $(B_1 > 0)$ . This may be carried out as follows. Call the orbital sequence  $X_0 = (x_n)$ ,  $(B_1 > 0)$ -ascending provided

$(x_n)$  is  $(B_1 > 0)$ -ascending:  $x_n (B_1 > 0) x_{n+1}$ , for all  $n$ .

Note that, by the very choice of this relation, our convention writes

$B_1(x_n, x_{n+1}) > 0$ , for all  $n$ ; referred to as:  $(x_n)$  is  $B$ -ascending.

We then say that the o-sequence  $X_0 = (x_n)$  has the *Ba-property*; or, equivalently, that  $X_0 = (x_n)$  is (Ba-o). The possibility of reaching such a property is a consequence of the reasoning below. Given the orbital sequence  $X_0 = (x_n)$ , we have two alternatives.

**Alt-1)** The orbital sequence  $X_0 = (x_n)$  is *telescopic*, in the sense there exists  $h \geq 0$ , such that  $d(x_h, x_{h+1}) = 0$ ; i.e.:  $x_h = x_{h+1}$ .

By the iterative definition of our sequence, one derives

$x_h = x_n$ , for all  $n \geq h$ ; whence,  $z := x_h$  is an element of  $\text{Fix}(T)$ .



Consequently, this case is completely clarified from the fixed point perspective.

**Alt-2)** The orbital sequence  $X_0 = (x_n)$  is *non-telescopic*, in the sense  $d(x_n, x_{n+1}) > 0$ , for all  $n$ .

Note that, in this case,

$$(B_1(x_n, x_{n+1}) > 0, \forall n); \text{ hence, } X_0 = (x_n) \text{ is (Ba-o).}$$

Summing up, the orbital sequences to be used in the sequel are endowed with the (Ba-o) property. In this case, the basic directions under which the investigations be conducted are described in the list below.

**(pic-1)** We say that the (Ba-o) sequence  $X_0 = (x_n)$  is *semi-Picard* (modulo  $(d, \mathcal{R}; T)$ ) when  $(x_n)$  is  $d$ -asymptotic

**(pic-2)** We say that the (Ba-o) sequence  $X_0 = (x_n)$  is *Picard* (modulo  $(d, \mathcal{R}; T)$ ) when  $(x_n)$  is  $d$ -Cauchy

**(pic-3)** We say that the (Ba-o) sequence  $X_0 = (x_n)$  is *strongly Picard* (modulo  $(d, \mathcal{R}; T)$ ) when  $x_\omega := \lim_n(x_n)$  exists with  $x_\omega \in \text{Fix}(T)$

**(pic-4)** Call the subset  $Y$  of  $X$ ,  $\mathcal{R}$ -almost-singleton (in short:  $\mathcal{R}$ -asingleton) provided  $y_1, y_2 \in Y, y_1 \mathcal{R} y_2 \implies y_1 = y_2$ ; and  $\mathcal{R}$ -singleton when, in addition,  $Y$  is nonempty. Then, let us say that

(fix-R-asing)  $T$  is *fix- $\mathcal{R}$ -asingleton*, if  $\text{Fix}(T)$  is  $\mathcal{R}$ -asingleton

(fix-R-sing)  $T$  is *fix- $\mathcal{R}$ -singleton*, in case  $\text{Fix}(T)$  is  $\mathcal{R}$ -singleton.

Likewise (cf. a previous convention), we say that

(fix-asing)  $T$  is *fix-asingleton*, if  $\text{Fix}(T)$  is asingleton

(fix-sing)  $T$  is *fix-singleton*, in case  $\text{Fix}(T)$  is singleton.

As a completion of these, we list the sufficient conditions to be used for getting such properties.

**(reg-1)** We say that  $X$  is *(Ba-o,d)-complete*, provided: for any (Ba-o) sequence  $Y_0 = (y_n; n \geq 0)$ , one has:  $(y_n)$  is  $d$ -Cauchy implies  $(y_n)$  is  $d$ -convergent

**(reg-2)** We say that  $T$  is *(Ba-o,d)-continuous*, if: for any (Ba-o) sequence  $Y_0 = (y_n; n \geq 0)$ , one has:  $y_n \xrightarrow{d} z$  implies  $Ty_n \xrightarrow{d} Tz$ .

**4-III)** To solve our posed problem along the precise directions, a lot of convergence type requirements is needed for the six-dimensional (geometric) contractive conditions we just introduced; these are strongly connected with the related developments in Turinici [40]. Letting  $X_0 = (x_n)$  be a (Ba-o) sequence in  $X$ , denote

$$[X_0] = \{x_n; n \geq 0\}, \quad [[X_0]] = \text{cl}([X_0]).$$

**(I)** The first condition upon our data is of *asymptotic* type. Two variants of it are of interest.

**I-a)** Call  $\Upsilon$ , *asymptotic* on  $[X_0]$  when

(asy) for each sequence  $(r_n)$  in  $R_+^0$  and each sequence  $(p_n)$  in  $R_+$  with  $((r_n, r_n, p_n, 0, r_{n+1}, r_{n+1}) \in \mathcal{Q}(B_1 > 0; [X_0]) \cap \mathcal{Y}$  and  $|p_n - r_n| \leq r_{n+1}, \forall n)$ , we have  $r_n \rightarrow 0$ ; hence,  $p_n \rightarrow 0$  as well.

**I-b)** Call  $\mathcal{Y}$ , *descending asymptotic* on  $[X_0]$  when

(desc-asy) for each sequence  $(r_n)$  in  $R_+^0$  and each sequence  $(p_n)$  in  $R_+$  with  $((r_n, r_n, p_n, 0, r_{n+1}, r_{n+1}) \in \mathcal{Q}(B_1 > 0; [X_0]) \cap \mathcal{Y}$  and  $|p_n - r_n| \leq r_{n+1}, \forall n)$ , we have that  $(r_n)$ =strictly descending and  $r_n \rightarrow 0$ ; hence (via  $(0 < r_n - r_{n+1} \leq p_n \leq r_n + r_{n+1}, \forall n)$ ),  $(p_n)$  is a sequence in  $R_+^0$  with (in addition)  $p_n \rightarrow 0$ .

Clearly, the inclusion below holds

$\mathcal{Y}$  is descending asymptotic on  $[X_0]$  implies  $\mathcal{Y}$  is asymptotic on  $[X_0]$ .

However, as we will see, the immense majority of contractive conditions is based on the descending asymptotic hypothesis; and not on the asymptotic one.

**(II)** The second condition to be considered is related to *right* properties. Take some point  $c = (c_1, \dots, c_6)$  in  $R_+^6$ . We say that the sequence  $(t^n := (t_1^n, \dots, t_6^n); n \geq 0)$  in  $R_+^6$  is *right* at  $c$ , if

$$(r-c) \quad (t_i^n \rightarrow c_i, \text{ as } n \rightarrow \infty, \forall i) \text{ and } (t_i^n \rightarrow c_i+, \text{ as } n \rightarrow \infty, \text{ when } c_i > 0).$$

Given  $b > 0$ , let us say that  $\mathcal{Y}$  is *nright* at  $b$  on  $[X_0]$ , if

(nright) for each sequence  $(t^n; n \geq 0)$  in  $\mathcal{Q}(B_3 > 0; [X_0]) \cap \mathcal{Y}$ , the right property at  $(0, b, b, b, b, 0)$  is not true.

The class of all these  $b > 0$  will be denoted as  $\text{nright}(\mathcal{Y}; [X_0])$ . In this case, define

(a-n-r)  $\mathcal{Y}$  is *almost nright* on  $[X_0]$ , if  $\Theta := \text{nright}(\mathcal{Y}; [X_0])$  is  $(>)$ -cofinal in  $R_+^0$  (for each  $\varepsilon \in R_+^0$  there exists  $\theta \in \Theta$  with  $\varepsilon > \theta$ )

(n-r)  $\mathcal{Y}$  is *nright* on  $[X_0]$ , if  $\Theta := \text{nright}(\mathcal{Y}; [X_0])$  is identical with  $R_+^0$ .

**(III)** The third condition involves *point* properties. Take some point  $c = (c_1, \dots, c_6)$  in  $R_+^6$ . We say that the sequence  $(t^n := (t_1^n, \dots, t_6^n); n \geq 0)$  in  $R_+^6$  is *point* at  $c$ , if

$$(p-c) \quad [t_i^n \rightarrow c_i \text{ as } n \rightarrow \infty, \forall i], \text{ and } [t_6^n = c_6, \forall n].$$

Given  $b > 0$ , let us say that  $\mathcal{Y}$  is *npoint* at  $b$  on  $[[X_0]]$ , if

(npoint) for each sequence  $(t^n; n \geq 0)$  in  $\mathcal{Q}(B_3 > 0; [[X_0]]) \cap \mathcal{Y}$ , the point property at  $(0, 0, b, 0, b, b)$  is not true.

The class of all these  $b > 0$  will be denoted as  $\text{npoint}(\mathcal{Y}; [[X_0]])$ . Then, define

(a-n-p)  $\mathcal{Y}$  is *almost npoint* on  $[[X_0]]$ , if  $\Theta := \text{npoint}(\mathcal{Y}; [[X_0]])$  is  $(>)$ -cofinal in  $R_+^0$

(n-p)  $\mathcal{Y}$  is *npoint* on  $[[X_0]]$ , if  $\Theta := \text{npoint}(\mathcal{Y}; [[X_0]])$  is identical with  $R_+^0$ .

(IV) The fourth condition to be posed is *normality*. Given  $a > 0$ , let us say that

- (norm)  $\Upsilon$  is *normal* at  $a$ , when  $(0, a, a, a, a, 0) \in \Upsilon$
- (n-norm)  $\Upsilon$  is *nnormal* at  $a$ , when  $(0, a, a, a, a, 0) \notin \Upsilon$ .

The class of all nnormal  $a > 0$  will be denoted as  $\text{nnorm}(\Upsilon)$ . Then, define

- (a-n-n)  $\Upsilon$  is *almost nnormal*, if  $\Theta := \text{nnorm}(\Upsilon)$  is  $(>)$ -cofinal in  $\mathbb{R}_+^0$
- (n-n)  $\Upsilon$  is *nnormal*, if  $\Theta := \text{nnorm}(\Upsilon)$  is identical with  $\mathbb{R}_+^0$ .

Some concrete examples of such objects will be given a bit further.

## 5 Main Result

Let  $(X, d)$  be a metric space. Further, let  $T \in \mathcal{F}(X)$  be a selfmap of  $X$ . As precise, we are interested to determine sufficient conditions for (uniqueness and) existence of elements in  $\text{Fix}(T)$ , via contractive type requirements involving iterative processes  $X_0 = (x_n)$  starting from an element  $x_0$  of  $X$ , and their associated sets

$$[X_0] = \{x_n; n \geq 0\}, [[X_0]] = \text{cl}([X_0]).$$

The basic directions and regularity conditions (relative to  $X_0$ ) under which the problem of determining fixed points of  $T$  is to be solved were already listed; and the contractive type framework (involving  $X_0$ ) was settled.

We are now in position to state our main result in this exposition.

**Theorem 4** *Let the subset  $\Upsilon \in \exp(\mathbb{R}_+^6)$  and its attached relation  $\mathcal{R}$  be such that  $(T$  is  $(d, \mathcal{R}; \Upsilon)$ -contractive, and)*

- (51-i)  $\mathcal{R}$  is  $B_1$ -admissible:  $(B_1 > 0) \subseteq \mathcal{R}$
- (51-ii)  $\Upsilon$  is asymptotic and almost nright on  $[Y_0]$ , for each  $(Ba-o)$  sequence  $Y_0 = (y_n)$ .

*Further, assume that  $X$  is  $(Ba-o, d)$ -complete; and let  $X_0 = (x_n)$  be a  $(Ba-o)$  sequence. Then,*

- (51-a)  $X_0 = (x_n)$  is Picard (modulo  $(d, \mathcal{R}; T)$ )
- (51-b)  $X_0 = (x_n)$  is strongly Picard (modulo  $(d, \mathcal{R}; T)$ ), whenever  $T$  is  $(Ba-o, d)$ -continuous
- (51-c)  $X_0 = (x_n)$  is strongly Picard (modulo  $(d, \mathcal{R}; T)$ ), whenever  $\Upsilon$  is npoint on  $[[X_0]]$
- (51-d)  $T$  is fix- $\mathcal{R}$ -asingleton (hence, fix- $\mathcal{R}$ -singleton, under any of these extra requirements) when, in addition to the conditions above,  $\Upsilon$  is nnormal  $((0, a, a, a, a, 0) \notin \Upsilon, \forall a > 0)$
- (51-e)  $T$  is fix-asingleton (hence, fix-singleton, under any of these extra requirements) when, in addition to the conditions above,  $\mathcal{R}$  is  $B_0$ -admissible and  $\Upsilon$  is nnormal.

**Proof** There are some steps to be passed.

**Step 1.** Denote for simplicity

$$(r_n = d(x_n, x_{n+1}); n \geq 0), (p_n := d(x_n, x_{n+2}); n \geq 0).$$

From the triangular inequality,

$$(rela-1) (\forall n) : |p_n - r_n| \leq r_{n+1}.$$

Moreover, we have by definition

$$(\forall n) : B_1(x_n, x_{n+1}) = \min\{r_n, r_{n+1}\} > 0; \text{ whence, } (x_n, x_{n+1}) \in \mathcal{R};$$

if we remember that  $\mathcal{R}$  is  $B_1$ -admissible. The contractive condition is applicable to the couples  $(x_n, x_{n+1})$ , for all  $n$ ; and gives (under these notations)

$$(rela-2) (r_n, r_n, p_n, 0, r_{n+1}, r_{n+1}) \in \mathcal{Q}(B_1 > 0; [X_0]) \subseteq \mathcal{Y}, \text{ for all } n.$$

By (rela-1)+(rela-2) (and the asymptotic property of  $\mathcal{Y}$ ), one derives

$r_n \rightarrow 0$  as  $n \rightarrow \infty$ ; so that,  
 $X_0 = (x_n)$  is semi-Picard (modulo  $(d, \mathcal{R}; T)$ ).

**Step 2.** Summing up,  $X_0 = (x_n)$  is (Ba-o) and  $d$ -asymptotic. On the other hand, as  $\mathcal{Y}$  is almost nright,  $\Theta := \text{nright}(\mathcal{Y}; [X_0])$  appears as  $(>)$ -cofinal in  $R_+^0$ . We show that, under these conditions,

$X_0 = (x_n)$  is  $d$ -Cauchy; hence,  
 $X_0 = (x_n)$  is Picard (modulo  $(d, \mathcal{R}; T)$ ).

Suppose that this is not true; and fix some index  $h \geq 1$ . By a preliminary statement, there exist a number  $\gamma \in \Theta$ , a rank  $J := J(\gamma, h) \geq 1$ , and a couple of rank-sequences  $(m(k); k \geq 0), (n(k); k \geq 0)$ , with

- (prop-1)  $J + k \leq m(k) < m(k) + 3h < n(k), \forall k$
- (prop-2)  $\forall s, t \in N[0, 3h]: V_k(s, t) := d(x_{m(k)+s}, x_{n(k)+t}) \rightarrow \gamma +$  as  $k \rightarrow \infty$ .

By the (Ba-o) and  $d$ -asymptotic properties of  $X_0 = (x_n)$ ,

- (prop-3)  $(t_1^k := r_{m(k)}; k \geq 0)$  and  $(t_6^k := r_{n(k)}; k \geq 0)$  are sequences in  $R_+^0$  with  $t_1^k, t_6^k \rightarrow 0$  as  $k \rightarrow \infty$ .

Moreover, taking (prop-2) into account, yields

$$(t_2^k := V_k(0, 0); k \geq 0), (t_3^k := V_k(0, 1); k \geq 0), \text{ and}$$

$$(t_4^k := V_k(1, 0); k \geq 0), (t_5^k := V_k(1, 1); k \geq 0)$$

are sequences in  $R_+^0$  with  $t_2^k, t_3^k, t_4^k, t_5^k \rightarrow \gamma +$  as  $k \rightarrow \infty$ ;

hence, putting these together,

(t-right) the vectorial sequence  $(t^k := (t_1^k, t_2^k, t_3^k, t_4^k, t_5^k, t_6^k); k \geq 0)$

is right at  $(0, \gamma, \gamma, \gamma, \gamma, 0)$ .

Concerning the contractive property of the same object  $(t^k; k \geq 0)$ , note that by (prop-2) and (prop-3),

$(\forall k) : B_3(x_{m(k)}, x_{n(k)}) = \min\{r_{m(k)}, V_k(0, 0), V_k(0, 1), V_k(1, 0), V_k(1, 1), r_{n(k)}\} > 0;$

whence,  $t^k \in \mathcal{Q}(B_3 > 0; [X_0]) \subseteq \mathcal{Y}$ .

This, via (t-right), contradicts the choice of  $\gamma$  as element of  $\Theta := \text{nrigh}(\mathcal{Y}; [X_0])$ . Hence, our working assumption is not acceptable; and the assertion follows.

**Step 3.** From these developments, we have, as  $X$  is (Ba-o,d)-complete

$X_0 = (x_n)$  is  $d$ -convergent:  $x_n \xrightarrow{d} z_0$  as  $n \rightarrow \infty$ , for some  $z_0 \in X$ .

We now claim that  $z_0$  is a fixed point of  $T$ . Two possible cases—treated in the steps below—are to be discussed.

**Step 4.** Suppose that  $T$  is (Ba-o,d)-continuous. Then,

$u_n := Tx_n \xrightarrow{d} Tz_0$  as  $n \rightarrow \infty$ .

On the other hand,  $(u_n = x_{n+1}; n \geq 0)$  is a subsequence of  $(x_n; n \geq 0)$ ; whence

$u_n \xrightarrow{d} z_0$  as  $n \rightarrow \infty$ .

Combining with  $d$ -separated, yields  $z_0 = Tz_0$ .

**Step 5.** Suppose that  $\mathcal{Y}$  is npoint on  $[[X_0]]$ . Three alternatives occur.

**Alter-1)** Suppose that

(Tz-rela-1)  $H_1 := \{n \in N; x_n = Tz_0\}$  is unbounded (in  $N$ ).

By a direct procedure (avoiding any use of (AC)) there may be obtained a strictly ascending sequence of ranks  $(i(n); n \geq 0)$ , such that

$a_n := x_{i(n)} = Tz_0$ , for all  $n$ .

But,  $(a_n; n \geq 0)$  is a subsequence of  $(x_n; n \geq 0)$ ; so that  $\lim_n(a_n) = z_0$ . Passing to limit in the relation above, gives  $z_0 = Tz_0$ .

**Alter-2)** Suppose that

(Tz-rela-2)  $H_2 := \{n \in N; Tx_n = Tz_0\}$  is unbounded (in  $N$ ).

By the same procedure (avoiding any use of (AC)) there may be obtained a strictly ascending sequence of ranks  $(j(n); n \geq 0)$ , such that

$$b_n := Tx_{j(n)} = Tz_0, \text{ for all } n.$$

But,  $(b_n = x_{j(n)+1}; n \geq 0)$  is a subsequence of  $(x_n; n \geq 0)$ ; so that  $\lim_n(b_n) = z_0$ . Passing to limit in the relation above, gives  $z_0 = Tz_0$ .

**Alter-3)** Suppose that

both subsets  $H_1$  and  $H_2$  are bounded (in  $N$ ).

This tells us that

$\exists i = i(z_0) \in N$ , such that:

$n \geq i$  implies  $Tx_n \neq Tz_0$  (hence,  $x_n \neq z_0$ ) and  $x_n \neq Tz_0$ .

Denote for simplicity  $(u_n = x_{n+i}; n \geq 0)$ ; clearly, by the preceding relation,

(non-id)  $(\forall n) : Tu_n \neq Tz_0$  (hence,  $u_n \neq z_0$ ) and  $u_n \neq Tz_0$ .

Again combining with the (Ba-o) property of  $X_0 = (x_n)$ , one derives

(posi-1)  $(\forall n) : Q_1(u_n, z_0) = d(u_n, Tu_n) > 0$ ,

$Q_2(u_n, z_0) = d(u_n, z_0) > 0$ ,  $Q_3(u_n, z_0) = d(u_n, Tz_0) > 0$ ,

$Q_4(u_n, z_0) = d(Tu_n, z_0) > 0$ ,  $Q_5(u_n, z_0) = d(Tu_n, Tz_0) > 0$ .

Suppose by contradiction that

(posi-2)  $b := d(z_0, Tz_0) > 0$  [whence,  $Q_6(x_n, z_0) = b > 0, \forall n$ ].

We show that this is not compatible with  $\gamma$  being npoint at  $b$ .

From the preceding observations, we have

$(\forall n) : B_3(u_n, z_0) = \min\{Q_1, \dots, Q_6\}(u_n, z_0) > 0$ ;

whence,  $\mathcal{Q}(u_n, z_0) \in \mathcal{Q}(B_3 > 0; [[X_0]]) \subseteq \gamma$ .

Let us evaluate the left part of this last relation. From the preceding facts

$(t_1^n := d(u_n, Tu_n); n \geq 0)$ ,  $(t_2^n := d(u_n, z_0); n \geq 0)$ ,

$(t_3^n := d(u_n, Tz_0); n \geq 0)$ ,  $(t_4^n := d(Tu_n, z_0); n \geq 0)$ ,

$$(t_5^n := d(Tu_n, Tz_0); n \geq 0), \text{ are sequences in } R_+^0 \text{ with}$$

$$(t_1^n, t_2^n, t_3^n, t_4^n, t_5^n) \rightarrow (0, 0, b, 0, b) \text{ as } n \rightarrow \infty.$$

At the same time,

$$(t_6^n := d(z, Tz) = b; n \geq 0) \text{ is a constant sequence (with } t_6^n \rightarrow b);$$

so, putting these together,

$$\text{(Rela-1) } (t^n := (t_1^n, \dots, t_6^n); n \geq 0) \text{ is point at } (0, 0, b, 0, b, b).$$

Finally, by simply replacing in the contractive condition above,

$$\text{(Rela-2) } (\forall n) : t^n \in \mathcal{Q}(B_3 > 0; [[X_0]]) \subseteq \mathcal{Y}.$$

This, via (Rela-1), contradicts the fact that  $b \in \text{npoint}(\mathcal{Y}; [[X_0]])$ . Hence, the assumption  $b > 0$  cannot be accepted; and then,  $b = 0$ ; that is:  $z_0 \in \text{Fix}(T)$ .

**Step 6.** Take the points  $z_1, z_2 \in \text{Fix}(T)$ , according to

$$z_1 \mathcal{R} z_2; \text{ and (by contradiction) } z_1 \neq z_2 \text{ (hence, } a := d(z_1, z_2) > 0).$$

As  $(z_1, z_2) \in \mathcal{R}$ , the contractive condition applies to  $(z_1, z_2)$ ; and gives

$$\mathcal{Q}(z_1, z_2) \in \mathcal{Y}; \text{ that is: } (0, a, a, a, a, 0) \in \mathcal{Y};$$

a contradiction with respect to the nnormal property of  $\mathcal{Y}$ . Hence, our working condition is not accepted; and the assertion follows.

**Step 7.** Take the points  $z_1, z_2 \in \text{Fix}(T)$ , according to

$$z_1 \neq z_2 \text{ (hence, } a := d(z_1, z_2) > 0).$$

This yields

$$B_0(z_1, z_2) = a > 0; \text{ whence, } (z_1, z_2) \in \mathcal{R};$$

if we note that  $\mathcal{R}$  is  $B_0$ -admissible. Consequently, the contractive condition applies to  $(z_1, z_2)$ ; and gives

$$\mathcal{Q}(z_1, z_2) \in \mathcal{Y}; \text{ that is: } (0, a, a, a, a, 0) \in \mathcal{Y};$$

a contradiction with respect to the nnormal property of  $\mathcal{Y}$ . Hence, our working condition is not accepted; and the assertion follows.

Note that, multivalued extensions of this result are possible, under the lines in Nadler [21]. On the other hand, an extended setting of these developments is possible, under the lines discussed in the 2001 PhD Thesis by Hitzler [13, Ch 1, Sect 1.2]; see also Pasicki [24]. Further aspects were delineated in Turinici [41].

## 6 Analytic Methods in RDC-Theory

As an application of the developments above, some analytic methods in the fixed point theory for Rhoades-Dutta-Choudhury contractions are being discussed.

Let  $(X, d)$  be a metric space. Further, let  $T \in \mathcal{F}(X)$  be a selfmap of  $X$ . As precise, we are interested to determine sufficient conditions for (uniqueness and) existence of elements in  $\text{Fix}(T)$ , via contractive type requirements involving iterative processes  $X_0 = (x_n)$  starting from an element  $x_0$  of  $X$ , and their associated sets

$$[X_0] = \{x_n; n \geq 0\}, [[X_0]] = \text{cl}([X_0]).$$

The basic directions and regularity conditions (relative to  $X_0$ ) under which the existence/uniqueness problem involving points of  $\text{Fix}(T)$  is to be solved were already listed; and the contractive type framework (involving  $X_0$ ) was settled. As a by-product of these, we established our main result in this exposition, Theorem 4. It is our aim in the following to show that, starting from this principle, it is possible to describe the basic lines for a kind of analytical approach in treating fixed points of contractions over metric spaces of the type introduced by Rhoades [31] and refined by Dutta and Choudhury [12].

(A) Roughly speaking, the approach to be considered requires an appropriate description of contractive conditions to be used. Denote, for  $x, y \in X$

$$\begin{aligned} Q_1(x, y) &= d(x, Tx), \quad Q_2(x, y) = d(x, y), \\ Q_3(x, y) &= d(x, Ty), \quad Q_4(x, y) = d(Tx, y), \\ Q_5(x, y) &= d(Tx, Ty), \quad Q_6(x, y) = d(y, Ty), \\ \mathcal{Q}(x, y) &= (Q_1(x, y), Q_2(x, y), Q_3(x, y), Q_4(x, y), Q_5(x, y), Q_6(x, y)). \end{aligned}$$

Then, let us construct the family of functions [for  $x, y \in X$ ]

$$\begin{aligned} B_0(x, y) &= \min\{Q_2, Q_5\}(x, y), \\ B_1(x, y) &= \min\{Q_1, Q_2, Q_5, Q_6\}(x, y), \\ B_2(x, y) &= \min\{Q_2, Q_3, Q_4, Q_5\}(x, y), \\ B_3(x, y) &= \min\{Q_1, Q_2, Q_3, Q_4, Q_5, Q_6\}(x, y), \end{aligned}$$



and put, for simplicity,

$$(B_i > 0) = \{(x, y) \in X \times X; B_i(x, y) > 0\}, i \in \{0, 1, 2, 3\}.$$

Let  $(\psi, \varphi)$  be a pair of functions over  $\mathcal{F}(R_+^0, R)$ . Further, let  $(u, v, w)$  be a triple of functions over  $\mathcal{F}(R_+^6, R_+)$ , and define the couple of subsets

$$\begin{aligned} \mathcal{T}_0 &= \{t \in R_+^6; u(t), v(t), w(t) > 0\}, \\ \mathcal{T} &= \{t \in \mathcal{T}_0; \psi(u(t)) \leq \psi(v(t)) - \varphi(w(t))\}. \end{aligned}$$

Finally, define the couple of relations over  $X$

$$R = Q^{-1}(\mathcal{T}_0), R^* = Q^{-1}(\mathcal{T}).$$

The obtained relations allow us defining two contractive conditions upon our data. The former of these is a functional contractive condition, written as:

$$\begin{aligned} \text{(fct-contr)} \quad T \text{ is } (d, \mathcal{R}; \psi, \varphi; u, v, w)\text{-contractive, if} \\ \psi(u \circ \mathcal{Q}(x, y)) \leq \psi(v \circ \mathcal{Q}(x, y)) - \varphi(w \circ \mathcal{Q}(x, y)), \text{ when } (x, y) \in \mathcal{R}. \end{aligned}$$

The latter of these appears as a set contractive condition and writes

$$\begin{aligned} \text{(set-contr)} \quad T \text{ is } (d, \mathcal{R}; \mathcal{T})\text{-contractive:} \\ \mathcal{Q}(x, y) \in \mathcal{T}, \text{ if } (x, y) \in \mathcal{R}; \text{ that is } \mathcal{Q}(\mathcal{R}) \subseteq \mathcal{T}. \end{aligned}$$

The connection between these is described in

**Proposition 5** *Under the above developments, one has*

*$T$  is  $(d, \mathcal{R}; \psi, \varphi; u, v, w)$ -contractive iff  $T$  is  $(d, \mathcal{R}; \mathcal{T})$ -contractive.*

**Proof** Two steps must be passed.

**Step 1.** Suppose that the functional contractive property (fct-contr) is holding; and let the couple  $(x, y) \in \mathcal{R}$  be arbitrary fixed. For the moment,

$$(x, y) \in \mathcal{R} \text{ implies } t := \mathcal{Q}(x, y) \in \mathcal{T}_0; \text{ that is: } u(t) > 0, v(t) > 0, w(t) > 0.$$

Moreover, by the very definition of contractive property,

$$\psi(u(t)) \leq \psi(v(t)) - \varphi(w(t)); \text{ that is: } t := \mathcal{Q}(x, y) \in \mathcal{T}.$$

This, by the arbitrariness of the couple  $(x, y) \in \mathcal{R}$ , shows that the set contractive property holds; and the left to right inclusion follows.

