

Springer Series in Supply Chain Management

Ming Hu *Editor*

Sharing Economy

Making Supply Meet Demand

 Springer

Springer Series in Supply Chain Management

Volume 6

Series Editor

Christopher S. Tang
University of California
Los Angeles, CA, USA

More information about this series at <http://www.springer.com/series/13081>

Ming Hu
Editor

Sharing Economy

Making Supply Meet Demand

 Springer

Editor

Ming Hu
Rotman School of Management
University of Toronto
Toronto, ON, Canada

ISSN 2365-6395 ISSN 2365-6409 (electronic)
Springer Series in Supply Chain Management
ISBN 978-3-030-01862-7 ISBN 978-3-030-01863-4 (eBook)
<https://doi.org/10.1007/978-3-030-01863-4>

Library of Congress Control Number: 2018964940

© Springer Nature Switzerland AG 2019

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, express 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

Sharing economy refers to a market model that enables and facilitates the sharing of access to goods and services. For example, Uber allows riders to share a car. Airbnb allows homeowners to share their extra rooms with renters. Groupon crowdsources demands, enabling customers to share the benefit of discounted goods and services, whereas Kickstarter crowdsources funds, enabling backers to fund a project jointly. Unlike the classic supply chain settings in which a firm makes inventory and supply decisions, in a sharing economy, supply is crowdsourced and can be modulated by a platform. The matching-supply-with-demand process in a sharing economy requires novel perspectives and tools to address challenges and identify opportunities.

This edited book examines the challenges and opportunities arising from today's sharing economy from an operations management perspective. Individual chapter authors present state-of-the-art research that examines the general impact of sharing economy on production and consumption, the intermediary role of a sharing platform, crowdsourcing management, and various context-based operational problems.

Toronto, Canada
March 2018

Ming Hu

Acknowledgments

This book cannot exist without the strong commitment of our colleagues. I am grateful to each of the contributing authors for sharing their cutting-edge research with us. I would like to thank Professor Chris Tang, the editor of the Springer Series in Supply Chain Management, who has encouraged me to work on this book. I am grateful to Mirko Janc for typesetting each chapter exceptionally carefully and rigorously.

This work was financially supported by the Natural Sciences and Engineering Research Council of Canada [Grant RGPIN-2015-06757]; the National Natural Science Foundation of China (NSFC) [Grant No. 71772006]; the Rotman School of Management, University of Toronto; and the program of Guanghua Thought Leadership at Guanghua School of Management, Peking University.

Contents

1	Introduction	1
	Ming Hu	
1.1	Overall Structure	1
1.2	Chapter Highlights	2
1.2.1	Part I: Impact of Sharing Economy	2
1.2.2	Part II: Intermediary Role of a Sharing Platform	3
1.2.3	Part III: Crowdsourcing Management	6
1.2.4	Part IV: Context-Based Operational Problems in Sharing Economy	7
	References	8
 Part I Impact of Sharing Economy		
2	Peer-to-Peer Product Sharing	11
	Saif Benjaafar, Guangwen Kong, Xiang Li, and Costas Courcoubetis	
2.1	Introduction	12
2.2	Literature Review	15
2.3	Model Description	17
2.3.1	Matching Supply with Demand	19
2.4	Equilibrium Analysis	21
2.4.1	Impact of Collaborative Consumption on Ownership and Usage	22
2.4.2	Impact of Collaborative Consumption on Consumers	25
2.5	The Platform's Problem	26
2.5.1	The For-Profit Platform	27
2.5.2	The Not-for-Profit Platform	29
2.5.3	Systems with Negative Externalities	31
2.5.4	The Impact of Extra Wear and Tear and Inconvenience Costs	33
2.6	Concluding Comments	33
	References	35

3	The Strategic and Economic Implications of Consumer-to-Consumer Product Sharing	37
	Baojun Jiang and Lin Tian	
3.1	Introduction	37
3.2	Modeling Framework	40
3.3	Effects of Sharing on Firm's Pricing Strategy, Profit, and Consumer Surplus	43
3.4	Effects of Sharing on Product Quality and Distribution Channel ..	48
	3.4.1 Effects of Sharing on Product Quality	48
	3.4.2 Effects of Sharing on Distribution Channel	49
3.5	Conclusions and Discussions	51
	References	53
4	Operational Factors in the Sharing Economy: A Framework	55
	Tunay I. Tunca	
4.1	Introduction	55
4.2	The Framework	56
4.3	Examples	60
	4.3.1 Ride Sharing	61
	4.3.2 Group Buying	64
4.4	Concluding Remarks	66
	References	68
5	Ride Sharing	73
	Siddhartha Banerjee and Ramesh Johari	
5.1	Introduction	73
5.2	Anatomy of a Modern Ridesharing Platform	75
	5.2.1 Timescales	75
	5.2.2 Strategic Choices	76
	5.2.3 Operation and Market Design	77
5.3	A Modeling Framework for Ridesharing Platforms	77
	5.3.1 Modeling Stochastic Dynamics of the Platform	78
	5.3.2 Platform Controls	81
	5.3.3 Platform Objectives	83
	5.3.4 Local Controls and Closed Queueing Models	84
	5.3.5 Modeling Endogenous Entry of Drivers	86
5.4	Analyzing the Model: Key Findings	87
	5.4.1 Fast-Timescale Control of Platform Dynamics	88
	5.4.2 The Slow Timescale: Pricing and Driver Entry	89
5.5	Related Literature	92
5.6	Conclusion	95
	References	96

Part II Intermediary Role of a Sharing Platform

6 The Role of Surge Pricing on a Service Platform with Self-Scheduling Capacity 101
 Gerard P. Cachon, Kaitlin M. Daniels, and Ruben Lobel

6.1 Introduction 102
 6.2 Literature Review 103
 6.3 Model 105
 6.4 Profitability of Commission Contract 108
 6.5 Impact of Dynamic Prices on Consumers 110
 6.6 Conclusion 111
 References 112

7 Time-Based Payout Ratio for Coordinating Supply and Demand on an On-Demand Service Platform 115
 Jiaru Bai, Kut C. So, Christopher S. Tang, Xiqun (Michael) Chen, and Hai Wang

7.1 Introduction 116
 7.2 Literature Review 117
 7.3 A Model of Wait-Time Sensitive Demand and Earnings Sensitive Supply 119
 7.3.1 Customer Request Rate λ and Price Rate p 120
 7.3.2 Number of Participating Providers k and Wage Rate w .. 120
 7.3.3 Problem Formulation 122
 7.4 The Base Model 122
 7.4.1 Special Case 1: When the Payout Ratio w/p Is Fixed 124
 7.4.2 Special Case 2: When the Service Level Is Exogenously Given 125
 7.5 Numerical Illustrations Based on Didi Data 127
 7.5.1 Background Information 127
 7.5.2 Number of Rides and Drivers Across Different Hours ... 128
 7.5.3 Travel Distance and Travel Speed 128
 7.5.4 Pricing and Wage Rates 129
 7.5.5 Strategic Factors and Their Implications 130
 7.5.6 Numerical Examples for Illustrative Purposes 130
 7.6 Conclusion 134
 References 135

8 Pricing and Matching in the Sharing Economy 137
 Yiwei Chen, Ming Hu, and Yun Zhou

8.1 Introduction 138
 8.1.1 Two-Sided Pricing 138
 8.1.2 Two-Sided Matching 138
 8.1.3 Pricing and Matching Under Strategic Behavior 140

8.2	Two-Sided Pricing and Fixed Commission	141
8.2.1	The Price and Wage Optimization Problem	141
8.2.2	The Fixed Commission Contract	143
8.2.3	Numerical Study	145
8.3	Dynamic Matching with Heterogeneous Types	146
8.3.1	Priority Properties of the Optimal Matching Policy	147
8.3.2	Bound and Heuristic	151
8.4	Pricing and Matching with Strategic Suppliers and Customers	153
8.4.1	Upper Bound of the Intermediary's Optimal Profit	157
8.4.2	A Simple Dynamic Policy: Asymptotic Optimality	158
8.5	Conclusion	163
	References	163
9	Large-Scale Service Marketplaces: The Role of the Moderating Firm	165
	Gad Allon, Achal Bassamboo, and Eren B. Çil	
9.1	Introduction	165
9.2	Literature Review	168
9.3	Model Formulation	170
9.4	No-Intervention Model	171
9.4.1	Characterization of SPNE	173
9.5	Operational Efficiency Model	175
9.5.1	Characterization of the Market Equilibrium	178
9.6	Communication Enabled Model	185
9.6.1	Characterization of the (δ, ϵ) -Market Equilibrium	186
9.7	A Marketplace with Non-identical Agents	188
9.8	Conclusion	189
	References	191
10	Inducing Exploration in Service Platforms	193
	Kostas Bimpikis and Yiangos Papanastasiou	
10.1	Introduction	193
10.2	Related Literature	194
10.3	Illustrative Example	197
10.4	Benchmark Model	198
10.5	Inducing Exploration	199
10.5.1	Strategic Information Disclosure	202
10.5.2	The Value of Information Obfuscation	205
10.5.3	Minimizing Regret	206
10.5.4	Incentivizing Customers Using Payments	208
10.6	Promising Directions	210
10.6.1	Learning in Dynamic Contests	210
10.6.2	Dealing with Misinformation	212
10.7	Concluding Remarks	213
	References	214

11 Design of an Aggregated Marketplace Under Congestion Effects: Asymptotic Analysis and Equilibrium Characterization..... 217
 Ying-Ju Chen, Costis Maglaras, and Gustavo Vulcano

11.1 Introduction..... 218

 11.1.1 Background and Motivation..... 218

 11.1.2 Overview of Results..... 220

 11.1.3 Literature Review..... 222

11.2 Model..... 224

 11.2.1 Description of the Market..... 224

 11.2.2 Problems to Address..... 226

11.3 Asymptotic Analysis of Marketplace Dynamics..... 227

 11.3.1 Background: Revenue Maximization for an $M/M/1$ Monopolistic Supplier..... 228

 11.3.2 Setup for Asymptotic Analysis..... 229

 11.3.3 Transient Dynamics via a Fluid Model Analysis..... 229

 11.3.4 State-Space Collapse and the Aggregate Marketplace Behavior..... 230

 11.3.5 Limit Model and Discussion..... 231

 11.3.6 A Numerical Example..... 232

11.4 Competitive Behavior and Market Efficiency..... 234

 11.4.1 Suppliers' First-Order Payoffs and the Capacity Game .. 234

 11.4.2 Suppliers' Second-Order Payoffs and the Pricing Game..... 235

 11.4.3 Centralized System Performance..... 236

 11.4.4 Competitive Equilibrium..... 237

 11.4.5 Coordination Scheme..... 243

 11.4.6 Simulation Results..... 245

11.5 Conclusions..... 247

References..... 248

12 Operations in the On-Demand Economy: Staffing Services with Self-Scheduling Capacity..... 249
 Itai Gurvich, Martin Lariviere, and Antonio Moreno

12.1 Introduction..... 250

12.2 Model..... 253

12.3 Analysis..... 256

 12.3.1 The Cost of Self Scheduling..... 256

 12.3.2 Earnings Constraint and Agent Flexibility..... 259

 12.3.3 Time-Varying Demand..... 259

 12.3.4 The Benefit of Flexible Capacity..... 260

12.4 Variants of the Base Model..... 261

 12.4.1 Volume-Dependent Compensation Schemes..... 263

 12.4.2 Price-Dependent Newsvendor..... 265

 12.4.3 When Maintaining a Larger Pool Costs More..... 267

 12.4.4 Period-Dependent Threshold Distributions..... 269

12.5 Concluding Remarks..... 270

References..... 277

13 On Queues with a Random Capacity: Some Theory, and an Application 279
 Rouba Ibrahim

13.1 Introduction 279

13.2 Theoretical Background: Queues with Uncertain Parameters 281

13.2.1 Self-Scheduling Servers: A Binomial Distribution 284

13.2.2 What Do the Asymptotic Results Mean? 285

13.3 Self-Scheduling Agents: A Long-Term Staffing Decision 289

13.3.1 The Model 289

13.3.2 Fluid Formulation 290

13.3.3 Optimal Staffing Policy 291

13.4 Short-Term Controls 294

13.4.1 Delay Announcements: Performance Impact 295

13.5 Joint Control of Compensation and Delay Announcements 299

13.6 Jointly Optimizing Long and Short-Term Controls 302

13.6.1 Low Minimum Wage 302

13.6.2 High Minimum Wage 303

13.7 Conclusions 303

Technical Appendix 304

References 315

Part III Crowdsourcing Management

14 Online Group Buying and Crowdfunding: Two Cases of All-or-Nothing Mechanisms 319
 Ming Hu, Mengze Shi, and Jiahua Wu

14.1 Introduction 319

14.2 Consumer Behavior Under All-or-Nothing Mechanisms 322

14.2.1 Empirical Model 323

14.2.2 Results 325

14.2.3 Potential Mechanisms Behind Threshold Effects 331

14.3 Coordination Under All-or-Nothing Mechanisms 333

14.3.1 Information Disclosure 333

14.3.2 Pricing 339

14.4 Conclusion 344

References 345

15 Threshold Discounting: Operational Benefits, Potential Drawbacks, and Optimal Design 347
 Simone Marinesi, Karan Girotra, and Serguei Netessine

15.1 Introduction 348

15.2 Literature Review 350

15.3 The Model 352

15.3.1 Preliminaries 352

15.3.2 The Traditional Approach: Seasonal Closure or Regular Discounting 353

15.3.3	Threshold Discounting	357
15.3.4	Comparing Threshold Discounting with the Traditional Approach	361
15.3.5	Impact of Strategic Customers on Threshold Discounting Performance	364
15.3.6	Mediated Threshold Discounting	366
15.3.7	Design Considerations in Threshold Discounting Offers	369
15.4	Discussion	375
	References	376
16	Innovation and Crowdsourcing Contests	379
	Laurence Ales, Soo-Haeng Cho, and Ersin Körpeoğlu	
16.1	Introduction	379
16.2	A General Model Framework for Innovation Contests	382
16.3	A Brief Taxonomy of Contest Literature	388
16.4	Contests with Uncertainty	390
	16.4.1 Optimal Award Scheme	390
	16.4.2 Open Innovation and Agents' Incentives	392
16.5	Contests with Heterogenous Agents	395
	16.5.1 Optimal Award Scheme	396
	16.5.2 Open Innovation and Agents' Incentives	397
16.6	Conclusion and Future Research	400
	References	405
 Part IV Context-Based Operational Problems in Sharing Economy		
17	Models for Effective Deployment and Redistribution of Shared Bicycles with Location Choices	409
	Mabel C. Chou, Qizhang Liu, Chung-Piaw Teo, and Deanna Yeo	
17.1	Introduction	410
	17.1.1 Review of the Bicycle-Sharing Systems	410
	17.1.2 Research Issues and Structure of the Chapter	412
17.2	The Stochastic Network Flow Model	413
	17.2.1 Equilibrium State in Time Invariant System	417
	17.2.2 Bicycle-Sharing System Design with Location Choice	419
17.3	Bicycle Sharing as Substitute for Train Rides	420
	17.3.1 Bicycle Deployment and Utilization	421
	17.3.2 Number of Bicycle Docks Needed	424
	17.3.3 Effectiveness of Bicycle Redistribution	425
17.4	Case Study on Bicycle Sharing with Location Decisions	427
17.5	Concluding Remarks	432
	References	434

18 Bike Sharing 435
Daniel Freund, Shane G. Henderson, and David B. Shmoys

18.1 Introduction 435

18.2 Data and Statistical Challenges 439

18.3 Motorized Rebalancing 442

 18.3.1 User Dissatisfaction Function 442

 18.3.2 Optimal Allocation Before the Rush 443

 18.3.3 Resulting Routing Problems 445

18.4 Allocating Capacity 448

 18.4.1 Model formulation 449

 18.4.2 Long-Run Average 450

 18.4.3 Measuring the Impact 451

18.5 Beyond Motorized Rebalancing 452

 18.5.1 Incentives 452

 18.5.2 Valets and Corrals 453

18.6 Expansion Planning 454

18.7 Conclusion 456

References 457

19 Operations Management of Vehicle Sharing Systems 461
Long He, Ho-Yin Mak, and Ying Rong

19.1 Introduction 461

19.2 Service Region Design 464

 19.2.1 Basic Model 464

 19.2.2 Customer Adoption 466

 19.2.3 Operational Profit 467

 19.2.4 Numerical Results 471

19.3 Fleet Sizing 472

 19.3.1 Two-Stage Stochastic Optimization Model 472

 19.3.2 Numerical Results 473

19.4 Fleet Repositioning 474

 19.4.1 Stochastic Dynamic Program Formulation 475

 19.4.2 The 2-Region System 477

 19.4.3 The N -Region System 479

19.5 Other Topics 479

 19.5.1 Dynamic Pricing 481

 19.5.2 Reservation Management 481

19.6 Discussion 482

References 483

20 Agent Pricing in the Sharing Economy: Evidence from Airbnb 485
 Jun Li, Antonio Moreno, and Dennis J. Zhang

20.1 Introduction 485

20.2 Literature Review and Hypothesis Development 487

 20.2.1 Literature Review 487

 20.2.2 Hypotheses Development 489

20.3 Empirical Setting and Data 491

 20.3.1 Empirical Setting: The Airbnb Platform 491

 20.3.2 Airbnb Data: Listings and Transactions 491

20.4 Performance of Professional vs. Nonprofessional Hosts:
 Econometric Specifications and Results 494

 20.4.1 Daily Revenue 494

 20.4.2 Occupancy Rate and Average Rent Price 496

 20.4.3 Exit Probability 498

20.5 Understanding the Differences in Performance 498

20.6 Conclusion 501

References 502

21 Intermediation in Online Advertising 505
 Santiago R. Balseiro, Ozan Candogan, and Huseyin Gurkan

21.1 Introduction 506

 21.1.1 Main Contributions 507

 21.1.2 Literature Review 507

21.2 Optimal Contracts for Intermediaries in Online Advertising 509

 21.2.1 Mechanism Design Problem 511

 21.2.2 Optimal Mechanism Characterization 514

 21.2.3 Economic Insights 517

21.3 Multi-stage Intermediation in Display Advertising 519

 21.3.1 Equilibrium Characterization 521

 21.3.2 Economic Insights 523

21.4 Concluding Remarks 527

References 527

Contributors

Laurence Ales Tepper School of Business, Carnegie Mellon University, Pittsburgh, PA, USA

Gad Allon University of Pennsylvania, Philadelphia, PA, USA

Jiaru Bai School of Management, Binghamton University, Binghamton, NY, USA

Santiago R. Balseiro Columbia University, New York, NY, USA

Siddhartha Banerjee Cornell University, Ithaca, NY, USA

Achal Bassamboo Northwestern University, Evanston, IL, USA

Saif Benjaafar Department of Industrial and Systems Engineering, University of Minnesota, Minneapolis, MN, USA

Kostas Bimpikis Graduate School of Business, Stanford University, Stanford, CA, USA

Gerard P. Cachon The Wharton School, University of Pennsylvania, Philadelphia, PA, USA

Ozan Candogan University of Chicago, Chicago, IL, USA

Xiqun (Michael) Chen College of Civil Engineering and Architecture, Zhejiang University, Hangzhou, China

Ying-Ju Chen School of Business and Management & School of Engineering, The Hong Kong University of Science and Technology, Kowloon, Hong Kong

Yiwei Chen Carl H. Lindner College of Business, University of Cincinnati, Cincinnati, OH, USA

Soo-Haeng Cho Tepper School of Business, Carnegie Mellon University, Pittsburgh, PA, USA

Mabel C. Chou Department of Analytics and Operations, NUS Business School, National University of Singapore, Singapore, Singapore

Eren B. Çil University of Oregon, Eugene, OR, USA

Costas Courcoubetis Engineering and Systems Design, Singapore University of Technology and Design, Singapore, Singapore

Kaitlin M. Daniels Olin Business School, Washington University in St. Louis, St. Louis, MO, USA

Daniel Freund Cornell University, Ithaca, NY, USA

Karan Girotra Cornell Tech, New York, NY, USA

Huseyin Gurkan Duke University, Durham, NC, USA

Itai Gurvich Cornell Tech, New York, NY, USA

Long He NUS Business School, National University of Singapore, Singapore, Singapore

Shane G. Henderson Cornell University, Ithaca, NY, USA

Ming Hu Rotman School of Management, University of Toronto, Toronto, ON, Canada

Rouba Ibrahim University College London, London, UK

Baojun Jiang Olin Business School, Washington University in St. Louis, St. Louis, MO, USA

Ramesh Johari Stanford University, Stanford, CA, USA

Guangwen Kong Department of Industrial and Systems Engineering, University of Minnesota, Minneapolis, MN, USA

Ersin Körpeoğlu School of Management, University College London, London, UK

Martin Lariviere Kellogg School of Management, Evanston, IL, USA

Jun Li Ross School of Business, University of Michigan, Ann Arbor, MI, USA

Xiang Li Department of Industrial and Systems Engineering, University of Minnesota, Minneapolis, MN, USA

Qizhang Liu Department of Analytics and Operations, NUS Business School, National University of Singapore, Singapore, Singapore

Ruben Lobel Airbnb, San Francisco, CA, USA

Costis Maglaras Columbia Business School, New York, NY, USA

Ho-Yin Mak Saïd Business School, University of Oxford, Oxford, UK

Simone Marinesi The Wharton School, University of Pennsylvania, Philadelphia, PA, USA

Antonio Moreno Harvard Business School, Boston, MA, USA

Serguei Netessine The Wharton School, University of Pennsylvania, Philadelphia, PA, USA

Yiangos Papanastasiou Haas School of Business, University of California, Berkeley, CA, USA

Ying Rong Antai College of Economics and Management, Shanghai Jiao Tong University, Shanghai, China

Mengze Shi Rotman School of Management, University of Toronto, Toronto, ON, Canada

David B. Shmoys Cornell University, Ithaca, NY, USA

Kut C. So The Paul Merage School of Business, University of California, Irvine, CA, USA

Christopher S. Tang Anderson School, University of California, Los Angeles, Los Angeles, CA, USA

Chung-Piaw Teo Department of Analytics and Operations, NUS Business School, National University of Singapore, Singapore, Singapore

Lin Tian School of Management, Fudan University, Shanghai, China

Tunay I. Tunca Robert H. Smith School of Business, University of Maryland, College Park, MD, USA

Gustavo Vulcano School of Business, Universidad Torcuato di Tella, Buenos Aires, Argentina

Hai Wang School of Information Systems, Singapore Management University, Singapore, Singapore

Jiahua Wu Imperial College Business School, Imperial College London, London, UK

Deanna Yeo Department of Analytics and Operations, NUS Business School, National University of Singapore, Singapore, Singapore
GE Healthcare, Singapore, Singapore

Dennis J. Zhang Olin Business School, Washington University in St. Louis, St. Louis, MO, USA

Yun Zhou DeGroote School of Business, McMaster University, Hamilton, ON, Canada

Chapter 1

Introduction



Ming Hu

Abstract This introduction provides an overview of the book with highlights of all chapters.

This book aims to address what it takes to be successful in today's sharing economy from an operational perspective. Sharing economy may refer to an online platform that enables individuals or small entities as buyers and sellers to “transact” effectively and efficiently or a market model that allows sharing of access to goods and services. Operations management has the tradition of coming from and going back to real-life applications. It deals with the management of the processes of matching supply with demand. The emerging business processes in a sharing economy call for active management, as well as adequate attention from operations researchers. However, as the business side of a sharing economy is still emerging and rapidly evolving, there is a lack of a comprehensive overview of ongoing academic efforts in addressing its operational problems. To fill the void, this book is, to the best of our knowledge, the first to present cutting-edge research on sharing economy from globally recognized field experts organized in one place. For future research directions, a good resource is Chen et al. (2018).

1.1 Overall Structure

This book is comprised of 21 chapters that are divided into four parts.

- The first part (Chaps. 2, 3, 4, and 5) explores the *general impact* of sharing economy on the production, consumption, and society. For example, with sharing dynamics taken into account, how the sharing economy affects the

M. Hu (✉)
Rotman School of Management, University of Toronto, Toronto, ON, Canada
e-mail: ming.hu@rotman.utoronto.ca

© Springer Nature Switzerland AG 2019
M. Hu (ed.), *Sharing Economy*, Springer Series in Supply Chain Management 6,
https://doi.org/10.1007/978-3-030-01863-4_1

consumption of goods and services and consumer welfare. Moreover, the section also highlights operational opportunities and challenges of a sharing economy.

- The second part (Chaps. 6, 7, 8, 9, 10, 11, 12, and 13) explores the *intermediary role* of a sharing platform that matches crowdsourced supply with demand. The decisions of the platform can be pricing decisions on the supply and demand sides, detailed matching decisions at the operational level, or decisions about capacity, information disclosure and payment schemes.
- The third part (Chaps. 14, 15, and 16) investigates the *crowdsourcing management* on a sharing platform with the goal to crowdsource both demand (group buying) and supply, such as funds (crowdfunding) and innovative ideas (tournament).
- The fourth part is (Chaps. 17, 18, 19, 20, and 21) dedicated to *context-based operational problems* of popular sharing economy applications, for example, how to dynamically rebalance bikes for a bike-sharing system, how to design service zones for one-way carsharing services such as Car2Go, and for homeowners how should they set prices on Airbnb.

Ultimately, the book introduces the reader to the fundamentals of operations in sharing economy and highlights the latest research on the topic.

1.2 Chapter Highlights

1.2.1 Part I: Impact of Sharing Economy

1.2.1.1 Economic Impact

In Chap. 2, Saif Benjaafar, Guangwen Kong, Xiang Li, and Costas Courcoubetis study an equilibrium model of peer-to-peer product sharing, or collaborative consumption, where individuals with different usage levels make decisions about whether to own or rent a homogenous product. Owners can generate income from renting their products to non-owners while non-owners can access these products through renting. The authors characterize equilibrium outcomes, including ownership and usage levels, consumer surplus, and social welfare. They compare these equilibrium outcomes in systems with and without collaborative consumption and examine the impact of various problem parameters. Their findings indicate that collaborative consumption can result in either lower or higher ownership and usage levels, with higher ownership and usage levels more likely when the cost of ownership is high.

In Chap. 3, Baojun Jiang and Lin Tian also examine the strategic and economic impact of product sharing among consumers. Consumers buy many products but end up not fully utilizing them. A product owner's self-use values can differ over time, and in a period of low self-use value, the owner may rent out her product in a product-sharing market. Transaction costs in the sharing market have a non-monotonic effect on the manufacturer's profits, consumer surplus, and social

welfare. When the manufacturer strategically chooses its retail price, consumers' sharing of products with high marginal costs is win-win both for the firm and the consumers, whereas their sharing of products with low marginal costs can be lose-lose. Moreover, in the presence of a sharing market, the firm will find it optimal to strategically increase its quality, leading to higher profits but lower consumer surplus. Lastly, within a distribution channel framework, product sharing can sometimes benefit the downstream retailer at the expense of the upstream manufacturer.

1.2.1.2 Operational Opportunity and Challenge

In Chap. 4, Tunay Tunca builds a framework for identifying, describing and analyzing operational factors that shape the efficiency of a sharing economy. In particular, these factors are: (1) utilization of sunk and fixed costs, (2) utilization of bit-sized resources, (3) utilization of human idle time, (4) utilization of networks to lower barriers to entry into workforce and markets and (5) assigning people new operational and economic roles. Then he discusses some potential downsides and pitfalls that arise as the side effects of these operational efficiencies of the sharing economy business models and foreseeable regulatory issues that may need attention.

In Chap. 5, Siddhartha Banerjee and Ramesh Johari outline the main challenges of ridesharing platforms in various aspects such as large-scale learning, real-time stochastic control, and market design. The authors present an approach to modeling, optimizing, and reasoning about such platforms, and describe how rigorous analysis has been used with great success in designing efficient algorithms for real-time decision making, in informing the market design aspects of these platforms, and in understanding the impact of these platforms in a broader societal context.

1.2.2 Part II: Intermediary Role of a Sharing Platform

1.2.2.1 Intermediation via Pricing and Matching

The following three chapters are primarily motivated by ridesharing platforms.

In Chap. 6, Gerard Cachon, Kaitlin Daniels, and Ruben Lobel focus on two-sided pricing as a moderating mechanism. The platforms may charge consumers prices and pay individual service providers wages, conditional on market conditions. The authors study several pricing schemes, with a specific focus on a contingent pricing policy that requires wages to be a fixed commission rate of dynamic prices. Although this heuristic policy is not optimal, it is shown to generally achieve nearly the optimal profit. As labor becomes more expensive, consumers are better off with the heuristic contingent pricing policy relative to fixed pricing, because they benefit both from lower prices during normal demand and expanded access to service during peak demand.

In Chap. 7, Jiaru Bai, Rick So, Chris Tang, Xiqun Chen, and Hai Wang also study two-sided pricing in a sharing economy. They adopt a queueing model with both the supply and demand endogenously dependent on the price the platform charges its customers and the wage the platform pays its independent providers. The authors use the steady-state performance in equilibrium to characterize the optimal price, optimal wage and optimal commission rate that maximize the profit of the platform. They find that it is optimal for the platform to offer time-dependent commission rates by providing a higher rate during peak hours and a lower rate otherwise.

In Chap. 8, Yiwei Chen, Ming Hu, and Yun Zhou study the pricing and matching decisions of a platform in simultaneously managing the supply and demand. First, the authors explore how the platform could optimally set the price and wage for a single service or product in different market conditions, and provide provable performance guarantee for the fixed commission contract. Second, even with determined pricing decisions, the platform still faces the task of matching customers with suppliers. Then they consider a stochastic, dynamic model with multiple demand types to be matched with multiple supply types over a planning horizon. They characterize the optimal matching policy by determining the priorities of the demand-supply pairs, under a sufficient condition on the reward structure. Finally, they study the joint pricing and matching decision by a platform for a single service or product and take into account suppliers' and customers' forward-looking behavior. They propose a simple heuristic policy and show it is asymptotically optimal when both sides of the market have sufficiently large volumes.

1.2.2.2 Intermediation via Information and Payment

In Chap. 9, Gad Allon, Achal Bassamboo, and Eren Çil study large-scale, web-based service marketplaces, where many small service providers compete among themselves in catering to customers with diverse needs. Customers who frequent these marketplaces seek quick resolutions and thus are usually willing to trade prices with waiting times. They discuss the role of the moderating platform in facilitating information gathering, operational efficiency, and communication among agents in such service marketplaces. Perhaps surprisingly, they show that operational efficiency may be detrimental to the overall efficiency of the marketplace. Then they establish that to reap the expected gains of operational efficiency for the marketplace, the moderating platform may need to complement the operational efficiency by enabling communication among its agents.

In Chap. 10, Kostas Bimpikis and Yiangos Papanastasiou focus on the information disclosure as a moderating scheme to incentivize customers to take system-optimal actions. Crowd-sourced content in the form of online product reviews or recommendations is an integral feature of most Internet-based service platforms and marketplaces. Customers may find such information useful when deciding among potential alternatives; at the same time, the process of generating such content is mainly driven by the customers' decisions themselves. The authors

focus on a platform that can potentially incentivize the actions of self-interested customers by appropriately designing an information provision policy or a payment scheme.

In Chap. 11, Ying-Ju Chen, Costis Maglaras, and Gustavo Vulcano study an aggregated marketplace where potential buyers arrive and submit requests-for-quotes. There are independent suppliers each modeled as a queueing system that competes for these requests. Each supplier offers a bid that comprises a fixed price and a dynamic target lead time, and the cheapest supplier wins the order as long as the quote meets the buyer's willingness to pay. The authors characterize the asymptotic performance of this system as the demand and the supplier capacities grow large and obtain insights into the equilibrium behavior of the suppliers. To overcome the efficiency loss from supplier competition, they propose a compensation-while-idling mechanism that the marketplace can impose: each supplier gets monetary compensation from other suppliers during his idle time. This mechanism induces suppliers to implement the centralized solution.

1.2.2.3 Intermediation in the Presence of Self-Scheduling Suppliers

Although the self-interested behavior of individual suppliers is an indispensable feature of most of the previous chapters in Part II, the following chapters build on classical operational models such as the newsvendor model and queueing systems and focus specifically on incorporating the self-scheduling behavior of individual suppliers.

In Chap. 12, Itai Gurvich, Martin Lariviere, and Antonio Moreno study capacity management of a service provider over a horizon when its workers have the flexibility to choose when they will (or will not) work and optimize their schedules based on the offered compensation and individual availability. The authors provide an augmented newsvendor formula to capture the tradeoffs for the firm and the agents. If the firm could keep the flexibility but have direct control of agents for the same wages, it would not only generate higher profit, as it is expected, but would also provide better service levels to its customers. If the agents require a "minimum wage" to remain in the agent pool, they will have to relinquish some of their flexibility. To pay a minimum wage, the firm must restrict the number of agents that can work in some time intervals. If the pool of agents is sufficiently large relative to peak demand, the firm benefits from self-scheduling behavior of individual suppliers.

In Chap. 13, Rouba Ibrahim also focuses on the self-scheduling behavior of individual workers. When such behavior is allowed, the number of workers available in any period is uncertain. She adopts a queueing-theoretic framework to study the effective management of service systems where the number of available agents is random. She begins by surveying some theoretical results on the control of queueing systems with uncertainty in the number of servers. Then, she illustrates how to apply those theoretical results to study the problems of staffing and controlling queueing systems with self-scheduling workers and impatient, time-sensitive, customers.

1.2.3 Part III: Crowdsourcing Management

1.2.3.1 Group Buying and Crowdfunding

In Chap. 14, Ming Hu, Mengze Shi, and Jiahua Wu investigate the two popular business models, namely, online group buying and crowdfunding. The former crowdsources demand, and the latter crowdsources funds. Both share the same unique feature of an all-or-nothing mechanism, where transactions will take place only if the total number of committed purchases or pledges exceeds a specified threshold within a specified period. The authors seek to understand the impact of the all-or-nothing mechanism on consumer behavior, as well as the optimal design of such mechanisms from the perspective of third-party platforms like Groupon and Kickstarter. First, using a dataset from the online group buying industry, they empirically identify two types of threshold-induced effects on consumer behavior. Next, they study the optimal design of all-or-nothing mechanisms from two different perspectives, namely, information disclosure and pricing.

In Chap. 15, Simone Marinesi, Karan Girotra, and Serguei Netessine study group buying and its impact on a service provider. They model a capacity-constrained firm offering service to a random-sized population of strategic customers in two representative time periods, a desirable hot period and a less desirable slow period. They show that strategic consumer behavior under group buying with an all-or-nothing threshold increases the firm's profits. When threshold discounts are offered through an intermediary platform, arrangements often used in practice distort the incentives of the intermediary, and typically result in a higher discount and a lower activation threshold relative to what would be optimal for the service firm. The authors consider alternative deal designs and find that the best designs compromise the service provider's flexibility to provide customers with clear offer terms.

1.2.3.2 Crowdsourcing Contest

In Chap. 16, Laurence Ales, Soo-Haeng Cho, and Ersin Körpeoğlu present a general model framework of innovation contests, in which an organizer crowdsources solutions to an innovation-related problem from a group of independent agents. Agents, who can be heterogeneous in their ability levels, exert efforts to improve their solutions, and their solution qualities are uncertain due to the innovation and evaluation processes. The framework captures main features of a contest and encompasses several existing models in the literature. Using this framework, the authors analyze two critical decisions of the organizer: a set of awards that will be distributed to agents and whether to restrict entry to a contest or to run an open contest. They provide a taxonomy of the contest literature and discuss past and current research on innovation contests as well as a set of exciting future research directions.

1.2.4 Part IV: Context-Based Operational Problems in Sharing Economy

1.2.4.1 Bike Sharing

In Chap. 17, Mabel Chou, Qizhang Liu, Chung-Piaw Teo, and Deanna Yeo develop practical operations models to support decision making in the design and management of public bicycle-sharing systems. They develop a network flow model with proportionality constraints to estimate the flow of bicycles within the network, and to estimate the number of trips and the number of docks needed at each station. The authors also examine the impact of periodic redistribution of bicycles in the network. The same approach can be extended to incorporate the decisions of station locations, by taking into account the proportional flow constraints into a mixed-integer programming formulation. Using a set of bus transit data, they implemented this approach to identify the ideal locations for the bicycle stations in a new town of Singapore.

In Chap. 18, Daniel Freund, Shane Henderson, and David Shmoys also discuss planning methods for bike-sharing systems. They study specific questions such as decisions related to the number of docks to allocate to each station, how to rebalance the system by moving bikes to match demand, and how to expand the network. They discuss linear integer programming models, specially-tailored optimization algorithms, and simulation methods. All of these methods rely on a careful statistical analysis of bike-sharing data, which they also briefly review. This chapter is based on their 4-year collaboration with Citi Bike in New York City, and its parent company Motivate.

1.2.4.2 Vehicle Sharing

In Chap. 19, Long He, Ho-Yin Mak, and Ying Rong study the free-float model of vehicle sharing, which allows users to start and end rentals at any location within a defined service region. Compared with conventional models of vehicle sharing, the free-float model offers its users the flexibility to make one-way, two-way and multi-stop trips, and as a result, provides a more viable alternative to individual vehicle ownership. On the other hand, the flexibility of the free-float model leads to many operations management challenges that must be overcome for such vehicle sharing systems to be economically sustainable. The authors review several operations management problems in vehicle sharing including system design, vehicle repositioning, fleet sizing, dynamic pricing and reservation policy. In particular, they discuss the optimization models for service region design and fleet repositioning.

1.2.4.3 Short-Term Rental

In Chap. 20, Jun Li, Antonio Moreno, and Dennis J. Zhang study Airbnb, the largest marketplace that allows people to rent short-term lodging from property owners. One of the distinct features of such a sharing-economy marketplace is that the supply side includes individual nonprofessional decision makers, in addition to firms and professional agents. Using a data set of prices and availability of listings on Airbnb, the authors find that there exist substantial differences in the operational and financial performance of professional and nonprofessional hosts. They provide empirical evidence to explain such performance differences between professionals and nonprofessionals: nonprofessional hosts are less likely to offer contingent rates across stay dates based on the underlying demand patterns.

1.2.4.4 Online Advertising

In Chap. 21, Santiago Balseiro, Ozan Candogan, and Huseyin Gurkan study online advertising, in which impressions are sold to advertisers via real-time auctions organized by central platforms referred to as ad exchanges. Advertisers participate in the auctions run by exchanges through intermediaries which acquire impressions on their behalf. Intermediaries are specialized entities that provide targeted services for a particular segment of the market, and typically there are multiple stages of intermediation. Moreover, an advertiser may have private information, e.g., budget, targeting criterion or value attributed to an impression. First, the authors study the mechanism design problem of an intermediary who offers a contract to an advertiser with a private budget and a private targeting criterion. They characterize the optimal mechanism and establish that the presence of the intermediary results in more straightforward bidding policies. Next, they study the strategic interaction among intermediaries organized in a chain network. They characterize a subgame perfect equilibrium of the resulting game among intermediaries and show that the most profitable position in the intermediation chain depends on the underlying value distribution of the advertiser.

References

Chen Y-J, Dai T, Körpeoğlu CG, Körpeoğlu E, Sahin O, Tang CS, Xiao S (2018, Forthcoming) Innovative online platforms: research opportunities. *Manuf Serv Oper Manag*

Part I
Impact of Sharing Economy

Chapter 2

Peer-to-Peer Product Sharing



Saif Benjaafar, Guangwen Kong, Xiang Li, and Costas Courcoubetis

Abstract We describe an equilibrium model of peer-to-peer product sharing, or collaborative consumption, where individuals with varying usage levels make decisions about whether or not to own a homogenous product. Owners are able to generate income from renting their products to non-owners while non-owners are able to access these products through renting on as needed basis. We characterize equilibrium outcomes, including ownership and usage levels, consumer surplus, and social welfare. We compare each outcome in systems with and without collaborative consumption and examine the impact of various problem parameters. Our findings indicate that collaborative consumption can result in either lower or higher ownership and usage levels, with higher ownership and usage levels more likely when the cost of ownership is high. Our findings also indicate that consumers always benefit from collaborative consumption, with individuals who, in the absence of collaborative consumption, are indifferent between owning and not owning benefitting the most. We study both profit maximizing and social welfare maximizing platforms and compare equilibrium outcomes under both in terms of ownership, usage, and social welfare. We find that the difference in social welfare between the profit maximizing and social welfare maximizing platforms is relatively modest.

This chapter is based on the paper “Peer-to-Peer Product Sharing: Implications for Ownership, Usage and Social Welfare in the Sharing Economy”, Published Online: 16 May 2018 in *Management Science*, <https://doi.org/10.1287/mnsc.2017.2970>

S. Benjaafar · G. Kong (✉) · X. Li
Department of Industrial and Systems Engineering, University of Minnesota, Minneapolis, MN,
USA
e-mail: saif@umn.edu; gkong@umn.edu; lix1315@umn.edu

C. Courcoubetis
Engineering and Systems Design, Singapore University of Technology and Design, Singapore,
Singapore
e-mail: costas@sutd.edu.sg

2.1 Introduction

We are witnessing, across a wide range of domains, a shift away from the exclusive ownership and consumption of resources to one of shared use and consumption. This shift is taking advantage of innovative new ways of peer-to-peer sharing that are voluntary and enabled by internet-based exchange markets and mediation platforms. Value is derived from the fact that many resources are acquired to satisfy infrequent demand but are otherwise poorly utilized (for example, the average car in the US is used less than 5% of the time). Several successful businesses in the US and elsewhere, such as Getaround for cars, Spinlister for bikes, 3D Hubs for 3D printers, LiquidSpace for office space, MachineryLink for farm equipment and JustPark for parking, provide a proof of concept and evidence for the viability of peer-to-peer product sharing or collaborative consumption (the term we use in the rest of the chapter). These businesses and others allow owners to rent on a short-term basis poorly utilized assets and non-owners to access these assets through renting on an as-needed basis. Collectively, these businesses and other manifestations of the collaborative consumption of products and services are giving rise to what is becoming known as the sharing economy.¹

The peer-to-peer sharing of products is not a new concept. However, recent technological advances in several areas have made it more feasible by lowering the associated search and transaction costs. These advances include the development of online marketplaces, mobile devices and platforms, electronic payments, and two-way reputation systems whereby users rate providers and providers rate users. Other drivers behind the rise of collaborative consumption are societal and include increased population density in urban areas around the world, increased concern about the environment (collaborative consumption is viewed as a more sustainable alternative to traditional modes of consumption), and increased desire for community and altruism among the young and educated.

Collaborative consumption has the potential of increasing access while reducing investments in resources and infrastructure. In turn, this could have the twin benefit of improving consumer welfare (individuals who may not otherwise afford a product now have an opportunity to use it) while reducing societal costs (externalities, such as pollution that may be associated with the production, distribution, use, and disposal of the product). It also has the potential of providing a source of net income for owners by monetizing poorly utilized assets, which are in some cases also expensive and rapidly depreciating. Take cars for example. The availability of a sharing option could lead some to forego car ownership in favor

¹The term sharing economy has been used to refer to businesses that enable the foregoing of ownership in favor of “on-demand” access. In several cases, this involves a single entity that owns the physical assets (e.g., Zipcar for short term car rentals). It also encompasses the peer-to-peer provisioning of services (e.g., Uber for transportation services, TaskRabbit for errands, and Postmates for small deliveries). For further discussion and additional examples, see Botsman and Rogers (2010), Malhotra and Van Alstyne (2014), Cusumano (2014), and Chase (2015).

of on-demand access. In turn, this could result in a corresponding reduction in congestion and emissions and, eventually, in reduced investments in roads and parking infrastructure. However, increased collaborative consumption may have other consequences, some of which may be undesirable. For example, greater access to cars could increase car usage and, therefore, lead to more congestion and pollution if it is not accompanied by a sufficient reduction in the numbers of cars.² It could also lead to speculative investments in cars and price inflation, or affect the availability and pricing of other modes of public transport, such as taxis, buses, and trains.

Collaborative consumption raises several important questions. How does collaborative consumption affect ownership and usage of resources? Is it necessarily the case that collaborative consumption leads to lower ownership, lower usage, or both (and therefore to improved environmental impact)? If not, what conditions would favor lower ownership, lower usage, or both? Who benefits the most from collaborative consumption among owners and renters? To what extent would a profit maximizing platform, through its choice of rental prices, improve social welfare? To what extent do frictions, such as extra wear and tear renters place on rented resources and inconvenience experienced by renters affect platform profit and social welfare?

In this chapter, we address these and other related questions. We describe an equilibrium model of peer-to-peer product sharing, where individuals with varying usage levels make decisions about whether or not to own a homogenous product. In the presence of collaborative consumption, owners are able to generate income from renting their products to non-owners while non-owners are able to access these products through renting. The matching of owners and renters is facilitated by a platform, which sets the rental price and charges a commission fee.³ Because supply and demand can fluctuate over the short run, we allow for the possibility that an owner may not always be able to find a renter when she puts her product up for rent. Similarly, we allow for the possibility that a renter may not always be able to find a product to rent when he needs one. We refer to the uncertainty regarding the availability of renters and products as matching friction and describe a model for this uncertainty. We also account for the cost incurred by owners due to the extra

²An article in the *New York Times* (2015) notes that “The average daytime speed of cars in Manhattan’s business districts has fallen to just under 8 miles per hour this year, from about 9.15 miles per hour in 2009. City officials say that car services like Uber and Lyft are partly to blame. So Mayor Bill de Blasio is proposing to cap their growth.” Note that, although the peer-to-peer product sharing we consider is different from the type of product sharing enabled by Uber (which requires the involvement of the owner as a service provider), the two share similarities in that they provide non-owners with access to a product without having to own it.

³A variety of pricing approaches are observed in practice. Some platforms allow owners to choose their own prices. Others (e.g., DriveMycar) determine the price. There are also cases where the approach is hybrid, with owners determining a minimum acceptable price but allowing the platform to adjust it higher (e.g., Turo), or with the platform suggesting a price (e.g., JustShareIt) but allowing owners to deviate. From conversations the authors had with several industry executives, there appears to be a push toward platform pricing, with several platforms investing in the development of sophisticated pricing engines to support owners.

wear and tear that a renter places on a rented product and for the inconvenience cost experienced by renters for using a product that is not their own.

For a given price and a commission rate, we characterize equilibrium ownership and usage levels, consumer surplus, and social welfare. We compare each in systems with and without collaborative consumption and examine the impact of various problem parameters including price, commission rate, cost of ownership, extra wear and tear cost, and inconvenience cost. We also do so when the price is a decision made by the platform to maximize either profit or social welfare. Our main findings include the following:

- Depending on the rental price, we show that collaborative consumption can result in either higher or lower ownership. In particular, we show that when the rental price is sufficiently high (above a well-specified threshold), collaborative consumption leads to higher ownership. We show that this threshold is decreasing in the cost of ownership. That is, collaborative consumption is more likely to lead to more ownership when the cost of ownership is high (this is because collaborative consumption allows individuals to offset the high ownership cost and pulls in a segment of the population that may not otherwise choose to own).
- Similarly, we show that collaborative consumption can lead to either higher or lower usage, with usage being higher when price is sufficiently high. Thus, it is possible for collaborative consumption to result in both higher ownership and higher usage (it is also possible for ownership to be lower but usage to be higher and for both ownership and usage to be lower).
- These results continue to hold in settings where the rental price is determined by a profit maximizing or a social welfare maximizing platform. In particular, collaborative consumption can still lead to either higher or lower ownership and usage with higher ownership and usage more likely when the cost of ownership is higher.
- We show that consumers always benefit from collaborative consumption, with individuals who, in the absence of collaborative consumption, are indifferent between owning and not owning benefitting the most. This is because among non-owners those with the most usage (and therefore end up renting the most) benefit the most from collaborative consumption. Similarly, among owners, those with the least usage (and therefore end up earning the most rental income) benefit the most.
- For a profit maximizing platform, we show that profit is not monotonic in the cost of ownership, implying that a platform is least profitable when the cost of ownership is either very high or very low (those two extremes lead to scenarios with either mostly renters and few owners or mostly owners and few renters). The platform is most profitable when owners and renters are sufficiently balanced. For similar reasons, social welfare is also highest when owners and renters are sufficiently balanced.
- We observe that profit is also not monotonic in the extra wear and tear renters place on a rented product, implying that a platform may not always have an incentive to reduce this cost. This is because the platform can leverage this cost to induce desirable ownership levels without resorting to extreme pricing, which can be detrimental to its revenue.

The rest of the chapter is organized as follows. In Sect. 2.2, we provide a review of related literature. In Sect. 2.3, we describe our model. In Sect. 2.4, we provide an analysis of the equilibrium. In Sect. 2.5, we consider the platform's problem. In Sect. 2.6, we offer concluding comments. Proofs and various extensions of the analysis can be found in the full length paper Benjaafar et al. (2018).