**Step 2.** Suppose that the set contractive property (set-contr) is holding; and let the couple  $(x, y) \in \mathcal{R}$  be arbitrary fixed. For the moment,

$$(x, y) \in \mathcal{R} \text{ implies } t := \mathcal{Q}(x, y) \in \mathcal{T}_0; \text{ that is: } u(t) > 0, v(t) > 0, w(t) > 0.$$

Moreover, by the very definition of  $\mathcal{Y}$ ,

$$\psi(u(t)) \leq \psi(v(t)) - \psi(w(t)); \text{ where } t := \mathcal{Q}(x, y).$$

This tells us that the couple  $(x, y)$  fulfills the functional contractive property; and proves the right to left inclusion.

In other words the functional contractive condition (fct-contr) is equivalent with the set-contractive condition (set-contr); to which the methods of the main result are applicable. Technically speaking, an effective application of the main result to the set contractive problem (set-contr) amounts to verifying the (sufficient) asymptotic, right, point, and normality properties upon  $\mathcal{Y}$ , in terms of its components; that is: the relation  $\mathcal{R}$ , the pair  $(\psi, \varphi)$ , and the triple  $(u, v, w)$ . Some preliminaries are needed.

For each function  $\psi \in \mathcal{F}(R_+^0, R)$ , and each  $a > 0$ , define

$\psi$  is right regulated at  $a$ :

$$\psi(a + 0) = \lim_{t \rightarrow a+} \psi(t) \text{ (the right limit of } \psi \text{ at } a \text{ exists (in } R \text{))}$$

$\psi$  is left regulated at  $a$ :

$$\psi(a - 0) = \lim_{t \rightarrow a-} \psi(t) \text{ (the left limit of } \psi \text{ at } a \text{ exists (in } R \text{))}$$

$\psi$  is regulated at  $a$ : both  $\psi(a + 0)$  and  $\psi(a - 0)$  exist (in  $R$ )

$$\psi \text{ is continuous at } a: \psi \text{ is regulated at } a \text{ and } \psi(a + 0) = \psi(a - 0) = \psi(a).$$

The class of all such  $a$  will be indicated, respectively, as

$\text{rreg}(\psi)$ = the right regulated domain of  $\psi$ ,

$\text{lreg}(\psi)$ = the left regulated domain of  $\psi$ ,

$\text{reg}(\psi)$ = the regulated domain of  $\psi$ ,

$\text{cont}(\psi)$ = the continuous domain of  $\psi$ .

In this case, we say that

(a-r-reg)  $\psi$  is *almost right regulated*, if  $\Theta := \text{rreg}(\psi)$  is ( $>$ )-cofinal in  $R_+^0$

(r-reg)  $\psi$  is *right regulated*, if  $\Theta := \text{rreg}(\psi)$  is identical with  $R_+^0$

(a-l-reg)  $\psi$  is *almost left regulated*, if  $\Theta := \text{lreg}(\psi)$  is ( $>$ )-cofinal in  $R_+^0$

(l-reg)  $\psi$  is *left regulated*, if  $\Theta := \text{lreg}(\psi)$  is identical with  $R_+^0$

(a-reg)  $\psi$  is *almost regulated*, if  $\Theta := \text{reg}(\psi)$  is ( $>$ )-cofinal in  $R_+^0$

(reg)  $\psi$  is *regulated*, if  $\Theta := \text{reg}(\psi)$  is identical with  $R_+^0$

(a-cont)  $\psi$  is *almost continuous*, if  $\Theta := \text{cont}(\psi)$  is ( $>$ )-cofinal in  $R_+^0$

(cont)  $\psi$  is *continuous*, if  $\Theta := \text{cont}(\psi)$  is identical with  $R_+^0$ .

*Remark 2* The class of almost right regulated functions is pretty large. This is shown in example below:

Let  $(a_n; n \geq 0)$  be a strictly descending sequence in  $R_+^0$  with

$$a_n \rightarrow 0 \text{ as } n \rightarrow \infty; \text{ whence: } \Theta := \{a_n; n \geq 0\} \text{ is } (>)\text{-cofinal in } R_+^0.$$

Then, let  $\psi : R_+^0 \rightarrow R$  be a function with

$$(r\text{-lim}) \psi(a_i + 0) \text{ exists, for all } i \in N.$$

In this case,  $\psi$  appears as an almost right regulated function. Note that, (r-lim) requires this behavior of  $\psi$  in the right neighborhoods of points in  $\Theta$ ; which means that, in the remaining points of  $R_+^0$ , the behavior of  $\psi$  is completely arbitrary.

In other words, the class of almost right regulated functions may contain many functions in  $\mathcal{F}(R_+^0, R)$  with an arbitrary behavior over large subsets of  $R_+^0$ .

*Remark 3* The class of almost continuous functions is pretty large too, in view of

(for each  $\psi \in \mathcal{F}(R_+^0, R)$ ):

$\psi$  is increasing or continuous implies  $\psi$  is almost continuous.

In fact, the affirmation concerning continuity is clear. On the other hand, the affirmation concerning increasing property is again clear, by a result in Natanson [23, Ch 8, Sect 1]. For, when  $\psi$  is increasing, there exists (by the quoted result) a denumerable part  $\Delta$  of  $R_+^0$ , with

$$\Theta := R_+^0 \setminus \Delta \subseteq \text{cont}(\psi);$$

and this along with  $\Theta$  being ( $>$ )-cofinal, proves the claim.

*Remark 4* Concerning the algebraic properties of such functions, the class of right regulated functions is invariant with respect to linear combinations and products:

$\psi_1, \psi_2 = \text{right regulated}$  imply  $\alpha_1\psi_1 + \alpha_2\psi_2$  (where  $\alpha_1, \alpha_2 \in R$ ) and  $\psi_1\psi_2$  are right regulated too.

This property is no longer true for the class of almost right regulated functions, in view of immediate fact

the ( $>$ )-cofinal in  $R_+^0$  property is not invariant to intersections: if  $\Theta_1$  and  $\Theta_2$  are ( $>$ )-cofinal in  $R_+^0$ , then  $\Theta_1 \cap \Theta_2$  need not be ( $>$ )-cofinal in  $R_+^0$ ;

just let  $\Theta_1$  and  $\Theta_2$  be the class of all rationals and irrationals of  $R_+^0$ , respectively. However, when one of these functions is continuous, this happens, in the sense

$\psi_1 = \text{almost right regulated}$  and  $\psi_2 = \text{continuous}$ , imply  $\alpha_1\psi_1 + \alpha_2\psi_2$  (where  $\alpha_1, \alpha_2 \in R$ ) and  $\psi_1\psi_2$  are almost right regulated.

Further aspects involving the concepts in question may be found in the 1960 monograph by Dieudonné [10, Ch VII, Sect 6].

Concerning these concepts, the following auxiliary fact is to be noted. Let  $\psi : R_+^0 \rightarrow R$  be a regulated function. Denote, for each  $b, c > 0$

$$(r\text{-osc-bc}) \quad \text{osc}(+)(\psi; b, c) = \psi(b + 0) - \psi(c + 0)$$

(the *right oscillation* of  $\psi$  at  $(b, c)$ )

$$(osc\text{-bc}) \quad \text{osc}(\psi; b, c) = \max\{\psi(u) - \psi(v); u \in \{b + 0, b, b - 0\}, v \in \{c + 0, c, c - 0\}\}$$

(the *oscillation* of  $\psi$  at  $(b, c)$ ).

Clearly, by the regulated hypothesis, these are finite real numbers. With an extra care, this definition may be extended to the case of  $b = 0$ . Precisely, assume that

(reg-ext)  $\psi$  is extended regulated:  $\psi$  is regulated, and  $\psi(0 + 0)$  exists.

A strong variant of this is

$\psi$  is strong extended regulated:

$\psi$  is regulated, and *zero abrupt*:  $\psi(0 + 0) = -\infty$ .

For example, the extended regulated property holds under  $\psi$ =increasing. However, the strong extended regulated property is not assured, in this general way; so, to get it, we must impose this property in a mandatory way.

Having these precise, suppose that  $\psi$  is extended regulated. We may then define for each  $c > 0$ ,

$$(r\text{-osc-0c}) \quad \text{osc}(+)(\psi; 0, c) = \psi(0 + 0) - \psi(c + 0); \text{ the } \textit{right oscillation} \text{ of } \psi \text{ at } (0, c)$$

$$(osc\text{-0c}) \quad \text{osc}(\psi; 0, c) = \max\{\psi(0 + 0) - \psi(v); v \in \{c + 0, c, c - 0\}\}; \text{ the } \textit{oscillation} \text{ of } \psi \text{ at } (0, c).$$

In the same general context, given the function  $\varphi : R_+^0 \rightarrow R$  and the point  $a \geq 0$ , let us introduce properties (where  $b \geq 0, c > 0$ ):

(r-bd-osc-bca)  $(\psi, \varphi)$  is *right bounded oscillating* at  $(b, c; a)$ : for each sequence  $(t_n)$  in  $R_+^0$  with  $t_n \rightarrow a+$ , we have  $\limsup_n(\varphi(t_n)) > \text{osc}(+)(\psi; b, c)$

(bd-osc-bca)  $(\psi, \varphi)$  is (*bilateral*) *bounded oscillating* at  $(b, c; a)$ : for each sequence  $(t_n)$  in  $R_+^0$  with  $t_n \rightarrow a$ , we have  $\limsup_n(\varphi(t_n)) > \text{osc}(\psi; b, c)$ .

**Proposition 6** *Let  $(\psi, \varphi)$  be a couple of functions over  $\mathcal{F}(R_+^0, R)$ , with*

*$\psi$  is extended regulated*

*and  $(a, b, c)$  be a triple of points with  $a, b \geq 0, c > 0$ . The following are valid*

**(62-1)** *If the couple  $(\psi, \varphi)$  is right bounded oscillating at  $(b, c; a)$ , then there are no sequences  $(t_n), (s_n)$  and  $(r_n)$  in  $R_+^0$ , such that*

$$t_n \rightarrow a+, s_n \rightarrow b+, r_n \rightarrow c+ \text{ and } \varphi(t_n) \leq \psi(s_n) - \psi(r_n), \text{ for all } n$$

**(62-2)** If the couple  $(\psi, \varphi)$  is bounded oscillating at  $(b, c; a)$  where  $b > 0$ , then there are no sequences  $(t_n), (s_n)$  and  $(r_n)$  in  $R_+^0$ , such that

$$t_n \rightarrow a, s_n \rightarrow b, r_n \rightarrow c \text{ and } \varphi(t_n) \leq \psi(s_n) - \psi(r_n), \text{ for all } n$$

**(62-3)** If the couple  $(\psi, \varphi)$  is bounded oscillating at  $(0, c; a)$ , then there are no sequences  $(t_n), (s_n)$  and  $(r_n)$  in  $R_+^0$ , such that

$$t_n \rightarrow a, s_n \rightarrow 0+, r_n \rightarrow c \text{ and } \varphi(t_n) \leq \psi(s_n) - \psi(r_n), \text{ for all } n.$$

**Proof** There are three parts to be passed.

**Part 1.** Suppose, by contradiction that there are sequences  $(t_n), (s_n)$  and  $(r_n)$  in  $R_+^0$ , such that

$$\text{(rela-1) } t_n \rightarrow a+, s_n \rightarrow b+, r_n \rightarrow c+ \text{ and } \varphi(t_n) \leq \psi(s_n) - \psi(r_n), \forall n.$$

Passing to  $\limsup$  as  $n \rightarrow \infty$  in this relation, gives

$$\limsup_n \varphi(t_n) \leq \psi(b + 0) - \psi(c + 0);$$

$$\text{whence (by the imposed conventions): } \limsup_n \varphi(t_n) \leq \text{osc}(+)(\psi; b, c);$$

absurd, by the oscillation type condition.

**Part 2.** Suppose, by contradiction that (under  $b > 0$ ) there are sequences  $(t_n), (s_n)$  and  $(r_n)$  in  $R_+^0$ , such that

$$\text{(rela-2) } t_n \rightarrow a, s_n \rightarrow b, r_n \rightarrow c \text{ and } \varphi(t_n) \leq \psi(s_n) - \psi(r_n), \text{ for all } n.$$

Let  $\mathcal{E} \in \{(>), (=), (<)\}$  be some relation. By definition, the unique limit of the sequence  $(\psi(c_n))$  when  $c_n \rightarrow c$  and  $(c_n \mathcal{E} c, \forall n)$  is denoted as  $\psi(c \mathcal{E})$ . Clearly,

$$\psi(c \mathcal{E}) = \psi(c + 0), \text{ when } \mathcal{E} \text{ is } (>),$$

$$\psi(c \mathcal{E}) = \psi(c), \text{ when } \mathcal{E} \text{ is } (=),$$

$$\psi(c \mathcal{E}) = \psi(c - 0), \text{ when } \mathcal{E} \text{ is } (<).$$

By the total property of the ordering in  $R$ , there exists a couple of relations  $\mathcal{E}_1, \mathcal{E}_2 \in \{(>), (=), (<)\}$  with

$$H_1 := \{n \in N; s_n \mathcal{E}_1 b\}, H_2 := \{n \in N; r_n \mathcal{E}_2 c\} \text{ are infinite.}$$

As  $H_1 := \{n \in N; s_n \mathcal{E}_1 b\}$  is infinite, there exists a strictly ascending sequence of ranks  $(i(n))$  such that  $s_{i(n)} \mathcal{E}_1 b, \forall n$ . Without loss—passing to a subsequence if necessary—one may assume  $(i(n) = n; n \geq 0)$ ; so that (by the above)

$$\text{(rela-3) } t_n \rightarrow a, s_n \rightarrow b, (s_n \mathcal{E}_1 b, \forall n), r_n \rightarrow c$$

$$\text{and } \varphi(t_n) \leq \psi(s_n) - \psi(r_n), \text{ for all } n.$$

As  $H_2 := \{n \in N; r_n \mathcal{E}_2 c\}$  is infinite, there exists a strictly ascending sequence of ranks  $(j(n))$  such that  $r_{j(n)} \mathcal{E}_2 c, \forall n$ . Without loss—passing to a subsequence if necessary—one may assume  $(j(n) = n; n \geq 0)$ ; so that (by the above)

$$\begin{aligned} & \text{(rela-4) } t_n \rightarrow a, s_n \rightarrow b, (s_n \mathcal{E}_1 b, \forall n), \\ & r_n \rightarrow c, (r_n \mathcal{E}_2 c, \forall n), \text{ and } \varphi(t_n) \leq \psi(s_n) - \psi(r_n), \text{ for all } n. \end{aligned}$$

Passing to  $\limsup$  as  $n \rightarrow \infty$  in this last relation, gives

$$\begin{aligned} & \limsup_n \varphi(t_n) \leq \psi(b \mathcal{E}_1) - \psi(c \mathcal{E}_2); \\ & \text{whence (by the imposed conventions): } \limsup_n \varphi(t_n) \leq \text{osc}(\psi; b, c); \end{aligned}$$

absurd, by the oscillation type condition.

**Part 3.** Suppose by contradiction that (under  $b = 0$ ) there are sequences  $(t_n), (s_n)$  and  $(r_n)$  in  $R_+^0$ , such that

$$\text{(rela-5) } t_n \rightarrow a, s_n \rightarrow 0+, r_n \rightarrow c \text{ and } \varphi(t_n) \leq \psi(s_n) - \psi(r_n), \text{ for all } n.$$

Let the preceding notations be in force. By the total property of the ordering in  $R$ , there exists a relation  $\mathcal{E} \in \{(>), (=), (<)\}$  with

$$H := \{n \in N; r_n \mathcal{E} c\} \text{ is infinite.}$$

There exists then a strictly ascending sequence of ranks  $(i(n))$  such that  $r_{i(n)} \mathcal{E} c, \forall n$ . Without loss—passing to a subsequence if necessary—one may assume  $(i(n) = n; n \geq 0)$ ; so that (by the above)

$$\begin{aligned} & \text{(rela-6) } t_n \rightarrow a, s_n \rightarrow 0+, r_n \rightarrow c (r_n \mathcal{E} c, \forall n), \\ & \text{and } \varphi(t_n) \leq \psi(s_n) - \psi(r_n), \text{ for all } n. \end{aligned}$$

Passing to  $\limsup$  as  $n \rightarrow \infty$  in this last relation, gives

$$\begin{aligned} & \limsup_n \varphi(t_n) \leq \psi(0 + 0) - \psi(c \mathcal{E}); \\ & \text{whence (by the imposed conventions): } \limsup_n \varphi(t_n) \leq \text{osc}(\psi; 0, c); \end{aligned}$$

absurd, by the oscillation type condition.

The proof is complete.

**(B)** We may now pass to the specific objective of our developments: to determine the (sufficient) asymptotic, right, point, and normality properties upon  $\Upsilon$ , in terms of its components; that is: the pair  $(\psi, \varphi)$  and the triple  $(u, v, w)$ . Two more basic conventions are in order.

Let  $(\xi, \eta, \zeta)$  be a triple of functions with  $\xi \in \mathcal{F}(R_+^0), \eta, \zeta \in \mathcal{F}(R_+)$ ; in particular, one may take  $\xi = I$  (the identity function  $(I(t) = t; t \in R_+^0)$ ); as well as  $\eta = 0$  and/or  $\zeta = 0$  (the identically zero function  $(0(t) = 0; t \in R_+)$ ).

**Conv-1)** We say that  $(u, v, w)$  is *right asymptotic of type*  $(\xi, \eta, \zeta)$ , when

(r-asy)  $\forall$  sequence  $(t^n; n \geq 0)$  in  $\mathcal{Q}(B_3 > 0; [X_0]) \cap \mathcal{Y}$ , and  $\forall b > 0: (t^n; n \geq 0)$  has the right property at  $(0, b, b, b, b, 0)$  implies  $u(t^n) \rightarrow \xi(b)+, v(t^n) \rightarrow \eta(b)+, w(t^n) \rightarrow \zeta(b)+.$

**Conv-2)** We say that  $(u, v, w)$  is *point asymptotic of type*  $(\xi, \eta, \zeta)$ , when

(pt-asy)  $\forall$  sequence  $(t^n; n \geq 0)$  in  $\mathcal{Q}(B_3 > 0; [[X_0]]) \cap \mathcal{Y}$ , and  $\forall b > 0 : (t^n; n \geq 0)$  has the point property at  $(0, 0, b, 0, b, b)$  implies  $u(t^n) \rightarrow \xi(b), v(t^n) \rightarrow \eta(b), w(t^n) \rightarrow \zeta(b).$

The following statement is an essential step towards the precise objective. Denote, for each couple  $\alpha, \beta \in R$ ,

$co(\alpha, \beta)$ =the convex cover of  $\{\alpha, \beta\}$ ;  
that is: the interval  $[\min\{\alpha, \beta\}, \max\{\alpha, \beta\}]$ .

**Theorem 5** Let  $(\psi, \varphi)$  be a couple of functions over  $\mathcal{F}(R_+^0, R)$ , with

(61-i)  $\varphi$  is strictly positive  $(\varphi(R_+^0) \subseteq R_+^0)$ .

Further, let  $(u, v, w)$  be a triple of functions over  $\mathcal{F}(R_+^6, R_+)$ , such that (under the precise notations)

(61-ii)  $\mathcal{R}$  is  $B_1$ -admissible:  $(B_1 > 0) \subseteq \mathcal{R}.$

Then,

**61-a)**  $\mathcal{Y}$  is descending asymptotic on  $[X_0]$ , provided

(61-a-i)  $\psi$  is right regulated,

and there exists a continuous strictly increasing  $\lambda \in \mathcal{F}(R_+^0)$ , with

(61-a-ii)  $(u, v, w)$  is strongly  $\lambda$ -iterative, in the sense: for each sequence  $(r_n)$  in  $R_+^0$  and each sequence  $(p_n)$  in  $R_+$  with  $(A_n := (r_n, r_n, p_n, 0, r_{n+1}, r_{n+1}) \in \mathcal{Q}(B_1 > 0; [X_0]) \cap \mathcal{Y}$ , and  $|p_n - r_n| \leq r_{n+1}$ , for all  $n$ ), we have  $u^n := u(A_n) = r_{n+1}, v^n := v(A_n) = \max\{r_n, r_{n+1}\}, w^n := w(A_n) \in co(\lambda(r_n), \lambda(r_{n+1}))$ ,  $\forall n$

(61-a-iii)  $(\psi, \varphi)$  is right bounded oscillating at  $(b, b; \lambda(b))$ , for each  $b > 0$

**61-b)**  $\mathcal{Y}$  is descending asymptotic on  $[X_0]$ , provided

(61-b-i)  $\psi$  is increasing (hence right regulated),

and there exists a continuous strictly increasing  $\lambda \in \mathcal{F}(R_+^0)$ , with

(61-b-ii)  $(u, v, w)$  is weakly  $\lambda$ -iterative, in the sense: for each sequence  $(r_n)$  in  $R_+^0$  and each sequence  $(p_n)$  in  $R_+$  with  $(A_n := (r_n, r_n, p_n, 0, r_{n+1}, r_{n+1}) \in \mathcal{Q}(B_1 > 0; [X_0]) \cap \mathcal{Y}$ , and  $|p_n - r_n| \leq r_{n+1}$ , for all  $n$ ), we have  $u^n := u(A_n) = r_{n+1}, v^n := v(A_n) \leq \max\{r_n, r_{n+1}\}, w^n := w(A_n) \in co(\lambda(r_n), \lambda(r_{n+1}))$ ,  $\forall n$

(61-b-iii)  $(\psi, \varphi)$  is right bounded oscillating at  $(b, b; \lambda(b))$ , for each  $b > 0$

**61-c)**  $\mathcal{Y}$  is almost nright on  $[X_0]$ , provided

(61-c-i)  $\psi$  is almost right regulated

and there exists  $\mu \in \mathcal{F}(R_+)$ , such that:

- (61-c-ii)  $(u, v, w)$  is right asymptotic of type  $(I, I, \mu)$
- (61-c-iii)  $(\psi, \varphi)$  is right bounded oscillating at  $(b, b; \mu(b))$ , for each  $b \in \text{rreg}(\psi)$
- 61-d**  $\mathcal{Y}$  is *npoint* (hence, almost *npoint*) on  $[[X_0]]$ , provided
- (61-d-i)  $\psi$  is extended regulated

and there exists a triple  $(\xi, \eta, \zeta)$  with  $\xi \in \mathcal{F}(R_+^0)$ ,  $\eta, \zeta \in \mathcal{F}(R_+)$ , such that:

- (61-d-ii)  $(u, v, w)$  is point asymptotic of type  $(\xi, \eta, \zeta)$
- (61-d-iii)  $(\psi, \varphi)$  is bounded oscillating at  $(\eta(b), \xi(b); \zeta(b))$ ,  $\forall b > 0$
- 61-e**  $\mathcal{Y}$  is *nnormal*, provided  $(u, v, w)$  is invariant, in the sense  $u(0, a, a, a, a, 0) = a, v(0, a, a, a, a, 0) = a, w(0, a, a, a, a, 0) = a, \forall a > 0$ .

**Proof** The argument consists in a number of parts.

**Part 1.** [ $\mathcal{Y}$  is descending asymptotic on  $[X_0]$  under the first lot of conditions].

Let the sequence  $(r_n)$  in  $R_+^0$  and the sequence  $(p_n)$  in  $R_+$  be such that

$$(\forall n) : A_n := (r_n, r_n, p_n, 0, r_{n+1}, r_{n+1}) \in \mathcal{Q}(B_1 > 0; [X_0]) \cap \mathcal{Y},$$

$$\text{and } |p_n - r_n| \leq r_{n+1}.$$

As  $(u, v, w)$  is strongly  $\lambda$ -iterative,

$$(\forall n) : u^n := u(A_n) = r_{n+1}, v^n := v(A_n) = \max\{r_n, r_{n+1}\},$$

$$w^n := w(A_n) \in \text{co}(\lambda(r_n), \lambda(r_{n+1})).$$

By the very representation of  $\mathcal{Y}$ , we must have

$$(\text{iter-1}) (\forall n) : \psi(r_{n+1}) \leq \psi(\max\{r_n, r_{n+1}\}) - \varphi(w^n),$$

where  $w^n \in \text{co}(\lambda(r_n), \lambda(r_{n+1}))$ .

If the alternative below holds

$$(\text{alter-1}) \quad r_n \leq r_{n+1}, \text{ for some } n$$

then, by the above relation

$$\varphi(w^n) \leq 0; \text{ in contradiction with } \varphi = \text{strictly positive}.$$

Hence, necessarily,

$$(\text{alter-2}) \quad r_n > r_{n+1}, \forall n; \text{ that is: } (r_n) \text{ is strictly descending.}$$

As a first consequence of this, one has that (iter-1) becomes

$$(\text{iter-2}) (\forall n) : \psi(r_{n+1}) \leq \psi(r_n) - \varphi(w^n), \text{ where } w^n \in [\lambda(r_{n+1}), \lambda(r_n)].$$



As a second consequence of this,  $r := \lim_n(r_n)$  exists in  $R_+$ . Suppose by contradiction that  $r > 0$ . As  $(r_n)$  is strictly descending,

$$r_n \rightarrow r+; \text{ hence } \lambda(r_n) \rightarrow \lambda(r)+, w^n \rightarrow \lambda(r)+$$

(cf. the choice of  $\lambda(\cdot)$  and  $(w^n)$ ). By (iter-2), we get

$$\text{(iter-3) } (0 <) \varphi(w^n) \leq \psi(r_n) - \psi(r_{n+1}), \forall n.$$

Passing to  $\limsup$  as  $n \rightarrow \infty$ , one gets

$$0 \leq \limsup_n \varphi(w^n) \leq \psi(r+0) - \psi(r+0) = 0; \text{ that is: } \lim_n \varphi(w^n) = 0;$$

in contradiction with  $(\psi, \varphi)$  being right bounded oscillating at  $(r, r; \lambda(r))$ . Hence,  $r = 0$ ; and the assertion follows.

**Part 2.** [ $\mathcal{Y}$  is descending asymptotic on  $[X_0]$  under the second lot of conditions].

Let the sequence  $(r_n)$  in  $R_+^0$  and the sequence  $(p_n)$  in  $R_+$  be such that

$$\begin{aligned} (\forall n) : A_n := (r_n, r_n, p_n, 0, r_{n+1}, r_{n+1}) \in \mathcal{Q}(B_1 > 0; [X_0]) \cap \mathcal{Y}, \\ \text{and } |p_n - r_n| \leq r_{n+1}. \end{aligned}$$

As  $(u, v, w)$  is  $\lambda$ -iterative,

$$\begin{aligned} (\forall n) : u^n := u(A_n) = r_{n+1}, v^n := v(A_n) \leq \max\{r_n, r_{n+1}\}, \\ w^n := w(A_n) \in \text{co}(\lambda(r_n), \lambda(r_{n+1})), \end{aligned}$$

By the very representation of  $\mathcal{Y}$  (and  $\psi$ =increasing)

$$\begin{aligned} \text{(iter-4) } (\forall n) : \psi(r_{n+1}) \leq \psi(v^n) - \varphi(w^n) \leq \psi(\max\{r_n, r_{n+1}\}) - \varphi(w^n), \\ \text{where } w^n \in \text{co}(\lambda(r_n), \lambda(r_{n+1})). \end{aligned}$$

By the strict positive condition upon  $\varphi$  and  $\psi$ =increasing, this yields

$$(\forall n) : \psi(r_{n+1}) < \psi(\max\{r_n, r_{n+1}\}); \text{ whence } r_{n+1} < r_n;$$

which tells us that  $(r_n)$  is strictly descending. As a first consequence of this, one has that (iter-4) becomes

$$\text{(iter-5) } (\forall n) : \psi(r_{n+1}) \leq \psi(r_n) - \varphi(w^n), \text{ where } w^n \in [\lambda(r_{n+1}), \lambda(r_n)].$$

As a second consequence of this,  $r := \lim_n(r_n)$  exists in  $R_+$ . Suppose by contradiction that  $r > 0$ . As  $(r_n)$  is strictly descending,

$$r_n \rightarrow r+; \text{ hence } \lambda(r_n) \rightarrow \lambda(r)+, w^n \rightarrow \lambda(r)+$$

(cf. the choice of  $\lambda(\cdot)$  and  $(w^n)$ ). By (iter-5), we get

$$(iter-6) \quad (0 <) \varphi(w^n) \leq \psi(r_n) - \psi(r_{n+1}), \forall n.$$

Passing to  $\limsup$  as  $n \rightarrow \infty$ , one gets

$$0 \leq \limsup_n \varphi(w^n) \leq \psi(r + 0) - \psi(r + 0) = 0; \text{ that is : } \lim_n \varphi(w^n) = 0;$$

in contradiction with  $(\psi, \varphi)$  being right bounded oscillating at  $(r, r; \lambda(r))$ . Hence,  $r = 0$ ; and the assertion follows.

**Part 3.** [ $\mathcal{Y}$  is almost nright on  $[X_0]$ , under the precise conditions].

Take some point  $c = (c_1, \dots, c_6)$  in  $R_+^6$ . We say that the sequence  $(t^n := (t_1^n, \dots, t_6^n); n \geq 0)$  in  $R_+^6$  is *right* at  $c$ , if

$$(r-c) \quad (t_i^n \rightarrow c_i, \forall i) \text{ and } (t_i^n \rightarrow c_i +, \text{ whenever } c_i > 0).$$

Given  $b > 0$ , let us say that  $\mathcal{Y}$  is *nright* at  $b$  on  $[X_0]$ , if

for each sequence  $(t^n; n \geq 0)$  in  $\mathcal{Q}(B_3 > 0; [X_0]) \cap \mathcal{Y}$ , the right property at  $(0, b, b, b, b, 0)$  is not true.

The class of all these  $b > 0$  will be denoted as  $\text{nright}(\mathcal{Y}; [X_0])$ . In this case, we say that  $\mathcal{Y}$  is

(a-n-r) *almost nright* on  $[X_0]$ , if  $\Theta := \text{nright}(\mathcal{Y}; [X_0])$  is  $(>)$ -cofinal in  $R_+^0$  (for each  $\varepsilon \in R_+^0$  there exists  $\theta \in \Theta$  with  $\varepsilon > \theta$ )

(n-r) *nright* on  $[X_0]$ , if  $\Theta := \text{nright}(\mathcal{Y}; [X_0])$  is identical with  $R_+^0$ .

We have to establish that the former property is retainable for our data. To do this, we start by noting that

$$\psi = \text{almost right regulated} \text{ implies } \Theta := \text{rreg}(\psi) \text{ is } (>)\text{-cofinal in } R_+^0.$$

We now claim that

$$\Theta \subseteq \text{nright}(\mathcal{Y}; [X_0]); \text{ wherefrom, } \text{nright}(\mathcal{Y}; [X_0]) \text{ is } (>)\text{-cofinal in } R_+^0.$$

This amounts to establish that

for each sequence  $(t^n := (t_1^n, \dots, t_6^n); n \geq 0)$  in  $\mathcal{Q}(B_3 > 0; [X_0]) \cap \mathcal{Y}$ , and each  $b \in \Theta$ , the right property at  $(0, b, b, b, b, 0)$  is not true.

Suppose by absurd that there exists a (vectorial) sequence  $(t^n := (t_1^n, \dots, t_6^n); n \geq 0)$  in  $R_+^6$  and some  $b \in \Theta$ , with

$$(right-b-1) \quad (\forall n): t^n := (t_1^n, \dots, t_6^n) \in \mathcal{Q}(B_3 > 0; [X_0]) \cap \mathcal{Y}; \text{ whence, } \psi(u(t^n)) \leq \psi(v(t^n)) - \varphi(w(t^n))$$

(right-b-2)  $(t^n := (t_1^n, \dots, t_6^n); n \geq 0)$  has the right property at  $(0, b, b, b, b, 0)$ ; i.e.:  $(t_i^n \rightarrow 0, \forall i \in \{1, 6\})$ , and  $(t_i^n \rightarrow b+, \forall i \in \{2, 3, 4, 5\})$ .

By the former of these properties

$$(\forall n) : (0 <) \varphi(w(t^n)) \leq \psi(v(t^n)) - \psi(u(t^n)).$$

As  $(u, v, w)$  is right asymptotic of type  $(I, I, \mu)$ ,

$$u(t^n) \rightarrow b+, v(t^n) \rightarrow b+, w(t^n) \rightarrow \mu(b) + .$$

Passing to  $\limsup$  as  $n \rightarrow \infty$  in this relation gives (by  $\varphi$ =strictly positive and  $\psi$ =right regular at  $b$ )

$$0 \leq \limsup_n \varphi(w(t^n)) \leq \psi(b + 0) - \psi(b + 0) = 0; \text{ that is: } \lim_n \varphi(w(t^n)) = 0;$$

in contradiction with  $(\psi, \varphi)$ =right bounded oscillating at  $(b, b; \mu(b))$ . Hence, our working assumption cannot be true; and the assertion follows.

**Part 4.**  $[\mathcal{Y}$  is npoint on  $[[X_0]]$  under the described assumptions].

Take some point  $c = (c_1, \dots, c_6)$  in  $R_+^6$ . We say that the sequence  $(t^n := (t_1^n, \dots, t_6^n); n \geq 0)$  in  $R_+^6$  is point at  $c$ , if

$$(\text{pt-c}) (t_i^n \rightarrow c_i, \forall i) \text{ and } [t_6^n = c_6, \forall n].$$

Given  $b > 0$ , let us say that  $\mathcal{Y}$  is npoint at  $b$  on  $[[X_0]]$ , if

(npoint) for each sequence  $(t^n; n \geq 0)$  in  $\mathcal{Q}(B_3 > 0; [[X_0]]) \cap \mathcal{Y}$ , the 6-point property at  $(0, 0, b, 0, b, b)$  is not true.

The class of all these  $b > 0$  will be denoted as  $\text{npoint}(\mathcal{Y}; [[X_0]])$ . In this case, we say that  $\mathcal{Y}$  is

(a-n-p) almost npoint on  $[[X_0]]$ , if  $\Theta := \text{npoint}(\mathcal{Y}; [[X_0]])$  is  $(>)$ -cofinal in  $R_+^0$

(n-p) npoint on  $[[X_0]]$ , if  $\Theta := \text{npoint}(\mathcal{Y}; [[X_0]])$  is identical with  $R_+^0$ .

We have to establish that the latter of these properties holds; that is,

for each sequence  $(t^n; n \geq 0)$  in  $R_+^6$  and each  $b > 0$ , the point property at  $(0, 0, b, 0, b, b)$  is not true.

Assume by contradiction that this assertion is false: there exists a (vectorial) sequence  $(t^n = (t_1^n, \dots, t_6^n); n \geq 0)$  in  $R_+^6$  and some  $b > 0$ , with

(point-b-1)  $(\forall n): t^n = (t_1^n, \dots, t_6^n) \in \mathcal{Q}(B_3 > 0; [[X_0]]) \cap \mathcal{Y}$ ; whence,  $\psi(u(t^n)) \leq \psi(v(t^n)) - \varphi(w(t^n))$

(point-b-2)  $(t^n = (t_1^n, \dots, t_6^n); n \geq 0)$  is point at  $(0, 0, b, 0, b, b)$ ; that is:  $(t_i^n \rightarrow 0, \forall i \in \{1, 2, 4\})$ ,  $(t_i^n \rightarrow b, \forall i \in \{3, 5, 6\})$ , and  $(t_6^n = b, \forall n)$ .

By the former of these properties

$$(\forall n) : \varphi(w(t^n)) \leq \psi(v(t^n)) - \psi(u(t^n)).$$

As  $(u, v, w)$  is point asymptotic of type  $(\xi, \eta, \zeta)$

$$u(t^n) \rightarrow \xi(b), v(t^n) \rightarrow \eta(b), w(t^n) \rightarrow \zeta(b).$$

Taking an auxiliary fact into account, we have that the obtained relations are impossible via  $(\psi, \varphi)$ -bounded oscillating at  $(\eta(b), \xi(b); \zeta(b))$ . Consequently, the working assumption is not acceptable; and conclusion follows.

**Part 5.** [ $\mathcal{Y}$  is normal under the posed condition].

Remember that, this property means:

$$(0, a, a, a, a, 0) \in \mathcal{Y} \text{ is impossible, for each } a > 0.$$

Suppose by contradiction that

$$\text{there exists } a > 0 \text{ with } (0, a, a, a, a, 0) \in \mathcal{Y}.$$

In view of  $(u, v, w)$ -invariant,

$$u(0, a, a, a, a, 0) = a, v(0, a, a, a, a, 0) = a, w(0, a, a, a, a, 0) = a.$$

This firstly means

$$(0, a, a, a, a, 0) \in \mathcal{Y}_0 \text{ (because } a > 0).$$

Secondly, by the complete definition of  $\mathcal{Y}$ , we derive

$$\psi(u(0, a, a, a, a, 0)) \leq \psi(v(0, a, a, a, a, 0)) - \varphi(w(0, a, a, a, a, 0)).$$

So, again combining with  $(u, v, w)$ -invariant, we must have

$$\psi(a) \leq \psi(a) - \varphi(a); \text{ that is: } \varphi(a) \leq 0;$$

absurd, by  $\varphi$ -strictly positive; and our claim follows.

Now, by simply combining the obtained fact with Theorem 4, one gets the following couple of fixed point statements with a practical value. Let  $(\psi, \varphi)$  be a pair of functions over  $\mathcal{F}(R_+^0, R)$ . Further, let  $(u, v, w)$  be a triple of functions over  $\mathcal{F}(R_+^6, R_+)$ , and define the couple of subsets

$$\mathcal{Y}_0 = \{t \in R_+^6; u(t), v(t), w(t) > 0\},$$

$$\mathcal{Y} = \{t \in \mathcal{Y}_0; \psi(u(t)) \leq \psi(v(t)) - \varphi(w(t))\}.$$

Finally, define the couple of relations over  $X$

$$\mathcal{R}_0 = \mathcal{Q}^{-1}(\mathcal{Y}_0), \mathcal{R}^* = \mathcal{Q}^{-1}(\mathcal{Y}).$$

The contractive property to be used here is

(fct-contr)  $T$  is  $(d, \mathcal{R}; \psi, \varphi; u, v, w)$ -contractive, if

$$\psi(u \circ \mathcal{Q}(x, y)) \leq \psi(v \circ \mathcal{Q}(x, y)) - \varphi(w \circ \mathcal{Q}(x, y)), \text{ when } (x, y) \in \mathcal{R}.$$

On the other hand, by the constructions above, one may consider the attached set contractive condition

(set-contr)  $T$  is  $(d, \mathcal{R}; \mathcal{Y})$ -contractive:  $\mathcal{Q}(x, y) \in \mathcal{Y}$ , if  $(x, y) \in \mathcal{R}$ .

The connection between these conditions is expressed as

(equi)  $T$  is  $(d, \mathcal{R}; \psi, \varphi; u, v, w)$ -contractive iff  $T$  is  $(d, \mathcal{R}; \mathcal{Y})$ -contractive.

In other words: our initial contractive condition in terms of  $\mathcal{R}$ ,  $(\psi, \varphi)$  and  $(u, v, w)$  is equivalent with a contractive condition in terms of  $\mathcal{R}$  and  $\mathcal{Y}$ .

The former of these statements (referred to as: Rhoades-Dutta-Choudhury principle for regulated functions; in short: (RDC-reg)) is based on a regulated condition upon  $\psi$ ; note that, in this case,  $\psi$  need not be increasing.

**Theorem 6** *Let the selfmap  $T$  be  $(d, \mathcal{R}; \psi, \varphi; u, v, w)$ -contractive, where  $(\psi, \varphi)$  is a couple of functions over  $\mathcal{F}(R_+^0, R)$ , with*

(62-i)  $\psi$  is extended regulated and  $\varphi$  is strictly positive,

and  $(u, v, w)$  is a triple of functions over  $\mathcal{F}(R_+^6, R_+)$ , such that (under the notations we just proposed)

(62-ii)  $\mathcal{R}$  is  $B_1$ -admissible:  $(B_1 > 0) \subseteq \mathcal{R}$ .

Further, let  $X$  be  $(Ba-o, d)$ -complete; and let  $X_0 = (x_n)$  be a  $(Ba-o)$  sequence.

Then,

**(62-a)**  $X_0 = (x_n)$  is semi-Picard (modulo  $(d, \mathcal{R}; T)$ ), provided there exists a continuous strictly increasing  $\lambda \in \mathcal{F}(R_+^0)$ , with

(62-a-i)  $(u, v, w)$  is strongly  $\lambda$ -iterative

(62-a-ii)  $(\psi, \varphi)$  is right bounded oscillating at  $(b, b; \lambda(b))$ , for each  $b > 0$

**(62-b)**  $X_0 = (x_n)$  is Picard (modulo  $(d, \mathcal{R}; T)$ ), provided (in addition to the above) there exists  $\mu \in \mathcal{F}(R_+)$ , such that

(62-b-i) the associated triple  $(u, v, w)$  is right asymptotic of type  $(I, I, \mu)$

(62-b-ii)  $(\psi, \varphi)$  is right bounded oscillating at  $(b, b; \mu(b))$ , for each  $b > 0$

**(62-c)**  $X_0 = (x_n)$  is strongly Picard (modulo  $(d, \mathcal{R}; T)$ ), whenever  $T$  is  $(Ba-o, d)$ -continuous

**(62-d)**  $X_0 = (x_n)$  is strongly Picard (modulo  $(d, \mathcal{R}; T)$ ), whenever (in addition to the above) there exists a triple  $(\xi, \eta, \zeta)$  with  $\xi \in \mathcal{F}(R_+^0)$ ,  $\eta, \zeta \in \mathcal{F}(R_+)$ , such that

- (62-d-i) the triple  $(u, v, w)$  is point asymptotic of type  $(\xi, \eta, \zeta)$
- (62-d-ii)  $(\psi, \varphi)$  is bounded oscillating at  $(\eta(b), \xi(b); \zeta(b))$ , for each  $b > 0$
- (62-e)**  $T$  is fix- $\mathcal{R}$ -asingleton (hence, fix- $\mathcal{R}$ -singleton, under any of these extra requirements) when, in addition to the conditions above, the triple  $(u, v, w)$  is invariant

$$(u(0, a, a, a, a, 0) = v(0, a, a, a, a, 0) = w(0, a, a, a, a, 0) = a, \forall a > 0)$$

- (62-f)**  $T$  is fix-asingleton (hence, fix-singleton, under any of these extra requirements) when, in addition to the conditions above,

$\mathcal{R}$  is  $B_0$ -admissible and the triple  $(u, v, w)$  is invariant.

The latter of these (referred to as: Rhoades-Dutta-Choudhury principle for increasing functions; in short: (RDC-incr)) is based on an increasing condition upon the ambient function  $\psi$ .

**Theorem 7** Let the selfmap  $T$  be  $(d, \mathcal{R}; \psi, \varphi; u, v, w)$ -contractive, where  $(\psi, \varphi)$  is a couple of functions over  $\mathcal{F}(R_+^0, R)$ , with

- (63-i)  $\psi$  is increasing and  $\varphi$  is strictly positive
- and  $(u, v, w)$  is a triple of functions over  $\mathcal{F}(R_+^6, R_+)$ , such that (under the notations we just proposed)

- (63-ii)  $\mathcal{R}$  is  $B_1$ -admissible:  $(B_1 > 0) \subseteq \mathcal{R}$ .

Further, let  $X$  be  $(Ba-o, d)$ -complete; and let  $X_0 = (x_n)$  be a  $(Ba-o)$  sequence.

Then,

- (63-a)**  $X_0 = (x_n)$  is semi-Picard (modulo  $(d, \mathcal{R}; T)$ ), provided there exists a continuous strictly increasing  $\lambda \in \mathcal{F}(R_+^0)$ , with

- (63-a-i)  $(u, v, w)$  is weakly  $\lambda$ -iterative
- (63-a-ii)  $(\psi, \varphi)$  is right bounded oscillating at  $(b, b; \lambda(b))$ , for each  $b > 0$

- (63-b)**  $X_0 = (x_n)$  is Picard (modulo  $(d, \mathcal{R}; T)$ ), provided (in addition to the above) there exists  $\mu \in \mathcal{F}(R_+)$ , such that

- (63-b-i) the associated triple  $(u, v, w)$  is right asymptotic of type  $(I, I, \mu)$
- (63-b-ii)  $(\psi, \varphi)$  is right bounded oscillating at  $(b, b; \mu(b))$ , for each  $b > 0$

- (63-c)**  $X_0 = (x_n)$  is strongly Picard (modulo  $(d, \mathcal{R}; T)$ ), whenever  $T$  is  $(Ba-o, d)$ -continuous

- (63-d)**  $X_0 = (x_n)$  is strongly Picard (modulo  $(d, \mathcal{R}; T)$ ), whenever (in addition to the above) there exists a triple  $(\xi, \eta, \zeta)$  with  $\xi \in \mathcal{F}(R_+^0)$ ,  $\eta, \zeta \in \mathcal{F}(R_+)$ , such that

- (63-d-i) the triple  $(u, v, w)$  is point asymptotic of type  $(\xi, \eta, \zeta)$
- (63-d-ii)  $(\psi, \varphi)$  is bounded oscillating at  $(\eta(b), \xi(b); \zeta(b))$ , for each  $b > 0$

- (63-e)**  $T$  is fix- $\mathcal{R}$ -asingleton (hence, fix- $\mathcal{R}$ -singleton, under any of these extra requirements) when, in addition to the conditions above, the triple  $(u, v, w)$  is invariant

$$(u(0, a, a, a, a, 0) = v(0, a, a, a, a, 0) = w(0, a, a, a, a, 0) = a, \forall a > 0)$$

**(63-f)**  $T$  is *fix-asingleton* (hence, *fix-singleton*, under any of these extra requirements) when, in addition to the conditions above,

$$\mathcal{R} \text{ is } B_0\text{-admissible and the triple } (u, v, w) \text{ is invariant.}$$

As already precise, an extended setting of these developments is possible, under the lines described in the 2001 PhD Thesis by Hitzler [13, Ch 1, Sect 1.2]; see also Turinici [41]. Further aspects will be discussed elsewhere.

## 7 Chandok-Choudhury Approach

Let  $(X, d)$  be a metric space. Further, let  $T \in \mathcal{F}(X)$  be a selfmap of  $X$ . As precise, we are interested to determine sufficient conditions for (uniqueness and) existence of elements in  $\text{Fix}(T)$ , via contractive type requirements involving iterative processes  $X_0 = (x_n)$  starting from an element  $x_0$  of  $X$ , and their associated sets

$$[X_0] = \{x_n; n \geq 0\}, [[X_0]] = \text{cl}([X_0]).$$

Let  $(\psi, \varphi)$  be a pair of functions over  $\mathcal{F}(R_+^0, R)$ . Further, let  $(u, v, w)$  be a triple of functions over  $\mathcal{F}(R_+^6, R_+)$ , and define the couple of subsets

$$\begin{aligned} \mathcal{Y}_0 &= \{t \in R_+^6; u(t), v(t), w(t) > 0\}, \\ \mathcal{Y} &= \{t \in \mathcal{Y}_0; \psi(u(t)) \leq \psi(v(t)) - \varphi(w(t))\}. \end{aligned}$$

Finally, define the couple of relations over  $X$

$$\mathcal{R} = \mathcal{Q}^{-1}(\mathcal{Y}_0), \mathcal{R}^* = \mathcal{Q}^{-1}(\mathcal{Y}).$$

The starting functional contractive property to be considered is

$$\begin{aligned} \text{(fct-contr)} \quad & T \text{ is } (d, \mathcal{R}; \psi, \varphi; u, v, w)\text{-contractive, if} \\ & \psi(u \circ \mathcal{Q}(x, y)) \leq \psi(v \circ \mathcal{Q}(x, y)) - \varphi(w \circ \mathcal{Q}(x, y)), \text{ when } (x, y) \in \mathcal{R}. \end{aligned}$$

On the other hand, the attached set contractive condition is

$$\text{(set-contr)} \quad T \text{ is } (d, \mathcal{R}; \mathcal{Y})\text{-contractive: } \mathcal{Q}(x, y) \in \mathcal{Y}, \text{ if } (x, y) \in \mathcal{R}.$$

As precise, the connection between these conditions writes

$$\text{(equi)} \quad T \text{ is } (d, \mathcal{R}; \psi, \varphi; u, v, w)\text{-contractive iff } T \text{ is } (d, \mathcal{R}; \mathcal{Y})\text{-contractive.}$$

For the set contractive condition, sufficient conditions were given upon  $\mathcal{R}$  and  $\mathcal{Y}$  so that the main result be applicable; these, naturally, are ultimately expressed in terms of  $\mathcal{R}$ ,  $(\psi, \varphi)$  and  $(u, v, w)$ . A by-product of these developments is the couple of fixed point results we just exposed. It is our aim in the following to discuss a lot of particular cases of these, with a practical meaning.

(A) For the moment, some direct constructions will be proposed.

**Constr-0)** Define the auxiliary system of functions  $(e_0, e_1, e_2, e_3, e_4)$  over the class  $\mathcal{F}(R_+^6, R_+)$ , as: for each  $t = (t_1, \dots, t_6) \in R_+^6$ ,

$$\begin{aligned} e_0(t) &= t_6(1 + t_4)/(1 + t_3), \\ e_1(t) &= t_6(1 + t_1)/(1 + t_2), e_2(t) = t_4(1 + t_3)/(1 + t_2), \\ e_3(t) &= \max\{t_1, t_6\}, e_4(t) = (1/2)(t_3 + t_4). \end{aligned}$$

**Constr-1)** Further, define the triple  $(u_1, u_2, u_3)$  over the class  $\mathcal{F}(R_+^6, R_+)$ , as: for each  $t = (t_1, \dots, t_6) \in R_+^6$ ,

$$(u\text{-def}) \quad u_1(t) = \max\{t_5, e_0(t)\}, u_2(t) = \max\{t_5, t_6\}, u_3(t) = t_5.$$

**Constr-2)** Then, let us construct the triple of functions  $(v_1, v_2, v_3)$  over the class  $\mathcal{F}(R_+^6, R_+)$ , as: for each  $t = (t_1, \dots, t_6) \in R_+^6$ ,

$$\begin{aligned} (v\text{-def}) \quad v_1(t) &= \max\{t_2, e_1(t), e_2(t)\}, v_2(t) = \max\{t_2, e_3(t), e_4(t)\}, \\ v_3(t) &= t_2. \end{aligned}$$

**Constr-3)** Finally, define the triple  $(w_1, w_2, w_3)$  of functions over  $\mathcal{F}(R_+^6, R_+)$ , as: for each  $t = (t_1, \dots, t_6) \in R_+^6$ ,

$$(w\text{-def}) \quad w_1(t) = \max\{t_2, e_1(t)\}, w_2(t) = \max\{t_2, t_6\}, w_3(t) = t_2.$$

Remember that, for the above explained technical reasons, we are forced to accept a regularity condition like

$$\mathcal{R} \text{ is } B_1\text{-admissible: } (B_1 > 0) \subseteq \mathcal{R}.$$

It is our aim to discuss of to what extent are the triples  $(u, v, w)$ , where

$$u \in \{u_1, u_2, u_3\}, v \in \{v_1, v_2, v_3\}, w \in \{w_1, w_2, w_3\}$$

compatible with the regularity conditions encountered in the couple of fixed point statements above. The response to this question is to be precise as below.

**R-1)** Given the sequence  $(r_n)$  in  $R_+^0$  and the sequence  $(p_n)$  in  $R_+$  with  $(A_n := (r_n, r_n, p_n, 0, r_{n+1}, r_{n+1}) \in \mathcal{Q}(B_1 > 0; [X_0]) \cap \mathcal{Y}, |p_n - r_n| \leq r_{n+1}, \text{ for all } n)$ , we have (by these definitions)



$$(\forall n) : e_0(A_n) = r_{n+1}/(1 + p_n) (\leq r_{n+1}), e_1(A_n) = r_{n+1}, e_2(A_n) = 0, \\ e_3(A_n) = \max\{r_n, r_{n+1}\}, e_4(A_n) = (1/2)p_n (\leq \max\{r_n, r_{n+1}\}).$$

This, along with our previous conventions, gives

$$\begin{aligned} \text{(u-iter)} \quad & u_1(A_n) = r_{n+1}, u_2(A_n) = r_{n+1}, u_3(A_n) = r_{n+1}, \\ \text{(v-iter)} \quad & v_1(A_n) = \max\{r_n, r_{n+1}\}, v_2(A_n) = \max\{r_n, r_{n+1}\}, v_3(A_n) = r_n, \\ \text{(w-iter)} \quad & w_1(A_n) = \max\{r_n, r_{n+1}\}, w_2(A_n) = \max\{r_n, r_{n+1}\}, w_3(A_n) = r_n. \end{aligned}$$

As a consequence of this, we have that

- (s-iter) any triple  $(u, v, w)$ , where  $u \in \{u_1, u_2, u_3\}$ ,  $v \in \{v_1, v_2\}$ ,  $w \in \{w_1, w_2, w_3\}$  is strongly  $I$ -iterative
- (iter) any triple  $(u, v_3, w)$ , where  $u \in \{u_1, u_2, u_3\}$ ,  $w \in \{w_1, w_2, w_3\}$  is weakly  $I$ -iterative.