2.2 Literature Review

Our work is related to the literature on peer-to-peer markets (see Einav et al. 2016 for a recent review). Within this literature, there is a small but growing stream that deals with peer-to-peer markets with collaborative consumption features. Fradkin et al. (2015) studies sources of inefficiency in matching buyers and suppliers in online market places. Using a counterfactual study, they show how changes to the ranking algorithm of Airbnb can improve the rate at which buyers are successfully matched with suppliers. Zervas et al. (2015) examine the relationship between Airbnb supply and hotel room revenue and find that an increase in Airbnb supply has only a modest negative impact on hotel revenue. Cullen and Farronato (2018) describe a model of peer-to-peer labor marketplaces. They calibrate the model using data from TaskRabbit and find that supply is highly elastic, with increases in demand matched by increases in supply per worker with little or no impact on price.

Papers that are closest in spirit to ours are Fraiberger and Sundararajan (2015) and Jiang and Tian (2018). Fraiberger and Sundararajan (2015) describe a dynamic programming model where individuals make decisions in each period regarding whether to purchase a new car, purchase a used a car, or not purchase anything. They model matching friction, as we do, but assume that the renter-owner matching probabilities are exogenously specified and not affected by the ratio of owners to renters (in our case, we allow for these to depend on the ratio of owners to renters which turns out to be critical in the decisions of individuals on whether to own or rent). They use the model to carry out a numerical study. For the parameter values they consider, they show that collaborative consumption leads to a reduction in new and used car ownership, an increase in the fraction of the population who do not own, and an increase in the usage intensity per vehicle. In this chapter, we show that ownership and usage can actually either increase or decrease with collaborative consumption and provide analytical results regarding conditions under which different combinations of outcomes can occur. We also study the decision of the platform regarding pricing and the impact of various parameters on platform profitability.

Jiang and Tian (2018) describe a two-period model, where individuals first decide on whether or not to own a product. This is followed by owners deciding in each period on whether to use the product themselves or rent it. They assume that demand always matches supply through a market clearing price and do not consider, as we do, the possibility of a mismatch, because of matching friction, between supply and demand. They focus on the decision of the product manufacturer. In particular, they

study how the manufacturer should choose its retail price and product quality in anticipation of sharing by consumers. In contrast, we focus on the decision of the platform which in our case decides on the rental price.

Empirical studies that examine the impact of peer-to-peer product sharing on ownership and usage are scarce. Clark et al. (2014) present results from a survey of British users of a peer-to-peer car sharing service. They find that peer-to-peer car sharing has led to a net increase in the number of miles driven by car renters. van der Linden and Franciscus (2016) examine differences in the prevalence of peer-to-peer car sharing among several European cities. They find that peer-to-peer car sharing is more prevalent in cities where a larger share of trips is taken by public transport and where there is a city center less suitable for car use. Ballus-Armet et al. (2014) report on a survey in San Francisco of public perception of peer-to-peer car sharing. They find that approximately 25% of surveyed car owners would be willing to share their personal vehicles through peer-to-peer car sharing, with liability and trust concerns being the primary deterrents. They also find that those who drive almost every day are less likely to rent through peer-to-peer car sharing, while those who use public transit at least once per week are more likely to do so. There are a few studies that consider car sharing that involves a third party service provider, such as a car rental company. For example, Nijland et al. (2015) (and also Martin and Shaheen 2011) find that car sharing would lead to a net decrease in car usage. On the other hand, a study by KPMG (Korosec 2015) projects a significant increase in miles driven by cars and attributes this to increased usage of on-demand transportation services. In general, there does not appear to be a consensus yet on the impact of car sharing on car usage and ownership. This chapter, by providing a framework for understanding how various factors may affect product sharing outcomes, could be useful in informing future empirical studies.

There is a growing body of literature that focuses on the concept of *servicization*. Servicization refers to a business model under which a firm that supplies a product to the market retains ownership of the product and instead charges customers per use (e.g., printer manufacturers charging customers per printed page instead of charging them for the purchase of a printer or car manufacturers renting cars on a short term basis instead of selling them or leasing them on a long term basis). Agrawal and Bellos (2017) examine the extent to which servicization affects ownership and usage and the associated environmental impact.⁴ Orsdemir et al. (2017) evaluate both the profitability and the environmental impact of servicization. Bellos et al. (2017) study the economic and environmental implications of an auto manufacturer, in addition to selling cars, offering a car sharing service. Additional discussion and examples of

⁴Under a servicization model, the firm can exert costly effort to improve certain characteristics of the product such as its energy efficiency during use or its durability. This could lower the corresponding operating costs, which in turn could result in higher usage. The phenomenon of higher efficiency leading to more usage is commonly referred to as the *rebound effect*. See Greening et al. [2000] for an overview and references. In our setting, the introduction of collaborative consumption can lead, under some conditions, to higher ownership because of the rental income owners derive from ownership.

servicization can be found in Agrawal and Bellos (2017) and the references therein. Peer-to-peer product sharing is different from servicization in that there is no single entity that owns the rental units, with owners being simultaneously consumers and suppliers of the market. As a result, the payoff of one side of the market depends on the availability of the other side. This, coupled with the fact that supply and demand are not guaranteed to be matched with each other, makes ownership and usage decisions more complicated than those under servicization.

Finally, we note that collaborative consumption has the features of a two sided-market (see for example Rochet and Tirole 2006; Weyl 2010; Hagiu and Wright 2015). Collaborative consumption is different from two-sided markets in several ways, the most important of which is that the two sides are not distinct. In collaborative consumption, being either an owner or a renter is a decision that users of the platform make, with more owners implying fewer renters, and vice-versa. Collaborative consumption shares also features of secondary markets for used goods (see for example Waldman 2003). Markets with collaborative consumption are different from those with a secondary market for used goods in that there is no permanent transfer of ownership from the seller to the buyer and the usage by the renter does not preclude usage by the owner.

2.3 Model Description

In this section, we describe our model of collaborative consumption. The model is applicable to the case of peer to peer product sharing where owners make their products available for rent when they are not using them and non-owners can rent from owners to fulfill their usage needs. We reference the case of car sharing. However, the model applies more broadly to the collaborative consumption of other products. We consider a population of individuals who are heterogeneous in their product usage, with their type characterized by their usage level ξ . We assume that usage is exogenously determined (i.e., the usage of each individual is mostly inflexible) and the utility derived by an individual with type ξ , $u(\xi)$ is linear in ξ with $u(\xi) = \xi$. We use a linear utility for ease of exposition and to allow for closed form expressions. A linear utility has constant returns to scale, and, without loss of generality, the utility derived from each unit of usage can be normalized to 1. Also without loss of generality, we normalize the usage level to $[0, 1]$, where $\xi = 0$ corresponds to no usage at all and $\xi = 1$ to full usage. We let $f(\xi)$ denote the density function of the usage distribution in the population.

We assume products are homogeneous in their features, quality, and cost of ownership. In the absence of collaborative consumption, each individual makes a decision about whether or not to own. In the presence of collaborative consumption, each individual decides on whether to own, rent from others who own, or neither. Owners incur the fixed cost of ownership but can now generate income by renting their products to others who choose not to own. Renters pay the rental fee but avoid the fixed cost of ownership.

We let p denote the rental price per unit of usage that renters pay (a uniform price is consistent with observed practices by certain peer-to-peer platforms when the goods are homogenous). This rental price may be set by a third party platform (an entity that may be motivated by profit, social welfare, or some other concern; see Sect. 2.5 for further discussion). The platform extracts a commission from successful transactions. We denote the commission rate as γ , where $0 \leq \gamma < 1$, so that the rental income seen by the owner per unit of usage is $(1 - \gamma)p$. We let α , where $0 \leq \alpha \leq 1$, denote the fraction of time in equilibrium that an owner, whenever she puts her product up for rent, is successful in finding a renter. Similarly, we denote by β , where $0 \leq \beta \leq 1$, the fraction of time that a renter, whenever he decides to rent, is successful in finding an available product (the parameters α and β are determined endogenously in equilibrium). A renter resorts to his outside option (e.g., public transport in the case of cars) whenever he is not successful in finding a product to rent. The owner incurs a fixed cost of ownership, denoted by c , which may include not just the purchase cost (if costs are expressed per unit time, this cost would be amortized accordingly) but also other ownership-related costs such as those related to storage and insurance. Whenever the product is rented, the owner incurs an additional cost, denoted by d_o , due to extra wear and tear the renter places on the product. Renters, on the other hand, incur an inconvenience cost, denoted by d_r (in addition to paying the rental fee), from using someone else's product and not their own. Without loss of generality, we assume that $c, p, d_o, d_r \in [0, 1]$ and normalize the value of the outside option (e.g., using public transport) to 0.

We assume that $p(1 - \gamma) \geq d_o$ so that an owner would always put her product out for rent when she is not using it. Note that usage corresponds to the portion of time an owner would like to have access to her product, regardless of whether or not she is actually using it. An owner has always priority in accessing her product. Hence her usage can always be fulfilled. We also assume that $p + d_r \leq 1$ so that a renter always prefers renting to the outside option. Otherwise, rentals would never take place as the outside option is assumed to be always available. There are of course settings where an individual would like to use a mix of options (e.g., different transportation methods). In that case, ξ corresponds to the portion of usage that an individual prefers to fulfill using the product (e.g., a car and not public transport).

The *payoff* of an owner with usage level ξ can now be expressed as

$$\pi_o(\xi) = \xi + (1 - \xi)\alpha[(1 - \gamma)p - d_o] - c, \quad (2.1)$$

while the payoff of a renter as

$$\pi_r(\xi) = \beta\xi - \beta(p + d_r)\xi. \quad (2.2)$$

The payoff of an owner has three terms: the utility derived from usage, the income derived from renting (net of the wear and tear cost), and the cost of ownership. The income from renting is realized only when the owner is able to find a renter. The

payoff of a renter is the difference between the utility derived from renting and the cost of renting (the sum of rental price and inconvenience cost). A renter derives utility and incurs costs whenever he is successful in renting a product.

An individual with type ξ would participate in collaborative consumption as an *owner* if the following conditions are satisfied

$$\pi_o(\xi) \geq \pi_r(\xi) \quad \text{and} \quad \pi_o(\xi) \geq 0.$$

The first constraint ensures that an individual who chooses to be an owner prefers to be an owner to being a renter. The second constraint is a participation constraint that ensures the individual participates in collaborative consumption. Similarly, an individual with type ξ would participate in collaborative consumption as a *renter* if the following conditions are satisfied

$$\pi_r(\xi) \geq \pi_o(\xi) \quad \text{and} \quad \pi_r(\xi) \geq 0.$$

Noting that, for any given pair of α and β in $[0, 1]$, $\pi_o(\xi) - \pi_r(\xi)$ is monotonically increasing and $\pi_r(\xi) \geq 0$ for $\xi \in [0, 1]$, collaborative consumption would take place if there exists $\theta \in (0, 1)$ such that

$$\pi_o(\theta) = \pi_r(\theta). \tag{2.3}$$

The parameter θ would then segment the population into owners and renters, where individuals with $\xi > \theta$ are owners and individuals with $\xi < \theta$ are renters (an individual with $\xi = \theta$ is indifferent between owning and renting).

We refer to $\omega = \int_{[\theta, 1]} f(\xi) d\xi$, the fraction of owners in the population, as the *ownership level* or simply *ownership*. In addition, we refer to $q(\theta) = \int_{[\theta, 1]} \xi f(\xi) d\xi + \beta \int_{[0, \theta]} \xi f(\xi) d\xi$, the total usage generated from the population, as the *usage level* or simply *usage*. Note that the first term is usage due to owners, and the second term is usage due to renters (and hence modulated by β).

2.3.1 Matching Supply with Demand

In the presence of collaborative consumption, let $D(\theta)$ denote the aggregate demand (for rentals) generated by renters and $S(\theta)$ the aggregate supply generated by owners, for given θ . Then, $D(\theta) = \int_{[0, \theta]} \xi f(\xi) d\xi$ and $S(\theta) = \int_{[\theta, 1]} (1 - \xi) f(\xi) d\xi$. Moreover, the amount of demand from renters that is fulfilled must equal the amount of supply from owners that is matched with renters. In other words, the following fundamental relationship must be satisfied

$$\alpha S(\theta) = \beta D(\theta). \tag{2.4}$$

The parameters α and β , along with θ , are determined endogenously in equilibrium.

As mentioned earlier, matching friction can arise because of short term fluctuations in supply and demand (even though overall supply and demand are constant in the long run). This short term fluctuation may be due to the inherent variability in the timing of individual rental requests or in the duration of individual rental periods. Consequently, an available product may not find an immediate renter and a renter may not always be able to find an available product. In constructing a model for α and β , the following are desirable properties: (i) α (β) increases (decreases) in θ ; (ii) α approaches 1 (0) when θ approaches 1 (0); (iii) β approaches 1 (0) when θ approaches 0 (1), and (iv) α and β must satisfy the supply-demand relationship in Eq. 2.4.

Below we describe a plausible model for the short term dynamics of matching owners and renters. The model takes the view that in the short term (e.g., over the course of a day) demand is not realized all at once but requests for rentals arise continuously over time with random interarrival times. The intensity of the arrival process is of course determined by the total demand (e.g., total demand per day). The supply translates into individual products available for rent (for simplicity assume that supply is realized all at once and does not fluctuate over the time during which rental requests arrive). Once a product is rented, it becomes unavailable for the duration of the rental time, which may also be random. Because of the randomness in the interarrival times between requests and rental times per request, a request may arrive and find all products rented out. Assuming renters do not wait for a product to become available, such a request would then go unfulfilled. Also, because of this randomness, a product may not be rented all the time even if total demand exceeds total supply.

The dynamics described above are similar to those of a *multi-server loss queueing system*.⁵ In such a system, $1 - \beta$ would correspond to the *blocking probability* (the probability that a rental request finds all products rented out, or, in queueing parlance, the arrival of a request finds all servers busy) while α would correspond to the *utilization of the servers* (the probability that a product is being rented out).

If we let m denote the mean rental time per rental, the arrival rate (in terms of rental requests per unit time) is given by $\lambda(\theta) = D(\theta)/m$, and service capacity (the number of rental requests that can be fulfilled per unit time) by $\mu(\theta) = S(\theta)/m$.⁶ Therefore, we can express the workload (the ratio of the arrival rate to the service

⁵In a multi-server loss queueing system, customers arrive over time to receive service from a set of identical servers. A customer who does not find an available server upon arrival leaves the system without getting service. A customer who finds one or more available servers proceeds to receive service from one of these servers. Service takes a specified amount of time. Upon completion of service, the corresponding server becomes available. Both the interarrival and service times can be stochastic (see Cooper 1981, for additional details).

⁶For example, suppose the aggregate demand for renting per unit time is $D(\theta) = 1000$ h and the aggregate supply for renting per unit time is $S(\theta) = 2000$ h. If the average rental period is $m = 5$ h, then the arrival rate and the service capacity of the system are respectively $\lambda(\theta) = D(\theta)/m = 200$ and $\mu(\theta) = S(\theta)/m = 400$ requests per unit time.

capacity) of the system as $\rho(\theta) = \lambda(\theta)/\mu(\theta) = D(\theta)/S(\theta)$ and the utilization as $\alpha = \beta\lambda(\theta)/\mu(\theta) = \beta D(\theta)/S(\theta)$ (these relationships are of course consistent with the supply-demand relationship in Eq. 2.4).

Assuming that we can approximate the arrival process by a Poisson process (this is a reasonable assumption given that the arrival process arises from a continuum of renters who make independent decisions about when to seek a rental), the blocking probability, $1 - \beta$, can be approximated by $1 - 1/(1 + \rho)$ (see Benjaafar et al. 2018, for details). This leads to the following approximation of α and β

$$\alpha = \frac{\rho}{1 + \rho} = \frac{D(\theta)}{S(\theta) + D(\theta)}, \quad (2.5)$$

and

$$\beta = \frac{1}{1 + \rho} = \frac{S(\theta)}{S(\theta) + D(\theta)}. \quad (2.6)$$

Note that α and β , as specified in the above expressions, satisfy the properties (i)–(iv) described above. Interestingly, these expressions can also be obtained directly from the supply-demand relationship in Eq. 2.4 if we require that $\alpha + \beta = 1$ (Eqs. 2.5 and 2.6 are in that case the unique solution to Eq. 2.4).

The expressions in Eqs. 2.5 and 2.6 allow for both α and β to be strictly less than one and for the possibility of matching friction for both owners and renters. In the rest of the chapter, we rely on this approximation for our analysis. The model for α and β specified by these expressions is not unique in satisfying properties (i)–(iv). We expect other plausible models that satisfy these properties to lead to results that are qualitatively similar to those we describe in the next two sections.

We are now ready to proceed with the analysis of the equilibrium. An equilibrium under collaborative consumption exists if there exists $(\theta, \alpha) \in (0, 1)^2$ that is a solution to Eqs. 2.3 and 2.5. When it exists, we denote this solution by (θ^*, α^*) . Knowing the equilibrium allows us to answer important questions regarding product ownership, usage, and social welfare, among others.

2.4 Equilibrium Analysis

In this section, we consider the case where the price is exogenously specified. In Sect. 2.5, we treat the case where the price is chosen optimally by the platform. As mentioned in Sect. 2.3, the rental price must satisfy $d_o/(1 - \gamma) \leq p \leq 1 - d_r$, since otherwise, either the owners or renters will not participate. We denote the set of admissible prices by $A = [d_o/(1 - \gamma), 1 - d_r]$. For ease of exposition and to allow for closed form expressions, we assume that ξ is uniformly distributed in $[0, 1]$.

Letting θ denote the solution to $\pi_o(\xi) = \pi_r(\xi)$ leads to

$$\theta = \frac{c - ((1 - \gamma)p - d_o)\alpha}{p + d_r + (1 - p - d_r)\alpha - ((1 - \gamma)p - d_o)\alpha}. \quad (2.7)$$

Given θ , the aggregate demand under collaborative consumption is given by $D(\theta) = \frac{1}{2}\theta^2$ and the aggregate supply by $S(\theta) = \frac{1}{2}(1 - \theta)^2$. This leads to $\rho(\theta) = \theta^2 \cdot (1 - \theta)^{-2}$, and by Eq. 2.5

$$\alpha = \frac{\theta^2}{(1 - \theta)^2 + \theta^2}. \quad (2.8)$$

An equilibrium exists if Eqs. 2.7 and 2.8 admit a solution (θ^*, α^*) in $(0, 1)^2$.

In the following theorem, we establish the existence and uniqueness of such an equilibrium. Let $\Omega = \{(p, \gamma, c, d_o, d_r) | c \in (0, 1), \gamma \in [0, 1), (d_o, d_r) \in [0, 1]^2, p \in A\}$.

Theorem 1 *A unique equilibrium (θ^*, α^*) exists for each $(p, \gamma, c, d_o, d_r) \in \Omega$. Moreover, θ^* and α^* both (i) strictly increase with the cost of ownership c , commission rate γ and extra wear and tear cost d_o , and (ii) strictly decrease with rental price p and inconvenience cost d_r .*

The existence of the equilibrium is guaranteed by the Intermediate Value Theorem. The uniqueness is due to the monotonicity of Eqs. 2.7 and 2.8; see Benjaafar et al. (2018) for a proof of this and all subsequent results.

Let ω^* and q^* denote the corresponding ownership and total usage in equilibrium. Then, $\omega^* = 1 - \theta^*$ and $q^* = \frac{1}{2}(1 - \alpha^*\theta^{*2})$, where the expression for q^* follows from noting that $q^* = \int_{[\theta^*, 1]} \xi \, d\xi + \beta \int_{[0, \theta^*]} \xi \, d\xi$ (note that total usage is the sum of usage from the owners and the fraction of usage from the non-owners that is satisfied through renting).

The following proposition describes how ownership and usage in equilibrium vary with the problems parameters.

Proposition 2 *In equilibrium, ownership ω^* and usage q^* both strictly increase in price p and inconvenience cost d_r , and strictly decrease in cost of ownership c , commission rate γ and extra wear and tear cost d_o .*

While the monotonicity results in Proposition 2 are perhaps expected, it is not clear how ownership and usage under collaborative consumption compare to those under no collaborative consumption. In the following subsection, we provide comparisons between systems with and without collaborative consumption, and address the questions of whether or not collaborative consumption reduces product ownership and usage.

2.4.1 Impact of Collaborative Consumption on Ownership and Usage

In the absence of collaborative consumption, an individual would own a product if $u(\xi) \geq c$ and would not otherwise. Let $\hat{\theta}$ denote the solution to $u(\xi) = c$. Then,

the fraction of the population that corresponds to owners (ownership) is given by $\hat{\omega} = \int_{[\hat{\theta}, 1]} f(\xi) d\xi = 1 - c$, with an associated usage given by $\hat{q} = \int_{[\hat{\theta}, 1]} \xi f(\xi) d\xi = \frac{1}{2}(1 - c^2)$.

In the following proposition, we compare ownership level with and without collaborative consumption. Without loss of generality, we assume here (and in the rest of the chapter) that $d_o/(1 - \gamma) < 1 - d_r$ so that the set of admissible prices consists of more than a single price.

Proposition 3 *There exists $p_\omega \in (d_o/(1 - \gamma), 1 - d_r)$ such that $\omega^* = \hat{\omega}$ if $p = p_\omega$, $\omega^* < \hat{\omega}$ if $p < p_\omega$, and $\omega^* > \hat{\omega}$ otherwise. Moreover, $\partial p_\omega / \partial \gamma > 0$, $\partial p_\omega / \partial c < 0$, $\partial p_\omega / \partial d_o > 0$, and $\partial p_\omega / \partial d_r < 0$.*

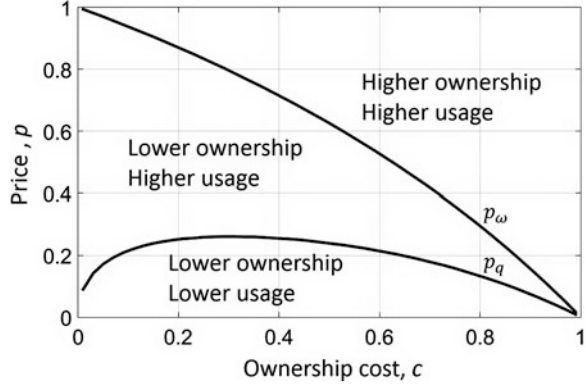
Proposition 3 shows that depending on the rental price p , collaborative consumption can result in either lower or higher ownership. In particular, when the rental price p is sufficiently high (above the threshold p_ω), collaborative consumption leads to higher ownership (e.g., more cars). Moreover, the threshold p_ω is decreasing in the cost of ownership c and renter's inconvenience d_r , and increasing in the commission rate γ and extra wear and tear cost d_o . The fact that p_ω is decreasing in c is perhaps surprising as it shows that collaborative consumption is more likely to lead to more ownership (and not less) when the cost of owning is high. This can be explained as follows. In the absence of collaborative consumption, when the cost of ownership is high, there are mostly non-owners. With the introduction of collaborative consumption, owning becomes more affordable as rental income subsidizes the high cost of ownership. In that case, even at low rental prices, there are individuals (those with high usage) who would switch to being owners. This switch is made more attractive by the high probability of finding a renter (given the high fraction of renters in the population). On the other hand, when the cost of ownership is low, only individuals with low usage are non-owners. For collaborative consumption to turn these non-owners into owners and lead to higher ownership, the rental price needs to be high. This is also needed to compensate for the low probability of finding a renter.

Similarly, usage can be either lower or higher with collaborative consumption than without it. In this case, there is again a price threshold p_q above which usage is higher with collaborative consumption, and below which usage is higher without collaborative consumption. When either d_o or d_r is sufficiently high, collaborative consumption always leads to higher usage. The result is formally stated in Proposition 4.

Proposition 4 *There exists $t \in (0, 1)$ such that (i) if $d_o/(1 - \gamma) + d_r < t$, then there exists $p_q \in (d_o/(1 - \gamma), 1 - d_r)$ such that $q^* = \hat{q}$ if $p = p_q$, $q^* < \hat{q}$ if $p < p_q$, and $q^* > \hat{q}$ if $p > p_q$; (ii) otherwise, $q^* \geq \hat{q}$ for all $p \in [d_o/(1 - \gamma), 1 - d_r]$.*

Unlike p_ω , the price threshold p_q is not monotonic in c (see Fig. 2.1). As c increases, p_q first increases then decreases. To understand the reason, note that collaborative consumption can lead to higher usage due to the new usage from non-owners. On the other hand, it can lead to lower usage if ownership decreases

Fig. 2.1 Ownership and usage for varying rental prices and ownership costs (higher/lower ownership/usage is relative to the case without collaborative consumption; $\gamma = 0.4$, $d_o = d_r = 0$)



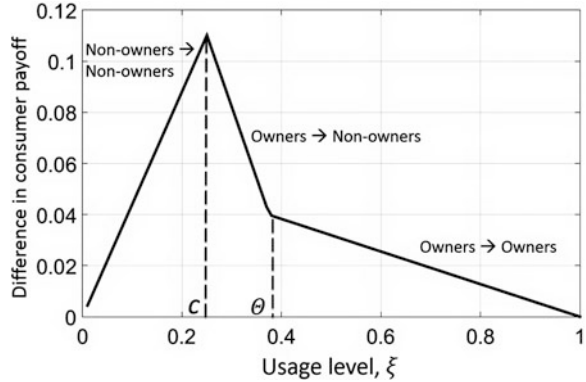
sufficiently (certainly to a level lower than that without collaborative consumption) such that the decrease in usage from those who switch from owning to renting is larger than the increase in usage from those who are non-owners. This implies that lower usage is less likely to happen if either (i) the individuals who switch from owning to renting can fulfill most of their usage via renting, or (ii) the usage from non-owners is high. The first scenario is true when the population of owners is high (i.e., the cost of ownership is low), whereas the second scenario is true when the population of non-owners is high (i.e., cost of ownership is high). Therefore, collaborative consumption is less likely to lead to lower usage when the cost of ownership is either very low or very high. Hence, the threshold p_q is first increasing then decreasing in c . When the cost of ownership is moderate, there is a balance of owners and non-owners without collaborative consumption, allowing for ownership to sufficiently decrease with relatively moderate rental prices, which in turn leads to lower usage and, correspondingly, a relatively higher threshold p_q .

The following corollary to Propositions 3 and 4 summarizes the joint impact of p and c on ownership and usage.

Corollary 5 *In settings where p_ω and p_q are well defined (per Propositions 3 and 4), collaborative consumption leads to higher ownership and higher usage when $p > p_\omega$, lower ownership but higher usage when $p_q < p \leq p_\omega$, and lower ownership and lower usage when $p \leq p_q$.*

Corollary 5, along with Propositions 3 and 4, show how price thresholds p_ω and p_q segment the full range of values of c and p into three regions, in which collaborative consumption leads to (i) lower ownership and lower usage, (ii) lower ownership but higher usage, and (iii) higher ownership and higher usage. These three regions are illustrated in Fig. 2.1. These results highlight the fact that the impact of collaborative consumption on ownership and usage is perhaps more nuanced than what is sometimes claimed by advocates of collaborative consumption. The results could have implications for public policy. For example, in regions where the cost of ownership is high, the results imply that, unless rental prices are kept sufficiently low or the commission extracted by the platform is made

Fig. 2.2 Impact of usage level on consumer payoff ($\theta^* > c$, $p = 0.4$, $\gamma = 0.4$, $c = 0.25$, $d_o = d_r = 0$)



sufficiently high, collaborative consumption would lead to more ownership and more usage. This could be an undesirable outcome if there are negative externalities associated with ownership and usage. Higher usage also implies less usage of the outside option (e.g., less use of public transport).

2.4.2 Impact of Collaborative Consumption on Consumers

Next, we examine the impact of collaborative consumption on consumer payoff. Consumer payoff is of course always higher with the introduction of collaborative consumption (consumers retain the option of either owning or not owning, but now enjoy the additional benefit of earning rental income if they decide to own, or of fulfilling some of their usage through renting if they decide not to own). What is less clear is who, among consumers with different usage levels, benefit more from collaborative consumption.

Proposition 6 *Let $\pi^*(\xi)$ and $\hat{\pi}(\xi)$ denote respectively the consumer payoff with and without collaborative consumption. Then, the difference in consumer payoff $\pi^*(\xi) - \hat{\pi}(\xi)$ is positive, piecewise linear, strictly increasing on $[0, c)$, and strictly decreasing on $[c, 1]$.*

An important implication from Proposition 6 (from the fact that the difference in consumer surplus $\pi^*(\xi) - \hat{\pi}(\xi)$ is strictly increasing on $[0, c)$ and strictly decreasing on $[c, 1]$) is that consumers who benefit the most from collaborative consumption are those who are indifferent between owning and not owning without collaborative consumption (recall that $[c, 1]$ corresponds to the population of owners in the absence of collaborative consumption). This can be explained by noting that there are always three segments of consumers. In the case where $\theta^* \geq c$ (see Fig. 2.2), which corresponds to the case where ownership decreases with collaborative consumption, the first segment corresponds to consumers who are non-owners in the absence of collaborative consumption and continue to be non-

owners with collaborative consumption (indicated by “non-owners→non-owners” in Fig. 2.2). The benefit these consumers derive from collaborative consumption is due to fulfilling part of their usage through accessing a rented product. This benefit is increasing in their usage.

The second segment corresponds to consumers who are owners in the absence of collaborative consumption and switch to being non-owners with collaborative consumption (indicated by “owners→non-owners”). These consumers have to give up the fulfillment of some usage (because a rental product may not always be available) and the amount they give up is increasing in their usage. Therefore, the amount of benefit they receive from renting decreases in their usage level. The third segment consists of consumers who are owners in the absence of collaborative consumption and continue to be owners with collaborative consumption (indicated by “owners→owners”). The benefit they experience is due to rental income. This income is decreasing in their usage (they have less capacity to rent when they have more usage). A similar explanation can be provided for the case where $\theta^* < c$.

2.5 The Platform’s Problem

In this section, we consider the problem faced by the platform. We first consider the case of a for-profit platform whose objective is to maximize the revenue from successful transactions. Then, we consider the case of a not-for-profit platform (e.g., a platform owned by a non-profit organization, government agency, or municipality) whose objective is to maximize social welfare.⁷ We compare the outcomes of these platforms in terms of ownership, usage and social welfare. We also benchmark the social welfare of these platforms against the maximum feasible social welfare.

A platform may decide, among others, on the price and commission rate. In this section, we focus on price as the primary decision made by the platform and treat other parameters as being exogenously specified (a survey of major peer-to-peer car sharing platforms worldwide reveals that commission rates fall mostly within a relatively narrow range, from 30% to 40% for those that include insurance, and do not typically vary across markets in which platforms operate). There are of course settings where the price is a decision made by the owners. Price may then be determined through a market clearing mechanism (i.e., the price under which supply equals demand; see Jiang and Tian 2018). In our case, because of the friction in matching supply and demand, the supply-demand balance equation in Eq. 2.4 can, per Theorem 1, be satisfied by any feasible price. Thus, the market clearing price is not unique and the system may settle on a price that maximizes neither social welfare nor platform revenue. Moreover, as we show in Sect. 2.5.1, platform

⁷An example of a not-for-profit platform is NeighborGoods, a peer-to-peer platform that facilitates the sharing of household goods. NeighborGoods allows owners to earn a rental fee but does not extract for itself a commission fee.

revenue (or social welfare) can be highly sensitive to price, giving the platform an incentive to optimize price. Platform pricing may also be beneficial to owners as it can serve as a coordinating tool and reduce competition among them. More significantly, and as we show in Sect. 2.5.2, the social welfare that results from a for-profit platform tends to be close to that resulting from a not-for-profit platform.

In what follows, we provide detailed analysis for the for-profit and not-for-profit platforms under the assumptions of Sect. 2.4. In Sects. 2.5.1, 2.5.2, and 2.5.3, we consider the case where $(d_o, d_r) = (0, 0)$. In Sect. 2.5.4, we discuss the case where $(d_o, d_r) \neq (0, 0)$.

2.5.1 The For-Profit Platform

For a for-profit platform, the objective is to maximize $\gamma p \alpha S(\theta)$, the commission income generated from the fraction of supply that is matched with demand. In particular, the platform's optimization problem can be stated as follows.

$$\max_p v_r(p) = \gamma p \alpha S(\theta) \quad (2.9)$$

$$\text{subject to } \pi_o(\theta) = \pi_r(\theta), \quad (2.10)$$

$$\alpha = \frac{D(\theta)}{D(\theta) + S(\theta)}, \quad (2.11)$$

$$\frac{d_o}{1 - \gamma} \leq p \leq 1 - d_r. \quad (2.12)$$

The constraints Eqs. 2.10 and 2.11 are the defining equations for the equilibrium (θ^*, α^*) . The constraint Eq. 2.12 ensures that price is in the feasible set A . In what follows, we assume that $\gamma > 0$ (the platform's revenue is otherwise always zero).

Under the assumptions of Sect. 2.4, the for-profit platform's problem can be restated as follows:

$$\max_p v_r(p) = \frac{1}{2} \gamma p \alpha (1 - \theta)^2 \quad (2.13)$$

subject to constraints Eqs. 2.7 and 2.8 and $p \in A$. It is difficult to analyze Eq. 2.13 directly. However, as the map between θ and p is bijective, we can use Eqs. 2.7 and 2.8 to express p in terms of θ as

$$p(\theta) = \frac{-\theta^3 + 2c\theta^2 - 2c\theta + c}{\theta(\theta - 1)(\gamma\theta - 1)}. \quad (2.14)$$

Hence, Eq. 2.13 can be expressed as

$$\max_{\theta} v_r(\theta) = \frac{\gamma}{2} \frac{(1 - \theta)\theta(\theta^3 - 2c\theta^2 + 2c\theta - c)}{((1 - \theta)^2 + \theta^2)(\gamma\theta - 1)} \quad \text{subject to } \theta \in [\underline{\theta}, \bar{\theta}] \quad (2.15)$$

where $\underline{\theta}$ is the solution to (2.7) and (2.8) at $p = 1$, $\bar{\theta}$ is the solution at $p = 0$, and $[\underline{\theta}, \bar{\theta}]$ is the set of solutions induced by $p \in [0, 1]$. We can use (2.14) to verify whether θ is in $[\underline{\theta}, \bar{\theta}]$. Specifically, $\theta < \underline{\theta}$ if $p(\theta) > 1$, $\theta \in [\underline{\theta}, \bar{\theta}]$ if $p(\theta) \in [0, 1]$, and $\theta > \bar{\theta}$ if $p(\theta) < 0$.

Proposition 7 *The function $v_r(\theta)$ is strictly quasiconcave in θ .*

Proposition 7 shows that the platform's problem is not difficult to solve. Depending on the value of γ and c , $v_r(\theta)$ is either decreasing or first increasing then decreasing on $[\underline{\theta}, \bar{\theta}]$. In both cases, the optimal solution to (2.15), which we denote by θ_r^* , is unique. We let p_r^* , ω_r^* , and q_r^* denote the corresponding price, ownership, and usage, respectively. We also use the notation v_r^* to denote the optimal revenue $v_r(\theta_r^*)$.

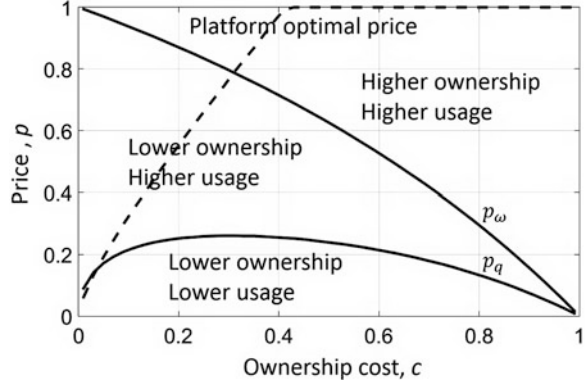
Proposition 8 *The platform's optimal revenue, v_r^* , is strictly quasiconcave in c , first strictly increasing and then strictly decreasing.*

Proposition 8 suggests that a platform would be most profitable when the cost of ownership is “moderate” and away from the extremes of being either very high or very low. In these extreme cases, not enough transactions take place because of either not enough renters (when the cost of ownership is low) or not enough owners (when the cost of ownership is high). This is perhaps consistent with the experience of iCarsclub, a peer-to-peer car sharing platform, that was first launched in Singapore, a country where the cost of ownership is exceptionally high and car ownership is low. iCarsclub struggled in Singapore and had to temporarily suspend operations. However, it is thriving in China where it operates under the name PPzuche and is present in several cities (Clifford Teo, CEO of iCarsclub, personal communication, 2015). This result also implies that a platform may have an incentive to affect the cost of ownership. For example, when the cost of ownership is low, a platform may find it beneficial to impose a fixed membership fee on owners, increasing the effective cost of ownership. On the other hand, when the cost of ownership is high, the platform may find it beneficial to lower the effective cost of ownership by offering, for example, subsidies (or assistance with financing) toward the purchase of new products.

Proposition 9 *There exists a threshold $c_{r,\omega} \in (0, 1)$ such that optimal ownership $\omega_r^* = \hat{\omega}$ if $c = c_{r,\omega}$, $\omega_r^* < \hat{\omega}$ if $c < c_{r,\omega}$, and $\omega_r^* > \hat{\omega}$ otherwise, with $c_{r,\omega}$ strictly increasing in γ .*

Proposition 9 shows that it continues to be possible, even when the price is chosen optimally by a revenue maximizing platform, for collaborative consumption to lead to either higher or lower ownership. In particular, collaborative consumption leads to higher ownership when the cost of ownership is sufficiently high (above the threshold $c_{r,\omega}$) and to lower ownership when the cost of ownership is sufficiently low (below the threshold $c_{r,\omega}$). This can be explained as follows. The platform has an incentive to somewhat balance supply and demand (otherwise few rentals will take place). When the cost of ownership is high, ownership is low in the absence of

Fig. 2.3 Impact of ownership cost on ownership and usage (higher/lower ownership/usage is relative to the case without collaborative consumption; $\gamma = 0.4$, $d_o = d_r = 0$)



collaborative consumption. In this case, the platform would try to induce, via higher prices, higher ownership, so as to generate more supply (hence, the result that a sufficiently high cost of ownership leads to higher ownership under collaborative consumption).⁸ Similarly, when the cost of ownership is low, the platform would try to induce lower ownership via lower prices, so as to generate more demand (hence, the result that a sufficiently low cost of ownership leads to low ownership under collaborative consumption).

We also observe that usage under platform pricing can be either higher or lower than that without collaborative consumption. Again, there exists a threshold $c_{r,q} < c_{r,\omega}$ in the cost of ownership, below which collaborative consumption leads to lower usage and above which collaborative consumption leads to higher usage. The impact of ownership cost on product ownership and usage under platform pricing is illustrated in Fig. 2.3 (the platform optimal price corresponds to the dashed curve).

2.5.2 The Not-for-Profit Platform

For a not-for-profit platform, the objective is to maximize social welfare (i.e., the sum of consumer surplus and platform revenue). Thus, the platform’s problem can be stated as

$$\max_p v_s(p) = \int_{[\theta, 1]} (\xi - c) f(\xi) d\xi + \int_{[0, \theta]} \beta \xi f(\xi) d\xi, \tag{2.16}$$

subject to constraints Eqs. 2.10, 2.11, and 2.12.

⁸This perhaps validates concerns expressed by the Singapore authorities that allowing peer-to-peer car sharing would increase car usage and road congestion and their initial decision to restrict peer-to-peer car rentals to evenings and weekends (Clifford Teo, CEO of iCarsclub, personal communication, 2015).

Under the assumptions of Sect. 2.4, the platform's problem can be restated as follows:

$$\max_p v_s(p) = \frac{1}{2}(1 - \alpha\theta^2) - (1 - \theta)c \quad (2.17)$$

subject to constraints Eqs. 2.7 and 2.8 and $p \in A$, or equivalently as

$$\max_{\theta} v_s(\theta) = \frac{1}{2} \left(1 - \frac{\theta^4}{(1 - \theta)^2 + \theta^2} \right) - (1 - \theta)c \quad \text{subject to } \theta \in [\underline{\theta}, \bar{\theta}]. \quad (2.18)$$

Analysis and results similar to those obtained for the for-profit platform can be obtained for the not-for-profit platform. In particular, we can show that the social welfare function, v_s , is strictly concave in θ , indicating that computing the optimal solution for the not-for-profit platform is also not difficult (we omit the details for the sake of brevity). The result also implies that Eq. 2.18 admits a unique optimal solution, which we denote by θ_s^* , with a resulting optimal social welfare which we denote by v_s^* .

The following proposition characterizes θ_s^* for varying values of γ .

Proposition 10 *There exists a strictly positive decreasing function $\gamma_s(c)$ such that $\theta_s^* \in (\underline{\theta}, \bar{\theta})$ if $\gamma < \gamma_s$, and $\theta_s^* = \underline{\theta}$ otherwise. Consequently, if $\gamma \leq \gamma_s(c)$, then*

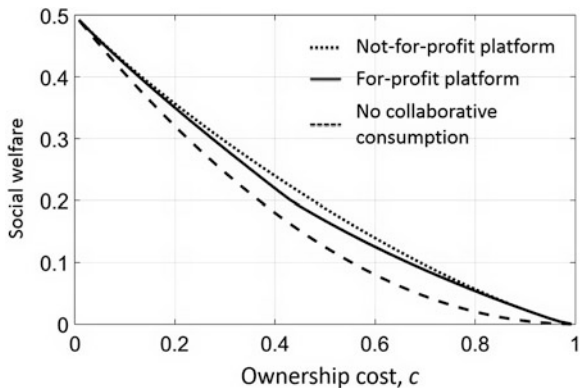
$$\max_{\theta \in [\underline{\theta}, \bar{\theta}]} v_s = \max_{\theta \in [0, 1]} v_s.$$

Proposition 10 shows that θ_s^* is an interior solution (satisfying $(\partial v_s / \partial \theta)(\theta_s^*) = 0$) if the commission rate is sufficiently low (below the threshold γ_s). Otherwise, it is the boundary solution $\underline{\theta}$. In particular, θ_s^* could never take the value of $\bar{\theta}$. An important implication of this result is that, when $\gamma < \gamma_s$, a not-for-profit platform that relies on price alone as a decision variable would be able to achieve the maximum feasible social welfare (i.e., the social welfare that would be realized by a social planner who can directly decide on the fraction of non-owners, θ). Note that this is especially true if the not-for-profit platform does not charge a commission fee (i.e., $\gamma = 0$).

Similar to the case of the for-profit platform, we can also show that a not-for-profit platform can lead to either higher or lower ownership or usage (relative to the case without collaborative consumption). Again, there are thresholds $c_{s,q} < c_{s,\omega}$ in the cost of ownership such that (i) ownership and usage are both lower if $c \leq c_{s,q}$, (ii) ownership is lower but usage is higher if $c_{s,q} < c \leq c_{s,\omega}$, and (iii) ownership and usage are both higher if $c > c_{s,\omega}$.

In the following proposition, we compare outcomes under the for-profit and not-for-profit platforms. In particular, we show that a not-for-profit platform would always charge a lower price than a for-profit platform. Therefore, it would also induce lower ownership and lower usage.

Fig. 2.4 Impact of ownership cost on social welfare ($\gamma = 0.2, d_o = d_r = 0$)



Proposition 11 Let p_s^* , ω_s^* and q_s^* denote the optimal price, ownership and usage levels under a not-for-profit platform, respectively. Then, $p_s^* \leq p_r^*$, $\omega_s^* \leq \omega_r^*$, and $q_s^* \leq q_r^*$.

A not-for-profit platform induces lower ownership by charging lower prices because it accounts for the negative impact of the cost of ownership on social welfare. In settings where there are negative externalities associated with ownership and usage, the result in Proposition 11 shows that, even without explicitly accounting for these costs, the not-for-profit platform would also lower such externalities (since both ownership and usage are lower). The fact that social welfare is maximized at prices lower than those that would be charged by a for-profit platform suggests that a regulator may be able to nudge a for-profit platform toward outcomes with higher social welfare by putting a cap on price.

Figure 2.4 illustrates the differences in social welfare between a system without collaborative consumption and systems with collaborative consumption under (a) a for-profit platform (a revenue-maximizing platform) and (b) a not-for-profit platform (a social welfare-maximizing platform). Systems with collaborative consumption can improve social welfare substantially, especially when the cost of ownership is neither too high nor too low (in those extreme cases, there are either mostly owners or mostly renters and, therefore, few transactions). However, the differences in social welfare between the for-profit and not-for-profit platforms are not very significant. This is because both platforms have a similar interest in maintaining a relative balance of renters and owners.

2.5.3 Systems with Negative Externalities

In this section, we consider settings where there are negative externalities associated with either usage or ownership. In that case, social welfare must account for the additional cost of these externalities. In particular, the following additional terms

must be subtracted from the expression of social welfare in Eq. 2.16

$$e_q q(\theta) + e_\omega \omega(\theta),$$

or equivalently

$$e_q \left(\int_{[\theta, 1]} \xi f(\xi) d\xi + \beta \int_{[0, \theta]} \xi f(\xi) d\xi \right) + e_\omega \int_{[\theta, 1]} f(\xi) d\xi,$$

where e_q and e_ω correspond to the social (or environmental) cost per unit of usage and per unit of ownership, respectively. This is consistent with the so-called *lifecycle* approach to assessing the social impact of using and owning products (see, for example Reap et al. 2008). The parameter e_q accounts for the social (or environmental) cost of using a product not captured by the utility (e.g., the cost of pollution associated with using a product), while e_ω would account for the social cost of product manufacturing, distribution, and end-of-life disposal.

For a not-for-profit platform, the optimization problem can then be restated as

$$\max_{\theta} v_e(\theta) = \frac{1}{2}(1 - e_q) \left(1 - \frac{\theta^4}{(1 - \theta)^2 + \theta^2} \right) - (c + e_\omega)(1 - \theta)$$

$$\text{subject to } \theta \in [\underline{\theta}, \bar{\theta}].$$

It is easy to show that the modified social welfare function v_e is still strictly quasiconcave in θ . Moreover, the optimal solution, which we denote by θ_e^* , is strictly increasing in both e_q and e_ω . As a result, the ownership and usage levels obtained under $e_q > 0$ and $e_\omega > 0$ are lower than those obtained under $e_q = e_\omega = 0$. Therefore, Proposition 11 continues to hold. However, Proposition 10 may no longer be valid if either e_q or e_ω is too large. That is, the platform may not be able to achieve the maximum feasible social welfare even if the commission rate is negligible. In this case, in order to achieve a higher social welfare, the platform may need to either subsidize non-owners (e.g., improve rental experience by reducing inconvenience cost) or penalize owners (e.g., make the ownership cost higher by charging extra tax on ownership), in addition to setting a low rental price.

We conclude this section by addressing the question of whether collaborative consumption reduces the total cost of negative externalities, $e_q q(\theta) + e_\omega \omega(\theta)$. Recall that collaborative consumption (under either a for-profit or a not-for-profit platform) leads to lower ownership and lower usage when the cost of ownership is sufficiently low, and it leads to higher ownership and higher usage when the cost of ownership is sufficiently high (see Fig. 2.3). This implies that collaborative consumption could either decrease or increase negative externalities, with a decrease more likely when the cost of ownership is low. In numerical experiments, we observe that there exists a threshold on the cost of ownership, which we denote by c_e , such that collaborative consumption reduces negative externalities if and only if $c < c_e$. We also observe that c_e is decreasing in e_q and increasing in e_ω , indicating

that collaborative consumption is more likely to reduce negative externalities if the social (or environmental) cost of using products is relatively low compared to that of owning.

2.5.4 *The Impact of Extra Wear and Tear and Inconvenience Costs*

In this section, we consider the case where $(d_o, d_r) \neq 0$. The extra wear and tear cost d_o reduces the payoff of owners and, therefore, places a lower bound on the set of admissible prices: $p \geq d_o/(1 - \gamma)$. Similarly, the inconvenience cost d_r reduces the payoff of renters and, consequently, places an upper bound on the price: $p \leq 1 - d_r$. Obtaining analytical results is difficult. However, we are able to confirm numerically that all the results obtained for $(d_o, d_r) = 0$ continue to hold (details are omitted for brevity).

Of additional interest is the impact of d_o and d_r on platform revenue and social welfare. For both the for-profit and not-for-profit platforms, we observe that social welfare is decreasing in both d_o and d_r . This is consistent with intuition. However, revenue for the for-profit platform can be non-monotonic in d_o . In particular, when the cost of ownership is low, platform revenue can first increase then decrease with d_o . This effect appears related to the fact that platform revenue is, per Proposition 8, non-monotonic in the cost of ownership. A higher value of d_o can be beneficial to the platform if it helps balance the amount of owners and renters (i.e., reduce ownership), leading to a greater amount of transactions. An implication of this result is that a for-profit platform may not always have an incentive to reduce the extra wear and tear cost.⁹ On the other hand, the inconvenience cost d_r does not have the same effect on platform revenue. An increase in d_r could lead to more transactions. However, it limits the price a platform could charge. The net effect is that the platform revenue is always decreasing in d_r .