**R-2)** Let the sequence  $(t^n = (t_1^n, \dots, t_6^n); n \geq 0)$  in  $\mathcal{Q}(B_3 > 0; [X_0]) \cap \mathcal{Y}$ , and the point  $b > 0$  be such that

$$\begin{aligned} (t^n = (t_1^n, \dots, t_6^n); n \geq 0) \text{ has the right property at } (0, b, b, b, b, 0); \\ \text{that is: } (t_i^n \rightarrow 0, \forall i \in \{1, 6\}), (t_j^n \rightarrow b+, \forall j \in \{2, 3, 4, 5\}). \end{aligned}$$

Then, by definition,

$$\begin{aligned} e_0(t^n) &= t_6^n(1 + t_4^n)/(1 + t_3^n) \rightarrow 0, \\ e_1(t^n) &= t_6^n(1 + t_1^n)/(1 + t_2^n) \rightarrow 0, \quad e_2(t^n) = t_4^n(1 + t_3^n)/(1 + t_2^n) \rightarrow b; \\ e_3(t^n) &= \max\{t_1^n, t_6^n\} \rightarrow 0, \quad e_4(t^n) = (1/2)[t_3^n + t_4^n] \rightarrow b. \end{aligned}$$

This, along with our preceding conventions, yields

$$\begin{aligned} \text{(u-pos)} \quad & u_1(t^n) = \max\{t_5^n, e_0(t^n)\} \rightarrow b+, \\ & u_2(t^n) = \max\{t_5^n, t_6^n\} \rightarrow b+, \quad u_3(t^n) = t_5^n \rightarrow b+, \\ \text{(v-pos)} \quad & v_1(t^n) = \max\{t_2^n, e_1(t^n), e_2(t^n)\} \rightarrow b+, \\ & v_2(t^n) = \max\{t_2^n, e_3(t^n), e_4(t^n)\} \rightarrow b+, \quad v_3(t^n) = t_2^n \rightarrow b+, \\ \text{(w-pos)} \quad & w_1(t^n) = \max\{t_2^n, e_1(t^n)\} \rightarrow b+, \\ & w_2(t^n) = \max\{t_2^n, t_6^n\} \rightarrow b+, \quad w_3(t) = t_2^n \rightarrow b+. \end{aligned}$$

Here, the convergence relations involving maximum type functions are obtained as

$$(u_1(t^n) \geq t_5^n, \forall n) \text{ and } u_1(t^n) \rightarrow b, t_5^n \rightarrow b+ \text{ imply } u_1(t^n) \rightarrow b+;$$

$$\begin{aligned}
&(u_2(t^n) \geq t_5^n, \forall n) \text{ and } u_2(t^n) \rightarrow b, t_5^n \rightarrow b + \text{ imply } u_2(t^n) \rightarrow b+; \\
&(v_1(t^n) \geq t_2^n, \forall n) \text{ and } v_1(t^n) \rightarrow b, t_2^n \rightarrow b + \text{ imply } v_1(t^n) \rightarrow b+; \\
&(v_2(t^n) \geq t_2^n, \forall n) \text{ and } v_2(t^n) \rightarrow b, t_2^n \rightarrow b + \text{ imply } v_2(t^n) \rightarrow b+; \\
&(w_1(t^n) \geq t_2^n, \forall n) \text{ and } w_1(t^n) \rightarrow b, t_2^n \rightarrow b + \text{ imply } w_1(t^n) \rightarrow b+; \\
&(w_2(t^n) \geq t_2^n, \forall n) \text{ and } w_2(t^n) \rightarrow b, t_2^n \rightarrow b + \text{ imply } w_2(t^n) \rightarrow b+.
\end{aligned}$$

As a consequence of this, we have that

(r-I-I) any triple  $(u, v, w)$ , where  $u \in \{u_1, u_2, u_3\}$ ,  $v \in \{v_1, v_2, v_3\}$ ,  $w \in \{w_1, w_2, w_3\}$  is right asymptotic of type  $(I, I, I)$ .

**R-3)** Let the sequence  $(t^n = (t_1^n, \dots, t_6^n); n \geq 0)$  in  $\mathcal{Q}(B_3 > 0; [[X_0]]) \cap \mathcal{Y}$ , and the point  $b > 0$  be such that

$$\begin{aligned}
&(t^n = (t_1^n, \dots, t_6^n); n \geq 0) \text{ has the point property at } (0, 0, b, 0, b, b); \\
&\text{that is: } (t_i^n \rightarrow 0, \forall i \in \{1, 2, 4\}), (t_j^n \rightarrow b, \forall j \in \{3, 5, 6\}), (t_6^n = b, \forall n).
\end{aligned}$$

Then, by definition,

$$\begin{aligned}
e_0(t^n) &= t_6^n(1 + t_4^n)/(1 + t_3^n) \rightarrow b/(1 + b), \\
e_1(t^n) &= t_6^n(1 + t_1^n)/(1 + t_2^n) \rightarrow b, \quad e_2(t^n) = t_4^n(1 + t_3^n)/(1 + t_2^n) \rightarrow 0, \\
e_3(t^n) &= \max\{t_1^n, t_6^n\} \rightarrow b, \quad e_4(t^n) = (1/2)[t_3^n + t_4^n] \rightarrow b/2.
\end{aligned}$$

This, along with our preceding conventions, yields

$$\begin{aligned}
&\text{(u-pt) } u_1(t^n) = \max\{t_5^n, e_0(t^n)\} \rightarrow b, \\
&u_2(t^n) = \max\{t_5^n, t_6^n\} \rightarrow b, \quad u_3(t^n) = t_5^n \rightarrow b, \\
&\text{(v-pt) } v_1(t^n) = \max\{t_2^n, e_1(t^n), e_2(t^n)\} \rightarrow b, \\
&v_2(t^n) = \max\{t_2^n, e_3(t^n), e_4(t^n)\} \rightarrow b, \quad v_3(t^n) = t_2^n \rightarrow 0 \\
&\text{(w-pt) } w_1(t^n) = \max\{t_2^n, e_1(t^n)\} \rightarrow b, \\
&w_2(t^n) = \max\{t_2^n, t_6^n\} \rightarrow b, \quad w_3(t^n) = t_2^n \rightarrow 0.
\end{aligned}$$

As a consequence of this, we have that

(pt-I-I-I) any triple  $(u, v, w)$ , where  $u \in \{u_1, u_2, u_3\}$ ,  $v \in \{v_1, v_2\}$ ,  $w \in \{w_1, w_2\}$ , is point asymptotic of type  $(I, I, I)$   
 (pt-I-I-0) any triple  $(u, v, w_3)$ , where  $u \in \{u_1, u_2, u_3\}$ ,  $v \in \{v_1, v_2\}$ , is point asymptotic of type  $(I, I, 0)$

- (pt-I-0-I) any triple  $(u, v_3, w)$ , where  $u \in \{u_1, u_2, u_3\}$ ,  $w \in \{w_1, w_2\}$ , is point asymptotic of type  $(I, 0, I)$
- (pt-I-0-0) any triple  $(u, v_3, w_3)$ , where  $u \in \{u_1, u_2, u_3\}$ , is point asymptotic of type  $(I, 0, 0)$ .

**R-4)** Let  $a > 0$  be arbitrary fixed; and put  $A = (0, a, a, a, a, 0)$ . By the very definition of these functions,

$$e_0(A) = 0, e_1(A) = 0, e_2(A) = a, e_3(A) = 0, e_4(A) = a.$$

This, along with our preceding conventions, yields

- (u-inv)  $u_1(A) = a, u_2(A) = a, u_3(A) = a,$
- (v-inv)  $v_1(A) = a, v_2(A) = a, v_3(A) = a$
- (w-inv)  $w_1(A) = a, w_2(A) = a, w_3(A) = a.$

As a consequence of this, we have that

- (inva) any triple  $(u, v, w)$ , where  $u \in \{u_1, u_2, u_3\}$ ,  $v \in \{v_1, v_2, v_3\}$ ,  $w \in \{w_1, w_2, w_3\}$  is invariant.

As a direct consequence of these facts and Rhoades-Dutta-Choudhury principles we just exposed, one may now derive a couple of fixed point statements with methodological value.

The former of these, starting from Rhoades-Dutta-Choudhury principle for regulated functions (RDC-reg) is based on a regulated condition upon  $\psi$ ; note that, in this case,  $\psi$  need not be increasing.

**Theorem 8** *Let the selfmap  $T$  be  $(d, \mathcal{R}; \psi, \varphi; u, v, w)$ -contractive, where  $(\psi, \varphi)$  is a couple of functions over  $\mathcal{F}(R_+^0, R)$ , with*

- (71-i)  $\psi$  is extended regulated and  $\varphi$  is strictly positive
- (71-ii)  $(\psi, \varphi)$  is right bounded oscillating at  $(b, b; b)$ , for each  $b > 0$ , and  $(u, v, w)$  is a triple of functions over  $\mathcal{F}(R_+^0, R_+)$ , with
- (71-iii)  $u \in \{u_1, u_2, u_3\}$ ,  $v \in \{v_1, v_2\}$ ,  $w \in \{w_1, w_2, w_3\}$ , such that (under the notations we just proposed)
- (71-iv)  $\mathcal{R}$  is  $B_1$ -admissible.

Further, let  $X$  be  $(Ba-o, d)$ -complete; and let  $X_0 = (x_n)$  be a  $(Ba-o)$  sequence.

Then,

- (71-a)**  $X_0 = (x_n)$  is Picard (modulo  $(d, \mathcal{R}; T)$ )
- (71-b)**  $X_0 = (x_n)$  is strongly Picard (modulo  $(d, \mathcal{R}; T)$ ), whenever  $T$  is  $(Ba-o, d)$ -continuous

**(71-c)**  $X_0 = (x_n)$  is strongly Picard (modulo  $(d, \mathcal{R}; T)$ ), whenever (in addition to the above) one of the alternatives below is holding

(71-c-i)  $v \neq v_3, w \neq w_3, (\psi, \varphi)$ =bounded oscillating at  $(b, b; b), \forall b > 0$

(71-c-ii)  $v \neq v_3, w = w_3, (\psi, \varphi)$ =bounded oscillating at  $(b, b; 0), \forall b > 0$

(71-c-iii)  $v = v_3, w \neq w_3, (\psi, \varphi)$ =bounded oscillating at  $(0, b; b), \forall b > 0$

(71-c-iv)  $v = v_3, w = w_3, (\psi, \varphi)$ =bounded oscillating at  $(0, b; 0), \forall b > 0$

**(71-d)**  $T$  is fix- $\mathcal{R}$ -asingleton (hence, necessarily, fix- $\mathcal{R}$ -singleton, under any of these extra requirements)

**(71-e)**  $T$  is fix-asingleton (hence, fix-singleton, under any of these extra requirements) when, in addition to the conditions above,

$\mathcal{R}$  is  $B_0$ -admissible.

The latter of these, starting from Rhoades-Dutta-Choudhury principle for increasing functions (RDC-incr) is based on an increasing assumption upon  $\psi$ . Note that, in this setting, we have

$(\forall i, j, k \in \{1, 2, 3\})$ :  $T$  is  $(d, \mathcal{R}; \psi, \varphi; u_i, v_j, w_k)$ -contractive implies

$T$  is  $(d, \mathcal{R}; \psi, \varphi; u_3, v_j, w_k)$ -contractive.

**Theorem 9** Let the selfmap  $T$  be  $(d, \mathcal{R}; \psi, \varphi; u_3, v, w)$ -contractive, where  $(\psi, \varphi)$  is a couple of functions over  $\mathcal{F}(R_+^0, R)$ , with

(72-i)  $\psi$  is increasing and  $\varphi$  is strictly positive

(72-ii)  $(\psi, \varphi)$  is right bounded oscillating at  $(b, b; b)$ , for each  $b > 0$ ,

and  $v, w$  are functions over  $\mathcal{F}(R_+^0, R_+)$ , with

(72-iii)  $v \in \{v_1, v_2, v_3\}, w \in \{w_1, w_2, w_3\}$ ,

such that (under the notations we just proposed)

(72-iv)  $\mathcal{R}$  is  $B_1$ -admissible.

Further, let  $X$  be  $(Ba-o, d)$ -complete; and let  $X_0 = (x_n)$  be a  $(Ba-o)$  sequence.

Then,

**(72-a)**  $X_0 = (x_n)$  is Picard (modulo  $(d, \mathcal{R}; T)$ )

**(72-b)**  $X_0 = (x_n)$  is strongly Picard (modulo  $(d, \mathcal{R}; T)$ ), whenever  $T$  is  $(Ba-o, d)$ -continuous

**(72-c)**  $X_0 = (x_n)$  is strongly Picard (modulo  $(d, \mathcal{R}; T)$ ) whenever (in addition to the above) one of the alternatives below is holding

(72-c-i)  $v \neq v_3, w \neq w_3, (\psi, \varphi)$ =bounded oscillating at  $(b, b; b), \forall b > 0$

(72-c-ii)  $v \neq v_3, w = w_3, (\psi, \varphi)$ =bounded oscillating at  $(b, b; 0), \forall b > 0$

(72-c-iii)  $v = v_3, w \in \{w_1, w_2, w_3\}$

**(72-d)**  $T$  is fix- $\mathcal{R}$ -asingleton (hence, necessarily, fix- $\mathcal{R}$ -singleton, under any of these extra requirements)

**(72-e)**  $T$  is fix-asingleton (hence, fix-singleton, under any of these extra requirements) when, in addition to the conditions above,

$\mathcal{R}$  is  $B_0$ -admissible.

**Proof** The only point to be clarified is that of the listed alternative related to  $v = v_3$  being sufficient for deriving the strongly Picard property. To do this, the same reasonings as the ones in the main result will be applied.

From the Picard property, we have, as  $X$  is (Ba-o,d)-complete

$$X_0 = (x_n) \text{ is } d\text{-convergent: } x_n \xrightarrow{d} z_0 \text{ as } n \rightarrow \infty, \text{ for some } z_0 \in X.$$

We now claim that  $z_0$  is a fixed point of  $T$ . Three alternatives are to be discussed.

**Alter-1)** Suppose that

$$(Tz\text{-rela-1}) \ H_1 := \{n \in N; x_n = Tz_0\} \text{ is unbounded (in } N).$$

By the same procedure as the one in the main result, we derive  $z_0 \in \text{Fix}(T)$ .

**Alter-2)** Suppose that

$$(Tz\text{-rela-2}) \ H_2 := \{n \in N; Tx_n = Tz_0\} \text{ is unbounded (in } N).$$

By the same procedure as the one in the main result, we derive  $z_0 \in \text{Fix}(T)$ .

**Alter-3)** Suppose that

both subsets  $H_1$  and  $H_2$  are bounded (in  $N$ ).

This tells us that

$$\exists i = i(z_0) \in N, \text{ such that:}$$

$$n \geq i \text{ implies } Tx_n \neq Tz_0 \text{ (hence, } x_n \neq z_0) \text{ and } x_n \neq Tx_0.$$

Denote for simplicity  $(u_n = x_{n+i}; n \geq 0)$ ; clearly, by the preceding relation,

$$\text{(non-id) } (\forall n) : Tu_n \neq Tz_0 \text{ (hence, } u_n \neq z_0) \text{ and } u_n \neq Tx_0.$$

Again combining with the (Ba-o) property of  $X_0 = (x_n)$ , one derives

$$\text{(posi-1) } (\forall n) : Q_1(u_n, z_0) = d(u_n, Tu_n) > 0,$$

$$Q_2(u_n, z_0) = d(u_n, z_0) > 0, \quad Q_3(u_n, z_0) = d(u_n, Tx_0) > 0,$$

$$Q_4(u_n, z_0) = d(Tu_n, z_0) > 0, \quad Q_5(u_n, z_0) = d(Tu_n, Tx_0) > 0.$$

Suppose by contradiction that

$$\text{(posi-2) } b := d(z_0, Tx_0) > 0 \text{ [whence, } Q_6(u_n, z_0) = b > 0, \forall n].$$

From the preceding observations, we have

$$(\forall n) : B_3(u_n, z_0) = \min\{Q_1, \dots, Q_6\}(u_n, z_0) > 0;$$

$$\text{so that, } (u_n, z_0) \in \mathcal{R}.$$

Consequence, the function contractive condition involving the systems  $(\psi, \varphi)$  and  $(u_3, v_3, w)$  applies; and gives

$$(\forall n) : \psi(d(Tu_n, Tz_0)) \leq \psi(d(u_n, z_0)) - \varphi(w(\mathcal{Q}(u_n, z_0))).$$

By the strict positivity assumption upon  $\varphi$ , one derives

$$(\forall n)\psi(d(Tu_n, Tz_0)) < \psi(d(u_n, z_0));$$

$$\text{whence } d(Tu_n, Tz_0) < d(u_n, z_0);$$

if we remember that  $\psi$  is increasing. Passing to limit as  $n \rightarrow \infty$ , we get (by a metrical property of  $d$ )

$$0 < d(z_0, Tz_0) \leq 0; \text{ a contradiction.}$$

Hence, our working assumption relative to  $z_0$  cannot be accepted; and then, we have  $b = 0$ ; that is:  $z_0 = Tz_0$ .

Some remarks are in order.

**Rem-1)** An extension of these statements to quasi-ordered spaces is immediate, by the developments in Turinici [43]. Further aspects will be delineated elsewhere.

**Rem-2)** Under the formulation above, Theorem 8 contains 27 fixed point statements. Among these, we have

(part-1) the (quasi-order) variant  $(u_3, v_1, w_1)$  of Theorem 8 is identical with the fixed point statement in Chandok et al. [6]

(part-2) the (quasi-order) variant  $(u_3, v_2, w_2)$  of Theorem 8 is identical with the fixed point statement in Choudhury et al. [7].

In order words: the quoted results admit extensions to regulated functions. Some other variants of Theorem 8 include a number of related fixed point results described in Radenović et al. [28]. But, most of these seem to be new, at least from the perspective of regulated (modulo  $\psi$ ) setting.

**Rem-3)** Under the formulation above, Theorem 9 contains 9 fixed point statements. The variant  $(u_3, v_3, w_3)$  of this statement includes the results in Wardowski [46] and Secelean [34]; see also Vujaković et al. [44]. For a different proof of it, we refer to the paper by Turinici [42]. The remaining variants  $(u_3, v_3, w)$  of the same where  $w \in \{w_1, w_2\}$ , seem to be new.

Finally, note that the proposed examples do not exhaust all possible variants described in the Rhoades-Dutta-Choudhury principles (RDC-reg) and (RDC-incr); for, e.g., the ones concerning strongly/weakly  $\lambda$ -iterative properties with  $\lambda \in \mathcal{F}(R_+^0)$  fulfilling  $\lambda \neq I$  are not present in this list of particular cases. Further technical aspects of this case will be delineated elsewhere.

### 8 Cosentino-Vetro Contractions

Let  $(X, d)$  be a metric space. Further, let  $T \in \mathcal{F}(X)$  be a selfmap of  $X$ . As precise, we are interested to determine sufficient conditions for (uniqueness and) existence of elements in  $\text{Fix}(T)$ , via contractive type requirements involving iterative processes  $X_0 = (x_n)$  starting from an element  $x_0$  of  $X$ , and their associated sets

$$[X_0] = \{x_n; n \geq 0\}, [[X_0]] = \text{cl}([X_0]).$$

Let  $(\psi, \varphi)$  be a pair of functions over  $\mathcal{F}(R_+^0, R)$ . Further, let  $(u, v, w)$  be a triple of functions over  $\mathcal{F}(R_+^6, R_+)$ , and define the couple of subsets

$$\begin{aligned} \mathcal{Y}_0 &= \{t \in R_+^6; u(t), v(t), w(t) > 0\}, \\ \mathcal{Y} &= \{t \in \mathcal{Y}_0; \psi(u(t)) \leq \psi(v(t)) - \varphi(w(t))\}. \end{aligned}$$

Finally, define the couple of relations over  $X$

$$\mathcal{R} = \mathcal{Q}^{-1}(\mathcal{Y}_0), \mathcal{R}^* = \mathcal{Q}^{-1}(\mathcal{Y}).$$

The functional and set contractive properties to be used here are

- (fct-contr)  $T$  is  $(d, \mathcal{R}; \psi, \varphi; u, v, w)$ -contractive, if
- $\psi(u \circ \mathcal{Q}(x, y)) \leq \psi(v \circ \mathcal{Q}(x, y)) - \varphi(w \circ \mathcal{Q}(x, y))$ , when  $(x, y) \in \mathcal{R}$
- (set-contr)  $T$  is  $(d, \mathcal{R}; \mathcal{Y})$ -contractive:  $\mathcal{Q}(x, y) \in \mathcal{Y}$ , if  $(x, y) \in \mathcal{R}$ .

The connection between these conditions is expressed as

- (equi)  $T$  is  $(d, \mathcal{R}; \psi, \varphi; u, v, w)$ -contractive iff  $T$  is  $(d, \mathcal{R}; \mathcal{Y})$ -contractive.

For the set contractive property, sufficient conditions were given upon  $\mathcal{R}$  and  $\mathcal{Y}$  so that the main result be applicable; these, as expected, were ultimately expressed in terms of  $\mathcal{R}$ ,  $(\psi, \varphi)$  and  $(u, v, w)$ . The by-product of these developments is the couple of Rhoades-Dutta-Choudhury fixed point results (RDC-reg) and (RDC-incr). Some particular cases of them were just exposed. It is our aim in the following to discuss some other particular cases of these, with a practical meaning.

(A) Let us say that  $v \in \mathcal{F}(R_+^6, R_+)$  is a *Cosentino-Vetro function*, when

- (CV-0)  $v$  is increasing and continuous in its variables
- (CV-1)  $v$  is (1, 2, 5, 6)-positive:  $v(a, a, 0, 0, a, a) > 0$ , for each  $a > 0$
- (CV-2)  $v$  is  $I$ -iterative:  $v(a, a, 2a, 0, a, a) \leq a$ , for each  $a > 0$
- (CV-3)  $v$  is invariant:  $v(0, b, b, b, b, 0) = b$ , for each  $b > 0$ .

Fix in the following such a function; and define  $\eta \in \mathcal{F}(R_+^0, R_+)$  as

$$(E\text{-def}) \eta(b) := v(0, 0, b, 0, b, b), b > 0.$$

Note that, by ( $v$ -increasing continuous) and the iterative condition:

(E-1)  $\eta(\cdot)$  is increasing and continuous (on  $R_+^0$ )

(E-2)  $0 \leq \eta(b) \leq v(b, b, 2b, 0, b, b) \leq b, \forall b > 0$ ; whence,  $\eta(0 + 0) = 0$ .

Unfortunately, the extremal inequalities above are not in general reducible to equalities. Then, define the couple  $(u, w)$  of functions over the class  $\mathcal{F}(R_+^6, R_+)$ , as: for each  $t = (t_1, \dots, t_6) \in R_+^6$ ,

$$(uw\text{-def}) \quad u(t) = t_5, w(t) = t_2;$$

in this case,  $(u, v, w)$  will be referred to as a *Cosentino-Vetro triple*.

Having this precise, let us introduce the class of contractive conditions

(C-V)  $T$  is Cosentino-Vetro  $(d, \mathcal{R}; \psi, \varphi; u, v, w)$ -contractive:

$$\psi(d(Tx, Ty)) \leq \psi(v(Q_1(x, y), \dots, Q_6(x, y)) - \varphi(d(x, y)), \forall (x, y) \in \mathcal{R}.$$

As usual, we assume that (under the notations we just proposed)

$$(B1\text{-adm}) \quad \mathcal{R} \text{ is } B_1\text{-admissible: } (B_1 > 0) \subseteq \mathcal{R}.$$

In this case, the following contractive property holds

(C-V-B1)  $T$  is Cosentino-Vetro  $(d, (B_1 > 0); \psi, \varphi; u, v, w)$ -contractive:

$$\psi(d(Tx, Ty)) \leq \psi(v(Q_1(x, y), \dots, Q_6(x, y)) - \varphi(d(x, y)), \forall (x, y) \in (B_1 > 0).$$

We are interested of to what extent is the triple  $(u, v, w)$  compatible with the regularity conditions encountered in Rhoades-Dutta-Choudhury principle for increasing functions (RDC-incr). The response to this question is contained in the developments below.

**R-1**) Given the sequence  $(r_n)$  in  $R_+^0$  and the sequence  $(p_n)$  in  $R_+$  with  $(A_n := (r_n, r_n, p_n, 0, r_{n+1}, r_{n+1}) \in \mathcal{Q}(B_1 > 0; [X_0]) \cap \mathcal{Y}, |p_n - r_n| \leq r_{n+1}$ , for all  $n$ ), we have (by these definitions)

$$(\forall n) : u(A_n) = r_{n+1}, w(A_n) = r_n.$$

On the other hand, under the convention

$$\alpha_n = \max\{r_n, r_{n+1}\}, n \geq 0$$

we have, by the iterative property

$$(\forall n) : v(A_n) = v(r_n, r_n, p_n, 0, r_{n+1}, r_{n+1}) \leq v(\alpha_n, \alpha_n, 2\alpha_n, 0, \alpha_n, \alpha_n) \leq \alpha_n.$$

This tells us that

the triple  $(u, v, w)$  is weakly  $I$ -iterative (see above);

and explains the choosing of Rhoades-Dutta-Choudhury principle for increasing functions (RDC-incr) as an appropriate tool for solving our problem.



**R-2)** Let the sequence  $(t^n = (t_1^n, \dots, t_6^n); n \geq 0)$  in  $\mathcal{Q}(B_3 > 0; [X_0]) \cap \mathcal{Y}$ , and the point  $b > 0$  be such that

$(t^n = (t_1^n, \dots, t_6^n); n \geq 0)$  has the right property at  $(0, b, b, b, b, 0)$ ; that is:  
 $(t_i^n \rightarrow 0, \forall i \in \{1, 6\}), (t_j^n \rightarrow b+, \forall j \in \{2, 3, 4, 5\})$ .

For the moment, we have by definition

$$u(t^n) = t_5^n \rightarrow b+, w(t^n) = t_2^n \rightarrow b+.$$

On the other hand, under the convention

$$\beta_n = \min\{t_2^n, t_3^n, t_4^n, t_5^n\}, n \geq 0 \text{ [hence, } \beta_n \rightarrow b+ \text{ as } n \rightarrow \infty\text{],}$$

we have (under the posed hypotheses upon  $v$ )

$$v(t^n) \geq v(0, \beta_n, \beta_n, \beta_n, \beta_n, 0) = \beta_n > b, \text{ for all } n,$$

$$\text{and } v(t^n) \rightarrow v(0, b, b, b, b, 0) = b \text{ as } n \rightarrow \infty;$$

$$\text{wherefrom } v(t^n) \rightarrow b+ \text{ as } n \rightarrow \infty.$$

Putting these together, it follows that

(r-asy) the triple  $(u, v, w)$  is right asymptotic of type  $(I, I, I)$ .

**R-3)** Let the sequence  $(t^n = (t_1^n, \dots, t_6^n); n \geq 0)$  in  $\mathcal{Q}(B_3 > 0; [[X_0]]) \cap \mathcal{Y}$ , and the point  $b > 0$  be such that

$(t^n = (t_1^n, \dots, t_6^n); n \geq 0)$  has the point property at  $(0, 0, b, 0, b, b)$ ;  
 that is:  $(t_i^n \rightarrow 0, \forall i \in \{1, 2, 4\}), (t_j^n \rightarrow b, \forall j \in \{3, 5, 6\}), (t_6^n = b, \forall n)$ .

Then, by definition,

$$u(t^n) = t_5^n \rightarrow b, v(t^n) \rightarrow v(0, 0, b, 0, b, b) = \eta(b), w(t^n) = t_2^n \rightarrow 0.$$

As a consequence of this, we have that

the triple  $(u, v, w)$  is point asymptotic of type  $(I, \eta, 0)$ .

**R-4)** Let  $a > 0$  be arbitrary fixed; and put  $A = (0, a, a, a, a, 0)$ . By the very definition of these functions,

$$u(A) = a, v(A) = v(0, a, a, a, a, 0) = a, w(A) = a.$$

And then, we have that

the triple  $(u, v, w)$  is invariant.

As a consequence of Rhoades-Dutta-Choudhury principle for increasing functions (RDC-incr), we derive the following fixed point statement (referred to as: *Cosentino-Vetro fixed point principle*; in short: (CV-fpp)).

**Theorem 10** *Let the selfmap  $T$  be Cosentino-Vetro  $(d, \mathcal{R}; \psi, \varphi; u, v, w)$ -contractive, where  $(\psi, \varphi)$  is a couple of functions over  $\mathcal{F}(R_+^0, R)$ , with*

- (81-i)  $\psi$  is increasing and  $\varphi$  is strictly positive
- (81-ii)  $(\psi, \varphi)$  is right bounded oscillating at  $(b, b; b)$ , for each  $b > 0$

and  $(u, v, w)$  is a Cosentino-Vetro triple over  $\mathcal{F}(R_+^6, R_+)$  such that (under the notations we just proposed)

- (81-iii)  $\mathcal{R}$  is  $B_1$ -admissible:  $(B_1 > 0) \subseteq \mathcal{R}$ .

Further, let  $X$  be  $(Ba-o,d)$ -complete; and let  $X_0 = (x_n)$  be a  $(Ba-o)$  sequence.

Then,

- (81-a)**  $X_0 = (x_n)$  is Picard (modulo  $(d, \mathcal{R}; T)$ )
- (81-b)**  $X_0 = (x_n)$  is strongly Picard (modulo  $(d, \mathcal{R}; T)$ ), whenever  $T$  is  $(Ba-o,d)$ -continuous
- (81-c)**  $X_0 = (x_n)$  is strongly Picard (modulo  $(d, \mathcal{R}; T)$ ), whenever (in addition to the above)

$$(\psi, \varphi) \text{ is bounded oscillating at } (\eta(b), b; 0), \text{ for all } b > 0$$

- (81-d)**  $T$  is fix- $\mathcal{R}$ -asingleton (hence, necessarily, fix- $\mathcal{R}$ -singleton, under any of these extra requirements)
- (81-6)**  $T$  is fix-asingleton (hence, fix-singleton, under any of these extra requirements) when, in addition to the conditions above,

$$\mathcal{R} \text{ is } B_0\text{-admissible: } (B_0 > 0) \subseteq \mathcal{R}.$$

A basic particular case of these developments corresponds to the choice

$$v(t) = \alpha_1 t_1 + \dots + \alpha_6 t_6, t = (t_1, \dots, t_6) \in R_+^6,$$

where  $(\alpha_1, \dots, \alpha_6)$  is a vector in  $R_+^6$  with the properties

- (Pro-1)  $\alpha_1 + \alpha_2 + \alpha_5 + \alpha_6 > 0$ ,
- (Pro-2)  $\alpha_1 + \alpha_2 + 2\alpha_3 + \alpha_5 + \alpha_6 \leq 1$ ,
- (Pro-3)  $\alpha_2 + \alpha_3 + \alpha_4 + \alpha_5 = 1$ .

Denote further

$$\eta := \alpha_3 + \alpha_5 + \alpha_6; \text{ hence, } 0 \leq \eta \leq 1.$$

The contractive condition upon  $T$  becomes a linear contraction like

(CV-lin)  $T$  is linear Cosentino-Vetro  $(d, \mathcal{R}; \psi, \varphi; u, v, w)$ -contractive:

$$\psi(d(Tx, Ty)) \leq \psi(\alpha_1 Q_1(x, y) + \dots + \alpha_6 Q_6(x, y)) - \varphi(d(x, y)), \forall (x, y) \in \mathcal{R}.$$

**Theorem 11** *Let  $T$  be linear Cosentino-Vetro  $(d, \mathcal{R}; \psi, \varphi; u, v, w)$ -contractive, where  $(\psi, \varphi)$  is a couple of functions over  $\mathcal{F}(R_+^0, R)$ , with*

(82-i)  $\psi$  is increasing and  $\varphi$  is strictly positive

(82-ii)  $(\psi, \varphi)$  is right bounded oscillating at  $(b, b; b)$  for each  $b > 0$ ,

and the constant vector  $\alpha = (\alpha_1, \dots, \alpha_6)$  in  $R_+^6$  fulfilling (Pro-1)-(Pro-3), be such that (under the notations we just proposed)

(B1-adm)  $\mathcal{R}$  is  $B_1$ -admissible:  $(B_1 > 0) \subseteq \mathcal{R}$ .

Further, let  $X$  be  $(Ba-o, d)$ -complete; and let  $X_0 = (x_n)$  be a  $(Ba-o)$  sequence.

Then,

(82-a)  $X_0 = (x_n)$  is Picard (modulo  $(d, \mathcal{R}; T)$ )

(82-b)  $X_0 = (x_n)$  is strongly Picard (modulo  $(d, \mathcal{R}; T)$ ), whenever  $T$  is  $(Ba-o, d)$ -continuous

(82-c)  $X_0 = (x_n)$  is strongly Picard (modulo  $(d, \mathcal{R}; T)$ ), whenever (in addition to the above)

$$(\psi, \varphi) \text{ is bounded oscillating at } (\eta b, b; 0), \text{ for all } b > 0$$

(82-d)  $T$  is fix- $\mathcal{R}$ -asingleton (hence, necessarily, fix- $\mathcal{R}$ -singleton, under any of these extra requirements)

(82-e)  $T$  is fix-asingleton (hence, fix-singleton, under any of these extra requirements) when, in addition to the conditions above,

$$\mathcal{R} \text{ is } B_0\text{-admissible: } (B_0 > 0) \subseteq \mathcal{R}.$$

Some remarks are in order.

**Rem-1)** An extension of these statements to quasi-ordered spaces is immediate, by the developments in Turinici [43]. Further aspects will be delineated elsewhere.

**Rem-2)** The nonlinear fixed point statement expressed via Theorem 10 seems to be new.

**Rem-3)** The variant  $(\alpha_2 = 1, \alpha_j = 0, j \neq 2)$  of Theorem 11 includes the statement in this area due to Wardowski [45, 46]; see also Secolean [34]. For a different proof of it, we refer to the paper by Turinici [42].

**Rem-4)** The remaining variants of Theorem 11 include a related statement due to Cosentino and Vetro [9], and its refinement in Vujaković et al. [44]. Further aspects may be found in Popescu and Stan [27].

### 9 Wardowski Type Contractions

Let  $(X, d)$  be a metric space. Further, let  $T \in \mathcal{F}(X)$  be a selfmap of  $X$ . As precise, we are interested to determine sufficient conditions for (uniqueness and) existence of elements in  $\text{Fix}(T)$ , via contractive type requirements involving iterative processes  $X_0 = (x_n)$  starting from an element  $x_0$  of  $X$ , and their associated sets

$$[X_0] = \{x_n; n \geq 0\}, [[X_0]] = \text{cl}([X_0]).$$

Let  $(\psi, \varphi)$  be a pair of functions over  $\mathcal{F}(R_+^0, R)$ . Further, let  $(u, v, w)$  be a triple of functions over  $\mathcal{F}(R_+^6, R_+)$ , and define the couple of subsets

$$\begin{aligned} \mathcal{Y}_0 &= \{t \in R_+^6; u(t), v(t), w(t) > 0\}, \\ \mathcal{Y} &= \{t \in \mathcal{Y}_0; \psi(u(t)) \leq \psi(v(t)) - \varphi(w(t))\}. \end{aligned}$$

Finally, define the couple of relations over  $X$

$$\mathcal{R} = \mathcal{Q}^{-1}(\mathcal{Y}_0), \mathcal{R}^* = \mathcal{Q}^{-1}(\mathcal{Y}).$$

The functional and set contractive properties to be used here are

- (fct-contr)  $T$  is  $(d, \mathcal{R}; \psi, \varphi; u, v, w)$ -contractive, if  $\psi(u \circ \mathcal{Q}(x, y)) \leq \psi(v \circ \mathcal{Q}(x, y)) - \varphi(w \circ \mathcal{Q}(x, y))$ , when  $(x, y) \in \mathcal{R}$
- (set-contr)  $T$  is  $(d, \mathcal{R}; \mathcal{Y})$ -contractive:  $\mathcal{Q}(x, y) \in \mathcal{Y}$ , if  $(x, y) \in \mathcal{R}$ .

And, the connection between these conditions is expressed as

$$\text{(equi)} \quad T \text{ is } (d, \mathcal{R}; \psi, \varphi; u, v, w)\text{-contractive iff } T \text{ is } (d, \mathcal{R}; \mathcal{Y})\text{-contractive.}$$

For the set contraction, sufficient conditions were given upon  $\mathcal{R}$  and  $\mathcal{Y}$  so that the main result be applicable; these, as expected, were ultimately expressed in terms of  $\mathcal{R}$ ,  $(\psi, \varphi)$  and  $(u, v, w)$ . The by-product of these developments is the couple of Rhoades-Dutta-Choudhury fixed point results (RDC-reg) and (RDC-incr). Some particular cases of them have been previously exposed.

As results from their proof, a common feature of all these is the fact that the only asymptotic property of  $\mathcal{Y}$  to be used there is the descending one; and not the general one. It is our aim in the following to show, by a standard example, that this general condition is ultimately applicable; and yields some interesting results.

To begin with, let  $(\psi, \varphi)$  be a pair of functions in  $\mathcal{F}(R_+^0, R)$ ; and  $(u, v, w)$  be the triple of functions introduced as: for each  $t = (t_1, \dots, t_6) \in R_+^6$

$$u(t) = t_5, v(t) = w(t) = t_2.$$

We have, according to our conventions

$$\mathcal{Y}_0 = \{t \in R_+^6; t_5, t_2 > 0\}, \mathcal{Y} = \{t \in \mathcal{Y}; \psi(t_5) \leq \psi(t_2) - \varphi(t_2)\}.$$

In addition, the attached relations to be used here are (see above)

$$\mathcal{R} = \mathcal{Q}^{-1}(\mathcal{Y}_0), \mathcal{R}^* = \mathcal{Q}^{-1}(\mathcal{Y}).$$

The contractive property attached to these data writes, in a shorter way

(fct-contr)  $T$  is  $(d, \mathcal{R}; \psi, \varphi)$ -contractive, if

$$\psi(d(Tx, Ty)) \leq \psi(d(x, y)) - \varphi(d(x, y)), \text{ when } (x, y) \in \mathcal{R}.$$

As before, we impose a regularity condition like

$$\mathcal{R} \text{ is } B_1\text{-admissible: } (B_1 > 0) \subseteq \mathcal{R}.$$

This, combined with the functional contractive condition, yields

(B1-contr)  $T$  is  $(d, (B_1 > 0); \psi, \varphi)$ -contractive :

$$\psi(d(Tx, Ty)) \leq \psi(d(x, y)) - \varphi(d(x, y)), \text{ when } (x, y) \in (B_1 > 0);$$

which is exactly the contractive condition we are dealing with. Note that under the stronger admissible assumption

$$\mathcal{R} \text{ is } B_0\text{-admissible: } (B_1 > 0) \subseteq \mathcal{R}$$

the underlying contractive condition is obtainable from the standard one

(B0-contr)  $T$  is  $(d, (B_0 > 0); \psi, \varphi)$ -contractive :

$$\psi(d(Tx, Ty)) \leq \psi(d(x, y)) - \varphi(d(x, y)), \text{ when } (x, y) \in (B_0 > 0).$$

But, the reciprocal implication is not in general valid.

To get the announced result about our introduced class of contractions, a couple of specific conditions is needed:

- (str-pos)  $\varphi$  is strictly positive:  $\varphi(R_+^0) \subseteq R_+^0$
- (r-as-pos)  $\varphi$  is right asymptotic positive: for each sequence  $(t_n; n \geq 0)$  in  $R_+^0$  and each  $b > 0$  with  $t_n \rightarrow b+$ , we must have  $\limsup_n \varphi(t_n) > 0$
- (tele-admi)  $(\psi, \varphi)$  is tele admissible: each sequence  $(t_n)$  in  $R_+^0$  with  $(\psi(t_{n+1}) \leq \psi(t_n) - \varphi(t_n), \forall n)$  fulfills  $\lim_n(t_n) = 0$
- (zero-her)  $(\psi, \varphi)$  is zero hereditary: there are no couple of sequences  $(t_n)$  and  $(s_n)$  in  $R_+^0$  and no points  $b > 0$  with  $(\psi(t_n) \leq \psi(s_n) - \varphi(s_n), \forall n)$  and  $t_n \rightarrow b, s_n \rightarrow 0$ .

Note that our asymptotic condition may be also written by means of the oscillation concepts we already introduced; but this is not important for us. Moreover, under the strict positive condition, this property is available under

(cont)  $\varphi$  is continuous (on  $R_+^0$ )

(loc-incr)  $\varphi$  is locally increasing: each  $b \in R_+^0$  has a neighborhood  $V_b$  such that  $\varphi$  is increasing on  $V_b$ .

The following auxiliary statement is the basic step towards our objective.

**Proposition 7** *Let  $(\psi, \varphi)$  be a couple of functions over  $\mathcal{F}(R_+^0, R)$ , with*

(91-I)  $\psi$  is almost right regulated ( $\Theta := rreg(\psi)$  is  $(>)$ -cofinal in  $R_+^0$ )

(91-II)  $\varphi$  is strictly positive and right asymptotic positive

(91-III)  $(\psi, \varphi)$  is tele admissible and zero hereditary.

Then,

(91-1)  $\Upsilon$  is asymptotic on  $[X_0]$

(91-2)  $\Upsilon$  is nright on  $[X_0]$

(91-3)  $\Upsilon$  is npoint on  $[[X_0]]$

(91-4)  $\Upsilon$  is nnormal.

**Proof** The argument consists in a number of parts.

**Part 1.** ( $\Upsilon$  is asymptotic on  $[X_0]$ ).

Let the sequence  $(r_n)$  in  $R_+^0$  and the sequence  $(p_n)$  in  $R_+$  be such that

$$(\forall n) : (r_n, r_n, p_n, 0, r_{n+1}, r_{n+1}) \in \mathcal{Q}(B_1 > 0; [X_0]) \cap \Upsilon, \text{ and } |p_n - r_n| \leq r_{n+1}.$$

By the very representation of  $\Upsilon$ , we must have

$$(\forall n) : \psi(r_{n+1}) \leq \psi(r_n) - \varphi(r_n).$$

This, along with  $(\psi, \varphi)$  being tele admissible, gives  $r_n \rightarrow 0$  (hence,  $p_n \rightarrow 0$ ); and proves the desired assertion.

**Part 2.** ( $\Upsilon$  is almost nright on  $[X_0]$ ).

Take some point  $c = (c_1, \dots, c_6)$  in  $R_+^6$ . We say that the (vectorial) sequence  $(t^n := (t_1^n, \dots, t_6^n); n \geq 0)$  in  $R_+^6$  is *right* at  $c$ , if

$$(r-c) (t_i^n \rightarrow c_i, \forall i) \text{ and } (t_i^n \rightarrow c_i +, \text{ whenever } c_i > 0).$$

Given  $b > 0$ , let us say that  $\Upsilon$  is *nright* at  $b$  on  $[X_0]$ , if

(nright) for each sequence  $(t^n; n \geq 0)$  in  $\mathcal{Q}(B_3 > 0; [X_0]) \cap \Upsilon$ , the right property at  $(0, b, b, b, b, 0)$  is not true.

The class of all these  $b > 0$  will be denoted as  $nright(\Upsilon; [X_0])$ . In this case, we say that  $\Upsilon$  is

(a-n-r) *almost nright* on  $[X_0]$ , if  $\Theta := nright(\Upsilon; [X_0])$  is  $(>)$ -cofinal in  $R_+^0$  (for each  $\varepsilon \in R_+^0$  there exists  $\theta \in \Theta$  with  $\varepsilon > \theta$ )

(n-r) *nright* on  $[X_0]$ , if  $\Theta := nright(\Upsilon; [X_0])$  is identical with  $R_+^0$ .

We have to establish that the former property is retainable for our data. To do this, we start by noting that

$$\psi = \text{almost right regulated implies } \Theta := \text{rreg}(\psi) \text{ is } (>)\text{-cofinal in } R_+^0.$$

We now claim that

$$\Theta \subseteq \text{nright}(\mathcal{Y}; [X_0]); \text{ wherefrom, } \text{nright}(\mathcal{Y}; [X_0]) \text{ is } (>)\text{-cofinal in } R_+^0.$$

This amounts to establish that

for each sequence  $(t^n := (t_1^n, \dots, t_6^n); n \geq 0)$  in  $\mathcal{Q}(B_3 > 0; [X_0]) \cap \mathcal{Y}$ , and each  $b \in \Theta$ , the right property at  $(0, b, b, b, b, 0)$  is not true.

Suppose—by *reductio ad absurdum*—that there exists a (vectorial) sequence  $(t^n := (t_1^n, \dots, t_6^n); n \geq 0)$  in  $R_+^6$  and some  $b \in \Theta$ , with

(right-b-1)  $(\forall n): t^n := (t_1^n, \dots, t_6^n) \in \mathcal{Q}(B_3 > 0; [X_0]) \cap \mathcal{Y}$ ; whence,  $\psi(t_5^n) \leq \psi(t_2^n) - \varphi(t_2^n)$

(right-b-2)  $(t^n := (t_1^n, \dots, t_6^n); n \geq 0)$  has the right property at  $(0, b, b, b, b, 0)$ ; that is:  $(t_i^n \rightarrow 0, \forall i \in \{1, 6\})$ , and  $(t_i^n \rightarrow b+, \forall i \in \{2, 3, 4, 5\})$ .

By the former of these properties

$$(\forall n) : (0 <) \varphi(t_2^n) \leq \psi(t_2^n) - \psi(t_5^n).$$

Passing to  $\limsup$  as  $n \rightarrow \infty$  in this relation gives (via  $\psi$ =right regulated at  $b$ )

$$0 \leq \limsup_n \varphi(t_2^n) \leq \psi(b + 0) - \psi(b + 0) = 0; \text{ that is: } \lim_n \varphi(t_2^n) = 0;$$

in contradiction with  $\varphi$ =right asymptotic positive. Hence, our working assumption cannot be true; and the assertion follows.

**Part 3.** ( $\mathcal{Y}$  is *npoint* on  $[[X_0]]$ ).

Take some point  $c = (c_1, \dots, c_6)$  in  $R_+^6$ . We say that the sequence  $(t^n := (t_1^n, \dots, t_6^n); n \geq 0)$  in  $R_+^6$  is *point* at  $c$ , if

$$(\text{pt-c}) \ (t_i^n \rightarrow c_i, \forall i) \text{ and } [t_6^n = c_6, \forall n].$$

Given  $b > 0$ , let us say that  $\mathcal{Y}$  is *npoint* at  $b$  on  $[[X_0]]$ , if

(*npoint*) for each sequence  $(t^n; n \geq 0)$  in  $\mathcal{Q}(B_3 > 0; [[X_0]]) \cap \mathcal{Y}$ , the *point* property at  $(0, 0, b, 0, b, b)$  is not true.

The class of all these  $b > 0$  will be denoted as  $\text{npoint}(\mathcal{Y}; [[X_0]])$ . In this case, we say that  $\mathcal{Y}$  is

(a-n-p) *almost npoint* on  $[[X_0]]$ , if  $\Theta := \text{npoint}(\mathcal{Y}; [[X_0]])$  is  $(>)$ -cofinal in  $R_+^0$

(n-p) *npoint* on  $[[X_0]]$ , if  $\Theta := \text{npoint}(\mathcal{Y}; [[X_0]])$  is identical with  $R_+^0$ .

We have to establish that the latter of these properties holds; that is, for each sequence  $(t^n; n \geq 0)$  in  $R_+^6$  and each  $b > 0$ , the point property at  $(0, 0, b, 0, b, b)$  is not true.

Assume by contradiction that this assertion is false: there exists a sequence  $(t^n = (t_1^n, \dots, t_6^n); n \geq 0)$  in  $R_+^6$  and some  $b > 0$ , with

(point-b-1)  $(\forall n): t^n = (t_1^n, \dots, t_6^n) \in \mathcal{Q}(B_3 > 0; [[X_0]]) \cap \mathcal{Y}$ ; whence,  $\psi(t_5^n) \leq \psi(t_2^n) - \varphi(t_2^n)$

(point-b-2)  $(t^n = (t_1^n, \dots, t_6^n); n \geq 0)$  is 6-point at  $(0, 0, b, 0, b, b)$ ; that is:  $(t_i^n \rightarrow 0, \forall i \in \{1, 2, 4\}), (t_i^n \rightarrow b, \forall i \in \{3, 5, 6\}),$  and  $(t_6^n = b, \forall n)$ .

The obtained inequality and convergence relations yield a contradiction with respect to  $(\psi, \varphi)$  being zero hereditary. Consequently, the working assumption is not acceptable; and conclusion follows.

**Part 4.** ( $\mathcal{Y}$  is nnormal).

Remember that, this property means:

$$(0, a, a, a, a, 0) \in \mathcal{Y} \text{ is impossible, for each } a > 0.$$

Suppose by contradiction that

$$\text{there exists } a > 0 \text{ with } (0, a, a, a, a, 0) \in \mathcal{Y}.$$

By the very definition of our triple,

$$u(0, a, a, a, a, 0) = a, v(0, a, a, a, a, 0) = a, w(0, a, a, a, a, 0) = a.$$

This, along with the definition of  $\mathcal{Y}$ , yields

$$\psi(a) \leq \psi(a) - \varphi(a); \text{ that is: } \varphi(a) \leq 0;$$

absurd, by  $\varphi$ =strictly positive; and our claim follows.

Now, by simply combining the obtained fact with our main result, one gets the following particular fixed point statement.

**Theorem 12** Assume that  $T$  is  $(d, \mathcal{R}; \psi, \varphi)$ -contractive, where  $(\psi, \varphi)$  is a couple over  $\mathcal{F}(R_+^0, R)$  with

(91-i)  $\psi$  is almost right regulated ( $\Theta := rreg(\psi)$  is  $(>)$ -cofinal in  $R_+^0$ )

(91-ii)  $\varphi$  is strictly positive and right asymptotic positive

(91-iii)  $(\psi, \varphi)$  is tele admissible and zero hereditary,

and the associated relation  $\mathcal{R}$  fulfills

(91-iv)  $\mathcal{R}$  is  $(B_1 > 0)$ -admissible:  $(B_1 > 0) \subseteq \mathcal{R}$ .

Further, assume that  $X$  is  $(Ba-o,d)$ -complete; and let  $X_0 = (x_n)$  be a  $(Ba-o)$  sequence. Then,



- (91-a)**  $X_0 = (x_n)$  is strongly Picard (modulo  $(d, \mathcal{R}; T)$ )
- (91-b)**  $T$  is fix- $\mathcal{R}$ -asingleton (hence, fix- $\mathcal{R}$ -singleton)
- (91-c)**  $T$  is fix-asingleton (hence, fix-singleton) when, in addition to the conditions above

$$\mathcal{R} \text{ is } (B_0 > 0)\text{-admissible: } (B_0 > 0) \subseteq \mathcal{R}.$$

A basic particular cases of these developments is the one of  $\varphi \in \mathcal{F}(R_+^0, R)$  being a constant function. For both practical and theoretical reasons, it will be useful, in the following, to discuss its technical aspects.

Let  $\psi \in \mathcal{F}(R_+^0, R)$  be a function; and  $C > 0$  be a (real) constant. Then, let  $(u, v, w)$  be the triple of functions introduced as before: for each  $t = (t_1, \dots, t_6) \in R_+^6$

$$u(t) = t_5, v(t) = w(t) = t_2.$$

We have, according to our conventions

$$\mathcal{Y}_0 = \{t \in R_+^6; t_5, t_2 > 0\}, \mathcal{Y} = \{t \in \mathcal{Y}_0; \psi(t_5) \leq \psi(t_2) - C\}.$$

In addition, the attached relations to be used here are (see above)

$$\mathcal{R} = \mathcal{D}^{-1}(\mathcal{Y}_0), \mathcal{R}^* = \mathcal{D}^{-1}(\mathcal{Y}).$$

The contractive property attached to these data writes, in a shorter way

(fct-contr)  $T$  is  $(d, \mathcal{R}; \psi, C)$ -contractive, if  $\psi(d(Tx, Ty)) \leq \psi(d(x, y)) - C$ , when  $(x, y) \in \mathcal{R}$ .

As before, we impose a regularity condition like

$$\mathcal{R} \text{ is } B_1\text{-admissible: } (B_1 > 0) \subseteq \mathcal{R}.$$

This, combined with the functional contractive condition, yields

$$\begin{aligned} \text{(B1-contr) } T \text{ is } (d, (B_1 > 0); \psi, C)\text{-contractive :} \\ \psi(d(Tx, Ty)) \leq \psi(d(x, y)) - C, \text{ when } (x, y) \in (B_1 > 0); \end{aligned}$$

which is exactly the contractive condition we are dealing with. Note that, under a stronger regularity condition like

$$\mathcal{R} \text{ is } B_0\text{-admissible: } (B_0 > 0) \subseteq \mathcal{R}.$$

the starting functional contractive condition yields

$$\begin{aligned} & \text{(B0-contr)} \quad T \text{ is } (d, (B_0 > 0); \psi, C)\text{-contractive} : \\ & \psi(d(Tx, Ty)) \leq \psi(d(x, y)) - C, \text{ when } (x, y) \in (B_0 > 0). \end{aligned}$$

Moreover, in this  $B_0$ -admissible setting, the following inclusion holds

$$\begin{aligned} & \text{(B0-B1)} \quad T \text{ is } (d, (B_0 > 0); \psi, C)\text{-contractive implies} \\ & T \text{ is } (d, (B_1 > 0); \psi, C)\text{-contractive.} \end{aligned}$$

But, the reverse inclusion does not hold, in general.

Formally, the conditions to be imposed upon the couple  $(\psi, C)$  are directly obtainable from the ones of Theorem 12, under the choice  $(\varphi(t) = C; t > 0)$ . In particular, the strictly positive and right asymptotic conditions upon  $\varphi$  are fulfilled here. Concerning the remaining ones, a basic instance when these hold writes

$$\text{(s-zero-abr)} \quad \psi \text{ is strongly zero abrupt: for each sequence } (t_n) \text{ in } R_+^0, (\psi(t_n) \rightarrow -\infty \text{ iff } t_n \rightarrow 0).$$

Putting these together, the following version of Theorem 12 is to be noted.

**Theorem 13** Assume that  $T$  is  $(d, \mathcal{R}; \psi, C)$ -contractive, where  $C > 0$  is a constant,  $\psi$  is a function in  $\mathcal{F}(R_+^0, R)$  with

- (92-i)  $\psi$  is almost right regulated and strongly zero abrupt, and the associated relation  $\mathcal{R}$  fulfills
- (92-ii)  $\mathcal{R}$  is  $(B_1 > 0)$ -admissible:  $(B_1 > 0) \subseteq \mathcal{R}$ .

Further, assume that  $X$  is  $(Ba-o, d)$ -complete; and let  $X_0 = (x_n)$  be a  $(Ba-o)$  sequence. Then,

- (92-a)  $X_0 = (x_n)$  is strongly Picard (modulo  $(d, \mathcal{R}; T)$ )
- (92-b)  $T$  is fix- $\mathcal{R}$ -asingleton (hence, fix- $\mathcal{R}$ -singleton)
- (92-c)  $T$  is fix-asingleton (hence, fix-singleton) when, in addition to the conditions above

$$\mathcal{R} \text{ is } (B_0 > 0)\text{-admissible: } (B_0 > 0) \subseteq \mathcal{R}.$$

**Proof** As precise,  $\varphi = C$  is strictly positive and right asymptotic positive. It will suffice establishing that the remaining conditions of Theorem 12 are working here to complete the argument.

**Part 1.** We claim that, necessarily,

$$\text{(tele-admi)} \quad (\psi, C) \text{ is tele admissible: each sequence } (t_n) \text{ in } R_+^0 \text{ with } (\psi(t_{n+1}) \leq \psi(t_n) - C, \forall n) \text{ fulfills } \lim_n(t_n) = 0$$

In fact, let  $(t_n)$  be as before. By a summation procedure, one arrives at

$$\psi(t_n) \leq \psi(t_0) - nC, \text{ for all } n.$$

Passing to limit as  $n \rightarrow \infty$ , yields

$$\psi(t_n) \rightarrow -\infty; \text{ hence, } t_n \rightarrow 0,$$

if we take the strongly zero abrupt property of  $\psi$  into account; hence the assertion.