2.6 Concluding Comments

In this chapter, we described an equilibrium model of collaborative consumption. We characterized equilibrium outcomes, including ownership and usage levels, consumer surplus, and social welfare. We compared each outcome in systems

⁹Note that, in some cases, a platform could exert costly effort to reduce this cost. For example, when extra wear and tear is, in part, due to renters' negligence, more effort could be invested in the vetting of would-be renters. Alternatively, the platform could provide more comprehensive insurance coverage or monitor more closely the usage behavior of a renter (such monitoring technology is already available for example in the case of cars).

Table 2.1 Notation

Symbol	Meaning
$\xi \in [0, 1]$	Individual usage level
$f(\xi)$	Density function for individual usage distribution
$p \in [0, 1]$	Rental price
$\gamma \in [0, 1)$	Commission fee
$c \in (0, 1)$	Cost of ownership
$d_o \in [0, 1]$	Moral hazard cost for owners
$d_r \in [0, 1]$	Inconvenience cost for renters
$\alpha \in [0, 1]$	Matching probability for owners
$\beta \in [0, 1]$	Matching probability for renters
$\pi_o(\xi)$	Payoff for an owner with usage level ξ
$\pi_r(\xi)$	Payoff for a renter with usage level ξ
θ	Individual usage level at which point $\pi_o(\xi) = \pi_r(\xi)$
$S(\theta)$	Aggregate rental supply generated from owners for a given θ
$D(\theta)$	Aggregate rental demand generated from renters for a given θ
θ^*	θ in equilibrium
α^*	α in equilibrium
ω^*	Ownership level in equilibrium in the presence of collaborative consumption
q^*	Total usage level in equilibrium in the presence of collaborative consumption
$\hat{\omega}$	Ownership level in the absence of collaborative consumption
\hat{q}	Total usage level in the absence of collaborative consumption
p_ω	Price threshold at which point $\omega^* = \hat{\omega}$
p_q	Price threshold at which point $q^* = \hat{q}$
v_r	Platform revenue and objective function for the for-profit platform
θ_r^*	Optimal solution to the for-profit platform's problem
v_r^*	Optimal platform revenue $v_r(\theta_r^*)$ and value function for the for-profit platform
p_r^*	Optimal price under the for-profit platform
ω_r^*	Optimal ownership under the for-profit platform
u_r^*	Optimal usage under the for-profit platform
v_s	Social welfare and objective function for the not-for-profit platform
θ_s^*	Optimal solution to the not-for-profit platform's problem
v_s^*	Optimal social welfare $v_s(\theta_s^*)$ and value function for the not-for-profit platform
ω_s^*	Optimal ownership under the not-for-profit platform
u_s^*	Optimal usage under the not-for-profit platform
v_e	Social welfare accounting for social costs on ownership and usage
e_ω	Social cost per unit of ownership
e_q	Social cost per unit of usage

with and without collaborative consumption and examined the impact of various problem parameters including rental price, platform's commission rate, cost of ownership, owner's extra wear and tear cost, and renter's inconvenience cost. Our findings indicate that collaborative consumption can result in either higher or lower ownership and usage levels, with higher ownership and usage levels more likely

when the cost of ownership is high. We showed that consumers always benefit from collaborative consumption, with individuals who, in the absence of collaborative consumption, are indifferent between owning and not owning benefitting the most. We also showed that the platform's profit is not monotonic in the cost of ownership (first increasing and then decreasing), implying that a platform is least profitable when the cost of ownership is either very high or very low (also suggesting that a platform may have an incentive to affect the cost of ownership by, for example, imposing membership fees or providing subsidies). In addition, we observed that, when the cost of ownership is low, platform profit can be increasing in the extra wear and tear cost, suggesting that a for-profit platform may not always have an incentive to eliminate this cost.

In Benjaafar et al. (2018), we consider several extensions to the model described in this chapter, including systems where (1) individuals are heterogeneous in their sensitivity to moral hazard and inconvenience, (2) usage is endogenous, (3) non-owners have the option of renting from a third party service provider, (4) usage has a general distribution, and (5) platforms may own assets of their own.

Possible avenues for future research are many. We mention a few examples. It would be of interest to consider settings where there is competition among multiple platforms, with owners and renters having the option of participating in one or more such platforms. Given that the effective demand the platform would face is non-monotonic in price, competition may not necessarily lead (as in a standard competitive setting) to lower prices. In this case, the competing platforms must account for the need to balance, via their choice of prices, the supply of owners and renters. It would also be interesting to investigate other forms of peer-to-peer product sharing, such as those that involve concurrent use of the product by multiple renters (as in car pooling) or by the owner and the renter (as in some forms of home sharing). In such cases, more usage by the owner may not necessarily imply less usage by non-owners.

Finally, in Table 2.1, we define all notations used throughout the chapter.

References

- Agrawal VV, Bellos I (2017) The potential of servicizing as a green business model. *Manag Sci* 63(5):1545–1562
- Ballus-Armet I, Shaheen S, Clonts K, Weinzimmer D (2014) Peer-to-peer carsharing exploring public perception and market characteristics in the San Francisco Bay area, California. *Transp Res Rec: J Transp Res Board* 2416:27–36
- Bellos I, Ferguson M, Toktay LB (2017) The car sharing economy: interaction of business model choice and product line design. *Manuf Serv Oper Manag* 19(2):185–201
- Benjaafar S, Kong G, Li X, Courcoubetis C (2018) Peer-to-peer product sharing: implications for ownership, usage and social welfare in the sharing economy. *Manag Sci.* 1–17. <http://pubsonline.informs.org/journal/mnsc/>
- Botsman R, Rogers R (2010) *What's mine is yours: the rise of collaborative consumption*. Harper Business, New York
- Chase R (2015) *Peers Inc: how people and platforms are inventing the collaborative economy and reinventing capitalism*. Public Affairs, New York

- Clark M, Gifford K, Le Vine S (2014) The usage and impacts of emerging carsharing business models: evidence from the peer-to-peer and business-to-business market segments. In: Transportation research board 93rd annual meeting 2014-1-12 to 2014-1-16, Washington DC, Paper No 14-1714
- Cooper RB (1981) Introduction to queueing theory, 2nd edn. North Holland, New York
- Cullen Z, Farronato C (2018) Outsourcing tasks online: matching supply and demand on peer-to-peer Internet platforms. Working paper, February 2018. Available at <https://www.hbs.edu/faculty/Pages/item.aspx?num=50051>. Accessed 14 Aug 2018
- Cusumano MA (2014) How traditional firms must compete in the sharing economy. *Commun ACM* 58(1):32–34
- Einav L, Farronato C, Levin J (2016) Peer-to-peer markets. *Annu Rev Econ* 8:615–635
- Fradkin A, Grewal E, Holtz D, Pearson M (2015) Bias and reciprocity in online reviews: evidence from field experiments on Airbnb. In: EC'15 proceedings of the sixteenth ACM conference on economics and computation. ACM, New York, p 641
- Fraiberger SP, Sundararajan A (2015) Peer-to-peer rental markets in the sharing economy. Available at SSRN: <http://ssrn.com/abstract=2574337>. Accessed 14 Aug 2018
- Hagiu A, Wright J (2015) Marketplace or reseller? *Manag Sci* 61(1):184–203
- Jiang B, Tian L (2018) Collaborative consumption: strategic and economic implications of product sharing. *Manag Sci* 64(3):1171–1188
- Korosec K (2015) The number of miles cars travel is about to explode. *Fortune*, 17 Nov 2015. Available at <http://fortune.com/2015/11/17/la-auto-show-vehicle-miles/>. Accessed 14 Aug 2018
- Malhotra A, Van Alstyne M (2014) The dark side of the sharing economy and how to lighten it. *Commun ACM* 57(11):24–27
- Martin E, Shaheen S (2011) Greenhouse gas emission impacts of carsharing in North America. *IEEE Trans Intell Transp Syst* 12(4):1074–1086
- New York Times (2015) Limiting Uber won't end congestion. Page A18, 18 July 2015. Available at <http://www.nytimes.com/2015/07/18/opinion/limiting-uber-wont-end-congestion.html>. Accessed 14 Aug 2018
- Nijland H, Van Meerkerk J, Hoen A (2015) Impact of car sharing on mobility and CO₂ emissions. Technical report. PBL Netherlands Environmental Assessment Agency
- Orsdemir A, Deshpande V, Parlakturk AK (2017) Is servicization a win-win strategy? Profitability and environmental implications of servicization. Available at SSRN: <http://ssrn.com/abstract=2679404>
- Reap J, Roman F, Duncan S, Bras B (2008) A survey of unresolved problems in life cycle assessment. *Int J Life Cycle Assess* 13(5):374–388
- Rochet J-C, Tirole J (2006) Two-sided markets: a progress report. *RAND J Econ* 37(3):645–667
- van der Linden DF, Franciscus D (2016) Explaining the differential growth of peer-to-peer car-sharing in European cities. Master's thesis. Utrecht University, Netherlands
- Waldman M (2003) Durable goods theory for real world markets. *J Econ Perspect* 17(1):131–154
- Weyl EG (2010) A price theory of multi-sided platforms. *Am Econ Rev* 100(4):1642–1672
- Zervas G, Proserpio D, Byers J (2015) A first look at online reputation on Airbnb, where every stay is above average. Available at SSRN: <http://ssrn.com/abstract=2554500>

Chapter 3

The Strategic and Economic Implications of Consumer-to-Consumer Product Sharing



Baojun Jiang and Lin Tian

Abstract In recent years, mobile communications technologies and online sharing platforms have made collaborative consumption among consumers a major trend in the economy. Consumers buy many products but end up not fully utilizing them. A product owners self-use values can differ over time, and in a period of low self-use value, the owner may rent out her product in a product-sharing market. This paper develops an analytical framework to examine the strategic and economic impact of product sharing among consumers. Our analysis shows that transaction costs in the sharing market have a non-monotonic effect on the firm's profits, consumer surplus, and social welfare. We find that when the firm strategically chooses its retail price, consumers sharing of products with high marginal costs is win-win for the firm and the consumers whereas their sharing of products with low marginal costs can be lose-lose. Further, in the presence of the sharing market, the firm will find it optimal to strategically increase its quality, leading to higher profits but lower consumer surplus. In addition, within a distribution channel framework, the existence of the sharing market is more likely to increase the downstream retailers profit than the upstream manufacturers profit, i.e., product sharing can sometimes benefit the downstream retailer at the expense of the upstream manufacturer.

3.1 Introduction

Many products that consumers buy or own are not fully utilized. In the sharing economy, these underutilized products can be put to use by other consumers through many sharing platforms. The recent economic recession and social concerns about consumption sustainability lead consumers and society as a whole to explore more

B. Jiang (✉)

Olin Business School, Washington University in St. Louis, St. Louis, MO, USA
e-mail: baojunjiang@wustl.edu

L. Tian

School of Management, Fudan University, Shanghai, China
e-mail: tianlin@fudan.edu.cn

© Springer Nature Switzerland AG 2019

M. Hu (ed.), *Sharing Economy*, Springer Series in Supply Chain Management 6,
https://doi.org/10.1007/978-3-030-01863-4_3

efficient use of resources and products. As a result, collaborative consumption has emerged as a global trend, enabled by technological advances in online, mobile communications. Many online platforms have helped to facilitate consumer-to-consumer sharing for a wide range of products and services such as bicycles (Spinlister), boats (Boatbound, GetMyBoat), cars (Turo, Getaround), working or parking spaces (Citizen Space, JustPark), car rides (Lyft, Uber, Zimride), short-term rental (Airbnb, Roomorama), gardens (Shared Earth, Landshare), clothing, portable tools/appliances, electronics, and household items (FriendsWithThings, Neighborgoods). In farmer communities in developing countries, the sharing of agricultural equipment is also common.

We can conceptually classify peer-to-peer sharing into two broad categories, depending on the types of resources being shared. First, if the shared underutilized resource is the consumers' own time and provision of services, we have peer-to-peer offering of services as on TaskRabbit, Uber, or Didi Chuxing. Second, if the shared underutilized resource is tangible assets or products that consumers own, we have consumer-to-consumer product sharing as on Turo or Neighborgoods. When consumers share their purchased products (e.g., sharing of cars on Turo), the manufacturers of the products are affected and may strategically change their decisions. By contrast, when consumers share their time or ability (e.g., a consumer on TaskRabbit assembles IKEA furniture for another consumer), typically there is no strategic manufacturer of time involved (i.e., each consumer's time resource is endowed from birth, not manufactured by some economic agent).

In this research we focus on the consumer-to-consumer sharing of products (e.g., on Turo or Neighborgoods), *not* the peer-to-peer offering of services (e.g., on Airbnb or Uber).¹ Many product-sharing transactions involve the renters paying a fee to the product owners through a sharing platform. From the consumer's perspective, sharing under-utilized products seems profitable and also environmentally responsible. How does product sharing affect the manufacturer? Though managers are wary of such sharing, anecdotal evidence shows that some firms are proactively responding to the emerging trend of collaborative consumption. For example, General Motors (GM) has worked with RelayRides (now Turo) to make it easier for owners to rent out their under-used OnStar-enabled GM vehicles to offset the cost of ownership, by introducing features such as remote unlocking of doors by authorized renters using their smartphones (General Motors 2012).

Our model captures the idea that a consumer's usage value for a product may vary over time. In each usage period, a product owner can decide whether to use the product herself or to rent it out to others through a third-party product-sharing platform, and consumers who did not purchase the product can decide whether to rent the product from the product-sharing platform. For each sharing transaction, the

¹Though we use consumer-to-consumer sharing as the context, our model applies equally to business-to-business sharing of products, e.g., the sharing of equipment among businesses or hospitals. The essence is that a firm/manufacturer's customers may rent out the product to its other potential customers during periods of low self-use value.

renting consumer pays a rental fee to the platform, which keeps a percentage of that rental fee as service fee and gives the rest to the product owner. The product owner will endure two types of transaction costs when renting out her product—one that is proportional to the sharing price (e.g., the platform’s fee) and one that is independent of the sharing price (e.g., costs of delivering and picking up the product). We develop an analytical framework with these key market features to study how a firm—the brand owner or manufacturer of the product—should strategically choose its retail price and product quality to respond to anticipated sharing by consumers.² We examine the impacts of product sharing on the firm’s pricing strategy, profits, consumer surplus, and social welfare.

We highlight a few major findings from our analysis. First, transaction costs in the sharing market have a non-monotonic effect on the firm’s profits, consumer surplus, and social welfare. Second, when the firm strategically chooses its retail price, consumers’ sharing of products with high marginal costs is win-win for the firm and the consumers whereas their sharing of products with low marginal costs can be lose-lose. Third, in the presence of the sharing market, the firm will find it optimal to strategically increase its quality, leading to higher profits but lower consumer surplus. Fourth, within a distribution channel framework, the existence of the sharing market is more likely to increase the downstream retailer’s profit than the upstream manufacturer’s profit. Further, we have analyzed the robustness of our results and insights to several alternative modeling assumptions, and provided some potential directions for future research.

Our research contributes to the emerging research literature on the consumer-to-consumer sharing or collaborative consumption. The fast growing trend in the sharing economy has recently received much attention in both practice and academia. Most extant work on consumer-to-consumer sharing has been conceptually depicting the phenomenon (e.g., Belk 2010, 2014). Both empirical and analytical researches are lacking in the published literature and are of great managerial and academic interest. Fraiberger and Sundararajan (2015) use the US automobile industry data and peer-to-peer car rental data from Getaround to study the welfare and distributional effects of a peer-to-peer rental market. Zervas et al. (2017) empirically estimate the impact of Airbnb’s entry on the incumbent hotels’ pricing and revenues, using a dataset collected on Airbnb listings in Texas and a decade-long panel of quarterly tax revenue for Texas hotels.

Weber (2014) shows that if the lender and the renter are risk-neutral, a sharing intermediary can eliminate the moral-hazard problem by providing optimal insurance to the lender and first-best incentives to the renter to exert care. Benjaafar et al. (2018) analyze how the platform should optimally set the rental price on the sharing platform; they find that, depending on the rental price, product ownership and usage levels can be higher or lower than the case without the sharing market. Bai et al. (2016) use a queueing model to study the coordination of supply and demand on an on-demand, peer-to-peer service platform. They show that when the

²For expositional convenience, we refer to the firm/manufacturer as “it” and a consumer as “she.”

potential customer demand increases, it is optimal for the platform to raise its price, wages, and the payout ratio. Weber (2016) uses an overlapping-generations model to show that a sharing market tends to increase the product price, and benefit the firm and the consumer for high-cost products. Bellos et al. (2017) examine the strategic interactions of business model choice and product line design in the car sharing economy. He et al. (2017) study how to design a geographical service region for car-sharing service providers to operate the service. They focus on the trade-offs between maximizing customer catchment by covering travel needs and controlling fleet operation costs.

Our research complements the stream of literature on secondary used-goods market (e.g., Hendel and Lizzeri 1999; Johnson 2011; Chen et al. 2013), where consumers sell their products rather than share their products. Note that a used-goods sale transaction involves a *permanent* transfer of product ownership from the seller to the buyer, whereas a product-sharing transaction involves merely a *temporary* transfer of use right from the product owner to the renter only for the particular sharing period and the product owner owns the continuation usage value of the product for the future periods. This distinction is shown to have a qualitatively different effect on the market outcome. Our research also complements the rent-or-buy and leasing literature (e.g., Desai and Purohit 1998; Desai 1999; Agrawal et al. 2012), which studies a firm's optimal selling or leasing strategy. By contrast, we examine the consumer's rent-or-buy decision in a marketplace where a firm's end *customers*, rather than the firm itself, can rent out their purchased products in a consumer-to-consumer sharing market. In essence, the firm's end customers can be its indirect competitors, because some potential buyers of the product may switch from buying to renting from the sharing market.

3.2 Modeling Framework

A monopolist firm produces a product of quality q at a constant marginal cost of c . The monopolist sells the product at price p to consumers, each of whom buys at most one unit and can derive usage value from the product in n time periods. Note that the consumer's product sharing is a short-term event. For example, car owners typically rent out their cars on Turo on a daily basis, but car manufacturers do not dynamically change their prices on a daily basis even when they respond to the existence of the product sharing market by pricing strategically or even changing their product quality (e.g., General Motors added new features to their cars to facilitate consumers' easier and more reliable car-sharing on RelayRides). So, to reduce analytical complexity, we focus on the fairly reasonable case where the firm will strategically choose its price but will not dynamically adjust that price from one sharing period to another.

At the end of the n usage periods, the product has some salvage value ε (e.g., it can be sold in a secondary used-goods market). Each consumer's per-period usage value from the product may vary over time. Consumer i 's usage values v_{ij} for

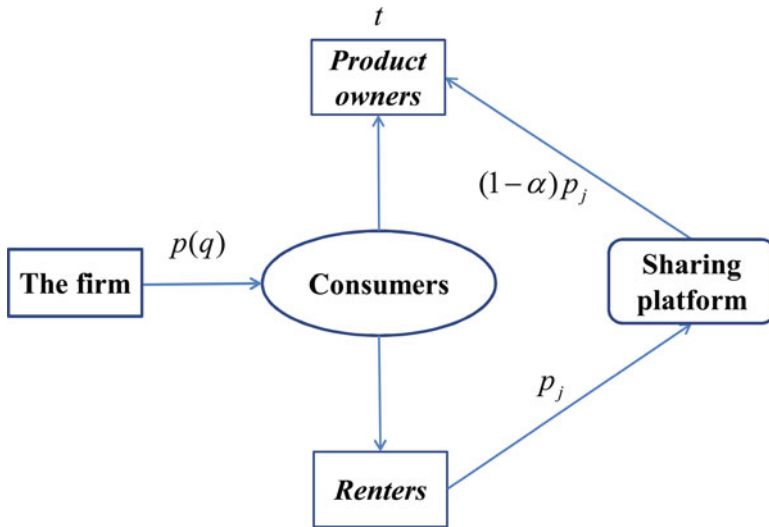


Fig. 3.1 Model structure

$j = 1, 2, \dots, n$ depend on the product’s quality (q) and her willingness to pay for quality (θ_{ij}); we assume $v_{ij} = q\theta_{ij}$, where θ_{ij} is uniformly distributed in the consumer population: $\theta_{ij} \sim U[0, 1]$. Without loss of generality, we normalize the total number of consumers to one. This type of model for quality and consumer heterogeneity has been widely adopted in the economics and marketing literature since Mussa and Rosen (1978). For expositional clarity/succinctness, we assume $n = 2$ for the main model. That is, the consumer can derive usage value from the product in two usage periods $j = 1, 2$; two periods are enough to capture the key market characteristic that a consumer’s usage value can vary across time and that during a period of low usage value she can forgo her self-use and earn some income by renting out her product through a sharing platform. If the consumer rents out her product, she will earn a rental fee for that period but needs to pay the platform a percentage fee, denoted by $\alpha \in [0, 1)$ fraction of the rental fee. Typically in practice, the sharing platform collects the rental fee from the renter, keeps a fixed α fraction of that fee as service charge, and will give the remaining fraction $(1 - \alpha)$ to the product owner. We model two types of transaction costs endured by the product owner when renting out her product—one that is proportional to the sharing price (e.g., the platform’s fee) and one that is independent of the sharing price (e.g., costs of delivering and picking up the product, or the accelerated product maintenance cost due to moral hazard of the renter). From here on, unless stated otherwise, we use the term “transaction cost” to refer to the latter and the platform fee to refer to the former. Let $t \geq 0$ denote the transaction cost for each sharing transaction. The market structure is illustrated in Fig. 3.1.

Consumer's Strategic Options. Consumers are forward-looking; at the time of their product-purchase decision, they rationally anticipate the possibility of sharing or renting the product in the sharing market. Each consumer i has eight (not-clearly-dominated) options: (i) buy the product and use it in both periods; (ii) buy the product, use it in period 1 and rent it out in the sharing market in period; (iii) buy the product, rent it out in the sharing market in period 1 and use it in period 2; (iv) do not buy the product but rent it in both periods; (v) do not buy the product but rent it only in period 1; (vi) do not buy the product but rent it only in period 2; (vii) buy the product (as a speculator) and rent it out in both periods; (viii) neither buy nor rent the product (i.e., the outside option).

Market Clearing Mechanism. With a sharing market, consumers may choose any of the above eight options. In each product-usage period, some consumers may rent out their purchased products while others may rent a product from the product-sharing market. In equilibrium, the supply and the demand for product sharing will be equal.³ In each period j , there will be a market-clearing price (p_j) that works to match the supply and demand; a consumer needs to pay p_j to rent the product from the market and a consumer who rents out her product will receive a fee of $(1 - \alpha)p_j$ while the platform keeps αp_j as its service fee.

Timing of Events. The timing of events in the core model is as follows. First, the firm chooses its retail price p . Second, consumers decide whether to buy the product. Third, in each usage period, consumers who bought the product before decide whether to use it themselves or to rent it out in the sharing market while consumers who did not buy the product decide whether to rent it from the sharing market, which clears at some endogenously determined equilibrium price p_j , at which there is no excess demand or supply for sharing. Note that after each sharing transaction, the product is returned from the renter to the original product owner, who will obtain the salvage value (ε) at the end of the last usage period. Note also that the platform's percentage fee (α) is taken as given; this is because in practice the sharing platform's percentage fee is the same across different products. If we endogenize α in our model, it would imply that the platform optimizes its percentage fee on an *individual* product basis, charging a different percentage for different products, which clearly does not reflect reality. In reality, the sharing platforms charge a fixed percentage fee across different products, and that percentage is typically between 10% (e.g., on Spinlister) to 25% (e.g., on Turo).

³The firm is assumed to play no direct role in the product-sharing market. In reality, in many markets, the firms (manufacturers) themselves do not offer hour-to-hour or day-to-day rentals of their products. This may be because the firm's transaction cost for managing renting of its products is much higher than that for consumers. For example, a consumer with an Xbox console can rent it to others in her local area on a daily or weekly basis much more efficiently than Microsoft, the producer of the Xbox, since the company would have many logistical issues (e.g., due to the lack of physical presence in the consumer's local area or city). Indeed, in reality, on these product-sharing platforms (or the firms' stores), we typically do not see the firms themselves offering to rent their products on a day-to-day basis; for example, we do not see General Motors offering hourly rental of their cars on car-sharing websites or at its own dealerships.

3.3 Effects of Sharing on Firm's Pricing Strategy, Profit, and Consumer Surplus

In this section, we assume that the firm has developed the product, which has a quality level of q with a marginal cost of production c . We examine the impact of consumer-to-consumer product sharing on the firm's price, profit, the consumer's surplus, and social welfare. In specific, we study how two key factors—transaction costs (t and the platform's fee α) and the firm's marginal cost (c)—affect the economic impact of collaborative consumption.

One may intuit that when the transaction cost for product sharing (i.e., t) increases, the firm should raise the price of its product since its customers will be less likely to offer the competing rental option to other consumers, making these consumers more likely to buy the product. However, we find that a higher transaction cost in the sharing market will actually lead to a price drop by the firm. This is because some of the product buyers, who have a high usage value in one period but a low usage value in the other period, will not be able to earn as much rental income from the sharing market and hence will no longer be willing to buy the product at the same price. To compensate and attract some of these buyers, the firm will find it optimal to reduce its price. In practice, the sharing platform often tries to reduce the transaction cost to the product owners by offering free insurance coverage or by enabling the product owners to rate the renters after sharing transactions, which will to some extent alleviate the moral hazard problem and reduce the transaction cost. Our result implies that increased incentives for the consumer's product sharing, i.e., lower friction or transaction costs in the sharing market, can actually induce the manufacturer to *raise* rather than lower its price.

Also we find, because of the firm's strategic pricing, a change in the transaction cost (t) can lead to a non-monotonic effect on the firm's profit, the social welfare and the total consumer surplus. So the sharing platform's efforts to reduce transaction costs may not always benefit the consumers or the manufacturer. As the transaction cost decreases, product owners are more likely to share their products to earn higher rental income (net of the transaction cost and the platform fee), but the firm's strategic increase of its price will not only reduce the customer's net sharing benefit but also force some customers who only use the product themselves to drop out of the market. Hence, the total consumer surplus in the market can drop.

Note that t is a product-sharing transaction cost that is independent of the sharing price whereas the sharing platform's percentage fee represents a transaction cost that is directly proportional to the sharing price. By similar analysis and intuition, we easily obtain the corollary that a decrease of the sharing platform's percentage fee (i.e., α) may not always benefit the consumers or the firm. In summary, frictions in the sharing market (e.g., the product-sharing transaction cost t and the platform percentage fee α) have a non-monotonic effect on the firm's profit and the total consumer surplus.

Some questions naturally arise. When is the firm likely to *benefit* from the consumers' product sharing? Under what situations will the sharing market *increase* consumer surplus and social welfare? Does product sharing necessarily affect the

consumer and the firm in opposite ways? Note that if the sharing platform's percentage fee is excessively high (i.e., $\alpha > 1 - t/q$), there will be no sharing transactions in equilibrium. For our main analysis, we assume that the platform's percentage fee is not too high; more specifically, $\alpha < \alpha^*$ for some $\alpha^* < 1 - t/q$. Our analysis shows that the firm's marginal cost and the transaction cost for sharing play an important role in determining the effects of the product-sharing market.

The conventional wisdom is that since the product-sharing market provides consumers with a sharing alternative to buying the product, it will lead to lower profitability for the firm, higher consumer surplus and social welfare. Our analysis shows that when the firm's unit cost and the transaction cost are low, product sharing among consumers will indeed hurt the firm. In the absence of the product-sharing market, the firm's optimal strategy for a low-cost product is to set a low price such that a large number of consumers will buy the product in equilibrium. However, when the product-sharing market exists, fewer consumers than before will buy the product even if the firm still sets the same low price, because some consumers can get the product from the sharing market, especially when the transaction cost is low (such that some product owners will share). Anticipating product sharing among consumers, the firm will find it optimal to raise its retail price, in an attempt to capture at least some of the value enhanced by the sharing market. However, the higher price is not enough to offset the loss in unit sales, resulting in a lower profit for the firm than when there is no sharing market.

Further, our analysis shows that when the firm's unit cost is low and the transaction cost is low, the existence of a product-sharing market not only reduces the firm's profit, but also lowers the total consumer surplus. The intuition hinges on the fact that the firm will significantly increase its price (from its low price in the absence of the sharing market), leading to fewer units of the product being sold and fewer consumers will use the product even with the sharing market. Contrary to the conventional wisdom, this finding suggests that product sharing among consumers may make them worse off if the firm anticipates such behaviors and strategically increases its price. This lose-lose scenario for product sharing happens for products with low marginal costs of production. Can the product-sharing market *benefit* the firm or the consumers?

Our analysis also shows that when the firm's unit cost of production is high, product sharing among consumers is win-win for the firm and the consumers. When the product-sharing market does not exist, the firm's optimal strategy for its high-cost product is to set a high price, and in equilibrium only a small number of consumers will buy the product. With the product-sharing market, the firm has an incentive to increase its retail price. However, this is quantitatively different from the case of a low-cost product—the firm's optimal price in the absence of sharing is already high, the magnitude of the price increase due to consumers' product sharing is not as dramatic for high-cost products. Furthermore, the firm can save a lot of marginal costs of production by selling fewer units at higher prices, which many consumers are willing to pay because of the potential earnings from renting out the product in the sharing market. Therefore the firm is better off. This finding suggests that a firm with high marginal costs of production may have incentives to promote

or improve the sharing market to encourage consumers to share the product even if the firm does not directly profit from the sharing market, because it can indirectly benefit by strategically raising its price to extract some value created by the sharing market.

Note that in the absence of the product-sharing market, only a small number of consumers will buy and use the high-cost product because of its high price. When a product-sharing market exists, some consumers with high usage values only in one period who otherwise will not buy the product will now buy the product even though the price is higher. This is because they anticipate the potential income from renting out the product during the period with low self-use value. The market-expansion effect is relatively stronger for products with high marginal costs—more consumers can use the product (either buy or rent) when the sharing market exists. As a result, consumers' product sharing will increase the total consumer surplus (and also the social welfare).

In summary, our analyses suggest that the consumer's sharing of high-cost products (such as high-tech products, cars, or agricultural equipment in developing countries) is overall beneficial for both the consumers and the manufacturer. In contrast, the sharing of products with very low marginal costs (such as digital products, information goods, or small tools) may be bad for both consumers and the firm. The findings are consistent with the anecdotal observations that firms in industries with high unit costs tend to encourage or facilitate sharing (e.g., GM) and firms selling information goods tend to discourage or curb consumers' sharing.

Next we analyze and discuss the robustness of our insights to several alternative modeling assumptions. First, we have extended our two-period product-sharing model to an n -period model. For analytical tractability, we assume that consumers learn their usage in each period at the beginning of that period. That is, when deciding whether to buy a product from the firm in the first period, consumers know their first-period usage value (i.e., v_{i1}), but for later periods $j = 2, \dots, n$ they know only the distribution of their usage value, i.e., $v_{ij} \sim U[0, q]$. We find that our main results remain qualitatively the same. Product sharing is win-win for the firm and the consumers when the firm's marginal cost is high, and lose-lose when the marginal cost is low.

Second, our core model has not explicitly considered any depreciation of the product over time. If we allow for product depreciation, e.g., the product quality is q for the first period but $q(1 - \Delta)$ for the second period, where Δ represents the rate of depreciation over time, we find that in equilibrium Δ will lower both the firm's retail price and the second-period sharing price. The analytical solutions for such a model become very cumbersome, but our main qualitative insights and intuitions remain the same as long as Δ is not too large.

Third, the sharing market has a salient moral-hazard problem—the consumer renting other's product may use it more abusively or carelessly than the product's owner does. For example, the renter may drive a rented car much less carefully with fast acceleration, hard braking or not slowing down on uneven or speed-bumped roads. In our core model, the transaction cost t can be interpreted as a reduced-form moral-hazard cost that is imposed on the product owner who rents out her

product. For example, we can set t to be the expected cost to the product owner, which is the damage or accelerated depreciation d multiplied by the probability w of such damage occurring. We have also analyzed a more explicit model of moral hazard in the sharing market. More specifically, we assume that for each period the product is rented out, its quality will decrease by δq ($\delta < 1$) and its salvage value will decrease by m ($< \varepsilon$). In reality, the renters may not be able to readily observe the quality-degradation of a previously rented product (i.e., whether a product has been rented out before). We analyze two cases. First, we analyze the case where the quality-degradation of a previously rented product (due to the renter's moral hazard) is observable. Second, we examine the case where the renter does not directly observe the quality-degradation of a previously rented product but will infer an expected degradation in quality of $\tilde{\delta}q$ with $\tilde{\delta} \leq \delta$, which in equilibrium will be fulfilled (from the early-period outcome). Our analysis shows that all our results remain qualitatively the same whether the renters observe the quality-degradation caused by moral hazard. The only difference with our core model is that no sharing transaction occurs in the first period, i.e., in equilibrium only consumers with high first-period usage value will buy the product. These consumers will use the product themselves in the first period and rent it out in the second period if their self-usage value is low. We acknowledge that this analysis is based on a two-period model for analytical tractability. In a general n -period model with moral hazard, the analysis for the sharing market becomes analytically intractable because there will be different quality variations of the products in the sharing market.

Fourth, our core model assumes that the product owner bears all transaction costs for sharing. But in reality, both the product owner and the renter have some transaction costs, for example, the renter may also have to incur some costs for picking up and returning the rented product. We have analyzed a product-sharing model with both parties having some transaction costs: the product owner incurs a cost t_1 and the renter incurs a cost t_2 for each sharing transaction. We show that this model extension is equivalent to our core model with the product owner's transaction cost t replaced by $t_1 + (1 - \alpha)t_2$. Note that the total transaction cost $t_1 + t_2$ is not perfectly or fully internalized through the sharing price—the effective total transaction cost is *smaller* than the direct sum. So, when the renter shares some of the total transaction cost, product-sharing transactions will be more likely to occur. The underlying reason for this effect hinges on the fact that the renter's transaction cost tends to reduce the product-sharing price, which lowers the platform fee paid by the product owner, making sharing more likely.

Besides analyzing the above four formal models, we would also like to briefly discuss how our results may be affected if some other model assumptions are relaxed. First, we have implicitly assumed that consumers know *ex ante* (at the time of purchase) their usage valuation for each period. Actually, this assumption is not necessary, e.g., if consumer i has her usage values v_{i1} and v_{i2} switched between two periods, it will make no difference in our analysis as long as the consumer learns her usage value for each period at the beginning of that period. The assumption we make is only that the consumer's usage value in the population is uniformly distributed in each period. In addition, we have assumed that the consumer's usage

values across different periods are not correlated. Note that if the consumers' usage values are perfectly positively correlated (e.g., each consumer has the same usage values across all periods), in equilibrium there will not be any product sharing among consumers. In a model in which the consumers' valuations are partially correlated across different periods, we expect our main results and intuitions to hold qualitatively the same as long as there is enough valuation heterogeneity across consumers and across periods such that product sharing will occur in the market. The effects of product sharing will be moderated by positive correlation and enhanced by negative correlation between the consumer's usage values.

Second, note that our model explicitly allows consumers to work as a third-party, rental company, which buys the product from the manufacturer/firm and rents it out in all periods. However, we have assumed that such a rental company acquires the product at the same price as consumers do, which leads to no speculators or pure rental agencies in equilibrium. In practice, a rental company might be able to buy the product at cheaper prices than consumers can, or perhaps has lower transaction costs than consumers. In that situation, the product-sharing price in the market will tend to be lower, and hence we expect that the impacts of consumer-to-consumer sharing will be moderated. The intuition and tradeoff from our analysis will still be relevant. As we observe in reality, even in markets with product rental agencies, consumer-to-consumer product sharing is still flourishing.

Third, our main model implicitly assumes an efficient sharing market, i.e., when a match of supply and demand for sharing at a sharing price is possible, the sharing transaction will occur with certainty. In reality, there can be inefficiency in the market; that is, there may be a positive probability that some sharing transactions will not take place at the theoretical market-clearing price, for random reasons such as severe weather. Conceptually, one can extend our model to incorporate that probability into the sharing market. In that case, when making the product purchase decision, the consumers will lower their expected revenue from the sharing market. However, we expect that barring extreme inefficiency in the sharing market, our key results and insights will remain qualitatively the same, albeit the parameter regions may change.

Lastly, our model assumes that all consumers are forward-looking and fully anticipate the possibility of product sharing. The opposite assumption is that consumers are all myopic, i.e., their purchase decisions are based only on their current-period utility and when deciding whether to buy the product they will not consider the potential income from product sharing in the future. In that extreme case, obviously, the firm's pricing and quality decisions will be the same as if the product-sharing market does not exist. However, since product owners can *ex post* decide to rent the product out during usage periods with low self-use values, the consumer surplus will be higher than in the case of forward-looking consumers—interestingly, consumers are better off being myopic than being strategic and forward-looking. In a model in which some consumers are myopic, we expect our main results and intuition to stay qualitatively the same as long as a large enough fraction of consumers are strategic, though the quantitative effects will be moderated.

3.4 Effects of Sharing on Product Quality and Distribution Channel

In this section, we examine the effects of consumer-to-consumer product sharing on firm's product quality decision, and its impacts on a distribution channel.

3.4.1 *Effects of Sharing on Product Quality*

With booming and maturing of sharing markets, one may expect that firms will over time become more strategic when they design their products in anticipation of the consumers' product sharing. In this section, we explore such a situation, where the firm responds to the anticipated product sharing among consumers by strategically choosing not only its price but also its product quality. How will such strategic behaviors by the firm influence the market outcome and the impact of the consumer's product sharing? We address this research question by extending the core model to allow for the firm's endogenous quality decision.

The firm's marginal cost of production typically depends on the quality level of the product. For example, a luxury model of a car will cost the manufacturer more to make than an economy model. For analytical tractability, we use the commonly adopted quadratic cost function: $c = k_1 q^2$. To simplify the later analytical expressions, instead of expressing the product's salvage value as some fraction of the cost (c) of producing the product, we write the product's salvage value as $\varepsilon = k_2 q^2$, where $k_2 < k_1$, i.e., the salvage value is k_2/k_1 fraction of the marginal cost of the product.

Note that the transaction cost (t) for sharing can be related to the product's value, which depends on the quality of the product. For example, other things being equal, the product owner will assess a higher (moral-hazard) cost for accelerated product depreciation or maintenance when sharing an expensive high-quality car than when sharing a low-quality economy car. This can be due to, for instance, the anticipated higher cost of maintenance services for the high-quality car (e.g., changing its high-performance tires or brakes) or other risks associated with sharing. For simplicity, we assume that $t = \tau q$, where τ represents the transactional friction of sharing. Note that the extended game builds on the core model analyzed in the earlier sections; the only difference is that the firm will now strategically choose both p and q to maximize its profit. Since no consumers will share the product if $\tau \geq 1 - \alpha$ (in which case whether the sharing market exists makes no difference), we will focus on the nontrivial parameter region of $\tau \in [0, 1 - \alpha)$.

Our analysis shows that the consumer's product sharing gives the firm a strategic incentive to increase product quality. With the sharing market, those consumers with a high usage value in one period and a low usage value in another period will be willing to pay more for the product since they can earn some rental income from the sharing market when their own usage value is low. This in effect increases those

consumers' willingness to pay for quality and hence gives the firm an incentive to raise its product quality in equilibrium. However, the increase in product quality does not lead to an increase in consumer surplus. In fact, because the firm will strategically raise its retail price to target a smaller number of customers, the existence of the sharing market will reduce the total consumer surplus even though some consumers with very high valuation for quality will become better off (due to the quality increase). By choosing its product quality and price strategically, the firm will make more profits when the sharing market exists. This result is different from the case where the firm strategically chooses only its price. The potential positive effect of product sharing on consumer surplus goes away and the firm is always better off when it strategically chooses both its product quality and price. This difference mainly comes from the fact that, in anticipating the consumers' product sharing, the firm's endogenous quality decision allows it to strategically select a price-quality pair (or equivalently a price-cost pair, since the firm's marginal cost is a function of quality) to ensure higher profitability by extracting more surplus from customers. We also find that the sharing market increases social welfare when the transaction cost is low but not when the transaction cost is high. This finding suggests that it can be overall socially beneficial to reduce the transaction cost in the sharing market.

3.4.2 Effects of Sharing on Distribution Channel

Most consumer products are sold through distribution channels or intermediaries rather than directly by manufacturers. In addition, manufacturers have to build production capacities long before the selling seasons. So, it is of both practical and academic interest to study how the fast-growing product-sharing phenomena affect distribution channels. In our current research (Tian and Jiang 2018), we address the following research questions. How will consumer-to-consumer product sharing affect different members (e.g., the upstream manufacturer and the downstream retailer) of the distribution channel? Will consumers' product sharing increase or decrease the manufacturer's optimal capacity? On one hand, one may intuit that because of higher utilization of the products, the manufacturer should reduce its capacity. On the other hand, the sharing economy allows firms better utilize any excess capacity during time of low demand, so firms should have more incentive to increase their capacity (PWC 2015). Ex ante, it is not clear how product sharing affects a manufacturer's optimal capacity or what factors are important in determining that effect. But clearly, as the sharing economy grows, a strategic firm with long-term vision should consider the effect of sharing in its capacity-planning decision.

We address these research questions by extending our earlier model to consider a distribution channel, in which an upstream manufacturer chooses the production capacity and sells its product at a wholesale price to a downstream retailer, which chooses its retail price to consumers. Using an n -period product-sharing model, we

analyze how the upstream manufacturer should strategically choose its production capacity and wholesale price, and how the retailer should strategically choose its retail price, to respond to anticipated product sharing by the consumers. For analytical tractability, we assume that consumers learn their usage in each period at the beginning of that period. Note that increasing production capacity often requires the manufacturer to make additional fixed-cost investments beforehand, e.g., expanding the production facilities and acquiring more machinery or equipment. So, in our model, the manufacturer first chooses what capacity level K to build—the maximum number of units of the product that the manufacturer can produce—subject to the commonly used quadratic fixed cost $C(K) = \varphi K^2$, where $\varphi > 0$ measures how costly it is for the manufacturer to build production capacity. A small φ means that the manufacturer is very cost-efficient at building up its production capacity.

Our analysis reveals that there exists a threshold for the capacity cost coefficient, above which product sharing will raise the manufacturer's optimal capacity and below which it will reduce the manufacturer's optimal capacity. The intuition lies in the tradeoff between the cannibalization effect and the value-enhancement effect of product sharing. The product-sharing market offers consumers a sharing alternative to buying the product, which lowers the demand intercept and tends to make fewer consumers buy the product—the cannibalization effect. However, anticipating the potential rental income from the sharing market, strategic consumers will ascribe a higher value to owning the product and will tend to be less price sensitive (i.e., the demand curve becomes less steep); thus, at a given retail price, consumers will be more likely to buy the product when a sharing market exists—the value-enhancement effect. When production capacity is not very costly to build, due to the high availability of products for sharing, the cannibalization effect of the sharing market will dominate its value-enhancement effect. But, if capacity is very costly to build, the relatively low availability of products available for sharing will lead to high sharing prices, making the value-enhancement effect of the sharing market dominate its cannibalization effect. Thus, in that case, consumer-to-consumer sharing can actually lead to higher production, increasing both unit sales and the efficiency of the product's utilization.

We find that when the manufacturer strategically chooses its capacity, consumer-to-consumer product sharing will increase the downstream retailer's equilibrium share of the total gross profit margin in the channel relative to the case of no sharing. Capacity cost tends to induce the manufacturer to raise its wholesale price, leaving less room for the retailer to price strategically, which effectively reduces double marginalization and the retailer's power in the channel. However, product sharing, relative to no sharing, gives consumers more value for owning the product, increasing their willingness to pay for the product and reducing their price sensitivity, which will therefore provide more room for the retailer's strategic pricing, increasing the retailer's power in the channel. In other words, the sharing market tends to exacerbate the double-marginalization problem in the channel.

We also show that product sharing can be lose-lose, lose-win, or win-win for the manufacturer and the retailer, depending on the manufacturer's capacity cost efficiency. When capacity is very costly, both the manufacturer and the retailer

will benefit from the consumer's product sharing; when capacity is not very costly, product sharing among consumers will make both the manufacturer and the retailer worse off. The intuition also lies in the tradeoff between the cannibalization effect and the value-enhancement effect of the sharing market. When production capacity is very costly to build, the value-enhancement effect of product sharing will dominate the cannibalization effect, making the firms better off; in contrast, when capacity is very inexpensive to build, the opposite is true, i.e., both firms are worse off. Furthermore, when the capacity cost is in the middle range, product sharing will make the manufacturer worse off but the retailer better off; this suggests that product sharing is more likely to increase the retailer's profit than the manufacturer's profit. This result is consistent with the result that in a channel with endogenous capacity, sharing among consumers tends to increase the downstream retailer's share of the total gross margin in the channel. In such a distribution channel, the retailer is more likely to have an incentive to promote or facilitate consumer-to-consumer sharing, especially when capacity is costly to build.

3.5 Conclusions and Discussions

Collaborative consumption has emerged as a major trend in recent years as the global economic recession has put financial pressure on consumers and as social concerns about consumption sustainability bring the society's attention to effective use of resources and products. Advances in mobile communication technologies and online product-sharing platforms have helped to facilitate product sharing among consumers on an unprecedented scale. Consumers share a wide range of products from bicycles, cars, videogame consoles, to clothing, portable tools, and household appliances. We have provided an analytical model that captures the idea that a consumer's own usage value for her purchased product may vary over time. In a period of low self-use value, the product owner can forgo her product use and rent it out to others through a third-party sharing platform. For each sharing transaction, the renting customer pays a rental fee and the product owner pays the platform a percentage fee. We have examined the consumer's purchasing and sharing decisions, and investigated how a brand owner or manufacturer of the product should strategically choose its price and product quality to respond to the anticipated product sharing among consumers. We have also examined how consumer-to-consumer product sharing affects a distribution channel.

We have shown several main findings. First, transaction costs in the sharing market, such as the transaction cost or the platform's percentage fee, have a non-monotonic impact on the firm's profit, consumer surplus, and social welfare. Second, if the firm strategically chooses its price (taking quality as given), then product sharing among consumers can be either lose-lose or win-win for the firm and the consumers. It is lose-lose when the firm's marginal cost is low and the transaction cost is not too high. In contrast, it is win-win if the firm's marginal cost is high. Third, if the firm strategically chooses both its price and product quality in

anticipation of the sharing market, consumer-to-consumer product sharing will lead to higher quality but even higher prices, increasing the firm's profit but lowering consumer surplus. Fourth, in a distribution channel, the existence of the sharing market is more likely to increase the retailer's profit than the manufacturer's profit, i.e., product sharing can sometimes benefit the downstream retailer at the expense of the upstream manufacturer.

We conclude by pointing out some caveats and potential directions for future research. First, we have analyzed only a monopoly market. We expect that if the firm has competitors, its ability to extract consumer surplus from the product's value enhanced by the sharing market will be moderated, depending on the level of competition and product differentiation among competitors. So, consumers will be more likely to benefit from sharing whereas the firms' gains from sharing may be very limited. Second, we have assumed that the firm plays no direct role in the sharing market. As the sharing economy grows, the firm itself may enter the sharing market. For example, BMW has recently entered into the U.S. car-sharing market and launched a new car-sharing service called ReachNow.⁴ Our future research will focus on the tradeoff that a firm faces when directly entering into the sharing market. Third, for analytically tractability, we have implicitly assumed that the firm keeps a uniform pricing strategy and does not adjust its retail price from one usage period to another. Clearly, given the already high complexity of the current analytical framework, one will have to make other modeling simplifications to be able to study dynamic pricing strategies in a market with consumer-to-consumer product sharing. We leave it to future research to explore the potential new insights from a framework with intertemporal price discriminations by the manufacturer and/or the retailer. Fourth, we have not explicitly modeled any uncertainty in the sharing market. In essence, the consumer is assumed to be risk neutral and makes her decision based on the average of the anticipated revenue from product-sharing transactions. We have also focused on search goods rather than experience goods, whose quality may not be fully observed by the consumers prior to purchase. We will leave it to future research to study the effects of uncertainty in the sharing market and uncertainty in the firm's product quality. Fifth, we have assumed an exogenous proportional fee by the sharing platform, in line with the observed reality, where the platform's percentage fee does *not* vary across different products or product categories. However, it may be of interest to examine what happens if the platform charges different fee percentages based on some product characteristics (e.g., a lower percentage for high-end products to encourage sharing of such products). Such studies may provide strategic recommendations different from the platforms' current strategies of not adjusting fee percentages based on products or product categories. Price discrimination issues by the sharing platform also deserve their own theoretical study in future research. Lastly, collaborative consumption in the sharing economy is a fast growing trend; we have studied consumer-to-

⁴<https://techcrunch.com/2016/04/08/bmw-just-jumped-into-the-u-s-car-sharing-biz-with-the-help-of-yc-alum-ridecell/>

consumer *product* sharing but not peer-to-peer *service* offerings such as Uber. When consumers share their purchased products (e.g., tools or cars), the manufacturers and the retailers of those products will be affected and can strategically change their decisions, which is the focus of our research. In contrast, when consumers offer their time or services (e.g., a consumer on TaskRabbit assembles IKEA furniture for another consumer), typically there is no strategic upstream supplier (manufacturer or retailer). In those situations, the labor market and the traditional service providers may be affected by the peer-to-peer service platforms. Theoretical and empirical research on both types of collaborative consumption are of great managerial and academic interest.

Acknowledgements This chapter is extensively based on the authors' paper in *Management Science* (Jiang and Tian 2018).

References

- Agrawal VV, Ferguson M, Toktay LB, Thomas VM (2012) Is leasing greener than selling? *Manag Sci* 58(3):523–533
- Bai J, So KC, Tang CS, Chen X, Wang H (2016) Coordinating supply and demand on an on-demand service platform: price, wage, and payout ratio. Working paper, University of California, Irvine
- Belk R (2010) Sharing. *J Consum Res* 36:715–734
- Belk R (2014) You are what you can access: sharing and collaborative consumption online. *J Bus Res* 67:1595–1600
- Bellos L, Ferguson M, Toktay LB (2017) The car sharing economy: interaction of business model choice and product line design. *Manuf Serv Oper Manag* 19(2):185–201
- Benjaafar S, Kong G, Li X, Courcoubetis C (2018, Forthcoming) Peer-to-peer product sharing: Implications for ownership, usage, and social welfare in the sharing economy. *Manag Sci*
- Chen J, Esteban S, Shum M (2013) Why do secondary markets harm firms. *Am Econ Rev* 103(7):2911–2934
- Desai P (1999) Competition in durable goods markets: the strategic consequences of leasing and selling. *Market Sci* 18(1):42–58
- Desai P, Purohit D (1998) Leasing and selling: optimal marketing strategies for a durable goods firm. *Manag Sci* 44(11):19–34
- Fraiberger S, Sundararajan A (2015) Peer-to-peer rental markets in the sharing economy. Working paper, New York University
- General Motors (2012) RelayRides and OnStar: baby, you can rent my car. GM news. http://media.gm.com/media/us/en/gm/news.detail.html/content/Pages/news/us/en/2012/Jul/0717_onstar.html
- He L, Mak HY, Rong Y, Shen ZJM (2017) Service region design for urban electric vehicle sharing systems. *Manuf Serv Oper Manag* 19(2):309–327
- Hendel I, Lizzeri A (1999) Interfering with secondary markets. *RAND J Econ* 30(1):1–21
- Jiang B, Tian L (2018) Collaborative consumption: strategic and economic implications of product sharing. *Manag Sci* 64(3):1171–1188
- Johnson JP (2011) Secondary markets with changing preferences. *RAND J Econ* 42(3):555–574
- Mussa M, Rosen S (1978) Monopoly and product quality. *J Econ Theory* 18:301–317
- PWC (2015) The sharing economy: consumer intelligence series. PWC research report. <http://www.pwc.com/us/en/technology/publications/sharing-economy.html>

- Tian L, Jiang B (2018) Effects of consumer-to-consumer product sharing on distribution channel. *Prod Oper Manag* 27(2):350–367
- Weber TA (2014) Intermediation in a sharing economy: insurance, moral hazard, and rent extraction. *J Manag Inf Syst* 31(3):35–71
- Weber TA (2016) Product pricing in a peer-to-peer economy. *J Manag Inf Syst* 33(2):573–596
- Zervas G, Proserpio D, Byers JW (2017) The rise of the sharing economy: estimating the impact of Airbnb on the hotel industry. *J Market Res* 54(5):687–705

Chapter 4

Operational Factors in the Sharing Economy: A Framework



Tunay I. Tunca

Abstract Applications of sharing economy, from ride and home sharing to crowd-funding to online freelance markets have been playing increasingly prominent roles in people's daily lives. An important reason for sharing economy's rise to prominence is the operational efficiencies it introduces, and the related savings and other economic benefits it unlocks. In this paper, we provide a framework identifying and describing the forces of sharing economy that fuel the success of the novel business models of the concept. In addition, we study two applications, namely ride sharing and group buying, in more detail, analyzing the operational efficiencies created by each business model with the framework we introduced and providing evidence from recent related literature for the efficiency gains they bring about. Finally we discuss some potential downsides and pitfalls that arise as the side effects of these operational efficiencies of sharing economy business models, and the related regulatory issues ahead that may need attention.

4.1 Introduction

Since the final decade of the twentieth century, the Internet has been changing the society and the global economy in a scale and intensity that can be considered unprecedented for any technology in human history. There are drastic and irreversible differences in the way people communicate, work, shop, search for and retrieve information, obtain news, consume entertainment, socially organize, travel, connect with friends, family and strangers, and carry out many other aspects of their lives today, compared to before Internet started entering daily life two just two decades ago. Among many ways the Internet has affected peoples' lives, perhaps one that had the most direct economic effect and that has generated one of the highest efficiency gains is its enabling of the sharing economy in the past decade.

T. I. Tunca (✉)

Robert H. Smith School of Business, University of Maryland, College Park, MD, USA
e-mail: ttunca@rhsmith.umd.edu

© Springer Nature Switzerland AG 2019

M. Hu (ed.), *Sharing Economy*, Springer Series in Supply Chain Management 6,
https://doi.org/10.1007/978-3-030-01863-4_4

Spanning many aspects and facets of the economy, today sharing economy touches hundreds of millions of peoples' lives daily both as consumers and providers. Its extent spans a wide variety of concepts and contexts from ride sharing to home sharing, group buying to crowdfunding, and peer-to-peer lending to online freelancing. In its foundation, the idea is utilizing distributed networks of people and resources, by bringing them together to increase efficiency of value creation by physical and financial assets, time, connections, expertise, and labor. The powerful and ever growing network effects on the Internet present an enormous potential for the growth of the sharing economy. The global sharing economy is expected to grow from \$14 billion in 2014 to \$335 billion dollars by 2025 (Yaraghi and Ravi 2017).

There are many ways the sharing economy creates and unlocks value. Many of the efficiencies it creates fall in the domain of operational efficiency since the force behind the engine of the sharing economy is essentially better matching of demand and supply through improved utilization of resources. This includes improved logistics, improved allocation of peoples' time and living spaces, efficient utilization of networks and incorporation of small scale scattered and dormant resources into the active economy. This article presents a framework for understanding the operational factors and advantages that shape and power the sharing economy. The goal is building a systematic and well-connected structure for identifying, understanding and analyzing parts and applications of the operations management on sharing economy under a unified umbrella. We further provide examples from recent applications and literature in the area to demonstrate how the elements of our framework appear in this research and how one can analyze the insights from these studies under the lens of this framework.

The rest of this article is organized as follows: Sect. 4.2 presents and discusses our framework. Section 4.3 presents two applications of sharing economy business models as examples, showing how the operational advantages presented in the framework yield efficiencies for these models, and discusses empirical evidence from recent literature on their value creation. Finally, Sect. 4.4 offers our concluding remarks.

4.2 The Framework

In this section, we present our framework of five sharing economy operational factors that shape the efficiency of the concept's business applications. In particular, these are:

1. Utilization of sunk and fixed costs
2. Utilization of bit sized resources
3. Utilization of human idle time
4. Utilization of networks to lower barriers to entry into workforce and markets
5. Assigning people new operational and economic roles

We next explain and analyze each one of these factors in detail.

1. Utilization of Sunk and Fixed Costs One of the most impactful ways the sharing economy unlocks value that would have been otherwise lost is its enabling increased usage of dormant resources. Consider the following: It is estimated that on average a given car is parked 95% of the time (Barter 2013; Morris 2016). This means that at any given time point, hundreds of millions of cars around the world are dormant resources that do not generate concrete economic value that fulfill their potential. Similarly, according to National Association of Home Builders (NAHB) there were approximately 7.5 million second homes in the United States (Zhao 2016). Including the empty rooms in the primary homes of more than 125 million households in the country (Statista 2016), tens of millions of rooms in American homes every day go unused. Therefore, it is not surprising that, once platforms were made available for cars to be utilized through ride sharing, within a few years of Uber's and Lyft's respective starts, the number of cars used for ride-sharing worldwide by Uber surpassed one million in 2015, and reached 700,000 for Lyft in 2017 (Lazo 2015; Weise and della Cava 2017). Similarly, in less than a decade from AirBnB's foundation, the number of listings available worldwide on that platform surpassed four million in 2017 (Hartmans 2017). In addition, another advantage these sharing economy providers have is that commercial transportation and hospitality services require not only provision of new replacement resources themselves, but also need additional fixed costs of operation, such as maintenance, insurance and registration fees for cars, and taxes, utilities, and other costs such as association dues for houses. Therefore, when one considers the competition between commercial entities and individual providers, this highlights another advantage for sharing economy business models: A significant portion of the fixed operating costs are sunk costs for sharing economy providers, and another significant portion of these costs can be better amortized and spread out over time and usage. This reduces the overall effective operating costs for the sharing economy providers and puts the traditional commercial providers such as taxi companies and hotels in a significant competitive cost disadvantage. This cost asymmetry further fuels the explosive increase of market share for the sharing economy providers versus traditional commercial entities.

2. Utilization of Bit-Sized Resources Through the power and reach of the Internet, the sharing economy has unprecedented ability to utilize small, distributed resources. Perhaps one of the best examples for this advantage is the rise of Peer-to-Peer (P2P) lending. In P2P lending sites, ordinary people come together to finance small size loans for each other. The loan amounts are usually limited, by rules, to a range of \$1,000 to \$40,000, and the average loan size is typically less than \$15,000 (Cunningham 2015; Singh 2016; Frankel 2017). The borrowers and the loan risk are rated by the online platform, and the lenders can view and choose by themselves whom and which loans to finance. In the past decade, P2P lending industry showed tremendous growth. The two largest P2P lending sites in the U.S., namely Lending Club and Prosper, both originated in mid 2000s, account for a combined \$36 billion dollars in loans processed as of 2017 (Investment Zen 2017; Lending Club 2017), and the market potential for P2P lending is estimated to be higher than \$350 billion

(Singh 2016). The driving force behind this explosive growth is the ability to bring together the demand and supply for very small resources with centralized, low-cost operations. Without an online business model like P2P lending, the loan matching has to be done through banks, and in some cases even through very high interest payday loan services, which add substantial amounts of financing costs and overhead, resulting in lower returns for the lenders and high interest rates for the borrowers. P2P lending, however, is a win-win for both lenders and borrowers. What is more, one can expect that with reduced interest rates on loans, many borrowers who would have been priced out of the market are now able to borrow money and use it, stimulating further economic activity, which means that P2P lending unlocks significant value in the economy. Parallel arguments can also be made for ride sharing and home sharing markets. Drivers who operate on ride sharing platforms like Uber and Lyft have full flexibility on when to drive and when to be inactive. That means that if they have a small amount of time in hand and prefer to drive at that point, they can drive to earn money and generate value. Similarly, without the home sharing business models like AirBnB, it would not be possible to organize and match hundreds of thousands of unused but potentially available rooms in private homes with the demand for those rooms. This is an unprecedented flexibility in utilization of such small sized resources, which, in aggregate constitute a substantial and otherwise untapped economic value.

3. Utilization of Human Idle Time According to United States Department of Labor, in 2015 average American, 15 years or older, spent 4 h and 59 min for leisure, relaxation and sports every day (Bureau of Labor Statistics 2015). According to the same report, unemployed individuals had nearly 7 h of free time while people with full-time employment had more than 4 h free time per day on average. Although the level for satisfactory amount of free time per day varies from person to person, especially at times of high unemployment, the amount of free time people have and prefer to be working can be substantial. Furthermore, this idle time comes with near zero opportunity (or shadow) cost and utilization of it is almost entirely a net upside. With the flexibility of the “work any time” concept, sharing economy business models enable people who have free time to participate in the economy, and can generate great amount of net value for the society. A very good example for this phenomenon is online freelance market places. Freelance marketplaces, such as *Upwork* (formerly *Odesk*, which had merged with *Elnance*, two pioneering freelance websites), *Fiverr* and *TaskRabbit* connect demand and supply for labor for a very broad spectrum of services, ranging from coding, web design, graphic design, writing, more recently to home improvement, consulting, marketing, and accounting. By doing so, these websites allow people from a broad range of professions to independently book jobs at a per task basis, which means that the providers have flexibility to organize their schedules and take a task during their idle times. This enables people to efficiently decide how to utilize potentially many different pieces of their free time, which they would not necessarily have been able to otherwise. A recent study has found that there are about 55 million freelancers in the U.S. corresponding to 35% of the U.S. workforce with total earnings of

approximately \$1 Trillion in 2016, and 66% of these freelancers declared that amount of work they found online had increased in that same year (Upwork 2016). Although some freelancers use these websites as their only source of employment, many use it for supplemental work as well. Similar arguments can be applied to ride sharing and to some ways home sharing. Overall the value generated by efficient utilization of human resources enabled by sharing economy business models is substantial, especially when one includes the additional societal and economic benefits of reduced unemployment brought about in the population.

4. Utilization of Networks to Lower Barriers to Entry into Workforce Many of the factors we discuss in our framework have effects on lowering the barriers to entry. However, through utilization of the global Internet infrastructure, sharing economy has another very powerful way in lowering the barriers to entry, specifically into labor markets. Once again, a very good example of this is ride sharing. For decades, the taxi industry had been tightly controlled through municipalities, with medallions (effectively licenses to operate individual taxicabs in a certain municipality) sold to a limited number of taxi companies and individuals at substantial prices. In many cases, the exorbitantly high prices of these medallions signalled the extent of imbalance in demand and supply. (Please see Sect. 4.3.1 below for detailed statistics in medallion prices and related factors.) Another indicator of this imbalance is how in the past few years the numbers of Uber and Lyft drivers skyrocketed. For instance in approximately 18 months following the launch of UberX, the number of drivers who were active providers (defined as those who completed at least four trips a month) increased to 160,000 (Hall and Krueger 2015). One of the most important reasons for this increase is reduction in search and matching costs. Before UberX enabled these drivers to connect potential customers, an individual effectively could not enter the market to provide rides for fares. This was not necessarily because of legal barriers, which proved to be not too difficult an obstacle to overcome (Lawler 2013; Geron 2013), especially compared to the much bigger obstacle of connecting of drivers and passengers in real time (U.S. Department of Transportation 2006; Chan and Shaheen 2012). The availability of smartphones with GPS technology finally solved the problem by 2010 and the biggest barrier to entry to the market was removed. This led to business models like Uber and Lyft, where the company only provides the platform for matching customers and drivers, hence solving the most challenging part, and having removed that barrier, every driver individually can choose to enter the market with ease, essentially by downloading an App and signing up with the company (and in some cases going through certain background checks). Similar arguments can be made for home sharing companies and online freelance market places as well. The effects of these lowered barriers on their respective industries can be tremendous, as we discuss in more detail for the ridesharing industry in Sect. 4.3.1.

5. Assigning People New Operational and Economic Roles One of the most unique ways the sharing economy affects the business environment is its ability to assign and leverage new business roles to people. Aside from the immediate

examples that come to mind like how with much ease ride-sharing allows people to become commercial drivers, and home sharing allows people to become providers in the hospitality industry, this effect in fact is much broader than it first catches the eye. Two interesting examples come from Crowdfunding and Group Buying. Originating in early 2000s Crowdfunding allows many small donors or investors to chip in funds for the realization of a project or getting a start up off the ground (Du et al. 2017; Miller 2017). There are two types of crowdfunding events: The first kind is the more philanthropic type, in which a number of donors effectively give their money for a cause or product they would like to see materialized without much expectations other than occasionally a token gift, which could be the funded product. This market, lead by websites such as Kickstarter and Indiegogo, with average funding amount per campaign about \$990, has reached a total estimated transaction volume of \$7.32 billion dollars in 2017, and is predicted to grow to a transaction volume of \$18.97 billion by 2021 (Statista 2017a). In this model, the participants are enabled to become philanthropic contributors or partners in projects, which otherwise would not have materialized. In the second form of Crowdfunding, which is also called *Crowdinvesting*, in which participants provide funds with the expectation that they will become equity holders in a start up company if the company successfully launches. Lead in the U.S. by sites such as EquityNet, CrowdCube and Seedrs, the global market size of this segment is \$5.69 billion with expected 2021 market size of \$19.33 billion, and the average project funding amount in this segment is about \$120,000 (Statista 2017b). Although being an inherently risky activity, Crowdinvesting allows individuals to effectively assume the role of small scale venture capitalists, an activity that became legal only recently in the United States by the “Jumpstart Our Business Startups” (JOBS) act of 2012, (Reiss 2016), while enabling the funding of thousands of potential start ups that can generate billions of dollars of value for the global economy. Finally, another striking example comes from Group Buying, which generates significant value for both companies and customers by incentivizing customers to assume the role of sales agents. We discuss Group Buying in more detail below in Sect. 4.3.2.

In different sharing economy business model applications, one may be able to see all of these five factors or a subset at work, in many cases the factors interacting and teaming up with each other to amplify and complement the effects of one another. We will next discuss two examples in detail to demonstrate the application of our framework.

4.3 Examples

In this section we discuss two sharing economy business models, namely Ride Sharing and Group Buying, more closely within the structure of our framework, highlighting the factors that affect the business landscape for these two particular applications.

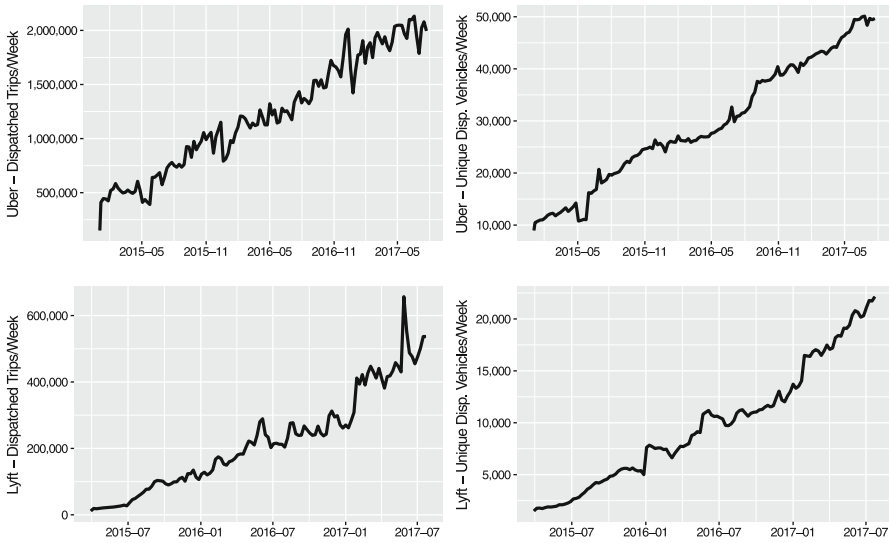


Fig. 4.1 The growth of Uber (the top two panels) and Lyft (the bottom two panels) in New York City from early 2015 to August 2017. For each company, the left panel illustrates weekly number of dispatches, and the right panel depicts the weekly number of unique vehicles dispatched for this time period. (Data source: New York City Open Data, <https://opendata.cityofnewyork.us/>)

4.3.1 Ride Sharing

Of all the many influential applications of the Sharing Economy, perhaps the most ubiquitous one is ride sharing. Considering that the concept materialized in practice less than a decade ago, there are few disruptive new business models in history that had industry-shaking effects as fast and strongly as ride sharing did. Today the industry is lead in the United States by Uber and Lyft, with an estimated total number of daily rides in the United States being more than 5.5 million and 1 million for these companies respectively (Locklear 2017), and their major global rival China’s Didi Chuxing accounting for upwards of 20 million daily rides (Millward 2016). Uber and Lyft are estimated to be worth more than \$60 billion and about \$7.5 billion dollars respectively (Russell 2017; Etherington 2017), while Didi Chuxing has an estimated value of \$50 billion (Macfarlane 2017).

Figure 4.1 demonstrates the growth of number of weekly dispatches and the number of unique vehicles dispatched by Uber and Lyft in New York City from early 2015 to August 2017. Note the growth in number of weekly dispatches for Uber from 151,235 in January 2015 to 1,995,459 in July 2017, a growth rate of approximately 9% a month or 181% a year. Similarly the number of weekly dispatches for Lyft grew from 11,588 in April 2015 to 536,966 in July 2017, by approximately 15% a month or 418% a year. Further, for Uber, the weekly number of unique vehicles dispatched increased from less than 8,980 in January 2015 to

49,665 in the beginning of August 2017, corresponding to an average monthly increase of close to 1,300 vehicles. For Lyft, weekly number of dispatched vehicles increased from 1,501 in April 2015 to 22,144 in August 2017, an average monthly increase of about 740. These numbers are even more striking when compared to the total number of taxi medallions in New York City, which was capped at 13,587 for 2016 and 2017 (Agovino 2017). According to Certify, the second largest expense management software provider in North America, in the fourth quarter of 2016, Uber captured the majority of ground transportation business transactions processed through their system, with a market share of 52% (Hagan 2017).

What accounts for this astonishing growth of ridesharing activity? When viewed through the lens of the framework we presented in Sect. 4.2, the components of the success of ride sharing companies like Uber and Didi emerge clearly. First, ride-sharing uses *existing resources*, namely private cars, which means that, fixed costs such as purchasing, periodic maintenance, registration and taxes are sunk. That is, other than drivers' opportunity cost of time, and the marginal operating costs of gas usage, amortization and added maintenance, which are similar or often lower compared to commercial taxi marginal operating costs, there are minimal additional cost considerations. This creates an enormous strategic advantage for ride sharing providers compared to taxi service providers at the time of entry decision since taxi service providers have to internalize the fixed costs of acquisition in their decisions, while drivers who join Uber or Lyft consider *most of the fixed costs as sunk*. Second, looking at the marginal costs, we see that for ride-sharing services, the main component of this type of cost that could differ from that of the taxi services is the driver's opportunity cost of time. Combining the availability of *human idle time* for unemployed and part-time employed individuals, and the ability to use *smaller bits of time* that enables people with full-time employment to provide ride-sharing services at times convenient to them, it can be claimed that individual opportunity cost for time for ride-share drivers is low. Considering the fact that many taxi drivers are third party agents that work for taxi owners, one can conclude that in large part, taxi marginal costs are further influenced by the requirement to provide the primary employment income for the drivers. Therefore, ride sharing not only has a fixed-cost advantage to commercial taxi model, but also a significant marginal cost advantage. Therefore, once the technology that could *utilize networks to enable real time peer-to-peer driver-customer matching* emerged, many people who could potentially be interested in providing rides for fares could be assigned *new roles as drivers* and easily enter the market. The resulting effect could be estimated as ride-sharing services taking significant market share away from taxi services.

One way to see the manifestation of this effect is looking at the revealed market valuation of license to operate a taxicab. Figure 4.2 displays the average monthly taxi medallion price in New York City from January 2010 to August 2017. Historically, due to limited number of availability medallions in major cities, medallion prices have been substantial. This can also be seen in the figure, where up until mid 2013, there is a steady upward trend in the price, with the average price peaking in June 2013 at \$1,050,625. After that point, the price becomes more volatile but still stays high until mid 2014, still exceeding \$1 million in July 2014.

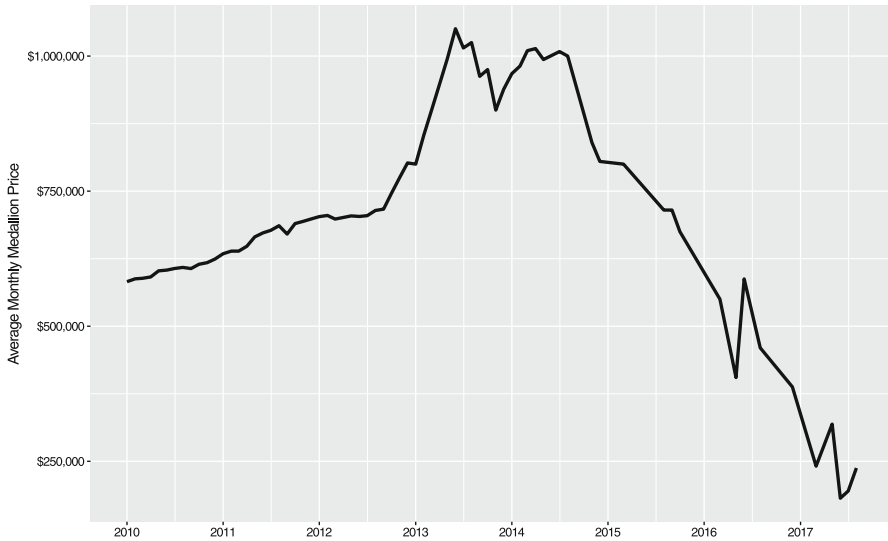


Fig. 4.2 Average monthly non-foreclosure taxi medallion price in New York City from January 2010 to August 2017 (for full medallion sales only). (Data source: New York City Taxi and Limousine Commission)

Beyond that, however, coinciding with the tremendous growth in ride-sharing as depicted in Fig. 4.1, there is a dramatic crash in the medallion price, with the average price falling as low as \$181,666 in June 2017. In August 2017 the average price for unrestricted, non-foreclosure medallion sales transactions in New York City was \$237,500, a 77.4% decline from the peak average monthly price within a time frame of about 4 years.

Another indicator of this sudden unexpected value loss in taxi licenses is the significant increase in the number of medallion foreclosures in the past few years. Figure 4.3 illustrates the rise of foreclosures in New York City, again on the 2010–2017 time frame. As can be seen from the figure, the number of foreclosures for each month between January 2010 to September 2014 was zero except for one foreclosure in August 2011. After September 2014, however, there is a steady and persistent rise in the number of foreclosures with 25 foreclosures in the first 8 months of 2017 and nine foreclosures in August 2017 alone. Overall, the totality of this evidence reflects the pattern of a strong disruptive technology with significant economic advantage, as pointed out by our discussion of the application of our framework, entering an existing market, fundamentally altering it and driving the cost-disadvantaged incumbents substantially out of market share and even the business itself.

Finally, one can also measure the magnitude of the welfare improving effects of the cost advantage coming from the ride sharing business model. Using transactional data obtained from Didi Chuxing, spanning December 2015 and January 2016, Ming et al. (2017) study demand, supply and price formation in the Beijing

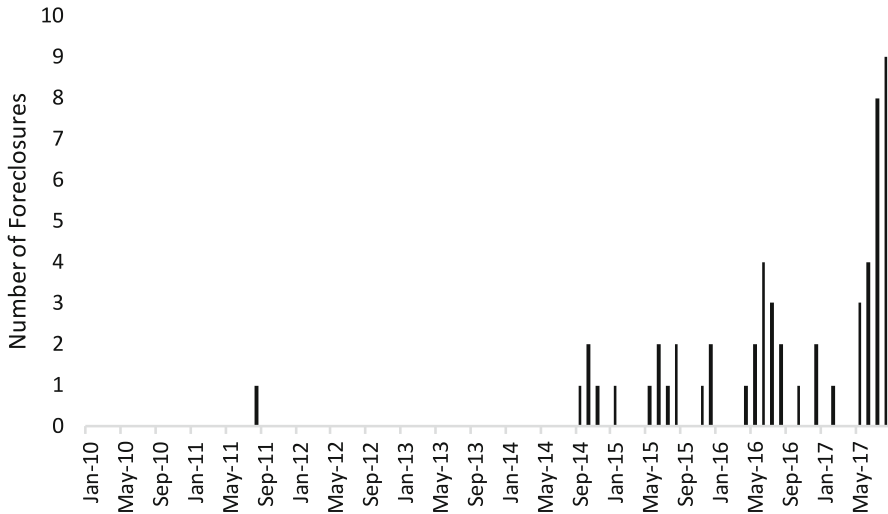


Fig. 4.3 Number of taxi medallion foreclosures in New York City per month from January 2010 to August 2017. (Data source: New York City Taxi and Limousine Commission)

ride-sharing market. Employing a two-sided multinomial logit preference based regression model, and controlling for factors such as time of the day, weather, air pollution, weekend and holiday effects, they estimate the consumer demand and driver supply for ridesharing services for 3-min intervals during a 24h cycle. They calculate that the estimated average ride-sharing price is lower than that of the average taxi price by approximately 29%, and using the demand and supply estimates from the regression analysis, through a counterfactual analysis, they show that increasing ride-sharing fares to the taxi fare levels would result in about 20% expected consumer welfare loss. That is, empirical evidence supports that lower consumer prices resulting from lower operating costs for the ride-sharing business model generate real economic advantage and create substantial value for customers.

4.3.2 Group Buying

Another application that demonstrates how sharing economy business models are able to extract economic value that would otherwise be very difficult or impossible to harness is group buying. Internet based group buying models initiated in early 2000s and brought into public conscience in the U.S. with companies like Groupon and LivingSocial. The main idea of group buying is that the higher the number of customers that join in the purchase the lower the unit price of the product or the service offered. This is usually implemented using discrete thresholds. That is, if the number of sign ups exceed a certain threshold, the unit price is discounted by a pre-

determined amount which is called the “deal discount”. In a group buying event, there is a fixed event window, which is announced in advance. During the event window, customers arrive at different times, (in most implementations) observe the current number of existing sign ups, and make a decision whether to join or not. In some implementations, there is a non-refundable deposit that a customer pays if she decides to join. At the end of the event window, depending on the number of sign ups, the unit price is announced, and at this point, some customers who signed up may leave if they choose to, forfeiting any deposit they paid at the time of sign up. The remaining customers pay the balance between the realized unit price and the deposit, and receive the good or the service from the provider.

Group Buying as a retailing method organized by e-commerce websites for particular retailers is gaining significant popularity recently around the world, especially in Asia. In 2016, Taobao, China’s largest e-commerce platform, hosted more than 200,000 individual group buying events, generating revenues exceeding 32 billion Chinese yuan or approximately \$4.8 billion US dollars for the retailers (www.sohu.com). Examining the benefits of group buying for customers and the retailers who sell through these events from the perspective of our framework helps reveal why. First, the main idea behind group buying is essentially *efficient use of bit sized resources*. Volume discounts is a procurement strategy traditionally very effectively used by large companies or entities that buy large quantities (Ovans 2000; Pei et al. 2011). Bulk buying reduces transaction costs that come from search and implementation, decreases overhead and *spreads fixed costs among a large number of units*, and thereby saves money to the seller, and the generated surplus can be shared with the buyer through volume discounts for a mutually beneficial transaction. One can easily see that there is great potential to unlock value if the strategy can be expanded to aggregate individual buyers’ transactions. From this perspective, each buyer’s unit purchase is a bit sized resource that can be put to efficient use only when aggregated with many other purchases to generate value for the entire system. Thus, as a sharing economy business model, group buying successfully makes use of existing scattered small resources, which otherwise would be either not used or used suboptimally. Before group buying mechanisms were enabled by the Internet, however, there had not been a method to efficiently bring this idea to life. With the mass communication and coordination capabilities that come through the Internet and utilization of networks that exist on the web, this *barrier to implementation and entry has been reduced significantly*.

Further, group buying, has another very important particular capability to generate value in a very creative way. Just as Crowdfunding and Crowdinvesting assign new roles to people to participate in the economy as philanthropists and small size venture capitalists, as we had discussed in Sect. 4.2, Group Buying creates *incentives for ordinary customers to assume the role of sales agents* for a retailer or a service provider. When customers participate or plan to participate in a group buying event, they benefit from a price reduction only if sufficient number of other customers sign up as well. Therefore, customers have an incentive to spread the word about the event in their social circles, informing and convincing others to join, including family, friends and acquaintances, as well as in some cases, utilizing social

networks on the Internet such as Facebook to advertise the potential deal and recruit potential customers. This can bring a significant demand boost and value for the retailer with minimal costs since the company gets all this extra promotion for free, and this surplus can be shared with the customers in the form of deal discounts. In addition, the mechanism also makes efficient economic use of *human idle time*, by allowing people to decide whether they would like to utilize their existing free time to spend on promoting the group buying event, thereby putting a low opportunity cost resource into use to create economic value.

The capability of group buying to incentivize customers to spend effort as sales agents had been recognized and studied in the academic literature, albeit only by theoretical studies (Jing and Xie 2011; Chen and Lu 2015). In a recent study Ming and Tunca (2017) provide empirical evidence for the existence of the value generated by group buying events and measure its magnitude. Ming and Tunca first develop a dynamic game-theoretical model to capture the sign up behavior of the customers during the group buying event window. They then use sign up data from 266 group buying events hosted by Taobao for a major Chinese appliance manufacturer in 2013 to estimate the customer arrival rates for these group buying events. Utilizing additional sales and product data for 2715 other instances where products were sold through traditional single-prices, and controlling for factors such as product review scores by customers, product types and characteristics, and date effects, they demonstrate that selling through a group buying event instead of traditional single-pricing boosted the demand on average by more than 15%. Further, through a counterfactual analysis, they compare the retailer's realized profits to its predicted profits for the case, where the products were sold through traditional single-pricing, and estimate that on average, selling through group buying improved retailer profits by approximately 11%. Overall this study provides empirical evidence for increased demand and the value unlocked by employing group buying events in the channel, and testifies for the effectiveness of this innovative sharing economy business model.

4.4 Concluding Remarks

In a very short amount of time, sharing economy has grown from virtual non-existence into a major economic force with a significant role in shaping the global business environment. In this article, we presented a framework to help dissect the forces that enable sharing economy business models to be so effectively disruptive in changing industries, generate substantial social and economic value, and distribute it broadly in the society. We then analyzed two sharing economy business models, namely Ride Sharing and Group Buying, in detail utilizing this framework, and discussed the empirical evidence from recent research supporting the predictions coming from the framework about the value generated by these business models. The insights from this analysis can be utilized for other examples of sharing economy applications, which can also be analyzed in a similar way.

An important thing to note however, is that the powerful forces that favor the sharing economy as we discussed in this article, also have downsides, which may limit or reduce the benefits the society can obtain from these innovative business models. One example is the problem of commercial leasing on AirBnB platform. In the home sharing market, the lowered barriers to entry not only allows individuals to make use of vacant rooms in their homes by renting them out on a short term basis, but also enables entities with commercial interests to take advantage of the platform to essentially run larger scale renting businesses by leasing out multiple units or entire buildings. Such entities in some cases are even claimed to acquire buildings in some neighborhoods solely for the purpose of renting them out on AirBnB. This in a way amounts to enabling businesses to run disguised hotel operations without going through regulated channels and in some cases in residential districts. According to a report by the American Hotel and Lodging Association (AH&LA), in the United States, AirBnB hosts renting two or more units accounted for 32.1% of listings and 89% of the platform's revenue growth in 2016 (CBRE 2017). One example of a city with this problem is Washington D.C. where, according to AH&LA president and CEO Katherine Lugar, 24% of AirBnB revenue comes from hosts with more than 20 units listed (Ting 2017). Although these figures are contested by AirBnB, the issue is certainly under the attention of the company as well as municipalities. In 2015 and 2016, AirBnB was accused of removing thousands of listings by commercial operators before making the data on its operations in New York City publicly available (Kulwin 2016; Rossi 2016). In response, cities like San Francisco and New York are implementing stricter laws on short-term rentals and as a result, being targeted with litigation by AirBnB. The commercial operators on AirBnB are not only a threat to traditional hospitality businesses like hotels, but also to cities and neighborhoods, and clandestine commercial entry to the market can steal business away from individuals, reducing the spread of social value generated by the business model.

Another issue that could become problematic in the future concerns individuals who make sharing economy platforms their main career or job outlet. Even though sharing economy is valued by providers and customers for enabling people to be hired for individual "one-off" tasks, thereby creating flexibility and improving utilization of resources, this also means that people who provide these services, such as drivers for ridesharing services or professionals who offer their services on online freelance marketplaces are by default, offer work for the sharing economy without benefits such as medical insurance or retirement accounts, or standard workers' protections. In fact two recent court decisions in London and New York declared that Uber must treat at least some of its drivers as employees, potentially requiring the company to pay millions of dollars in benefits such as healthcare contributions, and overtime and holiday pay (Kerr 2016; Furfaro 2017). On the worker's protection side, the company was sued in Los Angeles and in U.S. federal courts with allegations on committing fraud on driver compensation and for the rights of the drivers to unionize (Rubin 2017; Denton 2017). Concerns about benefits should also be kept in mind for freelancers who plan to use online platforms as their main source for work, since the income made by freelancers may be deceiving

without considering payments for health insurance and routine business expenses, especially also considering the fees paid to the platform. The issue of employee rights and benefits is still a significant concern that is not fully resolved for many sharing economy business models.

Finally, the powerful network and communication technologies that fuel the sharing economy by enabling millions of people around the globe to connect, coordinate and transact at will also raise important privacy concerns. For instance, in order to be able to match customers with drivers in real time and facilitate and keep track of payments, ride sharing platforms have to collect and record detailed GPS location records for millions of passengers and drivers alike. Similarly home sharing companies have records of travel destinations, dates and lengths of stay histories for millions of users. Peer-to-peer lending websites have access to information on the amounts and purposes for the loans their borrowers secure through their platforms, and the history of their creditworthiness, and personal payments or lack thereof. Collection of tracking of such large amounts of highly personal information has created significant personal and financial cybersecurity concerns, as well as concerns about the intrusive use of the data by the platforms themselves. For instance, Uber's tracking of customers through IP addresses even after they disabled their phones' GPS feature, its collection of data that goes beyond the basic trip information, and giving its employees improper access to that data caused significant debate and calls for new privacy legislation covering ride sharing companies' activities (Lyons 2015; Mueffelmann 2015). Unless carefully addressed, such concerns can create major setbacks and impede the growth and development of sharing economy business models.

Overall, the promise of the sharing economy in improving welfare of millions of people worldwide is immense. However, moving forward, industry self regulation as well as protective legislation is necessary to prevent abuse of the operational advantages offered by sharing economy business models, to reduce and control the negative externalities they may generate, and to protect their intended economic value creation for a broad base of the society.

Acknowledgements I am grateful to Taiming Cao and Sandeep Datta for help with data collection and processing, as well as Catherine Meyers and Liu Ming for helpful discussions.

References

- Agovino T (2017) How much is a NYC taxi medallion worth these days? Available at <https://www.cbsnews.com/news/how-much-is-a-nyc-taxi-medallion-worth-these-days/>. April, 2017. Accessed 30 Oct 2017
- Barter P (2013) "Cars are parked 95% of the time." Let's check. Available at <http://www.reinventingparking.org/2013/02/cars-are-parked-95-of-time-lets-check.html>. February, 2013. Accessed 30 Oct 2017
- Bureau of Labor Statistics (2015) American time use survey. Available at <https://www.bls.gov/tus/home.htm>. Accessed 30 Oct 2017

- CBRE (2017) Hosts with multiple units – a key driver of Airbnb growth. Available at https://www.ahla.com/sites/default/files/CBRE_AirbnbStudy_2017.pdf. March, 2017. Accessed 30 Oct 2017
- Chan ND, Shaheen SA (2012) Ridesharing in North America: past, present, and future. *Transp Rev* 32(1):93–112
- Chen YF, Lu HF (2015) We-commerce: exploring factors influencing online group-buying intention in Taiwan from a conformity perspective. *Asian J Soc Psychol* 18(1):62–75
- Cunningham S (2015) Peer to peer lending sites: an exhaustive review. Available at <http://www.lendingmemo.com/p2p-lending-sites/>. May, 2015. Accessed 30 Oct 2017
- Denton J (2017) Two federal lawsuits could spell big trouble for Uber. Available at <https://psmag.com/news/two-federal-lawsuits-could-spell-big-trouble-for-uber>. April, 2017. Accessed 30 Oct 2017
- Du L, Hu M, Wu J (2017) Contingent stimulus in crowdfunding. Working paper, University of Toronto
- Etherington D (2017) Lyft raises \$600M at \$7.5B valuation. Available at <https://techcrunch.com/2017/04/11/lyft-raises-600m-at-7-5b-valuation/>. April, 2017. Accessed 30 Oct 2017
- Frankel RS (2017) Prosper personal loans: 2017 comprehensive review. Available at <http://www.bankrate.com/loans/personal-loans/prosper-personal-loans-comprehensive-review/>. May, 2017. Accessed 30 Oct 2017
- Furfaro D (2017) Uber drivers should be legal employees with benefits: judge. Available at <http://nypost.com/2017/06/13/uber-drivers-should-be-legal-employees-with-benefits-judge/>. June, 2017. Accessed 30 Oct 2017
- Geron T (2013) California becomes first state to regulate ridesharing services Lyft, Sidecar, UberX. Available at <https://www.forbes.com/sites/tomiogeron/2013/09/19/california-becomes-first-state-to-regulate-ridesharing-services-lyft-sidecar-uber/>. September, 2013. Accessed 30 Oct 2017
- Hagan S (2017) Uber takes majority of ground transport market for U.S. business travelers. Available at <https://www.bloomberg.com/news/articles/2017-01-26/uber-takes-majority-of-ground-transport-market-for-u-s-business-travelers>. January, 2017. Accessed 30 Oct 2017
- Hall JV, Krueger AB (2015) An analysis of the labor market for Uber's driver-partners in the United States. Working paper, Princeton University
- Hartmans A (2017) Airbnb now has more listings worldwide than the top five hotel brands combined. Available at <http://www.businessinsider.com/airbnb-total-worldwide-listings-2017-8>. Accessed 30 Oct 2017
- Investment Zen (2017) Prosper review. Available at <http://www.investmentzen.com/peer-to-peer-lending-for-investors/prosper>. Accessed 30 Oct 2017
- Jing X, Xie J (2011) Group buying: a new mechanism for selling through social interactions. *Manag Sci* 57(8):1354–1372
- Kerr D (2016) UK court rules Uber drivers are employees, not contractors. Available at <https://www.cnet.com/news/uber-uk-court-ruling-drivers-employees-not-contractors/>. October, 2016. Accessed 30 Oct 2017
- Kulwin N (2016) Did Airbnb purge potentially illegal NYC listings ahead of its own data dump? Available at <https://www.recode.net/2016/2/10/11587752/did-airbnb-purge-potentially-illegal-nyc-listings-ahead-of-its-own>. February, 2016. Accessed 30 Oct 2017
- Lawler R (2013) A day after cutting a deal with Lyft, California regulator reaches an agreement with Uber as well. Available at <https://techcrunch.com/2013/01/31/a-day-after-cutting-a-deal-with-lyft-california-regulator-reaches-an-agreement-with-uber-as-well/>. January, 2013. Accessed 30 Oct 2017
- Lazo L (2015) Uber turns 5, reaches 1 million drivers and 300 cities worldwide. Now what? June, 2015. https://www.washingtonpost.com/news/dr-gridlock/wp/2015/06/04/uber-turns-5-reaches-1-million-drivers-and-300-cities-worldwide-now-what/?utm_term=.8c3981bcffa5. Accessed 30 Oct 2017
- Lending Club (2017) Lending club statistics. Available at <https://www.lendingclub.com/info/statistics.action>. June, 2017. Accessed 30 Oct 2017

- Locklear M (2017) Lyft reaches one million rides per day but is still well behind Uber. Available at <https://www.engadget.com/2017/07/05/lyft-million-rides-per-day-well-behind-uber/>. July, 2017. Accessed 30 Oct 2017
- Lyons K (2015) Surveillance society: Uber's use of customers' data raises concerns. Available at <http://www.post-gazette.com/business/tech-news/2015/07/12/Surveillance-Society-Uber-privacy-policy-allows-more-location-tracking/stories/201507120071>. July, 2015. Accessed 30 Oct 2017
- Macfarlane A (2017) China's Uber worth \$50 billion after raising more cash. Available at <http://money.cnn.com/2017/04/27/technology/didi-valuation-50-billion-uber/index.html>. April, 2017. Accessed 30 Oct 2017
- Miller Z (2017) What is crowdfunding? The main categories of crowdfunding and how you can get involved. Available at <https://www.thebalance.com/a-guide-what-is-crowdfunding-985100>. February, 2017. Accessed 30 Oct 2017
- Millward S (2016) In just six months, China's Didi doubles daily rides to 20 million. Available at <https://www.techinasia.com/china-didi-chuxing-20-million-daily-rides>. October, 2016. Accessed 30 Oct 2017
- Ming L, Tunca TI (2017) Consumer equilibrium, demand effects, and efficiency in group buying. Working paper, University of Maryland
- Ming L, Tunca TI, Xu Y, Zhu W (2017) An empirical analysis of price formation, utilization and value creation in ride sharing services. Working paper, University of Maryland
- Morris DZ (2016) Today's cars are parked 95% of the time. Available at <http://fortune.com/2016/03/13/cars-parked-95-percent-of-time/>. March, 2016. Accessed 30 Oct 2017
- Mueffelmann K (2015) Uber's privacy woes should serve as a cautionary tale for all companies. Available at <https://www.wired.com/insights/2015/01/uber-privacy-woes-cautionary-tale/>. Accessed 30 Oct 2017
- Ovans A (2000) E-procurement at Schlumberger. *Harv Bus Rev* 78:21–22. May–June 2000
- Pei PPE, Simchi-Levi D, Tunca TI (2011) Sourcing flexibility, spot trading, and procurement contract structure. *Oper Res* 59(3):578–601
- Reiss D (May, 2016) How to invest in equity crowdfunding: recent changes in legislation allow everyone to invest in an entry-level entrepreneurship. Available at <https://money.usnews.com/investing/articles/2016-05-27/how-to-invest-in-equity-crowdfunding>. Accessed 30 Oct 2017
- Rossi A (2016) Why Airbnb removed over 2,000 NYC listings. Available at <https://www.smartertravel.com/2016/07/15/airbnb-removed-2000-nyc-listings/>. July, 2016. Accessed 30 Oct 2017
- Rubin J (2017) Lawsuit accuses Uber of ripping off drivers, paying them smaller fares than what passengers pay. Available at <http://www.latimes.com/local/lanow/la-me-uber-drivers-lawsuit-20170429-story.html>. April, 2017. Accessed 30 Oct 2017
- Russell J (2017) SoftBank is reportedly keen to buy 'multi-billion dollar stake' in Uber. Available at <https://techcrunch.com/2017/07/25/softbank-is-reportedly-keen-to-buy-multi-billion-dollar-stake-in-uber/>. July, 2017. Accessed 30 Oct 2017
- Singh S (2016) The state of P2P lending. Available at <https://techcrunch.com/2016/01/30/the-state-of-p2p-lending/>. January, 2016. Accessed 30 Oct 2017
- Statista (2016) Number of households in the U.S. from 1960 to 2016 (in millions). Available at <https://www.statista.com/statistics/183635/number-of-households-in-the-us/>. Accessed 30 Oct 2017
- Statista (2017a) Crowdfunding. Available at <https://www.statista.com/outlook/335/100/crowdfunding/worldwide#>. Accessed 30 Oct 2017
- Statista (2017b) Crowdfunding. Available at <https://www.statista.com/outlook/377/100/crowdfunding/worldwide#>. Accessed 30 Oct 2017
- Ting D (2017) Airbnb's growth is being driven by commercial operators, report says. Available at <https://skift.com/2017/03/10/airbnbs-growth-is-being-driven-by-commercial-operators-report-says/>. March, 2017. Accessed 30 Oct 2017

- Upwork (2016) New study finds freelance economy grew to 55 million Americans this year, 35% of total U.S. workforce. Available at <https://www.upwork.com/press/2016/10/06/freelancing-in-america-2016/>. Accessed 30 Oct 2017
- US Department of Transportation (2006) Advanced public transportation systems: the state of the art update 2006. Available at <http://www.globaltelematics.com/pitf/FTA-dynamicRideSharingReview.pdf>. Accessed 30 Oct 2017
- Weise E, della Cava M (2017) Lyft faces its big moment to leap ahead of Uber. Available at <https://www.usatoday.com/story/tech/news/2017/03/08/lyft-uber-travis-kalanick-john-zimmer/98799804/>. March 2017. Accessed 30 Oct 2017
- Yaraghi N, Ravi S (2017) The current and future state of the sharing economy. Brookings India IMPACT Series No. 032017. March, 2017
- Zhao N (2016) Where are the nation's second homes? 2014 data. Available at <http://eyeonhousing.org/2016/03/where-are-the-nations-second-homes-2014-data/>. March, 2016. Accessed 30 Oct 2017

Chapter 5

Ride Sharing



Siddhartha Banerjee and Ramesh Johari

Abstract Ridesharing platforms such as Didi, Lyft, Ola and Uber are increasingly important components of the transportation infrastructure. However, our understanding of their design and operations, and their effect on society at large, is not yet well understood. From an academic perspective, these platforms present challenges in large-scale learning, real-time stochastic control, and market design. Their popularity has led to a growing body of academic work across several disciplines, with researchers addressing similar questions with vastly different tools and models. Our aim in this chapter is to outline the main challenges in ridesharing, and to present an approach to modeling, optimizing, and reasoning about such platforms. We describe how rigorous analysis has been used with great success in designing efficient algorithms for real-time decision making, in informing the market design aspects of these platforms, and in understanding the impact of these platforms in their larger societal context.