**Part 2.** We claim that under the admitted hypotheses

(zero-her)  $(\psi, C)$  is *zero hereditary*: there are no couple of sequences  $(t_n)$  and  $(s_n)$  in  $R_+^0$  and no points  $b > 0$  with  $(\psi(t_n) \leq \psi(s_n) - C, \forall n)$  and  $t_n \rightarrow b, s_n \rightarrow 0$ .

Suppose by contradiction that this is not true: there exists a couple of sequences  $(t_n)$  and  $(s_n)$  in  $R_+^0$  and some point  $b > 0$  with

(non-zh-1)  $\psi(t_n) \leq \psi(s_n) - C, \text{ for all } n \geq 0$

(non-zh-2)  $t_n \rightarrow b, s_n \rightarrow 0, \text{ as } n \rightarrow \infty.$

The second half of (non-zh-2) yields

$$\psi(s_n) \rightarrow -\infty \text{ (if we remember that } \psi = \text{strongly zero abrupt).}$$

But then, taking (non-zh-1) into account,

$$\psi(t_n) \rightarrow -\infty; \text{ hence, } t_n \rightarrow 0, \text{ in view of } \psi = \text{strongly zero abrupt.}$$

The obtained fact is in contradiction with the hypothesis  $t_n \rightarrow b > 0$ . Hence, our working assumption cannot be accepted; and this assertion is valid too.

Summing up, conditions of Theorem 12 are fulfilled by these data; and, from this, the conclusions in the statement follow.

A basic version of the obtained result is obtainable under the lines below. Let us say that  $\psi : R_+^0 \rightarrow R$  is a *Wardowski function*, if

$$\psi \text{ is increasing and } \psi(0 + 0) = -\infty.$$

By a previous observation, we have that any such function is almost right regulated. On the other hand, by an auxiliary statement in Turinici [39], any Wardowski function is strongly zero abrupt. By Theorem 13 we then have the following fixed point statement with a methodological meaning.

**Theorem 14** *Assume that  $T$  is  $(d, \mathcal{R}; \psi, C)$ -contractive, where  $C > 0$  is a (real) constant,  $\psi \in \mathcal{F}(R_+^0, R)$  is taken so as*

(93-i)  $\psi$  is Wardowski (see above)

*and the associated relation  $\mathcal{R}$  fulfills*

(93-ii)  $\mathcal{R}$  is  $(B_1 > 0)$ -admissible:  $(B_1 > 0) \subseteq \mathcal{R}$ .

*Further, assume that  $X$  is  $(Ba-o,d)$ -complete; and let  $X_0 = (x_n)$  be a  $(Ba-o)$  sequence. Then,*

- (93-a)  $X_0 = (x_n)$  is strongly Picard (modulo  $(d, \mathcal{R}; T)$ )  
 (93-b)  $T$  is fix- $\mathcal{R}$ -singleton (hence, fix- $\mathcal{R}$ -singleton)  
 (93-c)  $T$  is fix-singleton (hence, fix-singleton) when, in addition to the conditions above

$\mathcal{R}$  is  $(B_0 > 0)$ -admissible:  $(B_0 > 0) \subseteq \mathcal{R}$ .

The obtained statement may be viewed as a technical improvement of the 2012 fixed point principle in Wardowski [45]; based, among others, on the extra condition

$\psi$  is subunitary power compatible:  $\lim_{t \rightarrow 0^+} t^k \psi(t) = 0$ , for some  $k \in ]0, 1[$ .

This improvement is essentially related to

- (impr-1) the  $(B_1 > 0)$  setting of Theorem 14 is an effective extension of the  $(B_0 > 0)$  setting of the quoted statement  
 (impr-2) the subunitary power compatible condition appearing in the quoted article is avoided in our statement.

For a different proof of Theorem 7 we refer to the paper by Turinici [40]. Some other partial aspects were discussed in Secelean [34].

## References

1. M. Akkouchi, Common fixed points for weakly compatible maps satisfying implicit relations without continuity. *Dem. Math.* **44**, 151–158 (2011)
2. S. Banach, Sur les opérations dans les ensembles abstraits et leur application aux équations intégrales. *Fund. Math.* **3**, 133–181 (1922)
3. V. Berinde, F. Vetro, Common fixed points of mappings satisfying implicit contractive conditions. *Fixed Point Theory Algorithms Sci. Eng.* **2012**, Article ID 105 (2012)
4. P. Bernays, A system of axiomatic set theory: Part III. Infinity and enumerability analysis. *J. Symb. Log.* **7**, 65–89 (1942)
5. D.W. Boyd, J.S.W. Wong, On nonlinear contractions. *Proc. Amer. Math. Soc.* **20**, 458–464 (1969)
6. S. Chandok, B.S. Choudhury, N. Metiya, Fixed point results in ordered metric spaces for rational type expressions with auxiliary functions. *J. Egyptian Math. Soc.* **23**, 95–101 (2015)
7. B.S. Choudhury, P. Konar, B.E. Rhoades, N. Metiya, Fixed point theorems for generalized weakly contractive mapping. *Nonlinear Anal.* **74**, 2116–2126 (2011)
8. P.J. Cohen, *Set Theory and the Continuum Hypothesis* (Benjamin, New York, 1966)
9. M. Cosentino, P. Vetro, Fixed point results for  $F$ -contractive mappings of Hardy-Rogers type. *Filomat* **28**, 715–722 (2014)
10. J. Dieudonné, *Foundations of Modern Analysis* (Academic Press, New York, 1960)
11. N.V. Dung, V. I. Hang, A fixed point theorem for generalized  $F$ -contraction on complete metric space. *Vietnam J. Math.* **43**, 743–755 (2015)
12. P.N. Dutta, B.S. Choudhury, A generalisation of contraction principle in metric spaces. *Fixed Point Theory Algorithms Sci. Eng.* **2008**, Article ID 406368 (2008)
13. P. Hitzler, Generalized metrics and topology in logic programming semantics. PhD Thesis, National University Ireland, University College Cork (2001)

14. S. Kasahara, On some generalizations of the Banach contraction theorem. *Publ. Res. Inst. Math. Sci.* **12**, 427–437 (1976)
15. M.S. Khan, M. Swaleh, S. Sessa, Fixed point theorems by altering distances between the points. *Bull. Aust. Math. Soc.* **30**, 1–9 (1984)
16. S. Leader, Fixed points for general contractions in metric spaces. *Math. Japon.* **24**, 17–24 (1979)
17. J. Matkowski, *Integrable Solutions of Functional Equations*. *Dissertationes Mathematicae*, vol. 127 (Polish Sci. Publ., Warsaw, 1975)
18. A. Meir, E. Keeler, A theorem on contraction mappings. *J. Math. Anal. Appl.* **28**, 326–329 (1969)
19. G. H. Moore, *Zermelo's Axiom of Choice: its Origin, Development and Influence* (Springer, New York, 1982)
20. Y. Moskhovakis, *Notes on Set Theory* (Springer, New York, 2006)
21. S.B. Nadler Jr., Multi-valued contraction mappings. *Pacific J. Math.* **30**, 475–488 (1969)
22. H.K. Nashine, Z. Kadelburg, P. Kumam, Implicit-relation-type cyclic contractive mappings and applications to integral equations. *Abstr. Appl. Anal.* **2012**, Article ID 386253 (2012)
23. I.P. Natanson, *Theory of Functions of a Real Variable (Volume I)* (Frederick Ungar Publishing Co., New York, 1964)
24. L. Pasicki, A strong fixed point theorem *Topology Appl.* **282**, 107–130 (2020)
25. V. Popa, Some fixed point theorems for compatible mappings satisfying an implicit relation. *Dem. Math.* **32**, 157–163 (1999)
26. V. Popa, M. Mocanu, Altering distance and common fixed points under implicit relations. *Hacetetepe J. Math. Stat.* **38**, 329–337 (2009)
27. O. Popescu, G. Stan, Two fixed point theorems concerning  $F$ -contraction in complete metric spaces. *Symmetry* **2020**, 12, 58 (2019)
28. S. Radenović, Z. Kadelburg, D. Jandrić and A. Jandrić, Some results on weakly contractive maps. *Bull. Iranian Math. Soc.* **38**, 625–645 (2012)
29. S. Reich, Fixed points of contractive functions. *Boll. Un. Mat. Ital.* **5**, 26–42 (1972)
30. B.E. Rhoades, A comparison of various definitions of contractive mappings. *Trans. Amer. Math. Soc.* **226**, 257–290 (1977)
31. B.E. Rhoades, Some theorems on weakly contractive maps. *Nonlinear Anal.* **47**, 2683–2693 (2001)
32. I.A. Rus, *Generalized Contractions and Applications* (Cluj University Press, Cluj-Napoca, 2001)
33. E. Schechter, *Handbook of Analysis and its Foundation* (Academic Press, New York, 1997)
34. N.-A. Secelean, Weak  $F$ -contractions and some fixed point results. *Bull. Iran. Math. Soc.* **42**, 779–798 (2016)
35. A. Tarski, Axiomatic and algebraic aspects of two theorems on sums of cardinals. *Fund. Math.* **35**, 79–104 (1948)
36. M. Turinici, Fixed points of implicit contraction mappings. *An. Șt. Univ. A. I. Cuza Iași (S I-a, Mat)* **22**, 177–180 (1976)
37. M. Turinici, Fixed points of implicit contractions via Cantor's intersection theorem. *Bul. Inst. Polit. Iași (Sect I: Mat., Mec. Teor., Fiz.)* **26**(30), 65–68 (1980)
38. M. Turinici, Function pseudometric VP and applications. *Bul. Inst. Polit. Iași (Sect.: Mat., Mec. Teor., Fiz.)* **53**(57), 393–411 (2007)
39. M. Turinici, Wardowski implicit contractions in metric spaces. *Arxiv*, 1211–3164-v2 (2013)
40. M. Turinici, *Implicit Contractive Maps in Ordered Metric Spaces*. *Topics in Mathematical Analysis and Applications*, ed. by Th.M. Rassias, L. Toth (Springer International Publishing, Switzerland, 2014), pp. 715–746
41. M. Turinici, Contractive maps in locally transitive relational metric spaces. *The Sci. World J.* **2014**, Article ID 169358 (2014)
42. M. Turinici, Wardowski mappings and Meir-Keeler property, in *Modern Directions in Metrical Fixed Point Theory*, Paper 1–5 (Pim Editorial House, Iași, 2016)

43. M. Turinici, Rational implicit contractions in ordered metric spaces, in *Selected Topics in Metrical Fixed Point Theory* (Revised Ed.), Paper 3-1 (Pim Editorial House, Iași, 2017)
44. J. Vujaković, S. Mitrović, M. Pavlović, S. Radenović, On recent results concerning  $F$ -contraction in generalized metric spaces. *Mathematics* **8**, 767 (2020)
45. D. Wardowski, Fixed points of a new type of contractive mappings in complete metric spaces. *Fixed Point Theory Algorithms Sci. Eng.* **2012**, 94 (2012)
46. D. Wardowski, Solving existence problems via  $F$ -contractions. *Proc. Amer. Math. Soc.* **146**, 1585–1598 (2018)
47. E.S. Wolk, On the principle of dependent choices and some forms of Zorn's lemma. *Canad. Math. Bull.* **26**, 365–367 (1983)

# Nonlinear Dynamics of the KdV-B Equation and Its Biomedical Applications



Michail A. Xenos and Anastasios C. Felias

**Abstract** In recent years there is an incremental degree of bridging open questions in biomechanics with the help of applied mathematics and nonlinear analysis. Recent advancements concerning the cardiac dynamics pose important questions about the cardiac waveform. A governing equation, namely the KdV-B equation (Korteweg–de Vries–Burgers),

$$\frac{\partial u}{\partial t} + \gamma u \frac{\partial u}{\partial x} - \alpha \frac{\partial^2 u}{\partial x^2} + \beta \frac{\partial^3 u}{\partial x^3} = 0, \quad u = u(t, x), \quad \alpha, \beta, \gamma \in \mathbb{R}, \quad (1)$$

is a partial differential equation utilized to answer several of those questions. The cardiac dynamics mathematical model features both solitary and shock wave characteristics due to the dispersion and dissipation terms, as occurring in the arterial tree. In this chapter a focus is given on describing cardiac dynamics. It is customarily difficult to solve nonlinear problems, especially by analytical techniques. Therefore, seeking suitable solving methods, exact, approximate or numerical, is an active task in branches of applied mathematics. The phase plane of the KdV–B equation is analyzed and its qualitative behavior is derived. An asymptotic expansion is presented and traveling wave solutions under both shock and solitary profiles are sought. Numerical solutions are obtained for the equation, by means of the Spectral Fourier analysis and are evolved in time by the Runge–Kutta method. This whole analysis provides vital information about the KdV–B equation and its connection to cardiac hemodynamics. The applications of KdV–B, presented in this chapter, highlight its essence to human hemodynamics.

**Mathematics Subject Classification (2020)** 35Q35, 35Q53

---

M. A. Xenos (✉) · A. C. Felias

Department of Mathematics, University of Ioannina, Ioannina, Greece

e-mail: [mxenos@uoi.gr](mailto:mxenos@uoi.gr)

© Springer Nature Switzerland AG 2021

Th. M. Rassias (ed.), *Nonlinear Analysis, Differential Equations, and Applications*,

Springer Optimization and Its Applications 173,

[https://doi.org/10.1007/978-3-030-72563-1\\_26](https://doi.org/10.1007/978-3-030-72563-1_26)

## 1 Introduction

### 1.1 Background Information for KdV–B

In the last few decades, much attention from a rather diverse group of scientists such as physicists, engineers and applied mathematicians has been attracted to two contrasting themes: (a) the theory of dynamical systems, most popularly associated with the study of chaos, and (b) the theory of integrable (or nonintegrable) systems associated, among other things, with the study of solitary waves.

It is common knowledge that many physical phenomena, such as nonlinear shallow-water waves and wave motion in plasma, can be described by the Korteweg–de Vries (KdV) equation [29]. It is well known that solitons and solitary waves are the class of special solutions of the KdV equation. In order to study propagation of undular bores in shallow water [6, 27], liquid flow containing gas bubbles [54], fluid flow in elastic tubes [28], crystal lattice theory, nonlinear circuit theory and turbulence [20, 30, 51], the governing equation can be reduced to the so-called Korteweg–de Vries–Burgers equation (KdV–B) as follows [10],

$$\frac{\partial u}{\partial t} + \gamma u \frac{\partial u}{\partial x} - \alpha \frac{\partial^2 u}{\partial x^2} + \beta \frac{\partial^3 u}{\partial x^3} = 0, \quad u = u(t, x), \quad \alpha, \beta, \gamma \neq 0. \quad (2)$$

This is a nonintegrable equation in the sense that its spectral problem is nonexistent [19]. Multiplying  $t$ ,  $x$  and  $u$ , by constants can be used to make the coefficients of any of the above four terms equal to any given nonzero constant. Therefore, we focus on the case where  $\alpha \geq 0$ ,  $\beta > 0$  and  $\gamma \neq 0$ .

This equation is equivalent to the KdV equation with the addition of a viscous dissipation term ( $\alpha \frac{\partial^2 u}{\partial x^2}$ ). The studies of the KdV equation [29] ( $\alpha = 0$ ) and the Burgers equation [9] ( $\beta = 0$ ) have been undertaken, but the exact solution for the general case of equation (2) ( $\alpha \geq 0$ ,  $\beta > 0$ ,  $\gamma \neq 0$ ) has still not been completed.

### 1.2 Biomechanical Applications

Solitons are mathematical entities appearing as solutions of nonlinear wave equations [8]. They are waves of stable and steady form, although internal oscillations may occur, exhibiting unique characteristics when colliding with other solitary waves as described by Ablowitz and Segur [1]. During the last decade, soliton profiles are found when studying nonlinear optics, condensed matter Physics and quantum theory of matter and gravity [43]. Lately, an increasing number of studies focuses on describing the cardiac pulse as a soliton, due to the features those two seem to share. The pulsatility synchronization of the smooth arterial muscle allows the consideration of solitary profiles in cardiac hemodynamics [34].



Theoretical investigation for the blood waves have been developed by many researchers through the use of weakly nonlinear theories. The theoretical investigation of pulse wave propagation in human arteries has a long history starting from ancient times until today. Over the past decade, the scientific efforts have been concentrated on theoretical investigations of nonlinear wave propagation in arteries with a variable radius. The question “How local imperfections appeared in the artery can disturb the arterial wall deformation?” is important for understanding the nature and main features of various cardiovascular diseases, such as stenoses and aneurysms. Rowlands (1982) reported some extraordinary features of the cardiac pulse, leading to his conception of the arterial flows as a solitary motion [44]. A few years later, Otwinowski and collaborators presented a nonlinear differential equation whose solutions exhibited similar characteristics with those reported by Rowlands [38].

Based on those evolutionary theories, adding the inertial behavior of blood vessel in an one-dimensional cardiovascular model, researchers concluded that the KdV equation is a seemingly reliable tool in modelling cardiac dynamics. It was supported that the solitary wave formulation fits much better in describing the arterial pulse wave experimental results than the wave equation proposed by the majority of researchers [57]. An additional reason to support the above formulation is the peaking and steepening features of the pressure pulse, which coincide with the structure of soliton profiles of KdV [11].

The majority of studies on the wave propagation in blood flow is mainly based on linear waves. The linearized theories proposed by Resal, Witzig, Womersly, McDonald and others, consider the vessel as a straight, infinite, circular elastic tube filled with an isotropic and Newtonian fluid, blood [55]. Blood is studied as an incompressible fluid, a characterization justified by its compressibility being rather insignificant, compared to the dilation of the blood vessels. In 1958, Lambert based on the Euler equations of fluid motion, proposed the Method of Characteristics for the calculations concerning the nonlinear blood flow. All theories presented to model nonlinear blood flow are one-dimensional, meaning that both pressure and flow velocity are seen as functions of the axial distance along the vessel in time. Contributions in nonlinear modulation were done by Rudinger, Skalak, Rockwell, Hawley and Anliker. The suggested equations are basically the equations of continuity and motion coupled with an extra equation to describe the vessel wall distensibility [3, 43, 45, 49]. Sakanishi and Hasegawa proposed a soliton profile pulsatile wave modeling, based on the nonlinear elasticity of the vessel wall [46]. Yomosa and collaborators proposed a theory describing solitons in long arteries, where the viscous effects, the reflective effects caused by the arterial branch as well as the effects of the peripheral resistance are neglectible. For the above reasons, the latter modulation is unable to describe the pressure drop caused when moving away from the heart. Nevertheless, it points out that it does make some sense to attribute the special features of the pulsatile wave, including the “sudden steepening” and the change in the phase velocity, to the solitary profile [57]. While the pulsatile wave travels to arteries with smaller radius, viscosity seems to play a vital role in both the flow decay and the widening of the wave width [57].

Antar and Demiray studied the propagation of weak nonlinear waves in a thin elastic tube, under an initial stress distribution, due to the flow of an incompressible viscous fluid [4]. The propagation of pressure pulses in dilatable tubes has been studied by various researchers [22, 41]. Most of those studies, consider waves of small width, neglecting the nonlinear characteristics and focusing on their dispersive character [5, 12, 42]. It is widely accepted that a long-term evolution of weak nonlinear waves of either dispersion or dissipation, can be modeled by nonlinear dispersive equations. Two classical simplified and indicative examples are the Burgers equation and the KdV equation, exhibiting balance between nonlinearity and dissipation and nonlinearity and dispersion, respectively. On the other hand, when a balance is exhibited among nonlinearity, dissipation and dispersion, the simplest and most representative dispersive equation is the KdV–B equation, combining the KdV and Burgers equations. Via asymptotic methods, the propagation of small, but with finite width, waves in dilatable tubes has been studied sufficiently [4].

Hashizume and Yomosa showed that propagation, in the case of weak nonlinear waves in a thin and nonlinear elastic tube for incompressible flow, is determined by the KdV equation [57]. Erbay and collaborators, examining the propagation of weak nonlinear waves in a thin viscoelastic tube filled with fluid, were lead to the Burgers, KdV and KdV–B equations, depending on the parameters considered [15]. Demiray studied the propagation of slightly nonlinear waves in thin elastic and viscoelastic tubes for an incompressible fluid and finally concluded to the KdV and KdV–B equations, respectively. In all the above studies, an inviscid fluid was considered and the axial movement of the tube wall was neglected. However, regarding biological applications, blood is an incompressible and viscous fluid. So, Antar and Demiray formulated their mathematical model toward this direction [4].

In this chapter, an emphasis is given to the theoretical and numerical analysis of the KdV–B equation and its applications, providing vital information about the KdV–B equation and its connection to cardiac hemodynamics. More precisely, in the next section the phase plane of the KdV–B equation is analyzed and its qualitative behavior is derived. Furthermore, an asymptotic expansion is presented and traveling wave solutions under both shock and solitary profiles are derived. Finally, numerical solutions are obtained for the KdV–B equation, by means of spectral Fourier analysis and are evolved in time by the well known explicit 4th order Runge–Kutta method.

## 2 Phase Plane Analysis of KdV–B

In this section, the phase plane of the KdV–B equation is analyzed and its qualitative behavior is derived and further described. The wave variable  $\zeta$  is introduced as [25, 47],

$$\zeta = x - \lambda t, \quad (3)$$

with  $\lambda$  being the wave velocity. Then equation (2) using (3), can be written as,

$$(\gamma u - \lambda) \frac{du}{d\zeta} - \alpha \frac{d^2u}{d\zeta^2} + \beta \frac{d^3u}{d\zeta^3} = 0, \quad u = u(t, x) = u(x - \lambda t) = u(\zeta). \tag{4}$$

The so-called traveling-wave solution,  $u = u(\zeta)$ , shall be considered here. By integrating equation (4) with respect to  $\zeta$ , a nonlinear differential equation can be obtained as follows,

$$\frac{d^2u}{d\zeta^2} + c_1 \frac{du}{d\zeta} + c_2u^2 + c_3u = c_0, \tag{5}$$

where  $c_1 = -\frac{\alpha}{\beta}$ ,  $c_2 = \frac{\gamma}{2\beta}$ ,  $c_3 = -\frac{\lambda}{\beta}$  and the integral constant  $c_0 > -\frac{\lambda^2}{2\beta}$ .

In the case where  $c_0 \neq 0$ , a simple translation transformation,

$$u = u' + c'_0, \quad c'_0 = \frac{-c_3 \pm \sqrt{c_3^2 + 4c_0c_2}}{2c_2},$$

can be made, with  $u'$  satisfying the following equation,

$$\frac{d^2u'}{d\zeta^2} + c_1 \frac{du'}{d\zeta} + c_2u'^2 + (c_3 + 2c_2c'_0)u' = 0.$$

Without loss of generality, we shall confine ourselves to the consideration of  $c_0 = 0$  alone from now on. It can be further assumed that  $\lambda \geq 0$ , because the discussion on  $\lambda' = -\lambda$  can be made in the same manner for  $\lambda < 0$ .

Equation (5) can be written as an autonomous system of first-order equations,

$$\begin{cases} \frac{du}{d\zeta} = v, \\ \frac{dv}{d\zeta} = -\frac{u}{\beta}(\gamma \frac{u}{2} - \lambda) + \frac{\alpha}{\beta}v. \end{cases}$$

Now, we study the above system according to the qualitative theory of ordinary differential equations. Initially, we find the system's singular points, setting,

$$f_1(u, v) = v, \quad f_2(u, v) = -\frac{u}{\beta}(\gamma \frac{u}{2} - \lambda) + \frac{\alpha}{\beta}v.$$

The following conditions should be met,

$$f_1(u, v) = f_2(u, v) = 0$$

$$\Leftrightarrow \begin{cases} v = 0 \\ -\frac{u}{\beta}(\gamma \frac{u}{2} - \lambda) + \frac{\alpha}{\beta}v = 0 \end{cases}$$

$$\Leftrightarrow \begin{cases} v = 0 \\ u = 0, \quad u = \frac{2\lambda}{\gamma}. \end{cases}$$

Therefore, the singular points are,

$$\begin{cases} P_1 = (0, 0), \\ P_2 = (\frac{2\lambda}{\gamma}, 0). \end{cases}$$

Next, we are to find the eigenvalues of the linearization matrices, defined for our singular points, as follows,

$$A(P_1) = \begin{bmatrix} \frac{\partial f_1}{\partial u} & \frac{\partial f_1}{\partial v} \\ \frac{\partial f_2}{\partial u} & \frac{\partial f_2}{\partial v} \end{bmatrix}, \quad (P_1) = \begin{bmatrix} 0 & 1 \\ \frac{\lambda}{\beta} & \frac{\alpha}{\beta} \end{bmatrix},$$

so, for its eigenvalues we get,

$$\det(A_{P_1} - sI_2) = 0 \Leftrightarrow s^2 - \frac{\alpha}{\beta}s - \frac{\lambda}{\beta} = 0$$

$$\Leftrightarrow \begin{cases} s_1 = \frac{\frac{\alpha}{\beta} + \frac{1}{\beta}\sqrt{\alpha^2 + 4\lambda\beta}}{2} > 0 \\ s_2 = \frac{\frac{\alpha}{\beta} - \frac{1}{\beta}\sqrt{\alpha^2 + 4\lambda\beta}}{2} < 0 \end{cases}.$$

$$A(P_2) = \begin{bmatrix} \frac{\partial f_1}{\partial u} & \frac{\partial f_1}{\partial v} \\ \frac{\partial f_2}{\partial u} & \frac{\partial f_2}{\partial v} \end{bmatrix}, \quad (P_2) = \begin{bmatrix} 0 & 1 \\ -\frac{\lambda}{\beta} & \frac{\alpha}{\beta} \end{bmatrix},$$

so, for its eigenvalues we get,

$$\det(A_{P_2} - sI_2) = 0 \Leftrightarrow s^2 - \frac{\alpha}{\beta}s + \frac{\lambda}{\beta} = 0$$

$$\Leftrightarrow \begin{cases} s_1 = \frac{\frac{\alpha}{\beta} + \frac{1}{\beta}\sqrt{\alpha^2 - 4\lambda\beta}}{2} > 0 \\ s_2 = \frac{\frac{\alpha}{\beta} - \frac{1}{\beta}\sqrt{\alpha^2 - 4\lambda\beta}}{2} > 0 \end{cases}, \quad \alpha \geq 2\sqrt{\lambda\beta}.$$

$$\begin{cases} s_1 = \frac{\frac{\alpha}{\beta} + i\frac{1}{\beta}\sqrt{4\lambda\beta - \alpha^2}}{2} \\ s_2 = \frac{\frac{\alpha}{\beta} - i\frac{1}{\beta}\sqrt{4\lambda\beta - \alpha^2}}{2} \end{cases}, \quad \alpha \in (0, 2\sqrt{\lambda\beta}).$$

$$\begin{cases} s_1 = i\sqrt{\frac{\lambda}{\beta}} \\ s_2 = -i\sqrt{\frac{\lambda}{\beta}} \end{cases} \quad \alpha = 0.$$

We conclude that  $(0, 0)$  is invariably a saddle point, whereas  $(\frac{2\lambda}{\gamma}, 0)$  has three cases depending on the values of  $\alpha, \beta, \lambda$  [37],

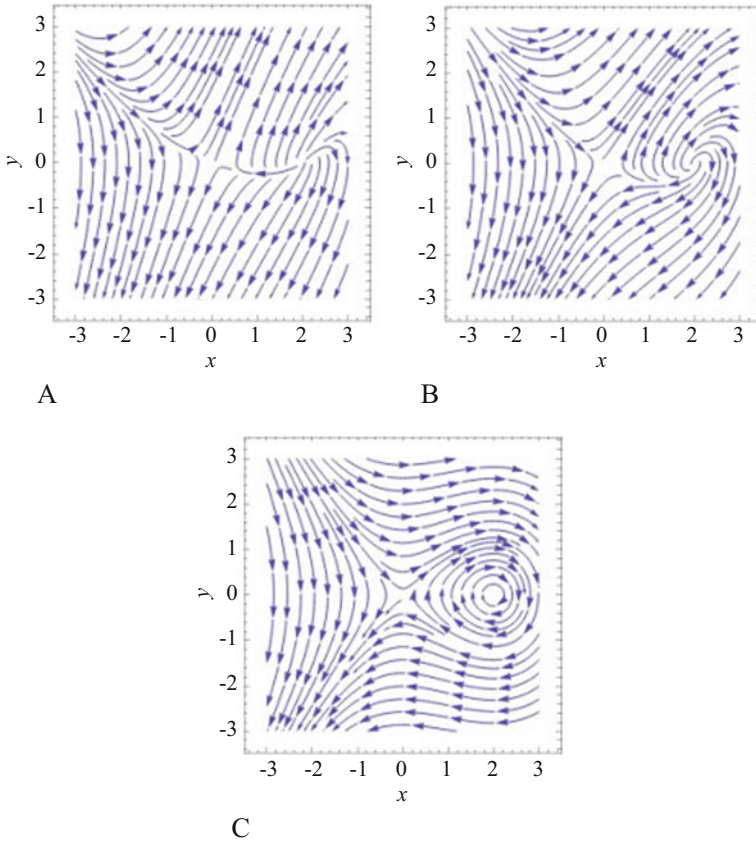
- A. a source for  $\alpha \geq 2\sqrt{\lambda\beta}$ ,
- B. a spiral source for  $\alpha \in (0, 2\sqrt{\lambda\beta})$ ,
- C. a central point for  $\alpha = 0$  (KdV).