5.1 Introduction

Since their founding over the last decade, ridesharing platforms have experienced extraordinary growth. At their core, these platforms reduce the friction in matching and dispatch for transportation. They do so based on a pair of matched driver and passenger mobile apps; a typical transaction starting with a potential passenger opening her app and requesting a ride, following which a centralized dispatcher matches her to a nearby driver if one is available. However, underlying this simple model are three features which fuel much of the success of these firms:

S. Banerjee (✉)
Cornell University, Ithaca, NY, USA
e-mail: sbanerjee@cornell.edu

R. Johari
Stanford University, Stanford, CA, USA
e-mail: ramesh.johari@stanford.edu

1. *Data Collection and Analytics*: The driver and passenger apps enable extremely high temporal and spatial resolution for data collection. Ridesharing platforms track the position of all drivers and passengers in the system. Moreover, modern graph analytics and predictive models allow the platform to leverage this data to obtain very good estimates of travel-times, instantaneous demand, and long-term driver and passenger engagement metrics.
2. *Real-time Operations and Control*: An important reason behind the success of ridesharing platforms is their reliability and lack of friction in requesting a ride. Critical to this is the ability of the platform to rebalance demand and supply over time and space. A key tool for this purpose is *dynamic pricing*: ridesharing platforms adjust prices in real-time; in addition, however, the platforms have several other real-time dispatch and rebalancing tools, as well as different means for regulating and/or pooling instantaneous demand.
3. *Market design*: Ridesharing platforms typically do not employ drivers, but rather, create a marketplace between passengers and freelance drivers. Drivers can choose when and where to work (or not), and earn a share of the earnings per ride. To deal with this uncertainty in supply, ridesharing platforms need to understand the longer term equilibrium impacts of their real-time control on driver and passenger decisions, as well as on the platform's overall performance.

The challenge in studying ridesharing platforms is that the above features interact with each other in fundamental ways. To analyze any particular aspect of the platform, one has to account for its effect on the others – for example, to test a new pricing strategy, the platform must control for short timescale spatio-temporal variations in demand and supply, as well as long-term effect on driver and passenger entry decisions. On the other hand, incorporating all aspects simultaneously may lead to intractable models.

Our main aim in this chapter is to outline a stochastic-network based micro-foundation for ridesharing platforms. The framework we present is adapted primarily from our prior work on these questions in Banerjee et al. (2015, 2017), which in turn were based on our experience in working on the design of pricing and matching algorithms at Lyft.¹ The popularity of ridesharing platforms has in recent years led to a growing body of work by researchers across several disciplines, including applied probability, optimization and network algorithms, economics and market design, transportation and urban planning, and even statistical physics. These works address similar problems, but using vastly differing models and tools, and this makes it difficult to translate the findings across different fields. While we do not in any way claim that the framework we present herein is the only way to model such platforms, we do believe that it captures the salient features of ridesharing, while being amenable to analysis and simulation. Our hope is that having such a common modeling framework will help unify the insights of researchers working in this exciting area.

¹SB worked there in 2014–2015, and was involved in designing their early pricing algorithms. www.lyft.com

The rest of the chapter is organized as follows. First, in Sect. 5.2, we provide a high-level description of the essential features of a ridesharing platform, and outline the different operational and market design challenges facing the platform. Next, in Sect. 5.3, building on our prior work in Banerjee et al. (2015, 2017), we describe how queueing-network models can be adapted to study ridesharing platforms. In particular, we focus on how they capture the critical features of such platforms: the high-resolution state description and two-sided nature, the real-time pricing and control tools, and the longer-term strategic interactions of drivers, passengers and the platform. In Sect. 5.4, we briefly summarize the operational and market design insights that can be derived from this modeling framework; in particular, we focus on how the model has been used to develop efficient control algorithms (based on results from Banerjee et al. 2017) and market mechanisms (following ideas outlined in Banerjee et al. 2015). Finally, in Sect. 5.5, we survey some of the related literature in ridesharing, and more generally, on control of stochastic networks and two-sided marketplaces.

5.2 Anatomy of a Modern Ridesharing Platform

In this section we describe the basic anatomy of a ridesharing platform. We divide our presentation in three parts: first, we discuss a fundamental separation of *timescales* that should guide any modeling of a ridesharing platform. Next, we discuss the *strategic choices* that guide the behavior of drivers and passengers on the platform. Finally, we discuss *operation and design* of the platform itself, taking both the timescale separation and incentives into account.

5.2.1 Timescales

There is an intrinsic timescale separation in the strategic interaction of drivers and passengers, as well as operation of the platform. In particular, ridesharing platforms have distinct behaviors on the following two timescales:

- (i) A fast timescale (roughly, intra-minute), which captures the instantaneous dynamics of cars and passengers in the network; and
- (ii) A slow timescale (roughly, intra-week), over which drivers make decisions as to how much and when to be on the platform.

The fast timescale provides the backdrop for the short-term operational and market design choices of the platform (especially pricing and matching), while the strategic consequences of these choices unfold over the slower timescale. While passengers primarily make entry decisions on the fast timescale (“Do I want to take this ride, given the current price and availability?”), drivers make such decisions on a longer timescale. This separation of the agent dynamics thus provides

a convenient separation between how the platform’s policies affect drivers and passengers: control policies influence instantaneous passenger-vehicle dynamics, while the aggregate effect of these policies affect the longer-term entry decisions of drivers. Both timescales are discussed in more detail in Sect. 5.2.2. We note that this viewpoint ignores other dynamics: for example, intra-hour changes in demand rates, or intra-year interactions between ridesharing firms and public transit providers. However, we argue that the two timescales we consider are crucial for understanding the first-order behavior of ridesharing platforms.

5.2.2 *Strategic Choices*

What are the main strategic choices made by participants in ridesharing platforms? We already alluded to the strategic modeling of one side of the platform: for the most part, passengers can be modeled by assuming they make an instantaneous decision of whether to participate based on price (and possibility also availability information) on the platform. Platforms refer to “app-opens” as opportunities to potentially engage a passenger; a subsequent “ride request” refers to the passenger actually choosing to request a ride. In what follows, we will typically assume that passengers choose to request a ride as long as the price of the ride is below a private reservation value.

Drivers exhibit far more complex strategic behavior. While there is some evidence that drivers will locally optimize on short timescales (e.g., perhaps moving to a nearby block if there is evidence that prices are higher there), for the most part it is reasonable to assume that drivers are relatively *inelastic* on short timescales, as noted above.

Instead, the key choice made by drivers on a longer timescale is *entry* – both where they choose to drive, as well as what days and times during the week they choose to do so. Drivers make these decisions in response to what they observe on shorter timescales, forming expectations based on their experiences while driving. These entry decisions can be quite sophisticated, reflecting spatio-temporal differences in the driver’s experience within the platform.

The incentive structure of platforms can be quite complex, in ways that we do not necessarily capture in the models discussed in this chapter. For example, both sides *rate* each other after a ride is complete; these rating systems play a key role in determining, for drivers in particular, whether they are allowed to stay in the platform. As another example, experienced drivers are relatively sophisticated about time-of-day effects (i.e., when demand is expected to be higher or lower), and they will choose to keep their app online or offline accordingly. Platforms can also provide longer-term incentives to drivers, particularly through the *fee* structure (i.e., what percentage is given to drivers as their pay); we do not study the optimal design of fees. Finally, drivers are also constantly making choices about how and whether to *multihome* – i.e., participate in multiple platforms at once. Multihoming has important consequences for the availability of drivers in each platform, and warrants further attention from academic researchers studying ridesharing.

5.2.3 Operation and Market Design

A platform's operation and market design can be roughly summarized by three pieces: *information revelation*, *pricing*, and *dispatch*.

First, platforms reveal information to both passengers and drivers about the state of the system. For passengers, this is in the form of ETA (estimated time of arrival), which captures the distribution of drivers locally near a passenger. For drivers, the platform provides information on the distribution of demand. There is also extensive information collected and visible on the ratings provided to each side of the platform.

Second, a crucial aspect of these platforms is that they constantly make choices about the fare that will be charged to passengers. The typical model is that platforms publish a *base fare* schedule, that determines a baseline rate for any trip. Next, they will modify this base fare by a *multiplier* that adjusts the base fare to account for local demand-supply imbalances: when supply is scarce, the multiplier increases. In the past, platforms would not publish a fare estimate, because passengers were not asked to enter a destination at the time of the ride request; in these settings, the platform simply displayed the price multiplier to the passenger at time of ride request. Now, platforms publish the expected fare for a ride to the passenger at the time of ride request, in response to solicitation of the destination. These fares incorporate the price multiplier.

Once ride requests are made, platforms must actually match drivers to passengers; this is *dispatch*. Platforms typically match passengers to their *closest* driver. A more recent development in ridesharing is the introduction of *carpooling* (UberPool, Lyft Line); these products match *multiple passengers* with a single driver. In addition, platforms are becoming more sophisticated in how they manage the dispatch problem; for example, while in the past occupied drivers were not considered in the dispatch problem, now platforms will anticipate the fact that a driver will free up before making the next match. Such policies reduce driver idle times.

We have only provided a brief overview of the operational aspects of these platforms, focusing on the elements that are most important for our models below. Of course, in reality there is a great deal more complexity. For example, platforms must work to develop a product interface that allows passengers and drivers to make good choices; they must develop marketing mechanisms and on-boarding mechanisms to attract driver supply; and they are working more and more to provide sophisticated long-term incentives to drivers, as noted above. These topics are important dimensions of the platforms, and may provide fruitful avenues for future study.

5.3 A Modeling Framework for Ridesharing Platforms

In this section, we outline a formal stochastic model of a ridesharing platform, and formulate the various associated control and market design problems. The basic framework we introduce below is adapted from the models proposed in Banerjee et al. (2015, 2017). It provides a rich modeling framework for ridesharing platforms,

allowing us to study many different features, controls and metrics. Moreover, despite its complexity, the model turns out to be surprisingly amenable to analysis. In subsequent sections, we describe how the framework can be used to study control policies for the system dynamics, market-design questions for the driver-passenger marketplace and inter-platform interactions. Moreover, we also describe how the framework can be extended to study other design aspects of such platforms.

As we discuss before, a key to modeling ridesharing platforms is identifying the appropriate timescales for different agent interactions. To this end, our framework combines a Markov chain model with time-invariant parameters for capturing the instantaneous dynamics of vehicles and passengers (the *fast-timescale*), with an equilibrium analysis that captures entry decisions of drivers and passengers as well as the objectives of the platform, based on the average system performance (the *slow-timescale*). Although the model below allows for modeling fairly complex agent behavior, we focus on a particular behavioral model, wherein we assume that passengers primarily react under the fast-timescale (i.e., to instantaneous vehicle availability and prices), while drivers react under the slow-timescale (i.e., based on long-term average earnings). This is a choice we make based on a combination of our experience on working on these platforms, as well as for pedagogical reasons: as described in Sect. 5.2, these interactions capture the first-order behavior of these platforms, while enabling a tractable analysis of the stochastic platform dynamics and long-term strategic interactions.

That said, we note that our modeling choices ignore four important dynamics: (i) short-term fluctuations in system parameters (e.g., changing demand or bursty arrivals), (ii) short-term strategic behavior of drivers (e.g., strategic repositioning and ride cancellations), (iii) long-term effects on passenger behavior (e.g., demand screening due to persistent low availability or high prices), and (iv) competitive interactions between platforms (e.g., price cuts, driver retention incentives). Understanding the impact these interactions is important, but are to some extent secondary to the questions we consider. Thus, though some of our results, as well as a growing body of work by others, apply to these questions, we choose not to dwell on them in this chapter.

5.3.1 Modeling Stochastic Dynamics of the Platform

We now define a stochastic model for the fast-timescale dynamics of a ridesharing system. The main elements of the model are summarized in Fig. 5.1.

State space and Markovian dynamics We model the fast-timescale dynamics of a ridesharing platform using a *stochastic processing network* framework (Kelly and Yudovina 2014): We consider a partition of a city into a set of n stations (corresponding to locations or neighborhoods in a city), and use a continuous-time Markov chain to track the positions of k units (i.e., vehicles) which are either idle at these stations or transiting (i.e., in a ride) between them. Each ride involves a driver picking up a passenger in one region, and dropping her off in another. These

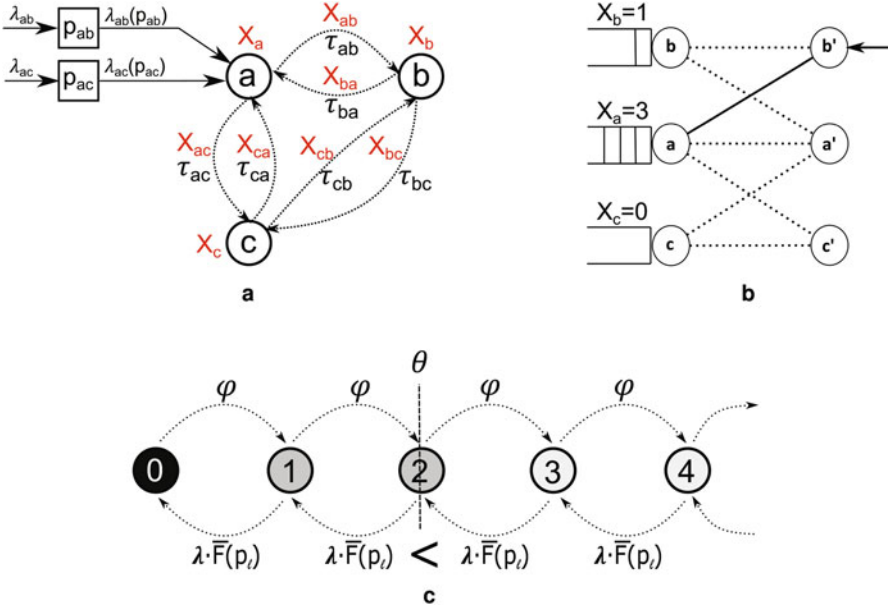


Fig. 5.1 Illustrating the stochastic dynamics of a ridesharing platform. (a) Summary of model parameters, (b) The dispatch graph, (c) Idle driver queue under state-dependent pricing. Figure a summarizes the primary components of our model for a ridesharing platform (in this case, with state-independent pricing). The platform depicted has 4 stations $V = \{a, b, c, d\}$, and m vehicles. Random variables are depicted in red; here, the random process $\{X_v\}$ tracks the number of idle units at stations, while $\{X_{ij}\}$ tracks the units in transit between stations i and j (with mean travel-time τ_{ij}). Zooming into station a , we see that passengers with destination $\{b, c\}$ arrive at a according to Poisson processes with rate $\{\lambda_{ab}, \lambda_{ac}\}$; these arrivals are then ‘thinned’ to $\lambda_{av}(p_{av})$ by setting (state-independent, but destination dependent) prices $\{p_{ab}, p_{ac}\}$. (Adapted from Banerjee et al. 2017). In figure b, we depict the bipartite graph for the dispatch problem, for the network in figure a under the assumption that station pairs (a, b) and (a, c) are close enough to use each other’s supply. An arriving passenger at station b can thus be matched to a vehicle at either station a or b – in the figure, we choose to match the arrival to a vehicle at station a . (Adapted from Banerjee et al. 2018). Figure c shows the birth-death chain for the number of idle drivers in a single station, under local state-dependent pricing policies. The arrival rate ϕ of vehicles to the station is determined by the overall network. The rate of departures (i.e., matched rides), however, is modulated by the pricing policy. Here, we have depicted a base arrival rate of passengers λ , and a simple *single-threshold* local pricing policy, where the platform uses a ‘base’ fare p_ℓ when the number of drivers is greater than a threshold θ , else charges a ‘primetime’ price $p_h > p_\ell$ (hence the queue drains slower when there are $\leq \theta$ drivers). (Adapted from Banerjee et al. 2015)

rides can be modulated via a set of controls – primarily pricing, but also dispatch and empty-unit rebalancing. To study the efficacy of different controls, we analyze the long-term average performance of the system; in other words, we study various metrics of interest in *steady state*. We henceforth use $V = \{1, 2, \dots, n\}$ to denote the set of stations, and $E = \{(i, j) \in V \times V\}$ to be the set of source-destination pairs.

Formally, at any time $t \geq 0$, the state of the ridesharing system is denoted as $\mathbf{X}(t) = \{\{X_i(t)\}_{i \in V}, \{X_{ij}(t)\}_{ij \in E}\}$, where $X_v(t)$ denotes the number of units which are idle at station i , and $X_{ij}(t)$ denotes the number of units in transit (either in a ride, or rebalancing) between stations i and j . As per our assumption, the sum over all states must at all times add up to k . Denoting $N = n + \binom{n}{2}$ to be the dimension of \mathbf{X} , we have that the state space of the Markov chain is given by $\mathcal{S}_{N,k} = \{\{x_i, x_{ij}\} \in \mathbb{N}_0^N \mid \sum_{i \in V} x_i + \sum_{ij \in E} x_{ij} = k\}$. Note that the state-space is finite. Since our focus is on the long-run average performance, i.e., under the steady state of the Markov chain, for ease of notation, we henceforth suppress the dependence on time t .

Passenger arrivals and ride requests Potential passengers who desire to travel between stations i and j (henceforth, type- ij passengers) arrive at station i following a stochastic process with average rate λ_{ij} (alternately, average inter-arrival time $1/\lambda_{ij}$). It is reasonable to assume the inter-arrival times to be independent, and hence, invoking the Palm-Khinchine theorem (cf. Chapter 14 in Kallenberg 2006), we assume that type- ij passengers arrive according to a Poisson process of rate λ_{ij} .

To model the ‘willingness-to-pay’ of the passengers, we assume that each type- ij passenger has a ride value drawn independently from a distribution $F_{ij}(\cdot)$. Upon arrival at i , a customer is quoted a *point-to-point* price p_{ij} (which may potentially depend on the current state $\mathbf{X}(t)$); she then requests a ride if her value exceeds this price, i.e. with probability $1 - F_{ij}(p_{ij})$. At this point, if at least one unit is available at station i (or more generally, at any sufficiently ‘nearby’ station), then she is matched to it. If on the other hand she is unwilling to pay the price, or is not matched to a vehicle, then she leaves the system *immediately*. We assume that F_{ij} has a density and that all values are positive with some probability, i.e. $F_{ij}(0) < 1$.

The passenger dynamics outlined above is referred to in the stochastic modeling community as a *loss-system* model. The possibility of immediate departure without a ride request captures the possibility that a passenger typically has outside options (walking/public transit/other ridesharing firms) which she turns to if the platform proves unattractive at the moment. In practice, with dynamic pricing, some passengers may tend to wait for a while to see if prices change (although ridesharing platforms may also freeze their prices for a given passenger, while still adjusting them for others), or cars become available. Such heterogeneity in passenger impatience can be accommodated in our model to some extent (for example, as a *negative* queue, akin to lost-sales models in inventory systems). However, the effects of such behavior is not well understood in practice, and currently, most ridesharing firms do not specifically account for passenger impatience in their policies.

Vehicle travel times Once a unit is dispatched to serve a passenger, it then needs to go pick up and drive the passenger to her destination station. We use the state variable $X_{ij}(t)$ to track the number of units in transition between stations i and j . When a customer engages a unit to travel from i to j , the state changes to $\mathbf{X} - e_i + e_j$

(i.e., $X_i \rightarrow X_i - 1$ and $X_{ij} \rightarrow X_{ij} + 1$). The unit remains in transit for a random time, drawn independently from some general distribution $G_{ij}(\cdot)$ with mean τ_{ij} . Upon reaching its destination, the unit drops off the passenger, and the system state changes to $\mathbf{X} - e_{ij} + e_j$.

For convenience, we will assume henceforth that transit times are exponentially distributed, with average transit-times τ_{ij} . This is primarily for ease of exposition, as it allows us to keep track of only the number of vehicles in transit between any two stations (as opposed to their exact time of arrival; this follows from the memoryless property of the exponential distribution). We note though that the predictions of the model remain essentially unchanged for any general (independent) travel time with mean τ_{ij} . We also assume that the demand characteristics and ride rewards are independent of the actual transit times (dependence on average transit times τ_{ij} can be embedded in the model parameters). Finally, note that we do not model stochastic correlation in travel times (e.g., that might arise because trips share a common road network) – a potentially interesting direction for future work.

We conclude with two additional observations about transit-times. First, we note that though the above discussion is primarily for vehicles dropping-off passengers, the transit times also apply to settings where empty vehicles move between stations to improve the demand-supply balance. Second, in many cases, the model and results greatly simplify if we assume that transit times are identically zero: in particular, note that in such a setting, we only need to keep track of idle vehicles at the stations. Introducing transit times tends to complicate analysis as it may lead to situations where availability is low as almost all vehicles are in transit (this corresponds to the so-called *heavy-traffic* regime in queueing models). Understanding the significance of transit times are in designing ridesharing policies is an under-explored question, and one which we will not deal with in this chapter in any significant detail.

5.3.2 Platform Controls

We now consider three primary ways in which the platform can intervene to affect the fast-timescale dynamics: (i) demand modulation via *pricing*, (ii) demand redirection via *dispatch*, and (iii) supply redirection via *empty-vehicle rebalancing*. We describe these in details below; note however that all these different controls are essentially linear transformations of the demand and supply flows, and moreover, can be combined together (and often are in practice).

Demand modulation (pricing) By adjusting the price p_{ij} for a ride from i to j , the platform can modulate the rate at which such rides are requested. To understand the effect of such a price, it is useful to define the *inverse demand* (or *quantile*)

function $q_{ij} = 1 - F_{ij}(p_{ij})$.² Now, for a fixed pricing policy \mathbf{p} with corresponding quantiles \mathbf{q} , the *effective demand rate* from i to j (i.e. type- ij passengers with value exceeding p_{ij}) follows a Poisson process with rate $\lambda_{ij}q_{ij}$ – this follows from the probabilistic thinning property of a Poisson process.

The most general model for pricing in the above model is that of *global state-dependent* prices, where the platform selects $p_{ij}(t)$ at time t as a *function of the overall state* $\mathbf{X}(t)$ – this induces a state-dependent Poisson process of type- ij passengers with rate $\lambda_{ij}q_{ij}(\mathbf{X})$. A natural relaxation of this is that of *local state-dependent* prices, where p_{ij} is a function of the *local state* $X_i(t)$ at the source. Finally, in *state-independent* pricing, p_{ij} is set to be independent of the instantaneous system state. The three pricing schemes decrease in complexity, and moreover, require decreasing levels of system engineering to enable – understanding their comparative behavior is thus of great importance. Banerjee et al. (2017) study the relation between global state-dependent and state-independent prices, while Banerjee et al. (2015) focus on local state-independent prices to understand the value of dynamic pricing in ridesharing platforms.

Demand redirection (dispatch) Though it is typically infeasible to redirect passengers to nearby stations (although this has been experimented with by Lyft and Uber in some markets), what is often possible is to match an incoming ride request station i to units which are idle at “nearby” stations. This is based on the underlying assumption that passengers are insensitive to small delays in pickup as compared to pickup time of the nearest unit. This is not strictly true in practice, as passengers are known to be sensitive to the pickup time (ETA); however, it is a convenient abstraction for our model, and moreover, can be refined by incorporating probabilistic ride cancellations due to longer pickup times. Moreover, longer dispatches may affect drivers, and this can be modeled by a cost for each possible dispatch decision.

To formally define a dispatch policy, we define a *compatibility graph* $G = (V, E)$ on the set of stations, with edges between pairs of stations that are near enough such that a passenger arriving at one can be served using a unit from the other (see Fig. 5.1b for an example). As with pricing, we can define a state-dependent dispatch policy $\mu(\mathbf{X})$ which, for each ride requested at station i , decides from which station in $\{i\} \cup \{j : (i, j) \in E\}$, the customer is served. Such a dispatch policy now induces a rate $f_{ij}(\mu)$ of customers arriving at i that travel to j using a unit from k , and a rate $z_{ik}(\mu)$ of customers arriving at i who are matched to a unit at k . Note that if μ is chosen in a state-independent manner (wherein a request is randomly routed to a neighboring station), then it may lead to a failed dispatch despite there being idle units; such policies however are more tractable to analyze, and hence have been considered in Ozkan and Ward (2016) and Banerjee et al. (2017). More recently,

²It is convenient to assume that the density of F_{ij} is positive everywhere in its domain, implying that there is a 1-1 mapping between prices and quantiles; this allows us to write $p_{ij} = F_{ij}^{-1}(1 - q_{ij})$. We note however that this is not necessary for the results we present.

state-dependent dispatch policies were analyzed by Banerjee et al. (2018) (albeit without costs for dispatch from neighboring stations).

Supply redirection (rebalancing) This is a catch-all for any control policy which allows the platform to affect the position of a unit at the end of a ride, i.e., whether the unit remains at the destination station or moves to another station without a passenger. Such a control is sometimes referred to as *empty-car rebalancing*, and such rebalancing typically is modeled as incurring a cost (for vehicle miles traveled/idle time of drivers). In practice, rebalancing is less common in current ridesharing systems (in comparison to bikesharing/carsharing systems), as drivers themselves choose where to go when idle – however, platforms do try to influence these decisions via information displays, incentive schemes, etc. With the potential introduction of autonomous vehicles, such control may become more prevalent.

We can model a rebalancing policy as a state-dependent control $\mathbf{r}(\mathbf{X})$ which, for each trip ending at a station i , redirects the unit to some station j (which could be i). This results in an increase in state X_{ij} , and has associated cost c_{ij} . Since redirection is costly for drivers, it is natural to assume that redirected units arriving at a station are not redirected again. Details of how to incorporate this in the above model are provided in Braverman et al. (2016) and Banerjee et al. (2017).

5.3.3 Platform Objectives

Given the above system dynamics (with *fixed* parameters $k, \lambda_{ij}, F_{ij}, \tau_{ij}$), our aim is to study the long-term average performance of various platform metrics. More precisely, we want to design control policies to maximize relevant performance metrics under the *stationary distribution* $\pi(\mathbf{x})$ of the Markov chain induced by our controls. Note that for given n, k and under any policy, the resulting Markov chain is finite-state (since the number of stations and units is fixed); furthermore, it is irreducible under weak assumptions on the prices and the demand (see Banerjee et al. 2017 for details). Now, using basic Markov chain theory, we have that our system has a unique steady-state distribution $\pi(\cdot)$ with $\pi(\mathbf{x}) \geq 0 \forall \mathbf{x} \in \mathcal{S}_{N,k}$ and $\sum_{\mathbf{x} \in \mathcal{S}_{N,k}} \pi(\mathbf{x}) = 1$.

Following Banerjee et al. (2017), we consider objective functions that decompose into per-ride reward functions $I_{ij}(p)$, which correspond to the reward obtained from a passenger riding between stations i and j at price p . In particular, such a structure admits three canonical objectives:

- *Volume of Trade or Throughput*: the total rate of rides in the system (setting $I_{ij}(p) = 1$).
- *Social welfare*: the contribution to social welfare from each $i \rightarrow j$ ride is given by $I_{ij}(p) = \mathbb{E}_{V \sim F_{ij}}[V \mid V \geq p]$.
- *Revenue*: to find the platform's revenue rate (assuming it keeps a fraction $1 - \gamma$ of the earnings), we set $I_{ij}(p) = (1 - \gamma) \cdot p$.

To formally define the objective, we focus on the case of pricing. Now, for a given objective $I_{ij}(\cdot)$, the aim of the platform is to select prices $\mathbf{p}(\mathbf{X})$ (with corresponding quantiles $\mathbf{q}(\mathbf{X})$) that maximizes the rate of reward accumulation under the stationary distribution. This can be written as:

$$\text{OBJ}_m(\mathbf{p}) = \sum_{\mathbf{x} \in \mathcal{S}_{N,k}} \pi(\mathbf{x}) \cdot \left(\sum_{i,j} \lambda_{ij} q_{ij}(\mathbf{x}) I_{ij}(p_{ij}(\mathbf{x})) \right), \quad (5.1)$$

where $\pi(\mathbf{x})$ is the stationary distribution of the Markov chain under pricing policy \mathbf{p} . Equation 5.1 can be understood as follows: at any station i , customers destined for j arrive via a Poisson process with rate λ_{ij} , and find the system in state $\mathbf{x} \in \mathcal{S}_{N,k}$ with probability $\pi(\mathbf{x})$ (this follows from the ‘‘Poisson Averages See Time Averages’’ or PASTA property; see Kelly and Yudovina 2014, for details). They are then quoted a price $p_{ij}(\mathbf{x})$, and engage a ride with probability $q_{ij}(\mathbf{x}) = 1 - F_{ij}(p_{ij}(\mathbf{x}))$. The resulting ride contributes $I_{ij}(q_{ij}(\mathbf{x}))$ to the expected objective. Recall that unavailability of units is captured by our assumption that $q_{ij}(\mathbf{x}) = 0$ whenever $x_i = 0$.

Though the above equation is most naturally written in terms of prices, it turns out to be non-concave even for a single station. However, a standard price-theoretic trick (for example, see Hartline 2013) in such cases is to instead write the objective in terms of quantiles, whereupon it turns out to be concave for most cases of interest. In particular, abusing notation to define $I_{ij}(q) := I_{ij}(F_{ij}^{-1}(1 - q))$, and defining *reward curves* $R_{ij}(q) := q \cdot I_{ij}(q)$, it can be shown that $R_{ij}(q)$ are concave in q for throughput and welfare under any distribution, and for revenue under *regular* distributions (a wide class of distributions which includes all increasing hazard-rate distributions; see Banerjee et al. 2017; Hartline 2013, for details).

However, the convexity of $R_{ij}(q)$ does not imply that the optimization problem in Eq. 5.1 admits a tractable solution, as we still have to determine the average under the stationary distribution. This involves solving a fixed-point constraint, which in general can be non-tractable. In fact, Banerjee et al. (2017) provide an example which shows that the problem is non-concave even for a setting with 3 stations and a single unit!

5.3.4 Local Controls and Closed Queueing Models

Although the stochastic dynamics described above is complex, it is still amenable to study via simulation. Moreover, in some special cases, its analysis can be greatly simplified using classical results from queueing theory (Serfoso 1999; Kelly 1979). In particular, a critical tool used in Banerjee et al. (2015, 2017) is the fact that *under state independent control policies (pricing, dispatch, rebalancing), as well as under local state-dependent pricing, the stationary distribution of the resulting Markov chains is known in closed form*. This now allows us to study the design of control policies in an analytic way. We now briefly provide some background behind this methodology.

The general Markov chain described in the previous sections (involving a fixed number k of units, located in one of N queues) is a special example of a *closed queueing network* (see Kelly 1979; Serfozo 1999). “Closed” here refers to the fact that the number of units remains constant; in contrast, in open networks, units may arrive and depart from the system. These networks are well-studied in applied probability, and in general may have complex stationary distributions. However, a critical property uniting the settings mentioned above (state-independent controls, local state-dependent pricing) is that the resulting Markov chain in each case is *quasireversible* (Kelly 1979). This is a particular structural property of Markov chains which generalizes the notion of reversibility. The exact definition is somewhat technical, but for our purposes, the crucial fact is that *quasi-reversibility is sufficient to ensure the stationary distribution is product form*, i.e.

$$\pi(\mathbf{x}) = \frac{1}{Z} \prod_i f_i(x_i) \prod_{ij} f_{ij}(x_{ij}),$$

where Z is the appropriate normalizing constant. The exact form of the local potentials f_i , f_{ij} depend on the precise nature of the system and controls; see Banerjee et al. (2015, 2017) for details. For illustration purposes, we develop this in more detail below for the special case of state-independent pricing and instantaneous transfers.

An important property of state-independent control prices is that the rate of units departing from any station i at any time t when $X_i(t) > 0$ is a constant, independent of the state of the network. The resulting model is a special case of a closed queueing model proposed by Gordon and Newell (1967).

Definition 1 A *Gordon-Newell network* is a continuous-time Markov chain on states $\mathbf{x} \in \mathcal{S}_{N,k}$, in which for any state \mathbf{x} and any $i, j \in [n]$, the chain transitions from \mathbf{x} to $\mathbf{x} - e_i + e_j$ at a rate $\mu_i r_{ij} \mathbf{1}_{\{x_i(t) > 0\}}$, where $\mu_i > 0$ is referred to as the *service rate* at station i , and $r_{ij} \geq 0$ are the *routing probabilities* that satisfy $\sum_j r_{ij} = 1$.

In other words, if units are present at a station i in state \mathbf{x} , then departures from that station occur according to a Poisson distribution with rate $\mu_i > 0$; conditioning on a departure, the destination j is chosen according to state-independent routing probabilities r_{ij} . For this network, the resulting steady-state distribution $\{\pi_{\mathbf{p}, m}(\mathbf{x})\}_{\mathbf{x} \in \mathcal{S}_{N,k}}$ was established to be product form via the celebrated Gordon-Newell theorem.

Theorem 1 (Gordon and Newell 1967) Consider a k -unit n -station Gordon-Newell network with transition rates μ_i and routing probabilities r_{ij} . Let $\{w_i\}_{i \in [n]}$ denote the invariant distribution associated with the routing probability matrix $\{r_{ij}\}_{i, j \in [n]}$, and define the traffic intensity at station i as $\rho_i = w_i / \sum_j r_{ij}$. Then the stationary distribution is given by:

$$\pi(\mathbf{x}) = \frac{1}{G_m} \prod_{j=1}^n (\rho_j)^{x_j}, \quad (5.2)$$

with normalization constant $G_m = \sum_{\mathbf{x} \in \mathcal{S}_{k,m}} \prod_{j=1}^n (\rho_j)^{x_j}$.

To see that the Markovian dynamics resulting from state-independent pricing policies fulfill the conditions of Gordon-Newell networks, observe that fixing a price p_{ij} (with corresponding q_{ij}) results in a Poisson process with rate $\lambda_{ij}q_{ij}$ of arriving customers *willing to pay price* p_{ij} . These customers engage a unit only if one is available, otherwise they leave the system. Thus, given quantiles \mathbf{q} , the time to a departure from station i is distributed exponentially with rate $\mu_i = \sum_j \lambda_{ij}q_{ij}$ when $X_i > 0$ and with rate 0 otherwise. Further, conditioned on an arriving customer having value at least equal to the quoted price, the probability that the customer's destination is j , is $r_{ij} = \lambda_{ij}q_{ij} / \sum_k \lambda_{ik}q_{ik}$, independent of system state. Now we can use the Gordon-Newell theorem to simplify the objective function in Eq. 5.1 to get an explicit function of the quantiles \mathbf{q} . The functions obtained are somewhat involved, and hence we omit them here; interested readers should refer Banerjee et al. (2017) for details.

5.3.5 Modeling Endogenous Entry of Drivers

Finally, we turn our attention to agent behavior in the slow-timescale – in particular, we discuss how the above model can be used to model the endogenous entry decisions made by drivers.

At a high level, the slow timescale allows us to capture the marketplace aspect of ridesharing platforms, by allowing us to specify the strategic aspects of agent-platform interactions. In particular, as we mention before, our primary use of the slow timescale is to model the endogenous entry decision of drivers, which thereby determines the number of vehicles k in equilibrium. We note though that a similar idea of determining the parameters of the fast-timescale model based on strategic considerations at a slower timescale can be used to model other agent interactions – in particular, an interesting open question is to model the effect of high ETAs or prices on passenger rates.

Our treatment here follows the model in Banerjee et al. (2015). The main assumption for the slow-timescale driver decisions is that each potential driver in the pool has a *reservation earning-rate* (or earning-rate target), and makes an *endogenous entry decision* (i.e., determines whether or not to work on the platform) by comparing their expected earning-rate on the platform to this target. We assume that the earning-rate target for each driver is drawn i.i.d. from some distribution G_d . On the other hand, the earning rate of a driver on the platform depends on the specific wage structure implemented by the platform – this is an aspect which different platforms have experimented with, and one whose effects are not yet well understood. For this chapter, as in Banerjee et al. (2015) (and as above), we assume that the driver gets a fraction γ of the price of each ride that he is matched to. Note that this is the policy currently followed by most ridesharing firms.

For convenience, we henceforth analyze the behavior of a single station under local state-dependent pricing; this can however be extended to the entire network using standard queueing theoretic tools. The setting is depicted in Fig. 5.1c. Let

X represent the instantaneous number of idle units at the station, and suppose the potential pool of drivers is of size \bar{k} . Given pricing policy $p(X)$, the passenger arrival rate is $\lambda \bar{F}(p(X_i))$ (where $\bar{F}(\cdot) = 1 - F(\cdot)$); the equilibrium number of drivers k is now determined by the system performance in the fast timescale (which in turn depends on k). Formally, the equilibrium rates of passenger requests and number of drivers must satisfy:

$$\lambda(X) = \lambda \bar{F}(p(X)), \quad k = \bar{k} G_d \left(\frac{\eta}{\iota + \tau} \right), \quad (5.3)$$

where η denotes the expected per-ride earnings, ι the expected waiting time for a driver between rides, and τ the expected ride time. Exact expressions for these can be computed using the product form characterization of the stationary distribution discussed in Sect. 5.3.4. Note though that η and ι depend on λ and μ , as well as the pricing policy $p(X)$. Details of these computations, and of the existence/uniqueness of the equilibrium, are given in Banerjee et al. (2015).

5.4 Analyzing the Model: Key Findings

We now briefly describe some of the insights that can we obtain by analyzing the queueing-theoretic model described laid out in the preceding sections. First, we summarize the results from Banerjee et al. (2017), where for any given number of units K , the authors show how we can design control policies for the fast-timescale dynamics, with strong performance guarantees. On the positive side, these policies surprisingly turn out to be state-independent. On the negative side, however, the results do not give insight into the number of units K that emerge in equilibrium. Moreover, the results critically depend on having full knowledge of the Markov chain parameters (in particular, passenger arrival rates λ_{ij} and willingness-to-pay distributions F_{ij}).

To characterize the equilibrium behavior of the system, as well as understand how to achieve good performance without perfect knowledge of system parameters, we turn to the use of state-dependent prices (in particular, local state-dependent prices). In Sect. 5.4.2, we summarize the results from Banerjee et al. (2015). Here, the authors show that while on the one hand state independent and dependent prices are asymptotically the same, the latter is much more robust to mis-specifications in system parameters. In more detail, they show that on the one hand, as the number of drivers and rate of passenger arrivals jointly scale to infinity, state-dependent pricing becomes asymptotically equal to state-independent pricing; on the other hand, they show that state-dependent prices are much less sensitive compared to state-independent prices under small perturbations in the system parameters.

5.4.1 *Fast-Timescale Control of Platform Dynamics*

We first turn to the question of choosing control policies for a given k that maximize the objective considered in Eq. 5.1. Note that these controls are extremely high-dimensional as they can in general depend on the instantaneous state of the system; moreover, as we discuss in Sect. 5.3.3, the problem is non-convex even in simple settings.

In Banerjee et al. (2017), the authors circumvent these problems via a novel technique for deriving control policies based on a convex relaxation which they term the *elevated flow relaxation*. The main idea behind this technique is to construct a concave pointwise upper bound for the objective, which is convex and hence admits a tractable optimization. This can be done essentially by assuming an infinite supply at each node, while simultaneously introducing additional flow-conservation constraints to capture the balance of units arriving to and exiting from any node in the stationary distribution. The authors prove that their new elevated objective is bounded below by the original objective, and thus optimal solutions in the elevated optimization problem are bounded below in value by optimal solutions in the original optimization problem. More importantly, they also prove lower bounds on the performance of natural state-independent policies provided by the relaxation, thereby establishing their approximation guarantees. They do so via the following three-step program:

1. First, they derive efficiently-computable upper bounds for the performance of any control policy, which encode essential ‘conservation laws’ of the system (in particular, flow balance at nodes and capacity constraints on number of vehicles on the road), while being amenable to optimization.
2. Next, they show that under an infinite-supply limit, where $k \nearrow \infty$ while all other parameters stay fixed, the achievable objective values under *state-independent control policies* exactly match the set of achievable upper bounds defined by the elevated flow relaxation.
3. Finally, using the product-form characterization of the stationary distribution under state-independent policies, they show that the performance of any policy in a setting with k units is within a factor of $1 + (n - 1)/k$ (assuming instantaneous transfers) of its performance in the infinite-supply setting.³

The versatility of the above framework allows it to be extended to more complex *multi-objective settings*, where the goal is to optimize some objective subject to a lower bound on another. A canonical example of this is the so-called *Ramsey pricing problem* (Ramsey 1927), where the platform aims to design a pricing policy to maximize system revenue subject to a lower bound on the system welfare; this is very relevant for most nascent ridesharing platforms, who are aiming to build

³Here we define the approximation factor as the ratio of the objective of the optimal policy to that of the proposed policy; this convention, which ensures the approximation ratio is always greater than 1, is common in the approximation algorithms literature.

a customer base. The authors also extend their analysis to include travel times, and show that increasing travel times may lead the approximation factor to degrade from $1 + O(1/k)$ under instantaneous transfers to $1 + O(1/\sqrt{k})$ in the worst case. One way to understand this phenomena is to note that assuming all other parameters stay fixed, increasing travel times leads to an increase in the amount of ‘work’ each incoming passenger needs from the k units; in the extreme case, most units are in transit at any time, and get engaged immediately as soon as they become idle.

The results in Banerjee et al. (2015) also recover and unify several other existing results in this area, and provide a general framework for deriving approximation algorithms for many other settings. In particular, the techniques provide an elementary proof of the so-called large-market (or ‘fluid limit’) optimality of the proposed state-independent policy; this form of limiting result was also obtained around the same time by Braverman et al. (2016) (for rebalancing) and Ozkan and Ward (2016) (for dispatch), via more technical limit interchange arguments.

Finally, the authors of Banerjee et al. (2017) also show that the bounds obtained via the above technique are *tight* compared to the optimal policy. This follows from an example earlier proposed in Wasserhole and Jost (2016), comprising of a ring of n stations, with equal request rate between every pair of adjacent stations except for one *bottleneck* link. The same example can also be modified to show similar bounds on the performance of *any* state-independent policy. Going beyond such policies is difficult; however, some recent work Banerjee et al. (2018) has proposed state-dependent algorithms which for large k , and under some additional conditions, can be shown to have $1 + e^{-O(k)}$ competitive ratio.

5.4.2 The Slow Timescale: Pricing and Driver Entry

The previous section provides a solution for the chief fast-timescale operational design questions in a ridesharing platform – given a supply of units, how best to use pricing, dispatch and rebalancing to balance the demand and supply. We next turn to understanding the slow-timescale response of drivers to such policies – in particular, we want to understand the impact of pricing policies on the overall marketplace equilibrium. To this end, we summarize the findings of Banerjee et al. (2015), who use the micro-foundations in Sect. 5.3 to compare the impact of two pricing schemes: (i) state-independent (i.e., *static*) pricing, where the price is fixed as a function of the fast-timescale system parameters, but not the instantaneous state⁴), and (ii) *dynamic* pricing policies, where the prices react to the system state. Following the discussion in Sect. 5.3.4, they focus on local state-dependent

⁴More specifically, these correspond to quasi-static policies, where the price remains fixed for blocks of time on the order of hours, but can be changed over slower timescales to reflect change in average demand/supply. Such policies are commonly used by traditional taxi firms.

pricing policies, as these admit closed-form stationary distributions; moreover, they focus on a simple class of *threshold* policies, wherein the platform raises the price whenever the number of available drivers in a region falls below a threshold.

The results in Banerjee et al. (2017) which we summarized in Sect. 5.4.1 indicate that when the supply is fixed, then state-independent policies, if correctly chosen, are very competitive compared to the optimal policy. The main contribution of Banerjee et al. (2015) is in understanding the effects of the pricing scheme on the equilibrium supply of drivers, and also, on characterizing the sensitivity of the two policies to parameter uncertainty. The major technical hurdle in doing so is that unlike the fast-timescale stationary distribution, the equilibrium of the system (as defined in Eq. 5.3) does not admit a closed-form expression for the driver/passenger arrival rates.

The Large-Market Scaling To circumvent this, the authors of Banerjee et al. (2015) study the system under a *large-market scaling*, wherein they consider a sequence of systems parametrized by ν , wherein $\bar{K}(\nu) = \bar{K}_0\nu$ and $\lambda(\nu) = \lambda\nu$, and all other system parameters, as well as the pricing policy \mathbf{p} , are held fixed. They then let ν approach ∞ , and study the *normalized* equilibrium supply state, i.e. $\lim_{\nu \rightarrow \infty} (K(\nu)/\nu)$, of the limiting system. For dynamic pricing policies, in addition to scaling \bar{K}_0 and μ_0 , they keep the set of prices fixed, but allow the threshold $\theta(\nu)$ to scale with ν .

Under the large-market scaling, the authors of Banerjee et al. (2015) characterize the equilibrium rates for the limiting system in closed form, for both static and dynamic pricing. An example of this convergence and limiting characterization can be seen graphically in Fig. 5.2a, which is similar to the plots given in Banerjee et al. (2015); here we have plotted the normalized equilibrium throughput (i.e., rate of matched rides) vs. static price p (the green curves), and also, for a class of dynamic pricing policies (the maroon curves) where we keep one price fixed at the red vertical dotted line. The dotted curves are numerically computed for $n \in \{1, 10, 100, 1000\}$, and can be seen to be monotonically converging up to the solid curves, which plot the theoretical large-market limits characterized in Banerjee et al. (2015).

Optimal Performance of Pricing Policies One surprising aspect of Fig. 5.2a is that the optimal throughput in the large-market limit over the dynamic pricing policies we consider appears to coincide with that obtained under static pricing. The authors of Banerjee et al. (2015) show, however, that this fact is true for *all* threshold dynamic pricing policies under fairly weak conditions: in particular, they prove that *as long as all passenger value distributions F_{ij} have an increasing hazard rate, then the optimal normalized throughput in the large-market limit under dynamic (i.e., local state-dependent) pricing collapses to that obtained under the optimal static (i.e., state-independent) pricing policy.* This combined with the results in Banerjee et al. (2017) shows that *the platform cannot improve performance much by employing state-dependent pricing.* Similar results are shown in Banerjee et al. (2015) for revenue and welfare, and also for multi-threshold pricing policies.

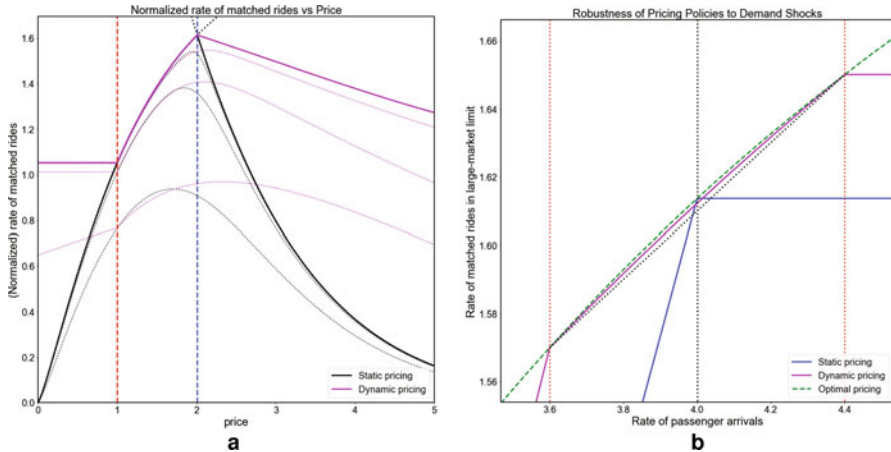


Fig. 5.2 Impact of pricing policies in ridesharing platforms: **(a)** Pricing policies in the large-market limit, **(b)** robustness of dynamic pricing. Figure **a** depicts the normalized equilibrium throughput in a ridesharing platform under static pricing (in black), and under dynamic pricing (in maroon) with one price fixed at the red vertical line. The dotted lines show the throughput curve for different values of the scaling parameter ν , with higher curves corresponding to higher values of ν . The solid curves plot the theoretical large-market limits. Note that in the large-market limit, the optimal throughput under both policies is the same (indicated by the black vertical dotted line). Figure **b** demonstrates the sensitivity of pricing policies to demand uncertainty: For a fixed \bar{K}_0 , we consider baseline passenger arrival rate $\lambda \in 4 \pm 10\%$, and compare the normalized throughput under (i) the optimal static policy with $\lambda = 4$ (indicated by the black vertical dotted line), and (ii) the dynamic-pricing policy which sets p_ℓ based on $\lambda = 3.6$, and p_h based on $\lambda = 4.4$ (indicated by the red vertical dotted lines). The dashed green curve shows the performance of the optimal static-pricing corresponding to the actual λ . We generated the plots for a single-node network, following the model described in Banerjee et al. (2015); in particular, we use demand value-distribution $F \sim \text{Exponential}(0.5)$, and driver reservation-value distribution $G_d \sim \text{Exponential}(0.8)$

We note here that the result given above is asymptotic; the plots in Fig. 5.2a clearly show that dynamic pricing does have gains over static pricing for smaller values of the scaling parameter ν . The non-trivial aspect is that the difference in performance vanishes in the limit. Note also that the performance of a dynamic pricing policy with prices (p_ℓ, p_h) is not identical to the performance with static price p_ℓ or p_h , and in fact, it can be shown that passengers experience both prices in the large-market limit.

Robustness of Pricing Policies More importantly, we note that the result that optimal static pricing and optimal dynamic pricing are asymptotically equivalent requires a key assumption: that the platform has knowledge of system parameters (e.g., the exogenous arrival rates of drivers and passengers, and distributions of reservation values). What should the platform do when these parameters are not well-known, and may even be highly variable?

To address this issue, the second main result in Banerjee et al. (2015) establishes a significant benefit that dynamic pricing holds over static pricing: *robustness*. Specifically, the authors of Banerjee et al. (2015) show that if the system operator

chooses the optimal threshold dynamic (resp., static) pricing policy assuming some predicted system parameters \bar{K}_0, λ , but the true parameters deviate from the predictions, then *dynamic pricing maintains a much higher share of the optimal throughput relative to the optimal static pricing*. This property is graphically depicted in Fig. 5.2b; for a more formal characterization of this property, we refer the reader to Banerjee et al. (2015).

5.5 Related Literature

In this section, we summarize the intellectual foundations our work builds on, as well as provide a brief survey of the growing literature on ridesharing across many fields. One of the great attractions of ridesharing is that the underlying questions lie at the intersection of several disciplines – economics, stochastic modeling and control, operations management, and network algorithms. The models and algorithms we have covered in this chapter borrows ideas from all these disciplines, and in a sense, we believe such a merging of ideas is critical to understanding these platforms. On the other hand, the diversity of disciplines makes it difficult to properly reference and comment on all the work related to this topic, and thus we acknowledge at the outset that our discussion below should be viewed as a survey of the main issues, rather than a comprehensive index of all relevant research.

Queueing networks and stochastic control Our model for the fast-timescale dynamics builds on the rich toolbox of *queueing network models*, starting from the seminal work of Jackson (1963) on open networks (i.e., where the number of units can change), and Gordon and Newell (1967) and Baskett et al. (1975) on closed networks (where the number of units is fixed, as in our setting). An excellent survey of these models is provided in the books by Kelly (1979), Kelly and Yudovina (2014) and Serfozo (1999). More recently, these models have found extensive use in several applied disciplines, including in the design of communications networks (Srikant and Ying 2013) and computer systems (Harchol-Balter 2013). A more recent line of work develops a similar theory for *matching queues*, obtaining surprising product form characterizations (Adan and Weiss 2012; Visschers et al. 2012; Moyal et al. 2017).

Optimal control of open queueing networks also has a long history, going back to the work of Whittle (1985), and more recently, these ideas have been extended to open matching networks (Bušić and Meyn 2014; Nazari and Stolyar 2016). However, there is much less work for closed networks. This in part due to the lack of a closed-form expression for the normalization constant though in our setting, it is computable in $O(nm)$ time via iterative techniques; see Buzen (1973) and Reiser and Lavenberg (1980). Consequently, most previous work on closed queueing networks used heuristics, with limited or no guarantees. In particular, heuristics based on ensuring ‘fairness’ properties (which are similar to the circulation constraints we discuss in Sect. 5.4.1) have been used in transportation setting to optimize weighted throughput (George 2012) and minimize rebalancing costs (Zhang and Pavone 2016).

The first formal approximation guarantees for control of closed networks was given by Wasserhole and Jost (2016), who derived a pricing policy for maximizing throughput; this result is a special case of the results of Banerjee et al. (2017) that we present in Sect. 5.4.1. Other recent works (Braverman et al. 2016; Ozkan and Ward 2016) have formally characterized the large-market limits for closed queueing networks, and proposed asymptotically optimal rebalancing and dispatch policies (these results can however be derived directly using the techniques in Banerjee et al. 2017).

Parallel to the work on control, there has also been a long line of research on strategic behavior in queueing systems; see Naor (1969) for early work in this area and Hassin and Haviv (2003) for an overview of these models. Typically, these works consider systems with a fixed number of servers, who serve arriving customers who are sensitive to price and delay. In contrast, our model considers strategic behavior on the part of the servers (i.e., the drivers). In this respect our treatment is closer to the recent work on queues with strategic servers (Gurvich et al. 2014; Gopalakrishnan et al. 2016).

Economics of two-sided platforms From an economics standpoint, the strategic aspects of our micro-foundations are in the spirit of the literature on the price theory of two-sided platforms (Rochet and Tirole 2006; Caillaud and Jullien 2003; Rysman 2009; Armstrong 2006; Armstrong and Wright 2007); refer to Weyl (2010) for an excellent summary and unification of this literature. This line of work typically studies the design two-sided markets under exogenously specified utility functions for agents. One critical difference in our approach is in building up the market model from the underlying stochastic dynamics, rather than specifying it exogenously. This is critical as having a dynamic model allows us to study dynamic pricing, which is one of the hallmarks of ridesharing platforms.

Stylized market models like the ones referenced above have been started being used for studying ridesharing as well. In particular, several recent works study the impact of spatial and temporal variations in prices on driver decisions in the fast timescale (Castillo et al. 2017; Bimpikis et al. 2016), a topic which we have not covered in this chapter. On the other hand, there is less work on using these to study inter-firm competition; one partial move in this direction is the work of Séjourné et al. (2017), which studies the additional societal costs of multi-firm ridesharing ecosystems under exogenous demand fragmentation.

An important difference in our treatment of the ridesharing marketplace is that we incorporate both stochastic dynamics and strategic interactions. Doing so is challenging as one needs to reason about market equilibria under a combination of dynamics and strategic interactions. One important approach that helps circumvent this is that of using *large-market limits*; this is the approach we adopt in Sect. 5.4.2, following the treatment in Banerjee et al. (2015). Large-market limits have grown in importance in recent years; see Kojima and Pathak (2009) and Azevedo and Budish (2012) for examples of this in the matching market literature, and Arnosti et al. (2014) for an application of this approach for dynamic matching markets.

Pricing and operations management A unique feature of ridesharing firms is that it is one of very few marketplaces where the platform can set the prices. Consequently, the study of such platforms shares commonalities to that of monopolist pricing. In particular, the comparison of static and dynamic pricing policies is a core topic of the literature on revenue management; refer to Talluri and Van Ryzin (2006) and Bitran and Caldentey (2003) for an overview of pricing approaches, and Gallego and Van Ryzin (1994) for an analysis of dynamic pricing based on current inventory levels. More recent work applies techniques from approximate dynamic programming to tackle problems in logistics with dynamic arrivals and pricing (Adelman 2007; Levi and Radovanović 2010; Hampshire et al. 2009). Though similar to the fast-timescale control problem, these approaches typically can deal only with small systems, as their dimensionality scales rapidly with the number of stations; moreover, many of the techniques have no provable guarantees.

We note that though the results in the revenue management literature are similar in spirit to ours (in particular, the optimality of static pricing in a large-system limit), there is a very significant difference in the underlying settings. In particular, while the primary concern of monopolist pricing in a one-sided platform is to regulate demand in the face of changing inventory levels, the role of prices in ridesharing is to simultaneously regulate both instantaneous demand-supply mismatches and the entry decisions of drivers.

Data-driven simulations and empirical studies Finally, in addition to the theoretical studies we mention above, there is also a parallel line of work which studies the same questions from a numerical perspective, either via data-driven simulation models, or experimental studies. Though such studies are of great importance, a big roadblock in this space is the lack of good data-sets and academic access to ridesharing platforms. An indicator of this state of affairs is the fact that several of the studies below were possible due to the New York City taxi dataset (http://www.nyc.gov/html/tlc/html/about/trip_record_data.shtml), which has played an important role as a public repository in this space. While we recognize the difficulties in data sharing, we emphasize that it is critical that the academic community works to cultivate publicly available data-sets that can be used as reference points for research in this area, across a range of platforms, geographies, and timescales.

In terms of data-driven simulation, a notable work is that of Santi et al. (2014), which used the NYC taxi dataset to study the benefits of *pooling* – combining multiple rides into one. The work introduced the idea of compatible rides based on a diversion threshold constraint, where two ride-requests are considered to be amenable to pooling if the total origin-to-destination time (with diversions) for each ride is within an additive constant of the trip time without pooling. This idea proved influential in the subsequent design of Lyft Line and Uber Pool; moreover, it also spurred further research into the algorithmic challenges of determining such pooled rides in real-time and at scale (Spieser et al. 2016; Alonso-Mora et al. 2017).

Empirical work on the internal dynamics of ridesharing platforms is challenging, as it depends on access to their proprietary data. As a consequence, most recent

work has involved collaborations between external researchers and data scientists. Several of these study the behavior of dynamic pricing. For example, Hall et al. (2015) presents a natural experiment that occurred when Uber’s surge pricing algorithm failed over the New Year’s Eve celebration in 2014–2015. In Chen and Sheldon (2016), the authors claim that dynamic pricing incentivizes more drivers to participate in the platform, and to meet supply shortages (analogous to the entry decisions described above in our analysis). The paper Hall et al. (2017) highlights the difference in timescales in market equilibration: in the short run driver supply is relatively inelastic, but in the long run driver supply enters in response to persistently higher prices.

5.6 Conclusion

In this chapter, we have outlined a stochastic modeling framework for ridesharing platforms, which is based on the ideas presented in our prior work on this topic (Banerjee et al. 2015, 2017), as well as work by several others on similar problems. In particular, in Sect. 5.3, we have presented this model in great detail, discussing all the assumptions that underlie this model, and arguing as to why these capture first-order phenomena in such platforms. We have also pointed out which aspects of these platforms lie outside our model.

One reason for our championing of this framework is its success in providing both theoretical insights and practical guidance into the design of pricing and dispatch policies on these platforms; we have summarized some of these results in Sect. 5.4. Our hope, however, is that this framework will go beyond the results presented to inspire and unify future studies into ridesharing platforms. To this end, we have outlined many open questions – in particular, there is a need for research into understanding the effect of driver strategic behavior on the fast timescale; the robustness of control policies to medium timescale changes in system parameters; the structure of efficient policies for ride pooling; the role of autonomous vehicles in the ridesharing landscape; the design of longer-term contracts for drivers and passengers; the interaction between multiple ridesharing platforms and between ridesharing and public transit; and the overall impact of ridesharing on society as a whole. These are difficult problems, and may not have any clear-cut answer – however, we hope the framework we have presented here will provide a benchmark for studying all these questions. We eagerly look forward to future research on these topics.

Acknowledgements The authors would like to thank the data science team at Lyft, particularly Chris Sholley; part of this work was carried out when SB was a technical consultant at Lyft. We gratefully acknowledge support from the National Science Foundation via grants CMMI-1234955 and CNS-1343253, the DARPA GRAPHS program, and Army Research Office grant W911NF-17-1-0094.

References

- Adan I, Weiss G (2012) Exact fcfs matching rates for two infinite multi-type sequences. *Oper Res* 60:475–489
- Adelman D (2007) Price-directed control of a closed logistics queueing network. *Oper Res* 55(6):1022–1038
- Alonso-Mora J, Samaranyake S, Wallar A, Frazzoli E, Rus D (2017) On-demand high-capacity ride-sharing via dynamic trip-vehicle assignment. *Proc Natl Acad Sci* 114(3):462–467
- Armstrong M (2006) Competition in two-sided markets. *RAND J Econ* 37(3):668–691
- Armstrong M, Wright J (2007) Two-sided markets, competitive bottlenecks and exclusive contracts. *Econ Theory* 32(2):353–380
- Arnosti N, Johari R, Kanoria Y (2014) Managing congestion in decentralized matching markets. Available at SSRN 2427960
- Azevedo EM, Budish E (2012) Strategy-proofness in the large as a desideratum for market design. In: *Proceedings of the 13th ACM conference on electronic commerce*. ACM, New York, p 55
- Banerjee S, Johari R, Riquelme C (2015) Pricing in ride-sharing platforms: a queueing-theoretic approach. In: *Proceedings of the 2015 ACM conference on economics and computation*. ACM, New York, p 517
- Banerjee S, Freund D, Lykouris T (2017) Pricing and optimization in shared vehicle systems: an approximation framework. In: *Proceedings of the 2017 ACM conference on economics and computation*. ACM, New York, p 517
- Banerjee S, Kanoria Y, Qian P (2018) The value of state dependent control in ride-sharing systems. arXiv preprint, arXiv:180304959
- Baskett F, Chandy KM, Muntz RR, Palacios FG (1975) Open, closed, and mixed networks of queues with different classes of customers. *J ACM (JACM)* 22(2):248–260
- Bimpikis K, Candogan O, Daniela S (2016) Spatial pricing in ride-sharing networks. Available at SSRN 2868080
- Bitran G, Caldentey R (2003) An overview of pricing models for revenue management. *Manuf Serv Oper Manag* 5(3):203–229
- Braverman A, Dai J, Liu X, Ying L (2016) Empty-car routing in ridesharing systems. arXiv preprint arXiv:160907219
- Bušić A, Meyn S (2014) Optimization of dynamic matching models. arXiv preprint, arXiv:14111044
- Buzen JP (1973) Computational algorithms for closed queueing networks with exponential servers. *Commun ACM* 16(9):527–531
- Caillaud B, Jullien B (2003) Chicken & egg: competition among intermediation service providers. *RAND J Econ* 34(2):309–328
- Castillo JC, Knoepfle D, Weyl G (2017) Surge pricing solves the wild goose chase. In: *Proceedings of the 2017 ACM conference on economics and computation*. ACM, New York, pp 241–242
- Chen MK, Sheldon M (2016) Dynamic pricing in a labor market: surge pricing and flexible work on the Uber platform. In: *Proceedings of the 2016 ACM conference on economics and computation*. ACM, New York, p 455
- Gallego G, Van Ryzin G (1994) Optimal dynamic pricing of inventories with stochastic demand over finite horizons. *Manag Sci* 40(8):999–1020
- George DK (2012) Stochastic modeling and decentralized control policies for large-scale vehicle sharing systems via closed queueing networks. PhD thesis, The Ohio State University
- Gopalakrishnan R, Doroudi S, Ward AR, Wierman A (2016) Routing and staffing when servers are strategic. *Oper Res* 64(4):1033–1050
- Gordon WJ, Newell GF (1967) Closed queueing systems with exponential servers. *Oper Res* 15(2):254–265
- Gurvich I, Lariviere M, Moreno A (2014) Staffing service systems when capacity has a mind of its own. Available at SSRN 2336514

- Hall J, Kendrick C, Nosko C (2015) The effects of Uber's surge pricing: a case study. The University of Chicago Booth School of Business
- Hall JV, Horton JJ, Knoepfle DT (2017) Labor market equilibration: evidence from Uber. Technical report, Working Paper, 1–42
- Hampshire RC, Massey WA, Wang Q (2009) Dynamic pricing to control loss systems with quality of service targets. *Probab Eng Inf Sci* 23(02):357–383
- Harchol-Balter M (2013) Performance modeling and design of computer systems: queueing theory in action. Cambridge University Press, Cambridge
- Hartline JD (2013) Mechanism design and approximation. Book draft October, p 122
- Hassin R, Haviv M (2003) To queue or not to queue: equilibrium behavior in queueing systems, vol 59. Springer, Boston
- Jackson JR (1963) Jobshop-like queueing systems. *Manag Sci* 10(1):131–142
- Kallenberg O (2006) Foundations of modern probability. Springer, New York
- Kelly FP (1979) Reversibility and stochastic networks. Cambridge University Press, Cambridge
- Kelly F, Yudovina E (2014) Stochastic networks, vol 2. Cambridge University Press, Cambridge
- Kojima F, Pathak PA (2009) Incentives and stability in large two-sided matching markets. *Am Econ Rev* 99(3):608–627
- Levi R, Radovanović A (2010) Provably near-optimal LP-based policies for revenue management in systems with reusable resources. *Oper Res* 58(2):503–507
- Moyal P, Basic A, Mairesse J (2017) A product form and a sub-additive theorem for the general stochastic matching model. arXiv preprint, arXiv:1711.02620
- Naor P (1969) The regulation of queue size by levying tolls. *Econometrica* 37(1):15–24
- Nazari M, Stolyar AL (2016) Optimal control of general dynamic matching systems. arXiv preprint, arXiv:1608.01646
- Ozkan E, Ward AR (2016) Dynamic matching for real-time ridesharing. Available at SSRN 2844451
- Ramsey FP (1927) A contribution to the theory of taxation. *Econ J* 37(145):47–61
- Reiser M, Lavenberg SS (1980) Mean-value analysis of closed multichain queueing networks. *J ACM (JACM)* 27(2):313–322
- Rochet JC, Tirole J (2006) Two-sided markets: a progress report. *RAND J Econ* 37(3):645–667
- Rysman M (2009) The economics of two-sided markets. *J Econ Perspect* 23(3):125–143
- Santi P, Resta G, Szell M, Sobolevsky S, Strogatz SH, Ratti C (2014) Quantifying the benefits of vehicle pooling with shareability networks. *Proc Natl Acad Sci* 111(37):13290–13294
- Séjourné T, Samaranyake S, Banerjee S (2017) The price of fragmentation in mobility-on-demand services. arXiv preprint, arXiv:1711.10963
- Serfozo R (1999) Introduction to stochastic networks, vol 44. Springer, New York
- Spieser K, Samaranyake S, Gruel W, Frazzoli E (2016) Shared-vehicle mobility-on-demand systems: a fleet operator's guide to rebalancing empty vehicles. In: Transportation Research Board 95th annual meeting, Washington, DC, 16–5987
- Srikant R, Ying L (2013) Communication networks: an optimization, control, and stochastic networks perspective. Cambridge University Press, Cambridge
- Talluri KT, Van Ryzin GJ (2006) The theory and practice of revenue management, vol 68. Springer
- Visschers J, Adan I, Weiss G (2012) A product form solution to a system with multi-type jobs and multi-type servers. *Queueing Syst* 70(3):269–298
- Waserhole A, Jost V (2016) Pricing in vehicle sharing systems: optimization in queueing networks with product forms. *EURO J Transp Logist* 5(3):293–320
- Weyl EG (2010) A price theory of multi-sided platforms. *Am Econ Rev* 100(4):1642–1672
- Whittle P (1985) Scheduling and characterization problems for stochastic networks. *J R Stat Soc Ser B (Methodol)* 47(3):407–428
- Zhang R, Pavone M (2016) Control of robotic mobility-on-demand systems: a queueing-theoretical perspective. *Int J Robot Res* 35(1–3):186–203

Part II
Intermediary Role of a Sharing Platform

Chapter 6

The Role of Surge Pricing on a Service Platform with Self-Scheduling Capacity



Gerard P. Cachon, Kaitlin M. Daniels, and Ruben Lobel

Abstract Recent platforms, like Uber and Lyft, offer service to consumers via “self-scheduling” providers who decide for themselves how often to work. These platforms may charge consumers prices and pay providers wages that adjust based on prevailing demand conditions. For example, Uber uses “surge pricing” which pays providers a fixed commission of its dynamic price. With a stylized model that yields analytical and numerical results, we study several pricing schemes that could be implemented on a service platform, including surge pricing. Our base model places no restrictions on the platform’s dynamic pricing and waging schemes, whereas our surge pricing analogue requires wages to be a fixed fraction of dynamic prices and our traditional taxi analogue requires prices to be fixed. We show that although surge pricing is not optimal, it generally achieves nearly the optimal profit, justifying its use in practice. Despite its merits for the platform, surge pricing has been criticized due to concerns for the welfare of consumers. In our model, as labor becomes more expensive, consumers are better off with surge pricing relative to fixed pricing because they benefit both from lower prices during normal demand

This chapter is adapted from Cachon et al. (2017), reprinted by permission, Gerard P Cachon, Kaitlin M Daniels & Ruben Lobel, *The Role of Surge Pricing on a Service Platform with Self-Scheduling Capacity*, Manufacturing & Service Operations Management, 2017. Copyright 2017, the Institute for Operations Research and the Management Sciences, 5521 Research Park Drive, Suite 200, Catonsville, MD 21228 USA.

G. P. Cachon
The Wharton School, University of Pennsylvania, Philadelphia, PA, USA
e-mail: cachon@wharton.upenn.edu

K. M. Daniels (✉)
Olin Business School, Washington University in St. Louis, St. Louis, MO, USA
e-mail: k.daniels@wustl.edu

R. Lobel
Airbnb, San Francisco, CA, USA
e-mail: ruben.lobel@gmail.com

and expanded access to service during peak demand. We conclude, in contrast to popular criticism, that both the platform and consumers can benefit from the use of surge pricing on a platform with self-scheduling capacity.

6.1 Introduction

The “gig-economy” has created new ways for firms to provide service to consumers. Instead of centrally scheduling the shifts of workers, gig-economy firms allow workers to independently set their hours with no advance notice required by the firm. The firm’s role becomes that of a platform that matches available providers to consumers demanding service. We use the term “self-scheduling” to describe this relationship.

Examples of platforms with self-scheduling capacity include ridesharing platforms Uber and Lyft, and delivery services Postmates and Instacart. A provider for one of these platforms first makes the long-term decision to join the platform (e.g. register as an Uber driver), which providers only do if they expect to earn more than their next best alternative. A provider then makes short term decisions about whether or not to work at each moment in time. This decision is influenced by the wage per service offered by the platform, the opportunity cost of the provider’s time (e.g. a provider’s opportunity cost is high at times when he has a medical appointment but low when he has nothing on his schedule), and the provider’s expectations about the availability of work. Note that providers respond not only to the volume of demand for service but also to the volume of competing providers. Specifically, a provider’s decision is based on the prevalence of demand relative to the amount of offered capacity.

Although the platform does not directly control the number of providers working at a time, it still takes an active role in managing its capacity to maximize its profit. The platform sets the price charged to consumers, the wage offered to providers, and regulates the maximum number of providers allowed to join the platform. Furthermore, the platform may dynamically set price and wage in response to changes in the state of the world. For example, Uber has become known for its dynamic policy called “surge pricing.”

While the platform is clearly motivated to maximize its profit, it must also consider the impact of its actions on other stakeholders. Dynamic pricing policies have come under fire from consumers who feel that the practice unfairly discriminates (Kosoff 2015; Stoller 2014). With potential regulation in mind, it is important to understand the extent to which there is tension between maximizing platform profit and maintaining consumer welfare.

Refer to a set of price, wage, and recruitment decisions as a contract. In this paper, we study the form of a gig-economy platform’s profit maximizing contract and its effect on the welfare of consumers. This “optimal” contract has state-dependent prices and independent, state-dependent wages. As a point of comparison, we introduce the “fixed” contract, in which the platform is restricted to offer the same

price in all states of the world. The comparison of consumer welfare elucidates the effect of dynamic prices on consumers. Echoing practice, we additionally consider the “commission” contract, in which the platform pays providers a fixed percentage of a dynamic price. Commission contracts are common (e.g. surge pricing), so analysis of this contract allows us to quantify the loss the platform incurs from imposing the fixed commission constraint.

While unconstrained dynamic incentives maximize the platform’s profit, in most cases the platform does not strongly prefer the optimal contract to the commission contract. Although not always optimal, the commission contract achieves near optimal profit and is simple to implement, which may explain its use in practice. Furthermore, we show that the improvement in platform profit from dynamic incentives is not necessarily to the detriment of consumer surplus. In markets where demand is rationed with the fixed contract, i.e., consumers do not always have access to service, dynamic incentives increase access to service and so serve more consumers. Consumers benefit from expanded access more than they suffer from increased prices during peak demand states, leading to a net benefit. However, in markets where the fixed contract provides sufficient capacity to serve all consumers in all states of the world, dynamic incentives provide the same level of access to service at a higher price during peak demand, making consumers worse off. Thus, if the lack of dynamic incentives leads to poor service, consumers benefit from the introduction of surge-pricing type policies.

6.2 Literature Review

Our work is primarily connected to three domains in the existing literature: research on capacity and pricing, revenue management models, and recent papers on peer-to-peer platforms and self-scheduling capacity. For simplicity and consistency, we refer to the various components in other papers using the terms relevant for our model. For example, the “platform” is the organization responsible for designing the market, “providers” generate capacity, “dynamic prices” are demand-contingent payments from consumers to the platform in exchange for service, and “dynamic wages” are demand-contingent payments from the platform to providers.

Several papers study competition among multiple providers and establish that competition can lead to excessive entry (e.g., Mankiw and Whinston 1986) and a platform should discourage competition to mitigate the losses in system value due to this issue (e.g., Bernstein and Federgruen 2005; Cachon and Lariviere 2005), but those papers do not consider dynamic wages or prices.

A set of papers considers peak-load pricing, the practice of charging higher prices during peak periods of demand (e.g., Gale and Holmes 1993). The primary motivation of peak-load pricing is to increase revenue by shifting demand from the peak period to the off-peak period. We do not incorporate this capability into our model. For example, consumers in need of transportation during a rainy evening are unable to postpone their need to a time with better weather.

There is work on the value of dynamic prices in systems that experience congestion, but with fixed capacity: e.g., Çelik and Maglaras (2008), Ata and Olsen (2009), and Kim and Randhawa (2017). Banerjee et al. (2015) considers the value of dynamic pricing in a model with random arrivals of consumers and providers. Unlike us, they find that dynamic pricing provides no benefit in terms of maximizing the platform's expected profit or system welfare, but they have a single demand regime whereas in our model some periods (importantly) have predictably higher demand than others for a given price.

There is a considerable literature on “two-sided markets” in which platforms earn rents by creating a market for buyers and sellers to transact (e.g., a game console maker as the platform between game developers and consumers) (e.g., Rochet and Tirole 2006). Daniels (2017) demonstrates the fundamental differences between classic two-sided markets models and models explicitly tailored to the gig-economy.

Peer-to-peer service platforms have attracted significant academic interest; e.g. Kabra et al. (2017), Hong and Pavlou (2014), Snir and Hitt (2003), Moreno and Terwiesch (2014), and Yoganarasimhan (2013). Those papers investigate how to subsidize different market players to accelerate the growth of a peer-to-peer platform, whether consumers have geographic preferences over providers, the influence of platform design on provider quality, and how provider reputation impacts the market. We do not explore those issues: our providers are ex-ante homogenous and do not build reputations. Fraiberger and Sundararajan (2015) investigate the interaction between ownership and sharing on a peer-to-peer marketplace, a dynamic that is not addressed in our model. Cohen et al. (2016) use Uber transaction data to measure the amount of consumer surplus generated given the implementation of surge pricing, but they do not estimate a counterfactual consumer surplus level with other contractual forms.

There is modeling and empirical work on the competition between peer-to-peer service marketplaces and existing markets: Einav et al. (2016), Zervas et al. (2017), Seamans and Zhu (2013), Cramer and Krueger (2016), and Kroft and Pope (2014). We do not directly consider the competition between the platform and incumbents.

Several papers (e.g., Hu and Zhou 2016; Allon et al. 2012) explore the process for matching providers to consumers when capacities are exogenous and all participants have preferences for the match they receive (e.g. a courier prefers to be matched to a nearby consumer). We do not consider matching because our consumers and providers are homogeneous, so careful matching does not provide a benefit.

Closest to our work are papers on self-scheduling capacity. Ibrahim and Arifoglu (2015) considers a model in which the platform chooses the number of providers and providers are either assigned by the platform to work in one of two different periods or they self select which of the two periods they work in. Unlike in our model, the platform can directly control the number of providers in the system. Taylor (2017) and Bai et al. (2016) study queuing systems in which a platform creates a market for service where arrivals of consumers and servers are endogenously determined based on decisions to seek and provide service respectively. Their models do not consider dynamic prices or wages, and the number of potential providers is exogenous (i.e., capacity decisions are made on a single, short-term, time scale). Gurvich et al.

(2016) studies a model in which a platform directly chooses the number of available providers, the wage for each provider who chooses to work, and a cap on the number of providers who are allowed to work: given the platform’s prevailing wage, more providers may want to work than the platform wants. They do not include dynamic pricing – in all versions of their model the platform selects a single price. They also do not impose an earnings constraint for providers. Instead, they impose an exogenous minimum wage. In our model providers decide whether to join the platform based on rational expectations of future earnings.

6.3 Model

To capture the relevant dynamics of the gig economy, we construct a two period model of provider behavior. At the outset, the platform sets the terms of its contract, which is comprised of prices charged to consumers, wages offered to providers, and a maximum number of recruited providers, N . In the first period, potential providers make the long-term decision about whether to join the platform (e.g. register as an Uber driver). We refer to this decision as the “joining decision.” This decision is relevant over a long horizon. In contrast, the second period represents a provider’s short term decisions about whether or not to offer service through the platform (e.g. go out on the road to drive for Uber). We refer to this as the “participation decision.” Only providers that join the platform in the first period can participate in the second period. While in practice providers make many participation decisions, for simplicity we condense these decisions into a single period (see Fig. 6.1).

When deciding to join, providers weigh their expected earnings from future participation in the platform’s market against the value of their next best alternative, denoted by c_1 . We assume that providers are homogeneous in the value of this alternative (e.g. everyone considering driving for Uber has similar qualifications and hence similar alternative employment opportunities) and that potential providers are numerous. Consequently, either all potential providers want to join to the platform, or none do. It follows that every profitable contract recruits exactly the

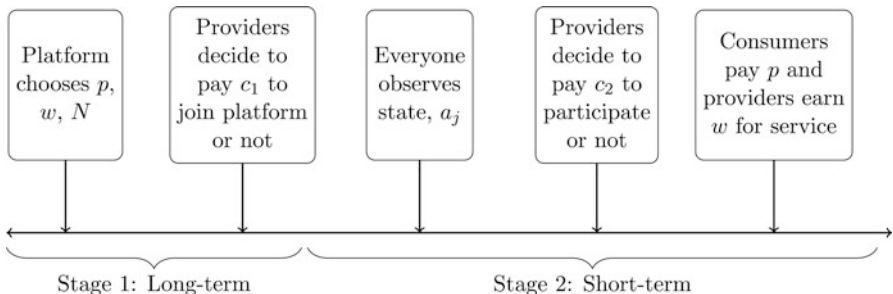


Fig. 6.1 Timeline of events

maximum number of allowed providers, N . Our model approximates a market with many potential providers and a highly elastic supply curve: with sufficiently high expected earnings the platform recruits ample capacity but few potential providers are interested if expected earnings fall below c_1 .

Providers join in expectation of a random shock that determines the state of the world in the second period. This shock determines “latent demand,” i.e., the number of consumers who would seek a free service, in the second period. Specifically, demand in state j is $D_j = (a_j - p_j)^+$, where p_j is the price the platform charges in state j and a_j is the latent demand in state j . As an illustration, consider the effect of rain on the operations of Uber. More consumers seek an Uber ride when it rains (Hall et al. 2015). Uber drivers cannot perfectly predict future weather, but the availability of demand affects earnings from participation. Hence providers evaluate joining based on expected future earnings. For simplicity we consider only two possible states, “high” and “low” with corresponding subscripts h and l , where the low state has smaller latent demand than the high state, i.e., $a_l < a_h$. Denote the probability of state j occurring by f_j .

In the second period, providers observe the state of the world (e.g. Uber drivers can see if it is raining) and decide whether to offer service. Providers must make themselves available to provide service before the platform matches them with consumers. For example, Uber matches drivers to nearby consumers, so drivers must be on the road to receive ride-requests. Hence providers decide to participate in advance of being matched to a customer. Suppose each provider has the capacity to serve at most one consumer in the second period. A provider then assesses his expected earnings from participation to be the product of the wage offered per service by the platform and the probability that the provider is matched with a consumer. Suppose consumers do not have preferences over providers and vice versa, so the platform is equally likely to match any consumer to any available provider. Then the probability that a provider is matched, ϕ_j , is the maximum of the ratio of demand to the number of competing providers and the provider’s capacity. Providers weigh these earnings against their opportunity cost, c_2 . This cost represents the value of the provider’s time; for example a provider with doctor’s appointment would have high c_2 while a provider without any plans would have low c_2 . A provider then participates only if the earnings from doing so exceed c_2 (for evidence that providers make profit maximizing decisions see Farber 2015; Chen and Sheldon 2015). For simplicity, the following analysis assumes that providers have homogeneous c_2 , though the results presented here generalize to heterogeneous c_2 (see Cachon et al. 2017, for details). Because providers must be available to receive a match, they incur their c_2 regardless of whether they are matched. A provider’s expected profit from participation in state j is

$$w_j \phi_j - c_2.$$

Because all providers are homogeneous, they play a symmetric mixed strategy. Let θ_j denote the probability that a provider participates in state j . If demand exceeds the number of recruited providers in the first stage, N , then all recruited

providers participate as long as their wage exceeds c_2 . If recruitment exceeds demand, then θ_j fraction of recruited providers participate, where

$$\theta_j = \min \left\{ \frac{w_j a_j - p_j}{c_2 N}, 1 \right\}.$$

Notice that if $c_2 \leq w_j(a_j - p_j)/N$, then all providers participate (i.e., $\theta_j = 1$) even though some are rationed. Otherwise, only some recruited providers participate and they all earn 0 profit. In the first stage providers expect to earn the corresponding profit,

$$\sum_j \left(w_j \min \left\{ \frac{D_j}{\theta_j N}, 1 \right\} - c_2 \right) f_j. \quad (6.1)$$

For any number of providers to join, the platform must ensure that (6.1) $\geq c_1$.

The platform anticipates the behavior of providers described above and correspondingly chooses its contract at the outset of the decision horizon. The contract is composed of prices, p_j , and wages w_j for each state of the world, and recruitment level N . The platform earns a margin $p_j - w_j$ on each service it successfully provides and may not provide service to more consumers than there are active providers. Hence, the platform solves:

$$\max_{w_l, w_h, p_l, p_h, N} \sum_j (p_j - w_j) \min \{a_j - p_j, \theta_j N\} f_j, \quad (6.2)$$

$$\text{s.t.} \quad \sum_j \left(w_j \min \left\{ \frac{D_j}{\theta_j N}, 1 \right\} - c_2 \right) f_j \geq c_1. \quad (6.3)$$

The platform may include additional constraints to ensure that prices are fixed, i.e., $p_l = p_h$, or that the platform offers a fixed commission, i.e., $w_l/p_l = w_h/p_h$. Notice that for a particular demand realization, price, and wage, it is possible that demand exceeds the platform's capacity to serve. In this case the probability that a provider gets a match is 1, and consumers are randomly rationed. Alternatively, it is possible that there are fewer available providers than there are consumers seeking service. In this case all consumers receive service but providers are randomly rationed, and $\phi_j < 1$. In either case, the mismatch in supply and demand creates a cost for the platform. If providers are not fully utilized, then they require an inflated wage to ensure that their expected earnings exceed c_1 . If instead consumers are rationed, the platform misses making a sale.

In addition to profit, we measure the consumer surplus resulting from the platform's contract choice. Similar to Cohen et al. (2015), we define consumer surplus with linear stochastic demand by:

$$S \doteq \sum_j \frac{a_j - p_j}{2} \min\{a_j - p_j, \theta_j N\} f_j.$$

Consumer surplus decreases in price and increases in the number of consumers served, which depends both on the number of recruited providers and on the fraction of recruits that participate. Hence the platform can increase consumer surplus by either decreasing price or expanding access to service.

Call the “optimal contract” the solution to the platform’s problem described by (1) and (2). The platform can achieve its optimal profit by offering the smallest payment required to induce participation in the low state (i.e., $w_l = c_2$), while compensating providers in the high state for both the cost of participating and the opportunity cost of joining (i.e., $w_h = c_2 + c_1/f_h$). The platform then recruits sufficiently many providers to ensure that all demand is met, so there is no demand rationing with the optimal contract. The platform’s optimal dynamic pricing scheme achieves platform profit and consumer surplus as outlined in the proposition below. (For proofs see the Technical Appendix of Cachon et al. 2017)

Proposition 1 *With the optimal contract, the platform earns*

$$U_o = \max \left\{ \frac{1}{4} \left((a_l - c_2)^2 f_l + \left(a_h - c_2 - \frac{c_1}{f_h} \right)^2 f_h \right), \frac{1}{4} (a_l f_l + a_h f_h - c_2 - c_1)^2 \right\}$$

and produces consumer surplus

$$S_o = \begin{cases} \frac{1}{8} ((a_l - c_2)^2 f_l + (a_h - c_2 - c_1/f_h)^2 f_h), & c_1 < (a_h - a_l) f_h, \\ \frac{1}{8} (a_l f_l + a_h f_h - c_2 - c_1)^2, & (a_h - a_l) f_h < c_1. \end{cases}$$

6.4 Profitability of Commission Contract

To quantify the cost of imposing a fixed commission, we construct the “commission contract.” In this setting, the platform is constrained to pay providers a fixed multiple, β , of the price in each state. The commission contract is the solution to:

$$\begin{aligned} & \max_{\beta, p_l, p_h, N} \sum_j (1 - \beta) p_j \min \{ a_j - p_j, \theta_j N \} f_j, \\ & \text{s.t.} \quad \sum_j \left(\beta p_j \min \left\{ \frac{D_j}{\theta_j N}, 1 \right\} - c_2 \right) f_j \geq c_1. \end{aligned}$$

This is equivalent to the platform’s original problem with an additional constraint requiring $w_l/p_l = w_h/p_h$. This constraint is irrelevant when joining cost is high: if $(a_h - a_l) f_h < c_1$ then the effective commission, w_j/p_j , with the optimal contract is the same in each state. In contrast, if the joining cost is “low,” i.e., $c_1 \leq (a_h - a_l) f_h$, then the platform uses a dynamic commission with the optimal contract. Hence, with low c_1 the commission contract forces platform offer a commission that is suboptimal in one or both states, reducing the platform’s profit relative to the optimal contract.

Proposition 2 *The following is a lower bound for the ratio of the platform's profit with the commission contract, U_β , and the platform's profit with the optimal contract, U_o : $\min\{U_\beta/U_o\} = (1 + \sqrt{f_h})/2$. This bound is achieved either when $c_1 = 0$ or $c_2 = 0$.*

The proposition above reports on a lower bound for the platform's profit with the commission contract. The commission contract performs poorly when one of the two costs is very low (either c_1 or c_2) and the probability of high demand is small. As $c_2 \rightarrow 0$, the optimal contract chooses a low commission when demand is low (to prevent too much participation) and, when demand is high, chooses a sufficiently high commission to give providers enough profit (c_1/f_h) to justify joining the platform. This disparity in the two commissions creates a challenge for the commission contract, which is required to choose a single commission. With the other extreme, $c_1 \rightarrow 0$, the joining constraint is not important. Instead, the focus is on the incentive for providers to participate. Because $p_l < p_h$, which implies $c_2/p_h < c_2/p_l$, the best commission with low demand is higher than with high demand (because both states must yield at least c_2 for the providers to participate). Again, the commission contract does not do well with this disparity in commissions. In the extreme, as $f_h \rightarrow 0$, the fixed commission contract earns only 1/2 of the profit of the optimal contract. However, when the two demand states are equally likely, the commission contract earns at least 85% of the optimal profit $(1/2)(1 + \sqrt{1/2})$.

To test the tightness of this bound, we numerically analyze 13,689 parameter combinations and evaluate the ratio U_β/U_o . According to the bound in Proposition 2, the commission contract has the most room for deviation from optimal profit when f_h is small. The tested parameter combinations all contain the extreme value $f_h = 0.05$ with the intention of illustrating the worst performance of the commission contract relative to the optimal contract. Without loss of generality, set the expected latent demand $\bar{a} \doteq a_l f_l + a_h f_h = 100$. We then vary the ratio of $\delta \doteq a_l/\bar{a}$, which is bounded between zero and one by definition and produces a corresponding $a_h = \bar{a}(1 - (1 - f_h)\delta)/f_h$. We consider only scenarios in which it is possible for the platform to serve low demand, i.e., $c_2/a_l \in [0, 1]$, and in which it is possible for the platform to recruit, i.e., $c_1/(a_h f_h + a_l f_l - c_2) \in [0, 1]$. The specific values of these parameters are summarized in Table 6.1.

Although there are cases in which the commission contract performs poorly relative to the optimal contract, this requires special parameters. For example, consider only the extreme cases in which $f_h = 0.05$, which yields a lower bound of $U_\beta/U_o = 0.612$. Evaluation of 13,689 evenly spaced observations throughout the feasible parameter space yields a minimum profit ratio close to the lower bound, $U_\beta/U_o = 0.6185$. (The lower bound is not achieved because the extreme border

Table 6.1 A summary of tested parameter values. All combinations of these values constitute 13,689 numerical experiments

Parameters	δ	c_2/a_l	$c_1/(a_h f_h + a_l f_l - c_2)$
Values	{0.1, 0.2, ..., 0.9}	{0.025, 0.050, ..., 0.975}	{0.025, 0.050, ..., 0.975}

Table 6.2 Quartiles and mean of U_β/U_o

	Minimum	Q1	Median	Mean	Q3	Maximum
U_β/U_o	0.6185	0.9871	0.9995	0.9755	1.000	1.000

conditions $c_1 = 0$ or $c_2 = 0$ are not included.) However, as illustrated in Table 6.2, the average ratio is $U_\beta/U_o = 0.9755$ and the median ratio is $U_\beta/U_o = 0.9995$. We conclude that for the majority of parameters, the commission contract yields nearly optimal profit.

6.5 Impact of Dynamic Prices on Consumers

The platform's pricing scheme affects not only the platform's profit but also the surplus accrued by consumers. Consumers have decried dynamic pricing of on-demand services as price gauging with such vehemence that regulators have considered limiting the practice (Kosoff 2015). Here we measure the effect of dynamic pricing on consumer surplus to determine whether dynamic pricing is as damaging as some consumers believe.

Call the platform's contract without dynamic pricing the "fixed contract." With this contract, the platform may adjust its payments to providers on a state-by-state basis, but prices remain fixed across states, i.e., $p_l = p_h$. With this contract, the platform is unable to match supply and demand in each state and so suffers from the resulting inefficiency. The fixed contract takes one of the following forms: full utilization of providers in all states but rationing of consumers in the high state, full service for consumers but rationing some providers in the low state, or service only for consumers in the high state. Refer to these possible options as the "poor service," "poor utilization," and "high demand only" outcomes, respectively.

Proposition 3 *The optimal contract has higher consumer surplus than the fixed contract if and only if "poor service" or "only high demand" is the best version of the fixed contract.*

Proposition 3 identifies the situations in which the optimal contract increases consumer surplus relative to the fixed contract. If providers are relatively expensive (high c_1) then the fixed contract involves demand rationing (poor service) and consumers benefit from switching from the fixed contract to the optimal contract. In these cases the fixed contract is unable to provide adequate supply and, even though consumers pay more in the high demand state with the optimal contract, the additional supply available with the optimal contract leads to higher consumer surplus. However, if providers are relatively cheap (low c_1) then the fixed contract fully serves consumers at the expense of provider utilization (e.g., the poor utilization version), and consumers are worse off with a switch to the optimal contract.

6.6 Conclusion

We study a platform that offers a service via a pool of independent providers. Providers self-schedule when they offer their service to the customers on the platform and decide whether or not to join the platform based on their earnings expectations. Demand varies over the long-term but is predictable in the short-term. Two inefficiencies can arise: (i) demand can be rationed either because too few providers join the platform or too few choose to participate; and (ii) capacity can be rationed because competition for a limited number of jobs leads too many providers to participate. Demand rationing is costly because some customers are unable to access the service that they value at the price charged, and the customers that do get the service might not be the ones that value it the most. Capacity rationing is costly because participating providers are not fully utilized but still incur their full opportunity cost of joining the platform. Both forms of rationing factor into the decision of providers as to whether to join the platform or not.

We study several contractual forms that vary in how prices respond to demand. The most basic contract, the fixed contract, sets a single price no matter what demand level occurs. The commission contract allows dynamic prices but requires prices to be a fixed multiple of the state-dependent wages offered to providers. The commission contract mimics pricing used in practice, such as Uber's surge pricing policy. Finally, we study an optimal contract, which imposes no restrictions on the platform's state-dependent prices and wages. Our main result is that even though the commission contract is not optimal, it yields nearly the optimal profit for the platform in the vast majority of plausible scenarios.

While maximizing profit is clearly an important objective for the platform, it isn't the only relevant one. A considerable amount of controversy has arisen over whether surge pricing gouges consumers. Hence, a platform should also be concerned with how it influences consumer surplus. The optimal contract leads to ambiguous welfare implications, which depend on how the fixed contract manages demand and capacity. If providers are relatively inexpensive (i.e., their opportunity cost to join the platform is low), then the fixed contract recruits an ample number of providers and underutilizes them during low demand periods. In this setting, switching to the optimal contract always works to the disadvantage of consumers because the platform recruits fewer providers and, in the high demand state, charges more and serves fewer customers. However, if providers have a high opportunity cost, then the fixed contract recruits a limited number of providers and forces customers during peak demand to suffer through poor service. In those cases, providers and consumers are better off with the introduction of the optimal contract: capacity expands to serve more customers in all demand states. To frame this in the context of ride-sharing, if with the fixed contract (e.g. taxis) it is hard to find service at peak demand times (e.g. a rainy evening), then Uber's introduction of surge pricing (i.e., dynamic pricing) is likely to make both Uber and consumers better off.

References

- Allon G, Bassamboo A, Çil EB (2012) Large-scale service marketplaces: the role of the moderating firm. *Manag Sci* 58(10):1854–1872
- Ata B, Olsen T (2009) Near-optimal dynamic lead-time quotation and scheduling under convex-concave customer delay costs. *Oper Res* 57(3):753–768
- Bai J, So KC, Tang C, Chen XM, Wang H (2016) Coordinating supply and demand on an on-demand service platform with impatient customers. Working paper, University of California at Irvine
- Banerjee S, Riquelme C, Johari R (2015) Pricing in ride-sharing platforms: a queueing-theoretic approach. In: Proceedings of the sixteenth ACM conference on economics and computation. ACM, New York, p 639
- Bernstein F, Federgruen A (2005) Decentralized supply chains with competing retailers under demand uncertainty. *Manag Sci* 51(1):18–29
- Cachon GP, Lariviere MA (2005) Supply chain coordination with revenue-sharing contracts: strengths and limitations. *Manag Sci* 51(1):30–44
- Cachon GP, Daniels KM, Lobel R (2017) The role of surge pricing on a service platform with self-scheduling capacity. *Manuf Serv Oper Manag* 19(3):368–384
- Çelik S, Maglaras C (2008) Dynamic pricing and lead-time quotation for a multiclass make-to-order queue. *Manag Sci* 54(6):1132–1146
- Chen MK, Sheldon M (2015) Dynamic pricing in a labor market: surge pricing and flexible work on the uber platform. Working paper, University of California Los Angeles
- Cohen MC, Lobel R, Perakis G (2015) The impact of demand uncertainty on consumer subsidies for green technology adoption. *Manag Sci* 62(5):1235–1258
- Cohen P, Hahn R, Hall J, Levitt S, Metcalfe R (2016) Using big data to estimate consumer surplus: the case of Uber (No. w22627). National Bureau of Economic Research
- Cramer J, Krueger AB (2016) Disruptive change in the taxi business: the case of Uber. Working paper, Harvard University. Available at <http://www.nber.org/papers/w22083>
- Daniels KM (2017) Distinguishing the gig-economy from two-sided markets. Working paper, Washington University in St. Louis
- Einav L, Farronato C, Levin J (2016) Peer-to-peer markets. *Annu Rev Econ* 8:615–635
- Farber HS (2015) Why you can't find a taxi in the rain and other labor supply lessons from cab drivers. *Q J Econ* 130(4):1975–2026
- Fraiberger SP, Sundararajan A (2015) Peer-to-peer rental markets in the sharing economy. Research paper, NYU Stern School of Business
- Gale IL, Holmes TJ (1993) Advance-purchase discounts and monopoly allocation of capacity. *Am Econ Rev* 83(1):135–146
- Gurvich I, Lariviere M, Moreno-Garcia A (2016) Operations in the on-demand economy: staffing services with self-scheduling capacity. Working paper, Kellogg School of Management
- Hall J, Kendrick C, Nosko C (2015) The effects of Uber's surge pricing: a case study. The University of Chicago Booth School of Business
- Hong Y, Pavlou PA (2014) Is the world truly "flat"? Empirical evidence from online labor markets. Working paper, Arizona State University
- Hu M, Zhou Y (2016) Dynamic type matching. Working paper, University of Toronto
- Ibrahim R, Arifoglu K (2015) Managing large service systems with self-scheduling agents. Working paper, University College London
- Kabra K, Belavina E, Girotra K (2017) The efficacy of incentives in scaling marketplaces. Working paper, INSEAD
- Kim J, Randhawa RS (2017) The value of dynamic pricing in large queueing systems. *Oper Res* 66(2):409–425
- Kosoff M (2015) A New York City politician wants to ban Uber's surge pricing—but that's a terrible idea. *Business Insider*. March 7, <http://www.businessinsider.com/banning-ubers-surge-pricing-is-a-terrible-idea-2015-2>

- Kroft K, Pope DG (2014) Does online search crowd out traditional search and improve matching efficiency? Evidence from Craigslist. *J Labor Econ* 32(2):259–303
- Mankiw NG, Whinston MD (1986) Free entry and social inefficiency. *RAND J Econ* 17(1):48–58
- Moreno A, Terwiesch C (2014) Doing business with strangers: reputation in online service marketplaces. *Inf Syst Res* 25(4):865–886
- Rochet JC, Tirole J (2006) Two-sided markets: a progress report. *RAND J Econ* 37(3):645–667
- Seamans R, Zhu F (2013) Responses to entry in multi-sided markets: the impact of Craigslist on local newspapers. *Manag Sci* 60(2):476–493
- Snir EM, Hitt LM (2003) Costly bidding in online markets for IT services. *Manag Sci* 49(11):1504–1520
- Stoller M (2014) Uber’s algorithmic monopoly. <http://mattstoller.tumblr.com/post/82233202309/ubers-algorithmic-monopoly-we-are-not-setting>. Date last accessed: 4 May 2017
- Taylor T (2018) On-demand service platforms. *Manuf Serv Oper Manag* 20(4):704–720
- Yoganarasimhan H (2013) The value of reputation in an online freelance marketplace. *Market Sci* 32(6):860–891
- Zervas G, Proserpio D, Byers J (2017) The rise of the sharing economy: estimating the impact of Airbnb on the hotel industry. *J Market Res* 54(5):687–705

Chapter 7

Time-Based Payout Ratio for Coordinating Supply and Demand on an On-Demand Service Platform



Jiaru Bai, Kut C. So, Christopher S. Tang, Xiquan (Michael) Chen, and Hai Wang

Abstract Many on-demand service platforms use a fixed payout ratio (i.e., the percentage of the platform’s revenue that is paid to the providers) regardless of the customer demand and the number of participating providers that tend to vary over time. In this chapter, we examine the implications of time-based payout ratios. To do so, we first present a queueing model with endogenous supply (number of participating providers) and endogenous demand (customer request rate) to model this on-demand service platform. In our model, earnings-sensitive independent providers have heterogeneous reservation price (for work participation) to serve wait-time and price-sensitive customers with heterogeneous valuation of the service. As such, both the supply and demand are “endogenously” dependent on the price the platform charges its customers and the wage the platform pays its independent providers. We use the steady state performance (associated with the M/M/1 queue) in equilibrium to characterize the optimal price, optimal wage and optimal payout ratio that maximize the profit of the platform. We find that it is optimal for the platform to offer time-based payout ratios by offering a higher payout ratio during peak hours and a lower payout ratio during non-peak hours.