Regarding the geometric nature of the above characterizations, we have the following [37],

1.  $(0, 0)$  being a saddle point, means that it's an unstable node and phase trajectories tend to move around it in hyperbolas, defined by the separatrices (i.e. straight lines directed along the two eigenvectors of the linearization matrix).
2.  $((\frac{2\lambda}{\gamma}, 0) : \alpha \geq 2\sqrt{\lambda\beta})$  being a source, means that it's an unstable node from where phase trajectories diverge away without any (or relatively little) rotation.
3.  $((\frac{2\lambda}{\gamma}, 0) : \alpha \in (0, 2\sqrt{\lambda\beta}))$  being a spiral source, means that it's an unstable focus where phase trajectories tend to spiral around before eventually diverge away from it.
4.  $((\frac{2\lambda}{\gamma}, 0) : \alpha = 0)$  being a central point, means that the phase trajectories tend to move in ellipses around the point, describing periodic motion of a point in the phase space.

The phase plots of Fig. 1 depict our three cases, where we have set for convenience  $\gamma = 1$ , since it is not related to any of  $\alpha, \beta, \lambda$  in effecting the stability statuses.

An essential tool in studying the phase portrait of nonlinear autonomous systems, like the above, is the Hartman–Grobman Theorem [22, 23, 37],



**Fig. 1** The point (0, 0) is invariably a saddle point whereas (2, 0) is a source point in **a**, a spiral source point in **b** and a central point in **c**. (**a**) ( $\alpha = 2, \beta = \lambda = 1$ ). (**b**) ( $\alpha = \beta = \lambda = 1$ ). (**c**) ( $\alpha = 0, \beta = \lambda = 1$ )

**Theorem 1 (Hartman-Grobman)** Consider a two-dimensional nonlinear autonomous system with a continuously differentiable field  $\bar{f}$ ,

$$\bar{x}' = \bar{f}(\bar{x})$$

and consider its linearization at a hyperbolic critical point  $\bar{x}_0$  (that is the Jacobian matrix has eigenvalues with non-zero real part),

$$\bar{u}' = (Df_0)(\bar{u}).$$

Then there is a neighborhood of the hyperbolic critical point where all the solutions of the linear system can be transformed into solutions of the nonlinear system by a continuous, invertible transformation.

*Remark 1* The above theorem implies that the phase portrait of the linear system in a neighborhood of the hyperbolic critical point can be transformed to the phase portrait of the nonlinear system by a continuous, invertible transformation. When that happens, we say that the two phase portraits are topologically equivalent.

Additional information about the phase plane of KdV–B equation can be found in [14].

### 3 Asymptotic Expansion for KdV–B

In the study of ordinary differential equations and their applications, an asymptotic expansion is of high importance. It would be very useful to understand thoroughly the property of the solution to the KdV–B equation. The asymptotic expansion would provide a reliable basis for estimating the advantages and disadvantages when seeking and applying numerical methods to our equation.

Here, by means of variable transformation and the qualitative theory of ordinary differential equations, the asymptotic behavior of the traveling wave solutions to the KdV–B equation is presented. The asymptotic expansion is real and continuous, if the argument is greater than a certain value.

The following variable transformation can be made [36, 48],

$$u = -e^{-\frac{c_1(1-k)\xi}{2} - \frac{c_1k^2}{c_2}y(\xi)}, \quad \xi = e^{-c_1k\xi}, \quad k = \sqrt{1 - \frac{4c_3}{c_1^2}} = \sqrt{1 + \frac{4\beta\lambda}{\alpha^2}} \geq 1. \tag{6}$$

Equation (5) ( $c_0 = 0$ ), can be reduced to the Emden–Fowler equation [35],

$$\frac{d^2y}{d\xi^2} = \xi^\sigma y^2, \quad \sigma = \frac{1 - 5k}{2k}. \tag{7}$$

It is obvious that,

$$\begin{cases} \sigma = -2, & \lambda = 0, \\ \sigma \in [-\frac{5}{2}, -2), & \lambda > 0. \end{cases} \tag{8}$$

Some characteristics of the KdV–B equation can be derived from equation (7).

Next we demonstrate some essential results that will help us in deriving the asymptotic expansion of KdV–B.

First, we show that the KdV–B equation, Equation (5) has finite isolated zero points only. Since  $y$  satisfies (7),  $y''$  does not change sign for  $\xi \in (0, \infty)$ , so Equation (7) has finite zero points only, except that it identically vanishes for some intervals. Equation (7) has finite zero points only. This indicates that the solution of

KdV–B is consistently positive negative or zero for large arguments, which depends upon the condition of infinite point.

Our next important tool is the Integral Rule of asymptotic formulae [31, 36, 48].

**Lemma 1 (Integral Rule of Asymptotic Formulae)** *Let,*

$$\phi(t) \sim f(t),$$

where  $f \neq 0$  and  $f$  does not change in sign. Then,

$$\begin{cases} \int_{t_0}^t \phi(t)dt \sim \int_{t_0}^t f(t)dt, & \text{if } \int_{t_0}^{\infty} |f(t)|dt = \infty, \\ \int_t^{\infty} \phi(t)dt \sim \int_t^{\infty} f(t)dt, & \text{if } \int_t^{\infty} |f(t)|dt < \infty. \end{cases}$$

Following the above result, we demonstrate the character of asymptotic expansion [31, 36, 48].

**Lemma 2 (Character of Asymptotic Expansion)** *If  $f(t) > 0$  and  $f'$  is continuous and non-negative as  $t \geq t_0$ , then,*

$$f' \leq f^{1+\epsilon}$$

for any  $t \geq t_0$  and for any  $\epsilon > 0$ , except perhaps in a set of intervals of finite total length, which depends upon  $\epsilon$ .

The final necessary result will be Hardy’s Theorem [31, 36, 48].

**Theorem 2 (Hardy)** *Any solution of an equation,*

$$\frac{df}{dt} = \frac{P(f, t)}{Q(f, t)},$$

which is continuous for  $t \geq t_0$ , is ultimately monotonic, together with all of its derivatives, and satisfies one of the following relations,

$$f \sim at^b e^{E(t)},$$

or

$$f \sim at^b (\ln t)^c,$$

where  $E(t)$  is a polynomial in time and  $a, b, c$  are constants.

Now, all the above three results can be adopted to derive the asymptotic expansion of KdV–B, for the different values of  $\lambda$ .

*Claim (Shu [48])* Let  $\lambda > 0$ . The negative asymptotic expansion of KdV–B has the following form,



$$u = -\frac{2k^2v^2U_\infty}{\delta}e^{-\frac{(k-1)v\zeta}{2\delta}} - \frac{8k^4v^2U_\infty^2}{(k-1)(3k-1)\delta}e^{-\frac{(k-1)v\zeta}{\delta}} [1+O(1)], \quad \zeta \rightarrow \infty, \tag{9}$$

where  $k = \sqrt{1 + \frac{4\lambda\delta}{v^2}}$  and  $U_\infty > 0$  is a constant.

In order to prove that, we consider  $\lambda > 0$  for which  $\sigma \in (-\frac{5}{2}, -2)$ . If  $u$  has a negative asymptotic expansion then  $y$  has a positive asymptotic expansion. Since,

$$\frac{d^2y}{d\xi^2} = \xi^\sigma y^2 > 0, \quad \xi > 0,$$

$y'$  must be strictly monotonically increasing for  $\xi > 0$  and  $y$  must be a monotone function for large  $\xi$ . Thus  $y'$  has three possible cases as  $\xi \rightarrow \infty$ ,

1.  $y' \rightarrow 0$
2.  $y' \rightarrow y'_0 = const > 0$
3.  $y' \rightarrow \infty$

Let us show that case (2) cannot hold.

If

$$y' \rightarrow y'_0 = const > 0,$$

then,

$$y \sim y'_0 \xi$$

and from equation (7),

$$y'' = y^2 \xi^\sigma \sim y_0'^2 \xi^{\sigma+2} > \frac{1}{2} y_0'^2 \xi^{\sigma+2},$$

whose integration yields,

$$y' > \frac{y_0'^2}{2(\sigma + 3)} \xi^{\sigma+3} \rightarrow \infty,$$

for large  $\xi$ , which leads to a contradiction. Then it will be shown that case (3) leads to a contradiction as well.

If

$$y' \rightarrow \infty,$$

then,

$$y' > M,$$

for large  $\xi$  and some  $M > 0$ , and hence,

$$y > M\xi .$$

Reverting to equation (7),

$$y'' = \xi^\sigma y^2 > M^2 \xi^{\sigma+2},$$

so,

$$y > \frac{M^2}{(\sigma + 3)(\sigma + 4)} \xi^{\sigma+4},$$

for large  $\xi$ . Continuing in this fashion,

$$y > y_0 \xi^5,$$

can be obtained for large  $\xi$  and the constant  $y_0$ . Hence, from equation (7),

$$y'' = \xi^\sigma y^2 > \sqrt{y_0} y^{\frac{3}{2}},$$

for large  $\xi$ . Since  $y'$  is positive,

$$y' y'' > \sqrt{y_0} y^{\frac{3}{2}} y',$$

whose integration yields,

$$y' > \frac{2y_0^{\frac{1}{4}}}{\sqrt{5}} y^{\frac{3}{4}},$$

which is impossible due to Lemma (2).

Consequently, we are left with case (1). Since  $y' < 0$  is strictly monotone increasing for  $\xi > 0$ , and  $y$  is strictly monotone decreasing for  $\xi > 0$ . Since  $y > 0$  for large  $\xi$ ,  $y$  has a finite limit  $U_\infty \geq 0$  as  $\xi \rightarrow \infty$ . Now, let us show that  $U_\infty \neq 0$ . If  $U_\infty = 0$ ,  $y(\xi_0) = \delta > 0$  is set to be small. Since  $y$  is strictly monotone decreasing,

$$\delta = y(\xi_0) = \int_{\xi_0}^\infty \left( \int_t^\infty \tau^\sigma y^2 d\tau \right) dt < \delta^2 \int_{\xi_0}^\infty \left( \int_t^\infty \tau^\sigma d\tau \right) dt$$

or

$$\delta > \frac{(\sigma + 1)(\sigma + 2)}{\xi_0^{\sigma+2}},$$

which leads to the contradiction for  $\delta$  sufficiently small.

Then let,

$$y(\infty) = U_\infty > 0, \quad y(\xi) = U_\infty + O(1), \quad \xi \rightarrow \infty.$$

Then

$$y'(\xi) = - \int_\xi^\infty y'' dt = - \int_\xi^\infty t^\sigma y^2 dt = \frac{U_\infty^2}{\sigma + 1} \xi^{\sigma+1} [1 + O(1)]$$

and thus,

$$y(\xi) = U_\infty - \int_\xi^\infty y' dt = U_\infty + \frac{U_\infty^2}{(\sigma + 1)(\sigma + 2)} \xi^{\sigma+2} [1 + O(1)].$$

The latter proves our claim.

*Claim (Shu [48])* Let  $\lambda = 0$ . The negative asymptotic expansion of the KdV-B equation has the following form,

$$u = -\frac{2kv}{\zeta} e^{\frac{-(k-1)v\zeta}{2\delta}}, \quad \zeta \rightarrow \infty, \tag{10}$$

where  $k = \sqrt{1 + \frac{4\lambda\delta}{v^2}}$ .

For the proof, we consider  $\lambda = 0$ , which gives us  $\sigma = -2$ . If  $u$  has a negative asymptotic expansion,  $y$  has a positive asymptotic expansion. Let  $\xi = e^s$ , obtaining from equation (7),

$$\frac{d^2y}{ds^2} - \frac{dy}{ds} - y^2 = 0. \tag{11}$$

If  $\frac{dy}{ds} = 0$  at  $s_0$ , then,

$$\frac{d^2y}{ds^2} = y^2 > 0$$

and  $y$  can only have a minimum at  $s_0$ . Hence,  $y$  is a monotone function for large  $\xi$ . Thus  $y$  has three possible cases as  $s \rightarrow \infty$ :

1.  $y \rightarrow 0$

2.  $y \rightarrow y_0 = \text{const} > 0$
3.  $y \rightarrow \infty$

Let us show that case (2) cannot hold.

If

$$y \rightarrow y_0 = \text{const} > 0,$$

then

$$\frac{d^2y}{ds^2} - \frac{dy}{ds} \sim y^2.$$

Integrating, we get,

$$\frac{dy}{ds} - y \sim y_0^2 s.$$

Since  $y \rightarrow y_0$ , this implies,

$$\frac{dy}{ds} \sim y_0^2 s$$

from which,

$$y \sim \frac{1}{2} y_0^2 s^2,$$

which contradicts  $y \rightarrow y_0$ . Next, we will show that case (3) is impossible.

If  $y \rightarrow \infty$ , let,

$$p = \frac{dy}{ds}.$$

Then equation (11) becomes,

$$p \frac{dp}{dy} - p - y^2 = 0. \tag{12}$$

Since  $y \rightarrow y_0$ , we have,

$$p \frac{dp}{dy} > 0.$$

Now, Theorem (2) indicates that  $p$  has two possible cases for large  $y$ ,

1.  $p \sim ay^b e^{E(y)}$ ,
2.  $p \sim ay^b (\ln y)^c$ ,

where  $E(y)$  is a polynomial in  $y$  and  $a > 0, b, c$  are constants. We also show that case (1) is impossible.

If  $E(y) \rightarrow -\infty$ , then,

$$p \rightarrow 0, \quad \frac{dp}{dy} \rightarrow 0$$

which leads to a contradiction by referring to equation (12).

If  $E(y) \rightarrow \infty$ , then,

$$p > y^2,$$

for large  $y$ , which contradicts Lemma (2). Hence,

$$E(y) = \text{const} .$$

If  $b > 1$ , then,

$$p > y^{\frac{b+1}{2}},$$

for large  $y$ , which is impossible due to Lemma (2).

If  $b \leq 1$ , then,

$$p \frac{dp}{dy} \sim y^2,$$

is obtained from equation (12). By integration,

$$\frac{1}{2} p^2 \sim \frac{1}{3} y^3,$$

is obtained, so that  $b = \frac{3}{2} > 1$ , which leads to a contradiction. Let us now show that case (2) is also impossible.

If  $b > 1$ , then,

$$p > y^{\frac{b+1}{2}},$$

for large  $y$ , which is impossible due to Lemma (2).

If  $b \leq 1$ , then,

$$\frac{1}{2} p^2 \sim \frac{1}{3} y^3,$$

is obtained, so that  $b = \frac{3}{2} > 1$ , which leads to a contradiction.

Consequently, we are left with case (1), where  $y \rightarrow 0$ . Let  $v = \frac{1}{y}$  and  $w = \frac{dv}{ds}$ , obtaining from equation (12),

$$w \frac{dw}{dv} - \frac{2w^2}{v} - w + 1 = 0. \quad (13)$$

Since  $y \rightarrow 0$ ,  $v \rightarrow \infty$  and  $\frac{dv}{ds} < 0$ , we have,

$$w = \frac{dv}{ds} = -\frac{1}{y^2} \frac{dy}{ds} > 0,$$

is obtained. Theorem (2) indicates that  $w$  has two possible cases for large  $v$ ,

1.  $w \sim av^b e^{E(v)}$ ,
2.  $w \sim av^b (\ln v)^c$ ,

where  $E(v)$  is a polynomial in  $v$  and  $a > 0$ ,  $b$ ,  $c$  are constants. It is now shown that if case (3) is satisfied,  $E(v) = \text{const}$  and  $b = 0$ . Similar to above,  $E(v) = \text{const}$  and  $b \leq 1$ .

If  $b = 1$ , then,

$$\frac{dw}{dv} \sim a > 0.$$

From equation (13),  $a = -1$  is obtained, which leads to a contradiction.

If  $b \in (0, 1)$ , then,

$$\frac{dw}{ds} \sim 1,$$

is obtained from equation (13). By integrating, we get,

$$w \sim v,$$

so that  $b = 1$ , which also leads to a contradiction.

If  $b < 0$ , then,

$$w \frac{dw}{dv} \sim -1,$$

is obtained from equation (13). By integrating, we get,

$$\frac{1}{2} w^2 \sim -v,$$

which leads to a contradiction.

Let us now show that if case (2) is satisfied, then  $b = c = 0$ . Similar to above, either  $b = 1, c \neq 0$  or  $b = 0$ .

If  $b = 1, c < 0$  or  $c > 0$ , then,

$$\frac{dw}{dv} \sim 1 \quad \text{or} \quad \frac{dw}{dv} \sim \frac{2w}{v},$$

is obtained from equation (13). By integrating, we get,

$$w \sim v \quad \text{or} \quad w \sim v^2,$$

is obtained, so that  $c = 0$ , which leads to a contradiction. Hence  $b = 0$ .

If  $c < 0$ , then,

$$w \frac{dw}{dv} \sim -1,$$

is obtained from equation (13). By integrating, we get,

$$\frac{1}{2}w^2 \sim -v,$$

is obtained, which leads to a contradiction.

If  $c > 0$ , then,

$$\frac{dw}{dv} \sim 1,$$

is obtained from equation (13). By integrating, we get,

$$w \sim v,$$

so that  $c = 0$ , which leads to a contradiction.

Summing up and from equation (13), we get,

$$w \sim 1,$$

so that,

$$\frac{dv}{ds} \sim 1, \quad s \rightarrow \infty.$$

By integrating, we finally obtain,

$$v \sim s,$$

as  $s \rightarrow \infty$ , so that,

$$y \sim \frac{1}{\ln \xi}, \quad \xi \rightarrow \infty.$$

That proves our claim.

*Claim* (Shu [48])

Let  $\lambda > 0$ . The negative asymptotic expansion of KdV-B can be written in the form (see Fig. 2),

$$u = -\frac{2k^2 v^2 U_\infty}{\delta} e^{-\frac{(k-1)v\zeta}{2\delta}} - \frac{2k^4 v^2}{\delta} \sum_{i=1}^{\infty} \frac{(2U_\infty)^{i+1} e^{-\frac{(i+1)(k-1)v\zeta}{2\delta}}}{\prod_{j=1}^i [j(k-1)+2k] j(k-1)}, \quad \zeta \rightarrow \infty, \tag{14}$$

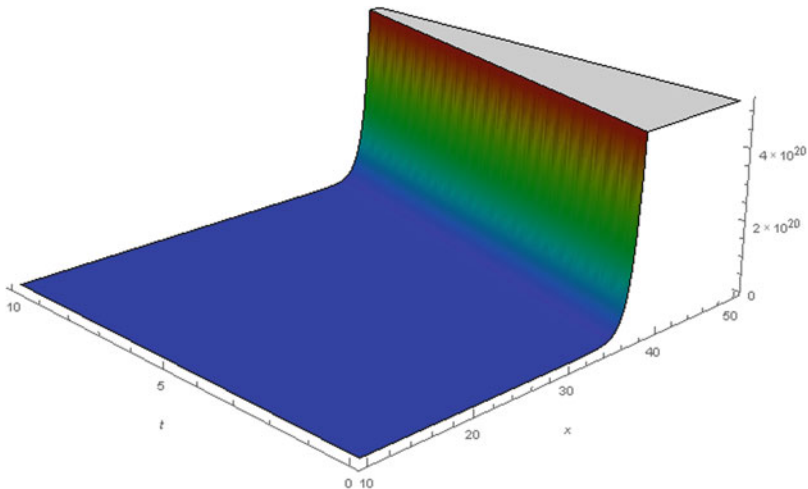
where  $k = \sqrt{1 + \frac{4\lambda\delta}{v^2}}$  and  $U_\infty > 0$  is a constant.

To prove the latter, we first notice that since,

$$e^{-\frac{(i+1)(k-1)v\zeta}{2\delta}}$$

exists, the infinite series converges. Let,

$$u_m = -\frac{2k^2 v^2 U_\infty}{\delta} e^{-\frac{(k-1)v\zeta}{2\delta}} - \frac{2k^4 v^2}{\delta} \sum_{i=1}^{\infty} \frac{(2U_\infty)^{i+1} e^{-\frac{(i+1)(k-1)v\zeta}{2\delta}}}{\prod_{j=1}^i [j(k-1)+2k] j(k-1)}$$



**Fig. 2** The graph shows the asymptotic expansion of the KdV-B equation (see Equation (14)) for the parameters,  $\alpha = 0.1$ ,  $\beta = 0.7$ ,  $\gamma = 1$ ,  $\lambda = 1$ ,  $U_\infty = 1$  and  $\sigma = -2.25$ . The shock wave characteristics can be observed



and

$$y_m = U_\infty + \sum_{i=1}^m \frac{2^{i-1} u_\infty^{1+i} \xi^{i(\sigma+2)}}{\prod_{j=1}^i [j(\sigma+2) - 1] j(\sigma+2)}.$$

Then,

$$y_{m+1}[1 + O(1)] = U_\infty + \int_\xi^\infty \left( \int_t^\infty \tau^\sigma y_m^2 [1 + O(1)]^2 d\tau \right) dt$$

can be obtained for an arbitrary integer  $m$ . Since  $u_m \rightarrow u$  as  $m \rightarrow \infty$ , we get,

$$y_m \rightarrow y_\infty, \quad m \rightarrow \infty,$$

so that,

$$y_\infty = U_\infty + \int_\xi^\infty \left( \int_t^\infty \tau^\sigma y_\infty^2 d\tau \right) dt$$

and  $y_\infty$  is the positive asymptotic expansion of equation (7).

### 4 Hyperbolic Methods for Traveling Wave Solutions of KdV-B

Since the late 1980s, various methods for seeking explicit exact solutions to the KdV-B equation have been independently proposed by many mathematicians, engineers and physicists. The first analytical traveling wave solution to the Burgers-KdV equation was obtained by Xiong [56] in 1989. Two different methods for the construction of exact solutions to the KdV-B equation were proposed by Jeffrey and Mohamad [26]. Wang [52] applied the homogeneous balance method to the study of exact solutions of the compound KdV-B equation. Demiray [13] proposed a so-called “hyperbolic tangent approach” for finding the exact solution to the KdV-B equation, which is actually the Parkes and Duffy’s automated method [39, 40]. Recently, Feng [16–18] introduced the first-integral method to study the exact solution of KdV-B, which is based on the ring theory of commutative algebra. The Cauchy problem for the KdV-B equation was investigated by Bona and Schonbek [7]. They proved the existence and uniqueness of bounded traveling wave solutions which tend to constant states at plus and minus infinity.

We focus on deriving traveling wave solutions for KdV-B, using the Tanh and Sech methods [21, 24, 32, 33, 53]. We start with the Tanh method, considering,

$$-\lambda u + \gamma \frac{u^2}{2} - \alpha \frac{du}{d\zeta} + \beta \frac{d^2u}{d\zeta^2} = 0, \tag{15}$$

where  $\lambda$  is the wave velocity, assuming that both our solution and its spatial derivatives vanish at either plus or minus infinity.

The Tanh method uses a finite series,

$$u(x, t) = u(\mu\zeta) = s(y) = \sum_{m=0}^M a_m y^m, \tag{16}$$

where  $\mu$  is the wave number, inversely proportional to the width of the wave, and  $M$  is a positive integer, in most cases, that will be determined. However if  $M$  is not an integer, a transformation formula is usually used. Substituting equation (16) into equation (15) yields an equation in powers of  $y$ .

To determine the parameter  $M$ , we usually balance the linear terms of highest order in the resulting equation with the highest order nonlinear terms. With  $M$  determined, we collect all coefficients of powers of  $y$  in the resulting equation where these coefficients have to vanish.