J. Bai (✉)

School of Management, Binghamton University, Binghamton, NY, USA
e-mail: jbai@binghamton.edu

K. C. So

The Paul Merage School of Business, University of California, Irvine, CA, USA
e-mail: rick.so@uci.edu

C. S. Tang

Anderson School, University of California, Los Angeles, Los Angeles, CA, USA
e-mail: chris.tang@anderson.ucla.edu

X. (Michael) Chen

College of Civil Engineering and Architecture, Zhejiang University, Hangzhou, China
e-mail: chenxiquan@zju.edu.cn

H. Wang

School of Information Systems, Singapore Management University, Singapore, Singapore
e-mail: haiwang@smu.edu.sg

7.1 Introduction

Recent advances in internet/mobile technologies have enabled the creation of various innovative on-demand service platforms for providing *on-demand* services anytime/anywhere. Examples include grocery delivery services (e.g., Instacart, Google Express), meal delivery services (e.g., Sprig, Blue Apron), and food delivery services directly from restaurants (e.g., DoorDash, Deliveroo (U.K.), Yelp’s Eat24), consumer goods delivery services (e.g., UberRush), dog-walking services (e.g., Wag), and taxi-style transportation (e.g., Uber, Didi).

Due to dynamic customer demand anytime/anywhere, most on-demand platforms use independent providers to fulfill customer demands. However, because independent agents do not get compensated for idle times, earnings depends on wage rate and utilization, which depends on customer demand. At the same time, the demand associated with wait-time and price-sensitive customers depends on two key factors: price and waiting time. Since customer’s waiting time depends on the number of participating agents, which depends on the wage and the customer demand. Therefore, the “supply” of participating agents and the “demand” of customer requests are endogenously dependent on the wage and the price specified by the firm. Hence, it is a big challenge for the platform to (1) set the right wage (i.e., compensation) to get the right supply (i.e., the right number of earnings-sensitive participating agents); and (2) charge the right price to control the right demand (i.e., the right amount of wait-time and price-sensitive customers).

In view of the intricate relationship between endogenous supply and demand through wage and price selections, we develop an analytical framework to examine how an on-demand service firm should set its price, wage and payout ratio (i.e., the ratio of wage over price). (Throughout this paper, we shall refer to “payout ratio” as the percentage of the price collected from the customers that is paid to the providers.) In our framework, we use a queueing model to study the situation where both supply (i.e., number of providers) and demand (i.e., customer arrival rate) are “endogenously” dependent on wage, price and other operating factors. Our model captures an operating environment where (1) wait-time and price-sensitive customers are “heterogeneous” in their *evaluation* of the service; and (2) earnings-sensitive independent providers are “heterogeneous” in their *reservation price* (i.e., the minimum wage for work participation).

By analyzing the steady state performance of our queueing model in equilibrium, we characterize the optimal price, wage and payout ratio (i.e., the ratio of wage over price) in the basic setting under which the objective is to maximize the firm’s profit. We then extend our analysis to a more general setting under which the objective is to maximize the firm’s profit plus the welfare of the consumers and providers. For both settings, we obtain two key findings:

1. When the potential customer demand becomes higher, it is optimal for the firm to charge a higher price, pay a higher wage, and offer a higher payout ratio.
2. When customers become more wait-time sensitive, it is optimal for the firm to pay a higher wage and offer a higher payout ratio; however, the firm may need to charge a lower price to sustain the demand of increasingly impatient customers.

Our findings have the following managerial implications. First, as both the optimal price and the optimal wage are increasing in the maximum potential customer demand rate, our result provides an additional explanation/justification for an on-demand service firm (such as Uber) to charge its customers a higher price and pay its independent providers a higher wage when demand is higher. Second, while it is simple to share a fixed percentage of its revenue with the independent agents (e.g., Uber shares 80% of its revenue with its drivers; see Damodaran 2014), we find that the firm can increase its own profit as well as the total (customer and provider) welfare by offering a higher payout ratio when demand is higher. We hope this result might motivate on-demand service firms to re-evaluate their current fixed revenue sharing scheme. For instance, the firm may offer a higher (lower) payout ratio during peak hours (non-peak hours). Third, we also find that it is optimal for the firm to reduce its payout ratio when the number of registered independent providers becomes larger. This analytical result provides an economic justification for explaining why Uber reduced its payout ratio from 0.8 (initial payout ratio for its first cohorts of drivers) to 0.75 (for its second cohorts of drivers in 2014). Fourth, for urgent on-demand services with highly wait-time sensitive customers, the firm may need to lower its price to sustain demand from increasingly impatient customers.

This chapter is organized as follows. We provide a brief review of related literature in Sect. 7.2. Section 7.3 presents our queueing model of endogenous supply and demand along with heterogeneous providers and customers. In Sect. 7.4, we analyze the equilibrium behavior of our queueing system to determine the optimal price, wage and payout ratio for maximizing the firm's profit. We also adapt our base model to two special cases when the firm uses a fixed payout ratio and when the firm sets a fixed service level. We construct some illustrative numerical examples in Sect. 7.5 based on actual data provided by Didi: the leading taxi-style transportation on-demand service in China. We conclude in Sect. 7.6. Due to page limits, all mathematical proofs for the results are available at: https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2831794.

7.2 Literature Review

Recent developments of various on-demand service platforms such as Uber and DoorDash (see Kokalitcheva 2015; Wirtz and Tang 2016; Shoot 2015) have motivated researchers to explore various operational issues. First, there is an on-going debate regarding the definition of independent contractors for various on-demand service platforms (e.g., see Roose 2014). At the same time, it is of interest to examine how dynamic wage affects supply especially when independent providers can freely choose whether and when to work. Chen and Sheldon (2015) examined transactional data associated with 25 million trips obtained from Uber and showed empirically that dynamic wage (due to surge pricing) could entice independent drivers to work for longer hours. Moreno and Terwiesch (2014) also examined

empirically the independent contractor's bidding behavior on freelancing platforms. Allon et al. (2012) explored the process for matching providers to consumers when capacities are exogenous.

A number of researchers have recently studied the impact of wage and price on supply and demand in the context of on-demand services. Specifically, they examined whether it would be beneficial for an on-demand service firm to adjust its prices and wages dynamically based on real-time system status including the current number of customers requesting service and the number of providers in the system. Riquelme et al. (2015) and Cachon et al. (2017) compared the impact of static versus dynamic prices and wages. By assuming that customers are heterogeneous in terms of valuation and the payout ratio is exogenously given, Riquelme et al. (2015) found that static pricing performs well. On the other hand, Cachon et al. (2017) found that surge pricing performs well by assuming that customers are homogeneous and the payout ratio is endogenously determined. Hu and Zhou (2017) developed a general model where supply purely depends on wage and demand purely depends on price, and derived the conditions under which the optimal revenue sharing ratio is a linear function of the demand rate. Gurvich et al. (2016) also developed a newsvendor-style model to examine the optimal price and wage decisions. This stream of research assumes that customer demand is independent of waiting time and supply (or capacity) is independent of system utilization over time. In contrast, our model captures the rational behavior of customers who are sensitive to wait-time (and price) and independent providers who are sensitive to earnings which depend on the system utilization.

One research stream in the queueing literature has studied pricing decisions for services where customers can incur waiting or delay costs. Of particular relevance to our paper, a number of research papers have examined an operating environment that uses a static uniform (non-discriminatory) pricing strategy for heterogeneous customers. Afeche and Mendelson (2004) analyzed the revenue-maximizing and socially optimal equilibria under uniform pricing for heterogeneous customers with different evaluations of their service, and found that the classical result that the revenue-maximizing admission price is higher than the socially-optimal price (see, e.g., Naor 1969) can be reversed under a more generalized delay cost structure. Zhou et al. (2014) analyzed the structure of the optimal uniform pricing strategies for two classes of customers with different service valuations and wait-time sensitivities. Armony and Haviv (2003) and Afanasyev and Mendelson (2010) studied the competition between two firms under uniform pricing for two classes of heterogeneous customers. All the above research papers, however, are based on the assumption that capacity is exogenously given. In contrast, our paper considers the case when the supply (capacity) depends on wage and system utilization, which needs to be determined endogenously.

Finally, our model is closely related to some recent work by Taylor (2017). To our knowledge, Taylor (2017) is the first to examine pre-committed price and wage based on customer demand and other operating factors in the context of on-demand services. He compared the optimal prices when the providers are independent contractors or regular employees, and examined the impact of wait-time sensitivity

on the optimal price and wage using a two-point distribution for both the customer valuation of the service and the provider's reservation price. Our model allows these two distributions to be continuous, and complements Taylor's work in two important ways. First, our focus is to examine the impact of demand rate, wait-time sensitivity, service rate, and the size of available providers (who are on-reserve) on the optimal price, wage and payout ratio (ratio between the optimal wage and the optimal price). Second, in addition to maximizing its profit, we also consider the case when the firm maximizes the sum of its own profit and the total consumer and provider surplus. We find that our key results continue to hold: the optimal price, the optimal wage and the optimal revenue sharing ratio are increasing in the potential customer demand rate.

7.3 A Model of Wait-Time Sensitive Demand and Earnings Sensitive Supply

We consider an on-demand service platform (e.g., Uber) that coordinates randomly arriving (price and wait-time sensitive) customers with (earning-sensitive) independent service providers. Customers arrive randomly at the platform to request for service, and each service request consists of a (random) amount of service units to be processed by a service provider (e.g., travel distance in km). The platform charges each customer a *fixed price rate* p per service unit (e.g., dollar per km), and offers a *fixed wage rate* w per service unit to each participating service provider. (Here, we use "wage rate" per service unit so that the payout ratio w/p is well defined. However, we shall compute "earnings rate" per unit time later for providers to decide whether to participate or not.) In the same spirit as in Taylor (2017), the price rate p and wage rate w are pre-committed. However, their values can vary across different time periods depending on the specific market characteristics such as the average customer demand rate and the expected number of available providers. Given p and w , each customer decides whether to use the platform to request for service, and each independent provider decides whether to participate.¹ The primary objective of the platform is to select the optimal price rate and wage rate, denoted by p^* and w^* , to maximize its average profit.

¹For each service request, the platform will assign one of the available participating providers to serve the customer. For instance, the service platform can assign an available participating provider based on certain specific criteria (e.g., Uber assigns an available driver closest to the pickup location), or can announce a service request to all available participating service providers and assign the request to the first respondent.

7.3.1 Customer Request Rate λ and Price Rate p

During a certain time period, the maximum potential customer demand rate for the service during this time period is given by $\bar{\lambda}$, each of which has a valuation of the service that is based on a value rate v per service unit, where v spreads over the range $[0, 1]$ according to a cumulative distribution function $F(\cdot)$ so that $F(0) = 0$ and $F(1) = 1$. To capture the notion of wait-time sensitivity, we assume that the utility function of a customer with value rate v is given by

$$U(v) = (v - p)d - cW_q, \quad (7.1)$$

where $v - p$ is the surplus per service unit and d represents the average service units dictated by a customer (not the provider), c denotes the cost of waiting per unit time, and W_q represents the expected wait-time for the service. In this case, a rational customer with valuation v will request for service only if $U(v) \geq 0$ i.e., only if $v \geq p + (c/d)W_q$. Hence, the platform can use p and w to indirectly control the effective demand (i.e., the customer request rate) λ so that

$$\lambda = \text{Prob}\{U(v) \geq 0\} \cdot \bar{\lambda} = \text{Prob}\{v \geq p + (c/d)W_q\} \cdot \bar{\lambda}.$$

By defining the “target” service level $s = \text{Prob}\{v \geq p + (c/d)W_q\}$, the effective customer request rate λ (i.e., demand) is given by:

$$\lambda = s\bar{\lambda}. \quad (7.2)$$

Because of the one-to-one correspondence between target service level s and the effective demand rate λ , we shall focus our analysis on s instead of λ throughout this paper for mathematical convenience. Using the fact that $s = \text{Prob}\{v \geq p + (c/d)W_q\}$ and that $v \sim F(\cdot)$, the price rate p satisfies the following equation:

$$p = F^{-1}(1 - s) - \frac{c}{d}W_q, \quad (7.3)$$

where the price rate p decreases in the expected wait-time W_q and the unit waiting cost c .

7.3.2 Number of Participating Providers k and Wage Rate w

Let K be the (maximum) number of potential earnings-sensitive providers who may decide to participate over the same time period. (Essentially, K represents the number of registered providers who are eligible to participate.) For any given (p, w) , let k be the actual number of providers participating on the platform, where $k \leq K$. Also, let μ denote the average service speed (number of service units

processed per unit time; e.g., travel speed measured in terms of km per hour) of the service providers so that μ/d represents the service rate of the providers (i.e., average number of customers served per hour).² Given the customer request rate λ and the number of participating providers k , the utilization of these k participating providers is equal to $\lambda/(k \cdot (\mu/d))$, where $\lambda d < k\mu$ to ensure system stability. The average wage per unit time of a participating provider (when working) is equal to the wage per service unit w multiplied by the average service speed μ . Accounting for the utilization, the average “earning rate” per unit time of a participating provider is equal to $w\mu \cdot (\lambda d/(k\mu)) = w(\lambda d/k)$.³

To model the notion of earnings-sensitivity, we assume that each potential provider has a reservation rate r per unit time (i.e., corresponding to his outside option), where r varies across different providers. To model the heterogeneity among providers, we assume that there is a continuum of provider types so that the reservation rate r spreads over the range $[0, 1]$ according to a cumulative distribution function $G(\cdot)$, where $G(\cdot)$ is a strictly increasing function with $G(0) = 0$ and $G(1) = 1$. For a (potential) provider with reservation rate r , he will participate to offer service only if his average earning rate $w(\lambda d/k)$ is at least equal to r .

Let β denote the proportion of providers who participate in the platform to offer service during this time period. Then, $\beta = \text{Prob}\{r \leq w(\lambda d)/k\} = G(w(\lambda d)/k)$, and the actual number of participating providers k (i.e., supply) is given by

$$k = \beta K. \quad (7.4)$$

Also, in equilibrium, $\beta = G(w(\lambda d)/k)$ so that:

$$G^{-1}(\beta) = w \frac{\lambda d}{k}. \quad (7.5)$$

From (7.4) and (7.5), we can express the wage rate w as a function of the number of participating providers k :

$$w = G^{-1}(\beta) \frac{k}{\lambda d} = G^{-1}\left(\frac{k}{K}\right) \frac{k}{\lambda d}. \quad (7.6)$$

²If the service units d are already measured in terms of time units, we can simply set $\mu = 1$ in this case.

³For independent service providers, utilization and wage rate are the two key factors for their participation. For example, DePhills (2016) reported that Uber drivers obtain higher earnings primarily because their utilization rate (measured in terms of percentage of miles driven with a passenger) is much higher than that for taxi drivers. For instance, Uber driver’s utilization is 64.2%, while taxi driver’s utilization is only 40.7% in Los Angeles.

7.3.3 Problem Formulation

Since the platform earns an average profit of $(p - w)d$ for each customer request, the platform's average total profit is then equal to $\pi = \lambda(p - w)d$. By substituting (7.3) and (7.6) into the profit function, we can express the profit function π as a function of (k, s) below:

$$\pi(k, s) = \lambda d \left[F^{-1}(1 - s) - \frac{c}{d} W_q - G^{-1}\left(\frac{k}{K}\right) \frac{k}{\lambda d} \right]. \quad (7.7)$$

Considering the system stability condition $\lambda d < k\mu$, the optimization problem of the platform can be formulated as:

$$\begin{aligned} \max_{k, s} \pi(k, s) &\equiv \lambda d \left[F^{-1}(1 - s) - \frac{c}{d} W_q - G^{-1}\left(\frac{k}{K}\right) \frac{k}{\lambda d} \right], \\ &\text{subject to } k > \lambda d / \mu. \end{aligned} \quad (7.8)$$

By using the optimal number of participating providers k^* and the optimal demand rate λ^* via optimal s^* through (7.2), we can use (7.3) and (7.6) to retrieve the corresponding optimal price rate p^* and optimal wage rate w^* from k^* and λ^* .

7.4 The Base Model

To enable us to characterize the optimal solution to problem (7.8), we shall assume that the distribution of value rate v and the reservation wage rate r are uniformly distributed over the range $[0, 1]$ so that $F(v) = v$ and $G(r) = r$. Also, we shall approximate the (expected) waiting time W_q given in the customer's utility function (7.1) based on an $M/M/1$ queue with service rate $k(\mu/d)$ so that the wait-time function W_q has the following simple closed-form expression:

$$W_q = \frac{\lambda}{(k \cdot \mu/d) \cdot (k \cdot \mu/d - \lambda)} = \frac{\lambda d^2}{k\mu(k\mu - \lambda d)}. \quad (7.9)$$

More formally, we shall assume that the following assumption holds for the remainder of this paper.

Assumption 1 $F(\cdot) \sim U[0, 1]$, $G(\cdot) \sim U[0, 1]$, and $W_q = \lambda d^2 / [k\mu(k\mu - \lambda d)]$.

Under Assumption 1, the price, wage and profit functions given in (7.3), (7.6) and (7.7), respectively, can be simplified as:

$$p = (1 - s) - c \left(\frac{\lambda d}{k\mu} \right) \frac{1}{k\mu - \lambda d} \quad (7.10)$$

$$w = \frac{k^2}{K\lambda d} \tag{7.11}$$

$$\pi(k, s) = \lambda d \left[(1 - s) - c \left(\frac{\lambda d}{k\mu} \right) \frac{1}{k\mu - \lambda d} - \frac{k^2}{K\lambda d} \right]. \tag{7.12}$$

By using the above expressions, we can maximize the expected profit $\pi(k, s)$ given in (7.7) subject to the system stability constraint: $k > \lambda d/\mu$, and obtain the following results:

Proposition 1 *The optimal price p^* , the optimal wage w^* , and the platform’s optimal profit π^* exhibit the following characteristics:*

1. When K or μ increases, w^* decreases, π^* increases, but p^* is not necessarily monotonic.
2. When c increases, w^* increases, π^* decreases, but p^* is not necessarily monotonic.
3. When $\bar{\lambda}$ or d increases, w^* , p^* and π^* increase.

As given in the proof of Proposition 1, we can also derive some monotonicity properties on how the different model parameters affect the optimal service level s^* , the optimal number of providers k^* , the optimal expected wait-time W_q^* , the optimal customer request rate λ^* , and the optimal system utilization $\rho^* = \lambda^*d/(k^*\mu)$. We summarize these monotonicity properties in Table 7.1.

Proposition 1 shows that when the maximum number of potential providers K (or when the service speed μ) increases, the potential capacity of the system becomes larger. As such, the platform can increase the number of providers k^* and increase the service rate s^* (or the corresponding customer request rate λ^*) by lowering its wage rate w^* , and can obtain a higher profit π^* . However, when k^* and s^* (as well as λ^*) increase, Eq. 7.10 reveals that the optimal price rate p^* is not necessarily monotonic. This explains the first statement. This result implies that it is beneficial for the platform to recruit more potential service providers K to join the platform, and help (if possible) to increase their average service speed μ .

Next, when the waiting cost c increases, the platform should lower the service level s^* so as to reduce the corresponding customer request rate λ^* and expected wait-time W_q as given in (7.9). Consequently, the platform earns less. However, as

Table 7.1 Impact of model parameters on s^* , k^* , W_q^* , λ^* and ρ^*

	s^*	k^*	W_q^*	λ^*	ρ^*
K	↑	↑	↓	↑	×
μ	↑	×	↓	↑	×
c	↓	×	↓	↓	↓
$\bar{\lambda}$	↓	↑	↑	↑	↑
d	↓	↑	↑	↓	↑

↑(increasing); ↓(decreasing); ×(non-monotonic)

the optimal number of providers k^* is not necessarily monotonic, Eq. 7.10 reveals that the optimal price rate p^* is also not necessarily monotonic. This explains the second statement.

Finally, when the potential customer demand rate $\bar{\lambda}$ increases, the third statement reveals that the platform should increase its price rate p^* to increase the customer request rate λ^* (even though the service level s^* is actually lower since $\lambda^* = s^* \cdot \bar{\lambda}$), and increase its wage rate w^* so as to attract more providers k^* to participate. Overall, the platform earns a higher profit π^* when the potential customer demand rate $\bar{\lambda}$ increases. Also, when the average amount of service units d increases, it increases the overall workload to the system for each customer request and essentially has the same effect as increasing the customer demand rate $\bar{\lambda}$. Consequently, the optimal price rate, the optimal wage rate and the optimal profit behave the same. This explains the third statement.

While optimal price rate is not necessarily monotonic with respect to K , μ and c , we can prove the following monotonicity property of the optimal payout ratio w^*/p^* as the model parameters change.

Proposition 2 *The optimal payout ratio w^*/p^* increases in c , $\bar{\lambda}$ and d , and decreases in K and μ .*

Proposition 2 shows that the platform should increase the payout ratio w^*/p^* to its providers when the customer's waiting cost c is higher, the maximum customer demand rate $\bar{\lambda}$ is higher or the average amount of service units d is higher. On the other hand, the platform should reduce the payout ratio when the maximum number of potential service providers K or the average service speed μ increases. One interesting implication of this result is that it would be more profitable for an on-demand transportation service platform (such as Uber) to increase the payout ratio to its participating drivers when the customer demand rate $\bar{\lambda}$ is higher and/or the travel speed μ is lower during rush hours.

Proposition 2 also indicates that it is more profitable for the platform to lower its payout ratio when the number of registered providers K increases. It is interesting to note that this result is consistent with Uber's strategy as reported by Huet (2014) that Uber offered a payout ratio of 0.8 for its first cohorts of drivers in San Francisco initially, but Uber lowered its payout ratio to 0.75 for its second cohorts of drivers in 2014 (i.e., as the number of registered drivers increases). Therefore, this result provides an economic justification for Uber to reduce its payout ratio as K increases.

7.4.1 Special Case 1: When the Payout Ratio w/p Is Fixed

As many on-demand service platforms (such as Uber and Didi) have adopted a fixed payout ratio to their service providers, we can adapt our base model to analyze this special case by imposing an additional constraint $w/p = \alpha$. We can use a similar analysis to establish the following result.

Proposition 3 *Under the additional constraint that $w/p = \alpha$, $0 < \alpha < 1$, both the optimal wage rate w^* and the optimal price rate p^* increase in $\bar{\lambda}$ and d .*

When the payout ratio is held constant so that $w/p = \alpha$, Proposition 3 implies that the optimal price rate p^* (and thus the optimal wage rate w^* due to a fixed payout ratio) should both be higher when customer demand rate for service $\bar{\lambda}$ is higher or when the average amount of service units d is higher. (We remark that the optimal price and wages rates, however, are not necessarily monotone in either the number of available service providers K , the average service speed μ , or the unit waiting cost c .) Our results thus suggest that an on-demand transportation service platform of using a fixed payout ratio (such as Uber) should charge a higher price (and thus provide a higher wage rate) during rush hours when the customer demand is high. This result is consistent with the notion of “surge pricing” as adopted by Uber and Lyft; see Cachon et al. (2017) for some recent discussions on the role of surge pricing.

We note that both Propositions 1 and 3 reveal that when the customer demand is higher, the platform should charge a higher price rate and offer a higher wage rate, regardless of whether the payout ratio is variable or fixed. In Sect. 7.6, we shall further compare the optimal profits of the platform between these two different settings using some numerical examples motivated by the sample data provided by Didi.

7.4.2 Special Case 2: When the Service Level Is Exogenously Given

As on-demand service platforms continue to emerge and innovate, a new start-up platform might need to target a very high service level to ensure high customer satisfaction and gain popularity, at the expense of a lower near-term profit, during the initial phase of its operations. We can adapt our base model to analyze this special case by imposing a fixed target service level s . In other words, when the parameter s (or equivalently, the customer request rate λ because $\lambda = s\bar{\lambda}$) is exogenously given, the optimization problem of the platform is now reduced to:

$$\max_k \pi(k) \equiv \lambda d \left[(1 - s) - c \left(\frac{\lambda d}{k\mu} \right) \frac{1}{k\mu - \lambda d} - \frac{k^2}{K\lambda d} \right], \text{ subject to } k > \frac{\lambda d}{\mu}.$$

It is straightforward to show that the above profit function $\pi(k)$ is concave and we can determine the optimal number of participating providers k^* using the first-order condition. Then, we can use (7.10) and (7.11) to retrieve the corresponding optimal price rate p^* and optimal wage rate w^* from the value of k^* . The following proposition summarizes the main results for this special case.⁴

⁴The results of Proposition 4 continue to hold under more general distributions $F(\cdot)$ or $G(\cdot)$ and general wait-time function W_q . For ease of exposition, we shall relegate the details to an appendix, which is available from the first author.

Proposition 4 *The optimal price p^* , the optimal wage w^* , and the optimal profit π^* exhibit the following characteristics:*

1. *When K or μ increases, p^* increases, w^* decreases, and π^* increases.*
2. *When c increases, p^* decreases, w^* increases, and π^* decreases.*
3. *When $\bar{\lambda}$ or d increases, p^* decreases and w^* increases.*

Proposition 4 can be interpreted as follows. When the maximum number of potential providers K or the service speed μ becomes higher, the potential capacity of the system increases. The first statement asserts that it is then optimal for the platform to charge a higher price p^* (because of lower wait-time due to higher capacity), offer a lower wage w^* (because there are plenty of potential providers), and earn a higher profit π^* . The second statement states that when customers become less patient (i.e., when c increases), the platform should lower its price p^* (to compensate for the higher waiting cost), offer a higher wage w^* (to entice more providers to offer service), and consequently, the platform earns less. Finally, when the customer request rate λ (or equivalent, $\bar{\lambda}$, as $\lambda = s\bar{\lambda}$ and s is fixed) or the average amount of service unit d increases, the average workload of the system increases. As such, the third statement reveals that the platform should lower its price p^* to compensate for the higher waiting cost and offer a higher wage w^* to entice more providers to participate.

By comparing the results of Propositions 1 and 4, one can observe that most of the results remain the same except for the characteristics of the optimal price rate p^* . When the service level s is endogenously determined, the optimal price rate p^* is not necessarily monotonic for certain model parameters, as stated in the first two statements of Proposition 1. However, when $\bar{\lambda}$ increases, the third statement of Proposition 1 reveals an opposite result, i.e., the optimal price rate p^* increases, versus p^* being decreasing in $\bar{\lambda}$ as given in Proposition 4. We can explain this opposite result as follows. When $\bar{\lambda}$ increases under a given s , the customer demand rate $\lambda = s\bar{\lambda}$ increases, and consequently, the platform has to offer a higher wage rate w^* to increase the number of providers k^* . Without the flexibility to adjust s , one can use (7.10) and the fact that $\lambda = s\bar{\lambda}$ to show that the corresponding optimal price rate p^* would decrease as stated in Proposition 4. On the other hand, when s (or thus the customer request rate λ) is endogenously determined, the platform has the flexibility to charge a higher price rate p^* and offer a higher wage rate w^* to better coordinate demand and supply as stated in Proposition 1.

We point out that Proposition 4 confirms the results (as also shown in Proposition 1) that it is always beneficial for the platform to recruit more potential service providers to join the platform (i.e., increase K), and to help (if possible) providers to increase their average service speed μ , regardless of whether the service level is fixed or endogenously determined.

We can use the results of Proposition 4 to characterize the optimal payout ratio w^*/p^* as follows:

Corollary 1 *The optimal payout ratio w^*/p^* increases in c , $\bar{\lambda}$ and d , but decreases in K and μ .*

The results in Corollary 1 are consistent with those given in Proposition 2. In particular, the platform should still increase the payout ratio w^*/p^* to its providers when the customer's waiting cost c is higher or the maximum customer demand rate $\bar{\lambda}$ is higher, but should reduce the payout ratio when the maximum number of potential service providers K or the average service speed μ increases, even if the platform intends to maintain a constant target service level s at any time.

7.5 Numerical Illustrations Based on Didi Data

To illustrate the implications of our analytical results presented in this paper, we have collected real data from Didi, the largest on-demand ride sharing service platform operating in over 360 cities in China.⁵

7.5.1 Background Information

Our data was based on rides that took place in the city of Hangzhou, the capital city of Zhejiang province with an urban population of over seven millions people, during the time periods between September 7–13 and November 1–30 in 2015.

In Hangzhou city, Didi has approximately 13,000 registered drivers offering different types of services including Taxi (traditional taxi service),⁶ Express/Private (equivalent to UberX/Black with on-demand drivers), and Hitch (passengers sharing similar routes). For our numerical illustrations, we shall focus on the data associated with the Express/Private service, which accounts for 60% of all rides provided by Didi in Hangzhou. There were 13,000 registered drivers for all services, but the exact number of Express/Private drivers was not known to us. We shall assume that 60% of Didi drivers were Express/Private drivers, i.e., the number of registered Express/Private drivers in Hangzhou city was estimated to be $K = 7,800$.

⁵<http://www.xiaojukeji.com/en/company.html>. Didi was founded in June 2012 and merged with Kuaidi (a major competitor) in February 2015 as a way to defend its market share when Uber officially launched its service in China in July 2014. In August 2016, Uber decided to retreat from China and its China operations merged with Didi.

⁶Unlike Uber's business model that aims to displace the traditional taxi services, Didi integrates taxi services into its business model by providing its mobile hailing service to taxi drivers free of charge.

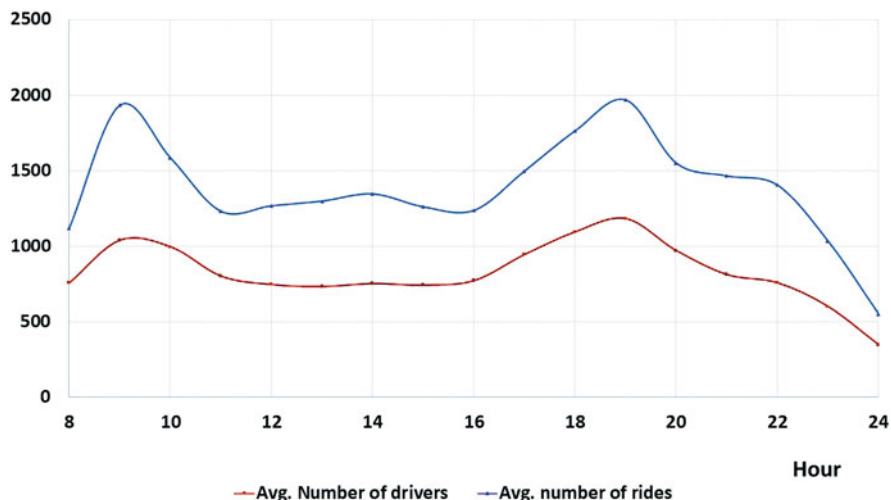


Fig. 7.1 Number of rides and drivers across different hours

7.5.2 Number of Rides and Drivers Across Different Hours

Figure 7.1 depicts the average number of Express/Private rides and drivers across different hours on any given day. (Here, Hour 8 represents 1-hour interval 8 a.m.–9 a.m., Hour 19 for 7 p.m.–8 p.m., and so on. Data for Hours 1–7 were omitted due to incomplete data in the database.) We observe from the Didi data that the pattern depicted in Fig. 7.1 is consistent throughout the weekdays, even though the average number of rides and drivers were slightly lower on Saturdays and Sundays, and that the peak hours are being Hours 9 and 19, and the slowest hours are being Hours 23 and 24. For instance, during the peak Hour 19, there were an average of 1,969 Express/Private rides and an average of 1,182 drivers in any given day. However, during the late night Hour 23, there were only an average of 1,033 rides and an average of 600 drivers.

7.5.3 Travel Distance and Travel Speed

While the average number of rides and drivers vary substantially across different hours of the day, it is interesting to note from Fig. 7.2 that the average travel distance for each Express/Private ride is rather stable across different hours. For example, the average travel distance d during the peak Hour 19 and during the late night Hour 23 were 6.3 and 6.6 km, respectively, while the average price per km p charged by Didi during these 2 hours were RMB 3.13 and RMB 2.76; respectively. We can also estimate the average travel speed across hours μ , and they were about 19

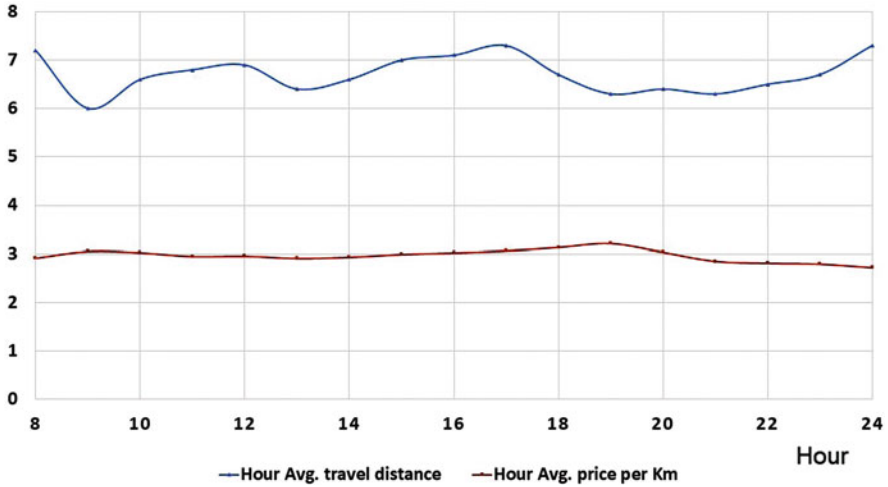


Fig. 7.2 Average travel distance and average price per kilometer across different hours

and 26 km/hour for Hour 19 and Hour 23, respectively. These numbers are thus consistent with the actual expected traffic conditions, which show that traffic is much less congested during late night hours.

7.5.4 Pricing and Wage Rates

Didi’s price p for its service consists of two components so that $p = p_1 + p_2$, where p_1 represents the fare that is primarily based on the travel distance, and p_2 represents surcharges (e.g., tolls). Accordingly, Didi paid its drivers according to the following scheme. When a passenger pays a total fee of p , the driver will receive $(p_1 \cdot 80\% - 0.5) \cdot (100\% - 1.77\%) + p_2 \cdot (100\% - 1.77\%)$ from Didi, but the driver needs to pay p_2 to cover various surcharges. Thus, the actual wage that Didi pays its drivers is approximately 80% of the total price; i.e., $w \approx 0.8p$. Figure 7.2 also shows that the average price per kilometer charged by Didi (excluding the surcharges) is relatively stable across different hours of the day. However, we observe that the average price per km p charged by Didi was RMB 3.13 during the peak hour (i.e., hour 19) and it was RMB 2.76 during the non-peak hour (i.e., hour 23). It is interesting to note that the observed price is higher during the peak hours, which corroborates with our results as stated in Proposition 3 that the optimal price rate p^* should be higher when customer demand rate for service is higher.

7.5.5 *Strategic Factors and Their Implications*

It is important to note that the observed price that Didi charged its passengers was heavily discounted or subsidized during the data collection periods for the following two strategic reasons: (a) Didi wanted to attract more passengers by pricing its service below the traditional taxi services⁷; and (b) Didi was engaged in a price war to compete with Uber by offering discount coupons to compete for market share. In addition to offering heavily discounted price to attract passengers, Didi also provided extra “side payments” to its drivers to entice more drivers to join their platform due to the intense market competition. In addition to the regular payments of approximately 80% of the fare collected from the passengers, Didi (and Uber) had offered extra payments and additional bonus (e.g., Didi offers an extra bonus if the number of rides provided by a driver exceeds a certain quota within a 7-day period).

While we were unable to obtain the details of the bonus scheme, BBC (2016) had reported that the extra payment can be as high as 110% of the fare paid by the passengers. With such generous payments, more drivers reported to work and there was no need for Didi to use surge pricing to attract more drivers to offer rides during peak hours. This explains why Didi was able to offer relatively stable pricing in Hangzhou as depicted in Fig. 7.2. Furthermore, the waiting time for Didi’s service was reasonably short with an adequate supply of drivers. Specifically, the average waiting time of all Express/Private rides over the aforementioned time periods was about 6 min, of which the waiting time for accepting a ride request was approximately 1 min, while the waiting time for picking up a passenger was approximately 5 min.

In view of the heavily discounted price due to the above strategic reasons, the price per km p as reported in Fig. 7.2 was biased and did not represent the regular prices p the firm should quote and the actual wages w should offer in equilibrium. Nevertheless, we shall use the data given in the Didi database to construct our numerical examples below to illustrate how our analytical results would compare with the actual prices/wages as reported in Didi’s data set.

7.5.6 *Numerical Examples for Illustrative Purposes*

We next provide some numerical examples based on the Didi data to illustrate our model results and discuss their implications. In all our numerical examples, we set the maximum number of drivers $K = 7,800$. We examined the average income for

⁷In Hangzhou, taxi charges RMB 11 initially and then RMB 2.6 per km. As a way to entice passengers to choose Didi over taxi service, Didi had priced its service below taxi rates to increase market share. Based on our discussions with passengers in China, there was an expectation that Didi’s price rate was lower than the taxi rate.

taxi drivers in Hangzhou city and the average major out-of-pocket expenses borne by the Didi drivers (including car insurance, license, fuel cost, etc.), and estimated that a minimum hourly wage of RMB 30 is required for a Didi driver to offer service. Thus, we assume that the hourly wage reservation r is distributed uniformly between RMB 30 to RMB 40 in our numerical examples.

As discussed earlier, the data were collected during the time when Didi (and Uber) was offering large fare discounts to attract riders, and so there was an expectation among riders that Didi price was less than the taxi rate in Hangzhou (which is RMB 2.6 per km). Thus, we used the taxi rate as a benchmark and assume that the customer value per km v is distributed uniformly between RMB 3 to RMB 4 in our numerical examples.

As shown in Fig. 7.2, the average travel distances did not vary significantly across hours, so we simply set the average travel distance $d = 6$ km across all hours in order to focus our discussions on how different demand and congesting levels would affect the optimal price and wage rates across different hours of the day. It is difficult to provide an accurate estimate of the waiting cost parameter c , and so we simply varied the value of c from RMB 200 to RMB 2,200 per hour to illustrate how the optimal price and wage rates would change with respect to the cost of customer waiting for service.⁸

We used data from two specific time periods to show our model results as illustrative examples for our discussions. In particular, we picked Hour 19 to represent peak-hour characteristics with high demand and travel congestion levels, and Hour 23 to represent non-peak hour characteristics with lower demand and congestion levels. Specifically, we set the average demand $\bar{\lambda} = 2,000$ with an average service speed $\mu = 19$ km/hour, and $\bar{\lambda} = 1,000$ with an average service speed $\mu = 26$ km/hour, respectively, in these two scenarios. In each scenario, we solved the base model as discussed in Sect. 7.4. The optimal number of participating drivers k^* , price rate p^* and wage rate w^* are given in Figs. 7.3 and 7.4 for the peak hour and non-peak hour scenarios, respectively.

These numerical results illustrate the properties as stated in Proposition 1. For example, the optimal wages w^* increases as the waiting cost c increases in both Figs. 7.3 and 7.4, which illustrates the results as stated in statement 2 of Proposition 1. By comparing the results in Figs. 7.3 and 7.4, we can also observe that the values of k^* (scale on the left), p^* and w^* (scale on the right) are all higher during the peak hour (Fig. 7.3) than those during the non-peak hour (Fig. 7.4). These properties illustrate the results as stated in statements 1 and 3 of Proposition 1 (and

⁸While it is difficult to estimate the waiting cost, Gómez-Ibáñez et al. (1999) reported that the waiting cost for a working class passenger in San Francisco is approximately 195% of the passenger's after-tax wages. Using this estimate and the fact that the average hourly wage of workers in Hangzhou is approximately RMB 40 per hour (Wu 2016), one can argue that the waiting cost for an average passenger in Hangzhou is approximately RMB 80 per hour. However, accounting for the income inequality and the impatient characteristics of most city dwellers in China (Li 2016), we simply choose to consider the range of c varying from RMB 200 to RMB 2200 for illustrative purposes.

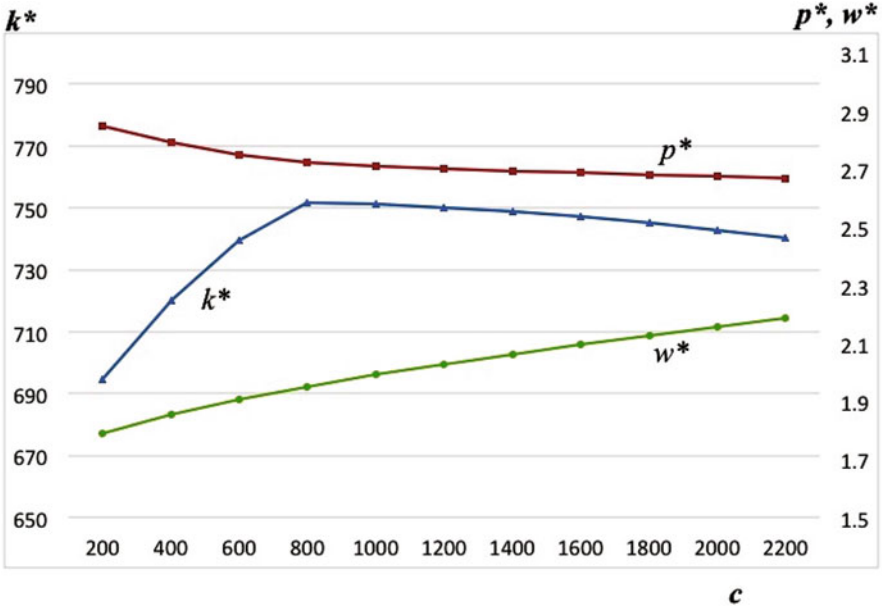


Fig. 7.3 Optimal number of participating drivers, optimal price and wage rates during peak hours ($\bar{\lambda} = 2000$ and $\mu = 19$ km/hour)

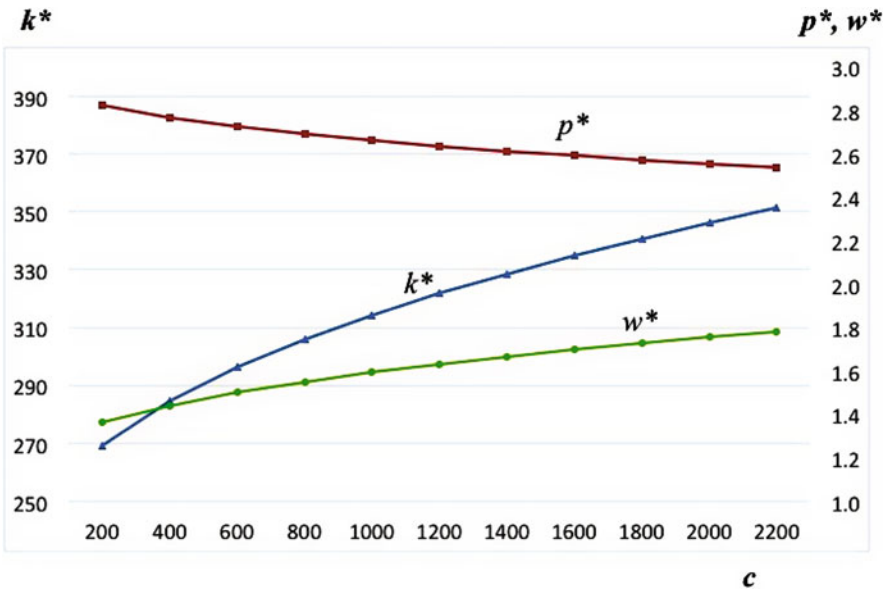


Fig. 7.4 Optimal number of participating drivers, optimal price and wage rates during non-peak hours ($\bar{\lambda} = 1000$ and $\mu = 26$ km/hour)

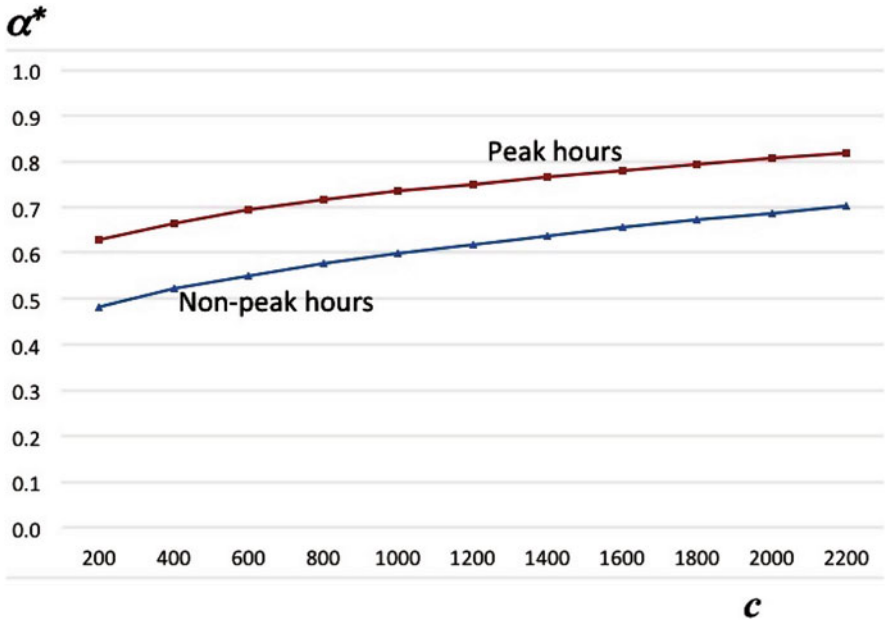


Fig. 7.5 Comparisons of the optimal dynamic payout ratios between peak and non-peak hours

Table 7.1), because the peak hour period has a higher customer demand rate $\bar{\lambda}$ and a slower service speed μ than that during non-peak hour period. However, the optimal number of participating driver k^* is not monotonic in the waiting cost c . In both scenarios, the optimal price rate p^* decreases as c increases. (However, the optimal price rate p^* is not necessarily monotonic in c (in general), as noted in statement 2 of Proposition 1.)

We also computed the optimal payout ratio $\alpha^* = w^*/p^*$; see Fig. 7.5. The optimal payout ratio α^* increases from 0.68 to 0.84 for the peak hour scenario and from 0.54 to 0.72 for the non-peak hour scenario, respectively, as c increases from 200 to 2,200. Also, observe from Fig. 7.5 that the optimal payout ratio is always higher for the peak hour scenario than that for the non-peak hour scenario for any fixed value of c . (We note that these monotonicity properties are proved in Proposition 2.)