This will give a system of algebraic equations involving the parameters  $a_m$ ,  $m = 0, \dots, M$ ,  $\mu$  and  $\lambda$ . Having determined these parameters, knowing that  $M$  is a positive integer in most cases, and using equation (16), we obtain an analytic solution in a closed form. We introduce,

$$y = \tanh(\mu\zeta), \tag{17}$$

that leads to the change of derivatives,

$$\begin{cases} \frac{d}{d\zeta} = \frac{d}{dy} \frac{dy}{dz} = \mu(1 - y^2) \frac{d}{dy}, \\ \frac{d^2}{d\zeta^2} = \frac{d}{d\zeta} \frac{d}{d\zeta} = \mu^2(1 - y^2) \left( -2y \frac{d}{dy} + (1 - y^2) \frac{d^2}{dy^2} \right). \end{cases} \tag{18}$$

Therefore, by replacing equation (16) in equation (15) and using equation (18), we derive an equation with respect to  $u$  as follows,

$$\begin{aligned} & -\lambda \left( \sum_{m=0}^M a_m y^m \right) + \frac{\gamma}{2} \left( \sum_{m=0}^M a_m y^m \right)^2 - \alpha \mu (1 - y^2) \frac{d}{dy} \left( \sum_{m=0}^M a_m y^m \right) \\ & + \beta \mu^2 (1 - y^2) \left( -2y \frac{d}{dy} \left( \sum_{m=0}^M a_m y^m \right) + (1 - y^2) \frac{d^2}{dy^2} \left( \sum_{m=0}^M a_m y^m \right) \right) = 0. \end{aligned}$$

To determine  $M$ , we follow the procedure described above to get  $M = 2$ . This gives the solution in the form,

$$s = \sum_{m=0}^M a_m y^m = a_0 + a_1 y + a_2 y^2, \quad a_2 \neq 0. \tag{19}$$

Substituting equation (19) into equation (15), we get,

$$\begin{aligned}
 & -\lambda(a_0 + a_1y + a_2y^2) + \frac{\gamma}{2}(a_0 + a_1y + a_2y^2)^2 \\
 & \quad - \alpha\mu(1 - y^2)(a_1 + 2a_2y) + \beta\mu^2(1 - y^2) \\
 & \quad \left( -2y(a_1 + 2a_2y) + 2a_2(1 - y^2) \right) = 0.
 \end{aligned}$$

Collecting the coefficients of different powers of  $y$ , gives the following system of algebraic equations for  $\lambda$ ,  $\mu$ ,  $a_0$ ,  $a_1$  and  $a_2$ ,

$$\begin{cases}
 a_2(\gamma a_2 + 12\beta\mu^2) = 0 \\
 \gamma a_1 a_2 + 2\beta\mu^2 a_1 + 2\alpha\mu a_2 = 0 \\
 -\lambda a_2 - 8\beta\mu^2 a_2 + \gamma a_0 a_2 + \frac{\gamma}{2} a_1^2 + \alpha\mu a_1 = 0 \\
 -\lambda a_1 - 2\beta\mu^2 a_1 + \gamma a_0 a_1 - 2\alpha\mu a_2 = 0 \\
 -\lambda a_0 + \frac{\gamma}{2} a_0^2 - \alpha\mu a_1 + 2\beta\mu^2 a_2 = 0
 \end{cases}$$

with solution,

$$\begin{cases}
 \lambda = \pm \frac{6}{25} \frac{\alpha^2}{\beta}, \\
 \mu = \pm \frac{\alpha}{10\beta}, \\
 a_0 = \frac{\lambda}{\gamma} + 12 \frac{\beta\mu^2}{\gamma}, \\
 a_1 = -\frac{12}{5} \frac{\alpha\mu}{\gamma}, \\
 a_2 = -12 \frac{\beta\mu^2}{\gamma}.
 \end{cases} \tag{20}$$

Using the trigonometric identities,

$$\begin{cases}
 \tanh^2(\theta) = 1 - \operatorname{sech}^2(\theta) \\
 \tanh(-\theta) = -\tanh(\theta)
 \end{cases}, \quad \theta \in \mathbb{R}$$

and requiring for both our solution and its spatial derivatives to vanish at plus infinity, we get the following traveling wave solution,

$$u_{1\infty}(\zeta) = \frac{3}{25} \frac{\alpha^2}{\beta\gamma} \left( \operatorname{sech}^2(\mu\zeta) - 2\tanh(\mu\zeta) + 2 \right), \quad \mu, \lambda > 0 \tag{21}$$

Requiring for both our solution and its spatial derivatives to vanish at minus infinity, we get the following traveling wave solution,

$$u_{2-\infty}(\zeta) = \frac{3}{25} \frac{\alpha^2}{\beta\gamma} \left( \operatorname{sech}^2(\mu\zeta) - 2\tanh(\mu\zeta) - 2 \right), \quad \mu > 0, \quad \lambda < 0 \tag{22}$$

*Remark 2* A notable result is that our traveling wave solutions are expressed as a composition of a bell-profile solitary wave (KdV) and a kink-profile solitary wave (Burgers’) with velocity  $\lambda = \pm \frac{6}{25} \frac{\alpha^2}{\beta}$ . The shock profile is dominant here. All those exact solutions, and others mentioned in literature, can be proved to be algebraically equivalent to each other [18]. That is, essentially only one explicit traveling solitary wave solution to the KdV–B equation is known which can be expressed as a composition of a bell-profile solitary wave and a kink-profile solitary wave. In other words, a feature of this solution is that is a linear combination of particular solutions of the KdV equation and the Burgers equation [18, 26].

By following similar steps as with the Tanh method for traveling wave solutions of KdV–B, the Sech method uses the variable transformation [21, 24, 53],

$$\begin{cases} y = \operatorname{sech}(\mu\zeta), & \mu \neq 0, \\ \frac{d}{d\zeta} = -\mu y \sqrt{1 - y^2} \frac{d}{dy}, \\ \left[ \frac{d^2}{d\zeta^2} = \mu^2 y [(1 - 2y^2) \frac{d}{dy} + (y - y^3) \frac{d^2}{dy^2}] \right] \end{cases} \tag{23}$$

and the ansatz,

$$u(x, t) = u(\mu\zeta) = s(y) = \sum_{m=0}^2 a_m y^m. \tag{24}$$

Then by replacing equation (24) in equation (15) and using equation (23), we derive an equation with respect to  $u$  as follows,

$$\begin{aligned} & [a_2(\gamma \frac{a_2}{2} - \mu\alpha - 6\beta\mu^2)]y^4 + [a_1(\gamma a_2 - \frac{\alpha\mu}{2} - 2\beta\mu^2)]y^3 \\ & + [-\lambda a_2 + \frac{\gamma}{2} a_1^2 + \gamma a_0 a_2 + 2\alpha\mu a_2 \\ & + 4\beta\mu^2 a_2]y^2 + [a_1(-\lambda + \gamma a_0 + \alpha\mu + \beta\mu^2)]y \\ & + \frac{\gamma}{2} a_0^2 - \lambda a_0 = 0. \end{aligned}$$

Above, we used the first order Taylor approximation,

$$\sqrt{1 - y^2} \sim 1 - \frac{y^2}{2}.$$

Collecting the coefficients of different powers of  $y$ , gives the following system of algebraic equations for  $\lambda$ ,  $\mu$ ,  $a_0$ ,  $a_1$  and  $a_2$ ,

$$\begin{cases} a_2(\gamma \frac{a_2}{2} - \mu\alpha - 6\beta\mu^2) = 0 \\ a_1(\gamma a_2 - \frac{\alpha\mu}{2} - 2\beta\mu^2) = 0 \\ -\lambda a_2 + \frac{\gamma}{2} a_1^2 + \gamma a_0 a_2 + 2\alpha\mu a_2 + 4\beta\mu^2 a_2 \\ a_1(-\lambda + \gamma a_0 + \alpha\mu + \beta\mu^2) = 0 \\ \frac{\gamma}{2} a_0^2 - \lambda a_0 = 0 \end{cases}$$

with a solution being,

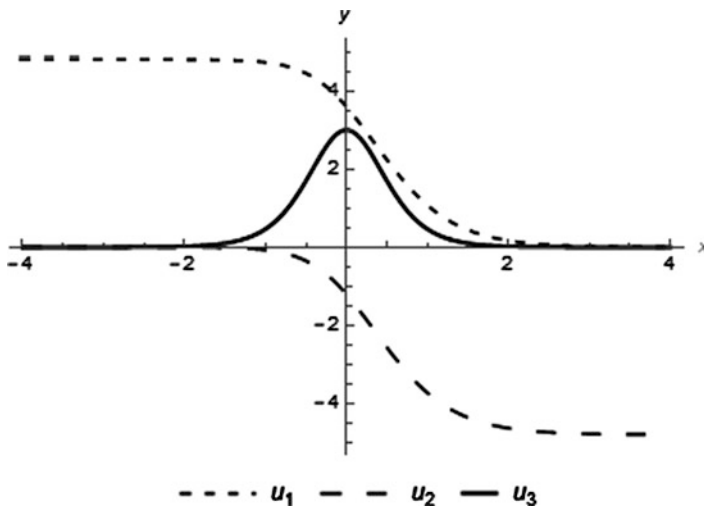
$$\begin{cases} \lambda = 2\mu(\alpha + 2\beta\mu) \\ a_0 = \frac{4\mu}{\gamma}(\alpha + 2\beta\mu) \\ a_1 = 0 \\ a_2 = \frac{2\mu}{\gamma}(\alpha + 6\beta\mu) \end{cases}, \quad \mu \neq 0 \tag{25}$$

giving a solitary profile traveling wave solution,

$$\frac{2\mu}{\gamma} \left( 2(\alpha + 2\beta\mu) + (\alpha + 6\beta\mu) \operatorname{sech}^2(\mu\zeta) \right), \quad \mu \neq 0. \tag{26}$$

*Remark 3* A notable result is that the Sech method can give “purely” solitary profile traveling wave solutions.

The following graph, Fig. 3, depicts the solutions studied in this section.



**Fig. 3**  $u_1$  and  $u_2$  are two shock profile traveling waves of the KdV-B equation vanishing at plus and minus infinity, respectively, whereas  $u_3$  is a solitary profile traveling wave solution of the KdV equation, for the parameters,  $\alpha = 1, \beta = 0.1, \gamma = 1, \lambda = 1$

## 5 Spectral Fourier Analysis for the Numerical Solution of KdV–B

We define the Fourier and Inverse Fourier Transform of a function, say  $f$ , in the sense that the following symbols make sense, to be [2, 31],

$$\begin{cases} F[f(x)] = \hat{f}(k) = \int_{-\infty}^{\infty} e^{-ikx} f(x) dx, \\ F^{-1}[\hat{f}(k)] = f(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{ikx} \hat{f}(k) dk. \end{cases} \quad (27)$$

It is easy to see, integrating by parts, that regarding the  $n$ th derivative of  $f$  and its Fourier Transform, the following results hold,

$$\begin{cases} \frac{d^n f}{dx^n} = (i^n) F^{-1} [k^n \hat{f}(k)], \\ \frac{d^n \hat{f}}{dk^n} = (-i)^n F [x^n f(x)]. \end{cases} \quad (28)$$

Now consider the KdV–B equation,

$$\frac{\partial u}{\partial t} + \gamma u \frac{\partial u}{\partial x} - \alpha \frac{\partial^2 u}{\partial x^2} + \beta \frac{\partial^3 u}{\partial x^3} = 0, \quad u = u(t, x). \quad (29)$$

Rearranging the terms of equation (29), we get,

$$\frac{\partial u}{\partial t} = -\gamma u \frac{\partial u}{\partial x} + \alpha \frac{\partial^2 u}{\partial x^2} - \beta \frac{\partial^3 u}{\partial x^3}. \quad (30)$$

By means of the Inverse Fourier Transform,  $F^{-1}$ , equation (30) can be written in the form [2, 50],

$$\frac{\partial u}{\partial t} = f(t, u), \quad (31)$$

where we have substituted the  $x$ -partial derivatives with,

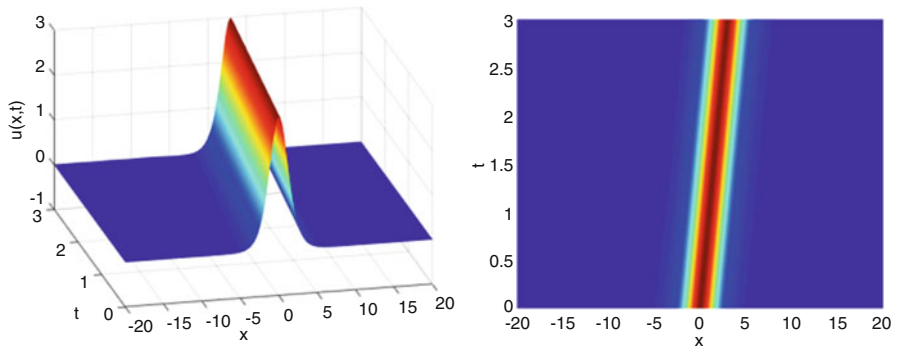
$$\begin{cases} \frac{\partial u}{\partial x} = i F^{-1}(\kappa \hat{u}), \\ \frac{\partial^2 u}{\partial x^2} = -F^{-1}(\kappa^2 \hat{u}), \\ \frac{\partial^3 u}{\partial x^3} = -i F^{-1}(\kappa^3 \hat{u}). \end{cases}$$

Now, equation (31) is suitable for applying the 4th order explicit Runge–Kutta method, giving us the following,

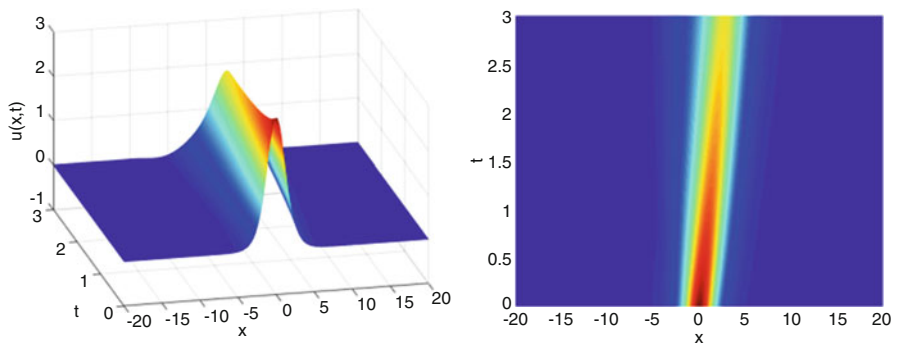
$$\begin{cases} u_{n+1} = u_n + \frac{h}{6}(k_1 + 2k_2 + 2k_3 + k_4) \\ t_{n+1} = t_n + h \\ k_1 = f(t_n, u_n) \\ k_2 = f(t_n + \frac{h}{2}, u_n + \frac{hk_1}{2}) \\ k_3 = f(t_n + \frac{h}{2}, u_n + \frac{hk_2}{2}) \\ k_4 = f(t_n + h, u_n + hk_3) \end{cases}, \quad n = 0, 1, \dots \quad (32)$$

For  $n = 0$ , we may choose either a soliton of the KdV equation or a similarity solution of the viscous Burgers equation or a traveling wave solution of KdV-B.

Below we exhibit the obtained numerical results for each case separately. In Fig.4, the evolution of an initial solitary profile solution of the KdV equation is depicted, where diffusive effects are absent. It can be observed that the solitary waveform is retained. In Fig.5, the evolution of a solitary profile solution of the

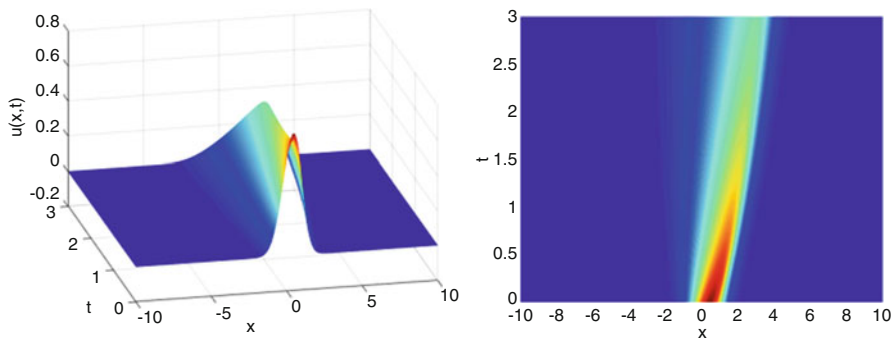


**Fig. 4** A solution of the KdV-B equation, evolving a solitary profile solution of the KdV equation, where diffusive effects are absent (KdV case), for the parameters,  $\lambda = 1, \alpha = 0, \beta = 0.7, \gamma = 1$

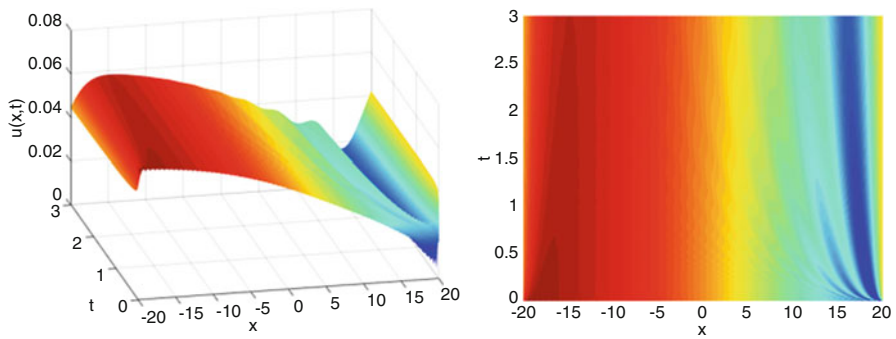


**Fig. 5** A solution of the KdV-B equation, evolving a solitary profile solution of the KdV equation, where both diffusive and dispersive effects coexist, for the parameters,  $\lambda = 1, \alpha = 0.5, \beta = 0.7, \gamma = 1$

KdV equation is presented, where both diffusive and dispersive effects coexist. It is observed that the profile loses energy and reduces in amplitude drastically. Additionally, in Fig. 6 we present the evolution of a similarity shock profile solution of the viscous Burgers equation, where both diffusive and dispersive effects coexist. In this case a wavefront can be observed revealing a shock-like behavior. Finally, in Fig. 7, the evolution of a traveling shock wave profile solution of the KdV–B equation is presented, where both diffusive and dispersive effects coexist. These numerical solutions clearly reveal both solitary and shock wave features of the KdV–B equation, revealing its connection to cardiac hemodynamics where all these phenomena, such as convection, diffusion and dispersion, can be observed.



**Fig. 6** A solution of the KdV–B equation, evolving a similarity shock profile solution of the viscous Burgers equation, where both diffusive and dispersive effects coexist, for the parameters,  $\lambda = 1.8, \alpha = 0.19, \beta = 0.01, \gamma = 3.4$



**Fig. 7** A solution of the KdV–B equation, evolving a traveling wave shock profile solution of the KdV–B equation, where both diffusive and dispersive effects coexist, for the parameters,  $\alpha = 0.3, \beta = 0.7, \gamma = 1$



## 6 Conclusions

Recent advancements concerning cardiac dynamics pose important questions about the cardiac waveform. A governing equation, namely the KdV–B equation (Korteweg–de Vries–Burgers), which is a partial differential equation can be utilized to answer several of those questions. The KdV–B equation features both solitary and shock wave characteristics due to the dispersion and dissipation terms, as also occurring in the arterial tree. This study focuses on describing cardiac dynamics with the applications of mathematics and nonlinear analysis. It is customarily difficult to solve nonlinear problems, especially by analytical techniques. Therefore, seeking suitable solving methods, such as, exact, approximate or numerical methods, is an active task in branches of applied mathematics and nonlinear analysis.

In this chapter, the phase plane of the KdV–B equation is analyzed and its qualitative behavior is derived, depicting the stability states of the equation contributing to the decisions made for further analytical and numerical consideration. The analysis reveals a saddle point  $(0, 0)$ , and an additional one that could be a source point or a spiral source point or a central point depending on the equation's parameters,  $\alpha$ ,  $\beta$  and  $\lambda$ .

Furthermore, an asymptotic expansion is presented, providing a reliable basis for estimating the advantages and disadvantages when seeking and applying numerical methods to KdV–B equation. Furthermore, traveling wave solutions under both solitary and shock profiles are obtained from the hyperbolic methods, whose strength is their ease of use to find which solitary wave structures and/or shock-wave (kinks) profiles satisfy nonlinear wave and evolution equations. These techniques allow to develop algorithms for symbolic software packages, so that nonlinear partial differential equations and difference equations, can be studied automatically whether (or not) they possess traveling wave solutions. Additionally, numerical solutions are obtained for the equation, by means of the Spectral Fourier analysis. Both these solutions and the latter traveling wave solutions are evolved in time by the Runge–Kutta method. These solutions clearly depict both solitary and shock wave characteristics of the KdV–B equation. This analysis provides vital information about the equation and its connection to cardiac hemodynamics.

## References

1. M.J. Ablowitz, H. Segur, *Solitons and the Inverse Scattering Transform*. (SIAM, Philadelphia, 1981)
2. M.J. Ablowitz, D.J. Kaup, A.C. Newell, H. Segur, The inverse scattering transform-fourier analysis for nonlinear problems. *Stud. Appl. Math.* **53**(4), 249–315 (1974)
3. M. Anliker, M.B. Hinstead, E. Ogden, Dispersion and attenuation of small artificial pressure waves in the canine aorta. *Circ. Res.* **23**(4), 539–551 (1968)
4. N. Antar, H. Demiray, Weakly nonlinear waves in a prestressed thin elastic tube containing a viscous fluid. *Int. J. Eng. Sci.* **37**(14), 1859–1876 (1999)

5. H. Atabek, Wave propagation through a viscous fluid contained in a tethered, initially stressed, orthotropic elastic tube. *Biophys. J.* **8**(5), 626–649 (1968)
6. D. Benney, Long waves on liquid films. *J. Math. Phys.* **45**(1–4), 150–155 (1966)
7. J. Bona, M. Schonbek, Travelling-wave solutions to the Korteweg-de Vries-Burgers equation. *Proc. Roy. Soc. Edinburgh Sect. A* **101**, 207–226 (1985)
8. R.K. Bullough, P.J. Caudrey, *Solitons*, vol. 17 (Springer Science & Business Media, Berlin, 2013)
9. J. Burgers, *Correlation Problems in a One-dimensional Model of Turbulence II* (North-Holland Publishing, Amsterdam, 1950)
10. J. Canosa, J. Gazdag, The Korteweg-de Vries-Burgers equation. *J. Comput. Phys.* **23**(4), 393–403 (1977)
11. E. Crépeau, M. Sorine, A reduced model of pulsatile flow in an arterial compartment. *Chaos Solitons Fractals* **34**(2), 594–605 (2007)
12. H. Demiray, Wave propagation through a viscous fluid contained in a prestressed thin elastic tube. *Int. J. Eng. Sci.* **30**(11), 1607–1620 (1992)
13. H. Demiray, A note on the exact travelling wave solution to the KdV–Burgers equation. *Wave motion* **38**(4), 367–369 (2003)
14. G.A. El, M.A.Hoefel, Dispersive shock waves and modulation theory. *Phys. D Nonlin. Phenom.* **333**, 11–65 (2016)
15. H. Erbay, S. Erbay, S. Dost, Wave propagation in fluid filled nonlinear viscoelastic tubes. *Acta Mech.* **95**(1–4), 87–102 (1992)
16. Z. Feng, The first-integral method to study the Burgers–Korteweg–de Vries equation. *J. Phys. A Math. Gen.* **35**(2), 343 (2002)
17. Z. Feng, On explicit exact solutions to the compound Burgers–KdV equation. *Phys. Lett. A* **293**(1–2), 57–66 (2002)
18. Z. Feng, Qualitative analysis and exact solutions to the Burgers–Korteweg–de Vries equation. *Dyn. Contin. Discret. Impul. Syst. Ser. A Math. Anal.* **9**(4), 563–580 (2002)
19. F. Feudel, H. Steudel, Non-existence of prolongation structure for the Korteweg–de Vries–Burgers equation. *Phys. Lett. A* **107**(1), 5–8 (1985)
20. G. Gao, A theory of interaction between dissipation and dispersion of turbulence. *Sci. Sin. (Ser. A) Math. Phys. Astronom. Tech. Sci.* **28**, 616–627 (1985)
21. Y.T. Gao, B. Tian, Generalized hyperbolic-function method with computerized symbolic computation to construct the solitonic solutions to nonlinear equations of mathematical physics. *Comput. Phys. Commun.* **133**(2–3), 158–164 (2001)
22. D.M. Grobman, Homeomorphism of systems of differential equations. *Dokl. Akad. Nauk SSSR* **128**(5), 880–881 (1959)
23. P. Hartman, A lemma in the theory of structural stability of differential equations. *Proc. Am. Math. Society* **11**(4), 610–620 (1960)
24. R. Hirota, Direct method of finding exact solutions of nonlinear evolution equations, in *Bäcklund Transformations, the Inverse Scattering Method, Solitons, and Their Applications* (Springer, Berlin, 1976), pp. 40–68
25. A. Jeffrey, T. Kakutani, Weak nonlinear dispersive waves: a discussion centered around the Korteweg–de Vries equation. *SIAM Rev.* **14**(4), 582–643 (1972)
26. A. Jeffrey, M. Mohamad, Exact solutions to the kdv-burgers’ equation. *Wave Motion* **14**(4), 369–375 (1991)
27. R. Johnson, Shallow water waves on a viscous fluid the undular bore. *Phys. Fluid.* **15**(10), 1693–1699 (1972)
28. T. Kawahara, Weak nonlinear magneto-acoustic waves in a cold plasma in the presence of effective electron-ion collisions. *J. Phys. Soc. Jpn.* **28**(5), 1321–1329 (1970)
29. D.J. Korteweg, G. De Vries, Xli. on the change of form of long waves advancing in a rectangular canal, and on a new type of long stationary waves. *London, Edinburgh, Dublin Philos. Mag. J. Sci.* **39**(240), 422–443 (1895)
30. S.D. Liu, S.K. Liu, KdV-burgers equation modelling of turbulence. *Sci. China Ser. A Math. Phys. Astron. Technol. Sci.* **35**(5), 576–586 (1992)

31. J.D. Logan, *Applied Mathematics* (John Wiley & Sons, Hoboken, 2013)
32. W. Malfliet, Solitary wave solutions of nonlinear wave equations. *Am. J. Phys.* **60**(7), 650–654 (1992)
33. W. Malfliet, W. Hereman, The tanh method: I exact solutions of nonlinear evolution and wave equations. *Phys. Scripta* **54**, 563–568 (1996)
34. A. Mangel, M. Fahim, C. Van Breemen, Control of vascular contractility by the cardiac pacemaker. *Science* **215**(4540), 1627–1629 (1982)
35. B. Mehta, R. Aris, A note on a form of the emden-fowler equation. *J. Math. Anal. Appl.* **36**(3), 611–621 (1971)
36. J.D. Murray, *Asymptotic Analysis*, vol. 48 (Springer Science & Business Media, Berlin, 2012)
37. G. Nagy, Ordinary Differential Equations. Online notes (2020)
38. M. Otwinowski, R. Paul, J. Tuszynski, An answer to the question: is the arterial pulse a soliton? *J. Biol. Phys.* **14**(2), 43–48 (1986)
39. E. Parkes, B. Duffy, An automated tanh-function method for finding solitary wave solutions to non-linear evolution equations. *Comput. Phys. Commun.* **98**(3), 288–300 (1996)
40. E. Parkes, B. Duffy, Travelling solitary wave solutions to a compound KdV-Burgers equation. *Phys. Lett. A* **229**(4), 217–220 (1997)
41. T.J. Pedley, Y. LX, *Fluid Mechanics of Large Blood Vessels* (Shaanxi People's Press, Beijing, 1995)
42. A.I. Rachev, Effects of transmural pressure and muscular activity on pulse waves in arteries. *J. Biomech. Eng.*, 119–123 (1980)
43. R.L. Rockwell, *Nonlinear Analysis of Pressure and Shock Waves in Blood Vessels*. Ph. D (Stanford Univ. CA, US, NASA Archives, Thesis, 1969)
44. S. Rowlands, Is the arterial pulse a soliton? *J. Biol. Phys.* **10**(4), 199–200 (1982)
45. G. Rudinger, Review of current mathematical methods for the analysis of blood flow, in *Biomedical Fluid Mechanics Symposium* (ASME, New York, 1966), pp. 1–33
46. Y. Shi, P. Lawford, R. Hose, Review of zero-d and 1-d models of blood flow in the cardiovascular system. *Biomed. Eng. Online* **10**(1), 33 (2011)
47. J.J. Shu, Exact n-envelope-soliton solutions of the Hirota equation. Preprint. arXiv:1403.3645 (2014)
48. J.J. Shu, The proper analytical solution of the Korteweg-de Vries-Burgers equation. Preprint. arXiv:1403.3636 (2014)
49. R. Skalak, Wave propagation in blood flow, in *Biomechanics Symposium* (ASME, New York, 1966), pp. 20–46
50. L.N. Trefethen, *Spectral Methods in Matlab, Volume 10 of Software, Environments, and Tools*, vol. 24 (Society for Industrial and Applied Mathematics (SIAM), Philadelphia, 2000)
51. M. Wadati, Wave propagation in nonlinear lattice. *I. J. Phys. Society Jpn.* **38**(3), 673–680 (1975)
52. M. Wang, Exact solutions for a compound KdV-Burgers equation. *Phys. Lett. A* **213**(5-6), 279–287 (1996)
53. A.M. Wazwaz, The tanh-coth and the sech methods for exact solutions of the Jaulent-miodek equation. *Phys. Lett. A* **366**(1–2), 85–90 (2007)
54. L.V. Wijngaarden, One-dimensional flow of liquids containing small gas bubbles. *Ann. Rev. Fluid Mech.* **4**(1), 369–396 (1972)
55. J.R. Womersley, WADC technical report tr 56-614 and *phil. Mag* **46**, 199–221 (1955)
56. S. Xiong, An analytic solution of Burgers-KdV equation. *Chin. Sci. Bull.* **34**(14), 1158–1162 (1989)
57. S. Yomosa, Solitary waves in large blood vessels. *J. Phys. Society Jpn.* **56**(2), 506–520 (1987)