As Didi used a fixed payout ratio $\alpha \approx 0.8$, it would be interesting to compare the corresponding optimal profit between using the dynamic payout ratio α^* as given in our model versus using a fixed payout ratio $\alpha = 0.8$ to examine the potential benefits of adopting the optimal dynamic payout ratio. We illustrate our results in Fig. 7.6 based on the peak hour scenario (i.e., Hour 19). Our numerical results show that, during these peak hour periods, using a dynamic payout ratio α^* can substantially increase the profit of the service platform over that using a fixed payout ratio of 0.8, especially when the waiting cost c is low when the optimal payout ratio is significantly different from 0.8 in our numerical examples here. For

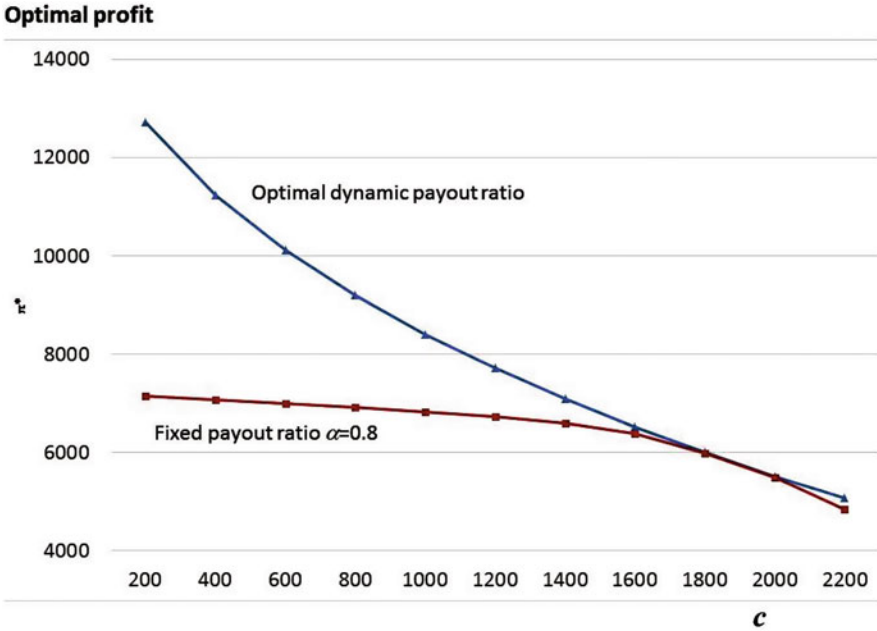


Fig. 7.6 Comparisons of optimal profit between the optimal dynamic payout ratio and a fixed payout ratio for the peak hour scenario

instance, when $c = 600$, the optimal profit is equal to 10,115 when the platform uses the optimal payout ratio $\alpha^* = 0.69$. However, if the platform uses a fixed payout ratio of 0.8, then the profit is equal to 7,001, which is much lower. However, it is important to point out that Didi (and Uber as well) used a fixed payout ratio due to various market considerations such as intense competition for driver participation and ridership as well as other practical implementation issues. Nevertheless, our results can serve as a guideline for understanding the magnitude of potential benefits for a hypothetical situation where such market considerations were no longer valid. Specifically, Fig. 7.6 suggests that, when the waiting cost c is low, using a dynamic payout ratio can enable the platform to earn a much higher profit.

7.6 Conclusion

In this chapter, we develop an analytical framework to understand how on-demand service platforms should set their optimal price and wage to match the needs of providers and customers taking into considerations the underlying supply and demand characteristics. The framework consists of a queueing model that captures some important market characteristics including time-sensitive customers and earning-sensitive suppliers. We analyze the steady state performance of a two-sided

queue in equilibrium and investigate the behavior of the optimal price and wage rates. We derive analytical results to show how different model parameters would affect these optimal price and wage rates. Our findings provide some interesting implications in managing prices and wages for on-demand service platforms.

Using some actual data collected from a major ride-sharing company in China, we construct some numerical examples to illustrate the results of our analytical model and discuss various implications on the optimal price and wage with respect to the underlying market characteristics. Although our model does not capture some important practical issues due to intense competition existed in China when the data were collected (and thus cannot be used to accurately predict the actual behavior of the players in the market), our analytical results can help to illustrate and explain some general observations that are consistent with the actual data provided by the company. More importantly, our model results can serve as a guideline for potentially increasing profitability when the underlying market conditions were to evolve to be consistent with the operating environment captured in our modeling framework. Specifically, we illustrate the potential benefits if the company were to adopt a time-based payout ratio versus their current practice of using a fixed payout ratio.

Acknowledgements The authors are grateful to Didi Chuxing (www.xiaojukeji.com) for providing us some sample data. The authors also thank Professors Ming Hu and Terry Taylor for their constructive comments on an earlier version of this paper. This paper was completed when the third author was serving as a visiting professor of the Institute for Advanced Study at the Hong Kong University of Science and Technology.

References

- Afanasyev M, Mendelson H (2010) Service provider competition: delay cost structure, segmentation and cost advantage. *Manuf Serv Oper Manag* 12(2):213–235
- Afeche P, Mendelson H (2004) Pricing and priority auctions in queueing systems with a generalized delay cost structure. *Manag Sci* 50(7):869–882
- Allon G, Bassamboo A, Çil EB (2012) Large-scale service marketplaces: the role of the moderating firm. *Manag Sci* 58(10):1854–1872
- Armony M, Haviv M (2003) Price and delay competition between two service providers. *Eur J Oper Res* 147(1):32–50
- BBC (2016) Uber ‘hide’ surge pricing notifications. <http://www.bbc.com/news/technology-36619414>. Accessed 24 June 2016
- Cachon GP, Daniels KM, Lobel R (2017) The role of surge pricing on a service platform with self-scheduling capacity. *Manuf Oper Serv Manag* 19(3):347–367
- Chen MK, Sheldon M (2015) Dynamic pricing in a labor market: surge pricing and flexible work on the Uber platform. Working paper, UCLA Anderson School
- Modaran A (2014) A disruptive cab ride to riches: the Uber payoff. *Forbes* (June 10) <http://www.forbes.com/sites/aswathmodaran/2014/06/10/a-disruptive-cab-ride-to-riches-the-uber-payoff/>
- DePhills L (2016) One reason you might be better off driving for Uber than in a taxi. *The Washington Post* (Mar 15)

- Gómez-Ibáñez J, Tye W, Winston C (1999) *Essays in transportation economics and policy*, vol 42. Brookings Institution Press, Washington, DC
- Gurvich I, Lariviere M, Moreno-Garcia A (2016) Operations in the on-demand economy: starting services with self-scheduling capacity. Technical report, Northwestern University. https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2336514
- Hu M, Zhou Y (2017) Price, wage and fixed commission in on-demand matching. Working paper, Rotman School of Management, University of Toronto. https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2949513
- Huet E (2014) Uber now taking its biggest UberX commission ever – 25 percent. *Forbes*, 22 Sept 2014
- Kokalitcheva K (2015) Uber and Lyft face a new challenger in Boston. *Fortune.com* (Oct 5)
- Li A (2016) Why are Chinese tourists so rude? A few insights. *South China Morning Post*, 10 Aug 2016
- Moreno A, Terwiesch C (2014) Doing business with strangers: reputation in online service marketplaces. *Inf Syst Res* 25(4):865–886
- Naor P (1969) The regulation of queue size by levying tolls. *Econometrica* 37(1):15–24
- Riquelme C, Banerjee S, Johari R (2015) Pricing in ride-share platforms: a queueing-theoretic approach. Working paper, Stanford University. https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2568258
- Roose K (2014) Does Silicon Valley have a contract-worker problem? *NYMag.com* (Sept 18)
- Shoot B (2015) Hot food, fast. *Entrepreneur* 68 (Aug)
- Taylor T (2017) On-demand service platforms. Working paper, University of California, Berkeley. https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2722308
- Wirtz J, Tang CS (2016) Uber: competing as market leader in the US versus being a distant second in China. Case study in: Wirtz J, Lovelock CH (eds) *Service marketing: people, technology and strategy*, 8th edn. World Scientific, Hackensack
- Wu Y (2016) Average salary in major Chinese cities is US\$900 and growing. *China Daily* (Jan 21). http://www.chinadaily.com.cn/china/2016-01/21/content_23183484.htm
- Zhou W, Chao X, Gong X (2014) Optimal uniform pricing strategy of a service firm when facing two classes of customers. *Prod Oper Manag* 23(4):676–688

Chapter 8

Pricing and Matching in the Sharing Economy



Yiwei Chen, Ming Hu, and Yun Zhou

Abstract Sharing economy platforms use crowdsourced suppliers to provide customers with services or goods. Their decision making often revolves around pricing and matching. Platforms like Uber charge the customers a price for using the services or goods and offer the crowdsourced suppliers a wage or pay for providing the services or goods. First, we study how the platform could optimally set the price and the wage for a single service or product in different market conditions, and investigate the performance of the fixed commission contract which uses a fixed commission percentage across all market conditions. Second, even with determined pricing decisions, the platform also faces the task of matching customers with suppliers. We consider a stochastic, dynamic model with multiple demand types to be matched with multiple supply types over a planning horizon. We characterize the optimal matching policy by determining the priorities of the demand-supply pairs, under a sufficient condition on the reward structure. Then, the results are applied to two cases with more specific reward structures; namely, the horizontal reward structure and the vertical reward structure, to better characterize the optimal policy. Finally, we study the joint pricing and matching decision by a platform for a single service or product and take into account suppliers' and customers' forward-looking behavior. We propose a simple heuristic policy: fixed price and wage plus waiting compensation, in conjunction with the greedy matching policy on a first-come-first-served basis. This heuristic policy induces forward-looking suppliers and customers to behave myopically and is shown to be asymptotically optimal.

Y. Chen

Carl H. Lindner College of Business, University of Cincinnati, Cincinnati, OH, USA
e-mail: chen3yw@ucmail.uc.edu

M. Hu (✉)

Rotman School of Management, University of Toronto, Toronto, ON, Canada
e-mail: ming.hu@rotman.utoronto.ca

Y. Zhou

DeGroote School of Business, McMaster University, Hamilton, ON, Canada
e-mail: zhouy185@mcmaster.ca

8.1 Introduction

The rise of the sharing economy has revolutionized industries such as taxi and hotel. Empowered with information technology, sharing economy platforms connect crowdsourced suppliers with customers. Uber, the venture-funded startup and one of the pioneers in the sharing economy, for example, matches a prospective passenger with a space in a nearby car owned by someone else. Uber's service has now expanded to over 300 cities in the world, recording its two billionth ride in 2016. Together with other ridesharing services such as Sidecar and Lyft, the company is leading a disruption to the taxi industry. The car rental company FlightCar offers travelers free parking at the airport, and then rent out their cars in exchange. In the logistics industry, firms are starting to hire independent couriers for last mile deliveries. One example is Amazon's "Amazon Flex" program, which uses independent contractor drivers to deliver packages.

8.1.1 *Two-Sided Pricing*

Pricing decisions are crucial for sharing economy platforms like Uber. Just like price-sensitive customers on the demand side, crowdsourced suppliers are also sensitive to their monetary rewards for providing services. Thus wage for suppliers and price for customers are the two critical controls for the platform such as Uber and Lyft to coordinate supply and demand. With time-varying market conditions with different supply and demand patterns, platforms often need to adjust its price and wage dynamically.

Sharing economy platforms in real life often use the fixed commission rate contract when they set their price and wage. Under this contract, all parties agree that the platform takes a fixed percentage of the revenue regardless of the market condition. Uber started its business taking a 20% commission on all rides and now has raised and lowered that rate in different cities depending on the supply of drivers and rider demand. Lyft currently charges a 25% commission. As supply and demand conditions change over time, the platform would ideally want to update both price and wage accordingly to match supply with demand. With the fixed commission contract, the pricing decision and its associated wage determined by the fixed commission contract is suboptimal. In this chapter, we investigate the performance of the best fixed-commission contract against the model in which the platform set price and wage optimally in each scenario to maximize its profit.

8.1.2 *Two-Sided Matching*

As intermediaries, platforms in the sharing economy also face the task of managing the matching between crowdsourced suppliers and customers. For example, carpooling platforms such as iCarpool and UberPool match a driver heading to a

Table 8.1 Comparisons between Monge sequence and modified Monge conditions

Monge sequence	Modified Monge condition
<i>Static, deterministic and balanced</i> transportation problem on a sequence sufficient and necessary a greedy algorithm: (1) priority property (2) <i>match as much as possible</i>	<i>Dynamic, stochastic and unbalanced</i> matching problem on pairs sufficient, and <i>robustly</i> necessary our result: (1) priority property (2) <i>match-down-to policy</i>

destination with several riders to the same destination (or in the same direction). Under “inventory commingling program,” its Amazon crowdsources inventories of an identical item from third-party merchants to its warehouses, to fulfill online orders.¹ A nonprofit organization, United Network for Organ Sharing (UNOS), allocates donated organs to patients in need of transplantation. In those examples, both supply and demand may have heterogeneous types. For instance, drivers and customers are differentiated by their locations; Organs and patients differ in health status. In addition to heterogeneous types, the arrivals of supply and demand are usually random and beyond the direct control of the platform. In this chapter, we formulate the intermediary firm’s dynamic matching problem as a discrete-time stochastic dynamic program. We derive structural properties of optimal matching policies and develop good heuristic policies.

More specifically, we establish the *modified Monge condition* that specifies a dominance relation between two pairs of demand and supply types. The modified Monge conditions are sufficient and robustly necessary for the optimal matching policy to satisfy the following priority properties in the dynamic matching problem. First, for *any* two pairs of demand and supply types with one strictly dominating the other, it is optimal to prioritize the matching of the dominating pair over the dominated pair. Second, it is optimal to greedily match a *perfect pair* of demand and supply types that dominates all other pairs sharing its demand or supply type. The modified Monge condition generalizes the condition of a Monge sequence, discovered by Gaspard Monge in 1781, which guarantees a static and balanced transportation problem to be solved by a greedy algorithm (see Table 8.1 for comparisons). As a result of the priority properties, the optimal matching policy boils down to a match-down-to structure (instead of matching as much as possible in the greedy algorithm) when considering a specific pair of demand and supply types, along with the priority hierarchy. In fact, in the optimal policy, if *some* pair of demand and supply types is not matched as much as possible, *all* pairs that are strictly dominated by this pair should not be matched at all.

¹A product ordered from Amazon or a third-party supplier may not have originated from the original supplier. The program gives Amazon the flexibility to ship products on the basis of their geographic proximity to customers, thus shortening delivery times and reducing shipping costs.

While two pairs of demand and supply types that share a common node may not be comparable under the modified Monge condition, the priority properties continue to hold for those pairs that indeed satisfy the modified Monge conditions, even when not all pairs are comparable. In addition, we provide bounds and heuristics for the general problem as follows.

As a heuristic method, we consider the deterministic counterpart of the stochastic dynamic problem for any period with the t amount of remaining time in the horizon and any given levels of demand and supply; this can be written as a linear program with $O(n \times m \times T)$ variables. We show that the deterministic model provides an upper bound on the optimal total surplus of the stochastic model, and that it is asymptotically optimal to re-solve the linear program for the current period and state and apply the solution as a heuristic policy, when the arrival rate of demand and supply of every type becomes increasingly large.

8.1.3 Pricing and Matching Under Strategic Behavior

Suppliers and customers can time their transactions based on market prices and their expectation of the likelihood of being matched. For example, freelance drivers can choose when to work, and luxury fashion customers may time their purchases. Chen et al. (2015) observe that Uber's customers often prefer to "wait out" the price surges. Then it is desirable to design market-making policies, i.e., pricing and matching policies that take into account the strategic behavior of customers and suppliers.

We propose a simple market-making policy in the market of a single product or service at which forward-looking customers and suppliers with stationary valuation and cost distributions arrive following Poisson processes with constant rates over a finite horizon. On the pricing side, the market maker posts *fixed* ask and bid prices plus a (time-dependent) price adjustment as compensation for the expected cost of waiting to be matched, and on the matching side, the market maker implements the *greedy* matching policy on a first-come-first-served basis. Those fixed prices balance demand and supply and can be computed efficiently. The commitment to the fixed base prices with a price adjustment for expected waiting cost induces the strategic customers and suppliers to behave myopically. They will submit a request for matching upon arrival without delay if the customer's valuation is no less than the fixed ask price or the supplier's cost is no more than the fixed bid price. The price adjustment on top of the fixed prices for a customer or a supplier arriving at a specific time is the expected cost of waiting to be matched after submitting his or her matching request. The effective prices offered after compensating for waiting are in general time-varying and tend to have opposite trends at the beginning and the end of the horizon. But we show that when the volumes of demand and supply are large, the compensation for waiting becomes negligible (more specifically, it is $O(1/n^{1/3})$ where n is the scaling of the arrival rates) and the effective price

trajectory tends to be stationary; moreover, when customers are willing to wait for a certain length of time without being compensated, the heuristic price trajectory can become stationary.

We use the mechanism design approach to establish a closed-form upper bound of the intermediary's optimal profit. We show that the relative profit loss of our simple market-making policy relative to the profit upper bound is $O(1/n^{1/2})$ as the arrival rates on both sides are scaled up by n . That is, our simple market-making policy is nearly optimal as the market size grows up.

8.2 Two-Sided Pricing and Fixed Commission

In this section, we consider a platform that coordinates the matching of customer demand with crowdsourced supply by optimally setting prices for both the demand and supply sides.

8.2.1 The Price and Wage Optimization Problem

Let $\mathcal{K} = \{1, 2, \dots, K\}$ be a set of possible scenarios of market conditions. Scenario k occurs with probability ρ_k . Each scenario is characterized by a demand curve and a supply curve. In Scenario k , the total amount of customers who are willing to pay for the service at a price of p is $d_k(p)$. We call $d_k(p)$ the *raw demand function* which naturally should satisfy the downward-sloping property; namely, $d_k(p)$ is a *decreasing* function in p .² As is consistent with the operations literature, the number of satisfied customers may be capped by the available supply. Given a posted wage w , the total amount of independent contractors or suppliers, who are willing to provide the service, is $s_k(w)$. This is the number of suppliers who would show up if they were guaranteed to be matched with a customer. We call $s_k(w)$ the *raw supply function*, which is assumed to be increasing in w . For simplicity, we assume $d_k(p)$ and $s_k(w)$ are continuous functions. One can also interpret $d_k(p)$ as the tail probability function of a customer's willingness-to-pay (up to a constant scalar), and $s_k(w)$ of the cumulative distribution function (c.d.f.) of a supplier's willingness-to-sell (up to a constant scalar). We assume that $s_k(0) = 0$ (i.e., no supplier would be willing to join the market if the wage is 0), and that $\lim_{p \rightarrow \infty} d_k(p) = 0$ (i.e., demand would be choked off if the price is set outrageously high).

The sequence of events from the perspective of the platform is as follows.

In Stage 1, the nature selects a scenario. Given a realized Scenario k , the platform decides on the price p_k , with the wage following from the pre-determined commission contract as $w_k = f(p_k)$.

²The monotonicity in this chapter is in its weaker sense unless otherwise specified.

In Stage 2, all customers and suppliers observe p and w and play a simultaneous game by deciding on whether to enter the marketplace. The platform clears the market by matching arriving demand and supply. If there are more suppliers than customers in the marketplace, suppliers are rationed to be matched with a customer, and vice versa. The rationing rule can be arbitrary and is announced upfront.

Given the observed price p and wage w , we characterize the behavior of customers and suppliers in a specific Scenario k . We allow customers or suppliers to anticipate their chances of being matched and assume that they know the demand curve $d_k(p)$ and supply curve $s_k(w)$. They join the market if their surplus of joining is nonnegative. Somewhat surprisingly, the following result shows that the matching quantity in equilibrium can be computed as if the customers or suppliers are naïve in the sense that they do not anticipate the likelihood of being matched. All proofs in this section can be found in Hu and Zhou (2018b).

Proposition 1 *For a given price p and wage w , the matching quantity when the suppliers and customers strategically anticipate the matching likelihood, is equal to $\min\{d_k(p), s_k(w)\}$.*

By Proposition 1, the profit of the platform in any Scenario k given price p and wage w , $\pi_k(p, w)$, is the product of the profit margin $p - w$ and the total matching quantity, i.e., $\pi_k(p, w) = (p - w) \min\{d_k(p), s_k(w)\}$. The platform's objective is to solve the following problem to maximize its profit in Scenario k .

$$\max_{p \geq 0, w \geq 0} \pi_k(p, w). \quad (8.1)$$

We proceed to solve the above problem. The following proposition shows that instead of maximizing the profit with respect to p and w , we can first find the optimal matching quantity z_k^* , from which the optimal price p_k^* and wage w_k^* can be recovered.

Theorem 1 (Maximization of platform profit) *Let $z_k^* \in \arg \max_{z \geq 0} [d_k^{-1}(z) - s_k^{-1}(z)]z$. Then the optimal price and wage are, respectively,*

$$p^* = d_k^{-1}(z_k^*) \quad \text{and} \quad w_k^* = s_k^{-1}(z_k^*),$$

where $d_k^{-1}(z) \equiv \max\{p \geq 0 \mid d_k(p) = z\}$ and $s_k^{-1}(z) \equiv \min\{w \geq 0 \mid s_k(w) = z\}$.

Once we solve (8.1), the optimal expected profit across all scenarios is $P^* = \sum_{k=1}^K \rho_k \pi_k(p_k^*, w_k^*)$.

Theorem 1 reduces the two-dimensional price and wage optimization problem to a one-dimensional problem. The intuition behind is that for a given scenario, in anticipation of the market formation, there is no incentive for the platform to set the price and wage such that there is more supply arriving than the demand or vice versa. That is, it is optimal for the platform to set the price and wage such that the arriving demand is equal to the arriving supply. Hence, the problem reduces to finding the optimal matching quantity, from which the optimal price and wage can be obtained.

8.2.2 The Fixed Commission Contract

Under a fixed commission contract, the following step precedes Steps 1 and 2 in the sequence of events.

Stage 0: the platform can decide and commit to a fixed commission rate, i.e., the wage is always a fixed fraction $\gamma \in (0, 1)$ of the price.

We present an illustrative example of the fixed commission contract.

Example 1 Suppose that in a scenario $k \in \mathcal{K}$, the supply function is $s_k(w) = s_{k0}F_k^s(w)$ and the demand function is $d_k(p) = d_{k0}[1 - F_k^d(p)]$. Here, F_k^s (resp. F_k^d) is the c.d.f. of the conditional normal distribution (conditioned on nonnegative values) with mean μ_k^s (resp. μ_k^d) and standard deviation σ_k^s (resp. σ_k^d). The supply and demand curves are obtained assuming that a supplier joins the market if and only if the wage exceeds his/her willingness-to-sell (opportunity cost for providing services), a customer joins the market if and only if the price is below his/her willingness-to-pay (valuation of the service), and both supplier’s willingness-to-sell and customer’s willingness-to-pay follow conditional normal distributions. The parameters s_{k0} and d_{k0} represent the numbers of the potential suppliers and customers, respectively. Recall that ρ_k is the probability for observing Scenario $k \in \mathcal{K}$.

We consider $K = 10$ scenarios. The parameters s_{k0} , d_{k0} , μ_k^s , σ_k^s , μ_k^d , σ_k^d and ρ_k are chosen as in Table 8.2. To isolate the effect of the market size, we hold the pool size of potential suppliers and customers fixed ($s_{k0} = 1$, $d_{k0} = 1.2$ for all $k \in \mathcal{K}$). The mean of the suppliers’ opportunity cost is increasing across the scenarios, and so is the mean of the customers’ valuation. Imagine as k increases, the weather condition worsens, and customers value more going by car more than say, using a bike sharing service, but at the same time, drivers also value staying at home more than driving. To focus on the first order effect, we set the standard deviation σ_k^s of the suppliers’ cost to be $\frac{1}{3}$ of the mean μ_k^s and the same on the demand side.

In each Scenario k , the optimal price p_k^* , optimal wage w_k^* , and the ratio $\gamma_k^* = w_k^*/p_k^*$ in the benchmark are displayed in Table 8.3. As expected, we see that both price p_k^* and wage w_k^* increase in k . Since both suppliers’ cost and customers’ valuation increase as the weather condition worsens, the platform needs

Table 8.2 Parameters in the illustrative example

Scenario k	1	2	3	4	5	6	7	8	9	10
s_{k0}	1	1	1	1	1	1	1	1	1	1
μ_k^s	15	16	17	18	19	20	21	22	23	24
σ_k^s	5	5.33	5.67	6	6.33	6.67	7	7.33	7.67	8
d_{k0}	1.2	1.2	1.2	1.2	1.2	1.2	1.2	1.2	1.2	1.2
μ_k^d	10	12	14	16	18	20	22	24	26	28
σ_k^d	3.33	4	4.67	5.33	6	6.67	7.33	8	8.67	9.33
ρ_k	0.05	0.05	0.1	0.1	0.2	0.2	0.1	0.1	0.05	0.05

Table 8.3 Optimal wages, prices and ratios in the illustrative example

Scenario k	1	2	3	4	5	6	7	8	9	10
p_k^*	14.20	16.52	18.79	25.90	27.11	28.32	29.53	30.74	31.95	33.16
w_k^*	9.24	10.59	11.87	13.10	14.28	15.44	16.57	17.69	18.79	19.88
γ_k^*	0.651	0.641	0.631	0.622	0.614	0.606	0.599	0.593	0.587	0.581

to increase the wage to attract more suppliers and raise the price to suppress the growing demand. We observe that the optimal wage/price ratio γ_k^* decreases in k . This implies that the platform does not need to increase the commission ratio to incentivize the suppliers. This is because even though γ_k^* decreases as k increases, the optimal wage $w_k^* = \gamma_k^* p_k^*$ is in fact increased and is sufficient for incentivizing the drivers to get on the street.

When the platform freely sets both wage and price in each scenario, the optimal expected profit of the platform in the benchmark is achieved by applying the wage-price pair (w_k^*, p_k^*) in Scenario k and is equal to $\sum_{k \in \mathcal{K}} \rho_k \pi_k(w_k^*, p_k^*) = 2.311$. Under the fixed commission contract $w = \gamma p$, the optimal expected profit of the platform is achieved by applying the price $\check{p}_k(\gamma) = \arg \max_{p \geq 0} \pi_k(\gamma p, p)$ for given γ in Scenario k . The optimal fixed commission contract for maximizing the platform's profit can be found by solving $\max_{\gamma \in [0, 1]} \sum_{k \in \mathcal{K}} \rho_k \pi_k(\gamma \check{p}_k(\gamma), \check{p}_k(\gamma))$. With the given parameters, the optimal ratio $\check{\gamma} = 0.6063$. This fixed commission contract achieves a surprisingly high profit 2.307, which is 99.82% of the optimal profit achieved by $\{(w_k^*, p_k^*)\}_{k \in \mathcal{K}}$. \square

Motivated by Example 1, we investigate the performance of the fixed commission contract. The following result shows that, under a mild condition on the supply curves, the optimal fixed commission contract can achieve a decent portion of optimality.

Theorem 2 (Concave supply curve: 3/4-optimality of fixed commission) *Suppose $s_k(w)$ is concave for all $k \in \mathcal{K}$.*

- (i) *For any Scenario k , $\gamma_k^* \equiv w_k^*/p_k^* \leq 50\%$.*
- (ii) *The expected profit achieved by a heuristic commission contract with $\gamma = (1 - \bar{\gamma})\bar{\gamma}/(1 - \underline{\gamma})$ ($\leq 50\%$) is at least $\frac{3}{4}P^*$, where $\bar{\gamma} = \max_{k \in \mathcal{K}} \gamma_k^*$ and $\underline{\gamma} = \min_{k \in \mathcal{K}} \gamma_k^*$.*

Theorem 2(ii) says that as long as the supply curve is concave, a fixed commission contract can achieve 75% of the optimality of the benchmark. Theorem 2(i) says that for any scenario the ratio of the optimal wage to price in the benchmark should be no more than 50% (i.e., the commission rate is over 50%). This ratio may seem overly low and may be caused by lack of competition in the model. Uber currently leaves the drivers a fraction of 70–80% of fares paid by the riders. Thus, Theorem 2(i) may imply that the Uber's current pricing practice is not profit maximizing, which is consistent with the news reports saying that Uber is not making a profit. Indeed, in practice, the platform may charge a lower commission due to fairness concerns and to improve supplier welfare. Nevertheless, platforms

Table 8.4 Statistics of the performance of the fixed commission contract: normal distributions

Maximum	Minimum	Mean	Median	Standard deviation
96.30%	82.54%	91.07%	91.33%	2.32%

Table 8.5 Statistics of the performance of the fixed commission contract: log-normal distributions

Maximum	Minimum	Mean	Median	Standard deviation
95.43%	76.52%	88.35%	88.59%	3.23%

like Uber have strong bargaining powers over their suppliers and may try to raise the commission to increase profit if it would not overly irritating the suppliers. For example, Uber and Lyft started with a 20% commission and later increased the rate to 25% in most cities. Currently, the Lyft fee in New York City is 31.4%.

8.2.3 Numerical Study

To close this section, we now present numerical experiments to further investigate the performance of the fixed commission contract against the optimal expected profit of the platform.

As in the illustrative example, we consider conditional normal distribution for the supplier's opportunity costs and customers' valuations. For the conditional-normally distributed supply cost, the supply curve is neither convex nor concave. We consider $K = 48$ scenarios in total. Following the same notation as in Example 1, we draw the parameters s_{k0} , μ_k^s , σ_k^s , d_{k0} , μ_k^d and σ_k^d independently and uniformly at random, in the following manner: $s_{k0} \sim \mathcal{U}[0, 1]$, $\mu_k^s \sim \mathcal{U}[10, 20]$, $\sigma_k^s \sim \mathcal{U}[0.1\mu_k^s, 0.4\mu_k^s]$, $d_{k0} \sim \mathcal{U}[0, 1]$, $\mu_k^d \sim \mathcal{U}[10, 20]$ and $\sigma_k^d \sim \mathcal{U}[0.1\mu_k^s, 0.4\mu_k^s]$, where $\mathcal{U}[A, B]$ denotes the uniform distribution over the interval $[A, B]$. To generate the probabilities ρ_k , $k \in \mathcal{K}$, we first draw $\tilde{\rho}_k \sim \mathcal{U}[0, 1]$ for every k , and then normalize the $\tilde{\rho}_k$'s, i.e., $\rho_k = \tilde{\rho}_k / \sum_{k \in \mathcal{K}} \tilde{\rho}_k$.

We generate a total number of 400 instances of a combination of parameters for the $K = 48$ number of scenarios and summarize the statistics on the performance of the best fixed commission contract in Table 8.4. The results show that the performance of the best fixed commission contract is consistently good, with the worst case achieving 82.54% of the optimality.

We further generate 400 instances where the suppliers' cost X and customers' valuation Y follow log-normal distributions, i.e., $\log(X)$ follows a normal distribution with mean μ_k^s and standard deviation σ_k^s , and $\log(Y)$ follows a normal distribution with mean μ_k^d and standard deviation σ_k^d . The number of scenarios is still $K = 48$. We randomly draw all the parameters in the same way as we did for the normal distributions. Table 8.5 shows the performance of the fixed commission contract. While the performance is slightly worse compared with that under the normal distributions, it is still consistently good, with even the worst case performing better than the performance guarantee 75% we obtained for the case of concave supply curves.

8.3 Dynamic Matching with Heterogeneous Types

We now study a dynamic matching model of the platform in which price and wage are either exogenous (e.g., price and wage have been fixed for a while and accepted by all parties when Uber assigns a rider to a driver) or irrelevant (e.g., for organ sharing, in the majority of countries except only a few it is not allowed to put a price tag on organs).

Consider a finite horizon with a total number of T periods. In practice, even though demand and supply arrive in continuous time, matching decisions are typically not made in real time. For example, Amazon periodically optimizes the way in which it matches customer orders and its warehouses (see Xu et al. 2009). At the beginning of each period, n types of demand and m types of supply arrive in *random* quantities. Let \mathcal{D} be the set of demand types and \mathcal{S} be the set of supply types. With a slight abuse of notation, we write $\mathcal{D} = \{1, 2, \dots, n\}$ and $\mathcal{S} = \{1, 2, \dots, m\}$, noting that \mathcal{D} and \mathcal{S} are disjoint sets. We use i to index a demand type and j to index a supply type. The pairs of demand and supply form a bipartite graph. An arc (i, j) represents a match between type i demand and type j supply. For simplicity, we consider a complete bipartite graph in the base model. In other words, any demand type can potentially be matched with any supply type, obviously with different rewards (or equivalently, mismatch costs). We denote the complete set of arcs by $\mathcal{A} = \{(i, j) \mid i \in \mathcal{D}, j \in \mathcal{S}\}$.

We denote, as system states, the demand vector by $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}_+^n$ and the supply vector by $\mathbf{y} = (y_1, \dots, y_m) \in \mathbb{R}_+^m$, where x_i and y_j are the quantity of type i demand and type j supply available to be matched. Although we assume that the states and the demand and supply arrivals are continuous quantities (and therefore so are the matching decisions), our results can be readily replicated if those quantities are discrete. On observing the state $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}_+^{n+m}$, the firm decides on the quantity q_{ij} of type i demand to be matched with type j supply, for any $i \in \mathcal{D}$ and $j \in \mathcal{S}$. For conciseness, we write the decision variables of matching quantities in a matrix form as $\mathbf{Q} = (q_{ij}) \in \mathbb{R}_+^{n \times m}$, with \mathbf{Q}_i its i -th row (as a row vector) and \mathbf{Q}^j its j -th column (as a column vector). We assume that there is a reward r_{ij} for matching one unit of type i demand and one unit of type j supply for all i, j . Similarly, we can write the rewards in a matrix form as $\mathbf{R} = (r_{ij}) \in \mathbb{R}^{n \times m}$. Thus the total reward from matching is linear in the matching quantities. That is, $\mathbf{R} \circ \mathbf{Q} \equiv \sum_{i=1}^n \sum_{j=1}^m r_{ij} q_{ij}$, where “ \circ ” gives the sum of elements of the Hadamard product of two matrices. The *post-matching levels* of type i demand and type j supply are given by $u_i = x_i - \mathbf{1}^m \mathbf{Q}_i^T = x_i - \sum_{j=1}^m q_{ij}$ and $v_j = y_j - \mathbf{1}^n \mathbf{Q}^j = y_j - \sum_{i=1}^n q_{ij}$, respectively. That is, $\mathbf{u} = \mathbf{x} - \mathbf{1}^m \mathbf{Q}^T$ and $\mathbf{v} = \mathbf{y} - \mathbf{1}^n \mathbf{Q}$. The post-matching levels cannot be negative; i.e., $\mathbf{u} \geq \mathbf{0}, \mathbf{v} \geq \mathbf{0}$.

After the matching is done in each period, each unit of unmatched demand and supply incurs a holding cost of c and h respectively. The cost on the demand side could be the loss of goodwill or customers’ waiting costs. Consequently, the total holding cost amounts to $c \mathbf{1}^n \mathbf{u}^T + h \mathbf{1}^m \mathbf{v}^T = c \sum_{i=1}^n u_i + h \sum_{j=1}^m v_j$. The unmatched demand and supply carry over to the next period with carry-over rates α and β ,

respectively. In other words, $(1 - \alpha)$ fraction of demand and $(1 - \beta)$ fraction of supply leave the system. Without loss of generality, we assume they leave the system with zero surpluses.

The firm's goal is to determine a matching policy $\mathbf{Q}^* = (q_{ij}^*)$ that maximizes the expected total discounted surplus. (Our perspective is the maximizing of social welfare. Alternatively, the formulation can account for profit maximization if r_{ij} is interpreted as the revenue collected from a matching, and c and h are interpreted as the penalty paid to demand and supply for showing up but without a successful match in a period.) Let $V_t(\mathbf{x}, \mathbf{y})$ be the optimal expected total discounted surplus given that it is in period t and the current state is (\mathbf{x}, \mathbf{y}) . We formulate the finite-horizon problem by using the following stochastic dynamic program:

$$V_t(\mathbf{x}, \mathbf{y}) = \max_{\mathbf{Q} \in \{\mathbf{Q} \geq \mathbf{0} | \mathbf{u} \geq \mathbf{0}, \mathbf{v} \geq \mathbf{0}\}} H_t(\mathbf{Q}, \mathbf{x}, \mathbf{y}),$$

$$H_t(\mathbf{Q}, \mathbf{x}, \mathbf{y}) = \mathbf{R} \circ \mathbf{Q} - c \mathbf{1}^n \mathbf{u}^T - h \mathbf{1}^m \mathbf{v}^T + \gamma E V_{t+1}(\alpha \mathbf{u} + \mathbf{D}_t, \beta \mathbf{v} + \mathbf{S}_t), \quad (8.2)$$

where $\gamma \leq 1$ is the discount factor. The boundary conditions are $V_{T+1}(\mathbf{x}, \mathbf{y}) = 0$ for all (\mathbf{x}, \mathbf{y}) , without loss of generality. In other words, at the end of the horizon, all unmatched demand and supply leave the system with zero surpluses.

8.3.1 Priority Properties of the Optimal Matching Policy

We provide sufficient conditions for prioritizing the demand-supply pairs (i.e., the arcs of the bipartite network) under the optimal matching policy. The conditions we provide are imposed on the reward matrix and independent of any other system parameters. Those conditions guarantee that specific priority structural properties hold for the dynamic problem at *any* time and with *any* realized demand and supply. For succinctness, we may only present the definitions and results on one side of the market, analogous definitions and results can be easily stated and obtained for the other side by symmetry.

We define a *relation* " \succeq " between arcs as follows and show it is a *partial order*. We refer readers to Hu and Zhou (2018a) for proofs in this section. First, we consider neighboring arcs in the bipartite graph.

Definition 1 (Modified Monge condition for arcs with a common node) $(i, j) \succeq (i, j')$, if

- (i) $r_{ij} \geq r_{ij'}$ and
- (ii) $r_{ij} + r_{i'j'} \geq r_{ij'} + r_{i'j}$ for all $i' \in \mathcal{D}$. (D)

(When $i' = i$, condition (D) holds automatically. It is easy to see that $(i, j) \succeq (i, j')$ holds automatically for $j' = j$.)

We further define a relation between arcs that do not share any node but can be connected through a sequence of neighboring arcs regulated by the relation “ \succeq ”.

Definition 2 (Modified Monge condition for arcs without common nodes) For $i \neq i'$ and $j \neq j'$, we say $(i, j) \succeq (i', j')$ if there exists a sequence of arcs $(i_1, j_1), (i_2, j_2), \dots, (i_k, j_k)$ such that either $i_k = i_{k+1}$ or $j_k = j_{k+1}$ for $k = 1, \dots, n - 1$, and $(i, j) = (i_1, j_1) \succeq (i_2, j_2) \succeq \dots \succeq (i_k, j_k) = (i', j')$.

Further, we define the equivalence and strict dominance of two pairs (i, j) and (i', j') as follows.

Definition 3 We say (i, j) is equivalent to (i', j') , denoted by $(i, j) \sim (i', j')$, if $(i, j) \succeq (i', j')$ and $(i', j') \succeq (i, j)$. We say (i, j) strictly dominates (i', j') , denoted by $(i, j) \succ (i', j')$, if $(i, j) \succeq (i', j')$ but (i, j) is not equivalent to (i', j') .

For any arc $(i, j) \in \mathcal{A}$, we define a set of neighboring arcs that are *strictly* dominated by (i, j) : $\mathcal{L}_{ij} \equiv \{(i'', j) \mid (i, j) \succ (i'', j)\} \cup \{(i, j'') \mid (i, j) \succ (i, j'')\}$. We also define

$$w_{ij} = w_{ij}(\mathbf{Q}, \mathbf{x}, \mathbf{y}) \equiv \min \left\{ x_i - \sum_{j': (i, j') \notin \mathcal{L}_{ij}} q_{ij'}, y_j - \sum_{i': (i', j) \notin \mathcal{L}_{ij}} q_{i'j} \right\}.$$

If $w_{ij} = 0$, type i or j is exhausted by the matching over arcs outside the set \mathcal{L}_{ij} .

The following theorem establishes a sufficient condition for a demand-supply pair (i, j) to be prioritized over another pair (i', j') .

Theorem 3 (Partial order implies priority) *Without loss of generality, assume $\mathbf{x} > \mathbf{0}$ and $\mathbf{y} > \mathbf{0}$ in period t .³ There exists an optimal decision \mathbf{Q}^* such that for any $(i, j) \succ (i', j')$, $\min\{w_{ij}^*, q_{i'j'}^*\} = 0$, i.e., either \mathbf{Q}^* exhausts type i or j over arcs outside \mathcal{L}_{ij} , or $q_{i'j'}^* = 0$.*

The northwest corner rule under the assumption of a Monge sequence can completely solve the deterministic and balanced version of the problem in a greedy fashion. For the stochastic version, we show in Theorem 3 that the priority structure preserves under the modified Monge conditions, a somewhat stronger set of assumptions than the Monge sequence.⁴ However, even a pair has higher priority in the optimal matching, they are not necessarily matched in a greedy fashion; when they are not exhausted, all pairs that have strictly lower priority should not be matched.

Next, we provide a sufficient condition for a pair of demand and supply types to be matched greedily over all other possible matching options.

³If $x_i = 0$ (or $y_j = 0$), we can delete demand node i (or supply node j) and all its connected arcs, on which matching quantities are set to zero.

⁴If all arcs are comparable under our partial order along the sequence, then it is a Monge sequence. But we do not require all arcs to be comparable.

Theorem 4 (When greedy matching is optimal) *If $(i, j) \succeq (i, j')$ for all $j' \in \mathcal{S}$ and $(i, j) \succeq (i', j)$ for all $i' \in \mathcal{D}$, then $q_{ij}^* = \min \{x_i, y_j\}$.*

Theorem 4 is *not* a direct consequence of Theorem 3. By directly applying Theorem 3, we can only say that under the conditions in Theorem 4, it is optimal for the firm to prioritize the matching of type i demand and type j supply over any other possibilities. However, it may still be possible that the firm has reserved some type i demand and type j supply without greedily matching them.

As an immediate application of Theorem 4, consider demand and supply types that are specified by their locations in a Euclidean space. The reward of matching supply with demand is a fixed prize minus the disutility proportional to the Euclidean distance between the supply location and the demand location. It is easy to verify that a demand type and a supply type from the *same* location forms a perfect pair, and by Theorem 4, they should be matched as much as possible. To see why they are a perfect pair, we have $r_{ii} + r_{i'j'} \geq r_{ij'} + r_{i'i}$ because $d_{i'j'} \leq d_{ij'} + d_{i'i}$, where d_{ij} is the Euclidean distance between the locations of type i demand and type j supply. The latter inequality is simply the triangle inequality. We summarize this result as follows.

Corollary 1 *In a Euclidean space with horizontally differentiated types as locations, it is optimal to match the demand and supply from the same location greedily.*

Corollary 1 suggests that with geographic locations as types, the intermediary firm such as Uber and Amazon should always match a demand with a supply if they are originated from the same geographic region, or practically speaking if they are sufficiently close to each other.

In the followings, we present two reward structures, for which Theorem 3 can be applied to characterize the optimal matching policy.

8.3.1.1 Unidirectional Horizontal Types

Let the n demand and m supply types be distributed on a fixed route C (e.g., a line segment) with a given direction. All the demand types have distinct locations (otherwise we can treat two demand types sharing the same location as the same type) and so do the supply types. For any two types t_1, t_2 , we write $t_1 \rightarrow t_2$ to denote that t_1 is located before t_2 , along the given direction. We denote by $\mathbf{d}(t_1, t_2)$ the travel distance from the location of t_1 to that of t_2 along the given direction. The unit matching reward r_{ij} between type i demand and type j supply is a nonincreasing function f of the distance between the two types, which is measured as follows. Distance-based reward implies horizontal preferences since a demand/supply type would “prefer” (in the eyes of the platform) a supply/demand type that is “closer” to herself.

For $i \in \mathcal{D}$ and $j \in \mathcal{S}$ such that $j \rightarrow i$, we define $d_{ij} = \mathbf{d}(j, i)$. For $i \in \mathcal{D}$ and $j \in \mathcal{S}$ such that $i \rightarrow j$, we consider one of the following definitions:

- (i) (Directed line segment) $d_{ij} = N$, where N is an arbitrarily large number;
- (ii) (Directed circle) $d_{ij} = |C| - \mathbf{d}(i, j)$, where $|C|$ is the length of the route C ;
- (iii) (Undirected line segment) $d_{ij} = \mathbf{d}(i, j)$.

In case (i), a supply type is not allowed to travel counter to the given direction. In this case, C is a directed line segment, on which $i \in \mathcal{D}$ and $j \in \mathcal{S}$ can be matched with each other if and only if $j \rightarrow i$. The product upgrading model has the structure of a directed line segment (see, e.g., Yu et al. 2015).

In case (ii), a type j supply can travel along the given direction to reach a type i demand if $j \rightarrow i$. If $i \rightarrow j$, j needs to travel to the end of the route along the given direction, then “reappears” at the origin of the route and continues along the direction to reach i . This is equivalent to the case in which C is a directed (say, clockwise) circle and a supply type always needs to go clockwise on the circle to reach a demand type.

In case (iii), a supply type can go along or counter to the direction to reach a demand type. Thus the direction no longer plays a role, and C is equivalent to an undirected line segment.

For $i \in \mathcal{D}$ and $j \in \mathcal{S}$, let $\text{seg}(i \leftarrow j)$ denote the segment of the route traveled by j to reach i .

Provided that the unit reward $r_{ij} = f(d_{ij})$ is linear in d_{ij} , the following priority structure can be inferred from Theorem 3.

Theorem 5 (Distance-based priority) *Suppose f is a linear and decreasing function.*

- (i) *If $\text{seg}(i \leftarrow j) \subseteq \text{seg}(i \leftarrow j')$, then $(i, j) \succeq (i, j')$. Similarly, if $\text{seg}(i \leftarrow j) \subseteq \text{seg}(i' \leftarrow j)$, then $(i, j) \succeq (i', j)$.*
- (ii) *In the case of undirected line segment, if $\text{seg}(i \leftarrow j) \subseteq \text{seg}(i' \leftarrow j')$ and $\text{seg}(i \leftarrow j)$ has the same direction with $\text{seg}(i' \leftarrow j')$, then $(i, j) \succeq (i', j')$.*
- (iii) *In the case of directed line segment and circle, $\text{seg}(i \leftarrow j) \subseteq \text{seg}(i' \leftarrow j')$ is equivalent to $(i, j) \succeq (i', j')$.*

Theorem 5 (i) implies that a shorter distance leads to a higher priority when we compare two pairs of demand and supply with a common node. Moreover, as suggested by the following theorem, for the directed line segment and the directed circle, each demand or supply type should be matched greedily with its closest match.

Theorem 6 (Greedy match of perfect pairs) *Consider the directed line segment case or the directed circle case. Suppose that $\text{seg}(i \leftarrow j)$ does not contain any other types than themselves. If f is nonincreasing and convex, $q_{ij}^* = \min\{x_i, y_j\}$.*

8.3.1.2 Vertical Types

We consider the case where the demand and supply types have “quality” differences. The unit reward takes the form $r_{ij} = r_i^d + r_j^s$, where r_i^d and r_j^s are contributions made by type i demand and type j supply, respectively. Without loss of generality, we index the types such that $r_1^d > \dots > r_n^d$ and $r_1^s > \dots > r_m^s$.

With the additive reward structure, $r_{ij} + r_{i'j'} = r_{ij'} + r_{i'j}$ for all $i, i' \in \mathcal{D}$ and $j, j' \in \mathcal{S}$. This implies that for two neighboring arcs, $(i, j) \succeq (i', j)$ if and only if $r_i^d \geq r_{i'}^d$, and $(i, j) \succeq (i, j')$ if and only if $r_j^s \geq r_{j'}^s$. This observation can easily be

generalized as $(i, j) \succeq (i', j')$ if and only if $i < i'$ and $j < j'$. By Theorem 4, it is optimal to match type 1 demand and type 1 supply greedily. From Theorem 3, the arc (i, j) has priority over (i, j') and (i', j) for all $j' > j$ and $i' > i$.

This leads to an optimal policy that follows a top-down matching procedure.

Corollary 2 (Top-down matching) *Line up demand types and supply types separately in increasing order of their indices. Match from the top, down to some level. The optimal matching decision \mathbf{Q} in a period is fully determined by a total matching quantity $Q \equiv \sum_{i'=1}^n \sum_{j'=1}^m q_{i'j'}$.*

8.3.2 Bound and Heuristic

We study the deterministic counterpart of the stochastic problem in its general form. We show that the heuristic suggested by the deterministic model can be computed efficiently and is asymptotically optimal for the stochastic problem.

We consider the deterministic model by ignoring the uncertainty and assume that the mean demand quantity $\lambda_{it} = \mathbb{E} D_{it}$ and mean supply quantity $\mu_{jt} = \mathbb{E} S_{jt}$ arrive in each period. The following linear program gives the formulation of the problem from period τ to period T :

$$\begin{aligned}
 (\mathbf{P}_\tau^{\mathbf{x}, \mathbf{y}}) \quad & \max_{q_{ijt}, x_{it}, y_{jt}} \sum_{t=\tau}^T \gamma^{t-1} \left[\sum_{i=1}^n \sum_{j=1}^m r_{ij} q_{ijt} - c \left(\sum_{i=1}^n x_{it} - \sum_{i=1}^n \sum_{j=1}^m q_{ijt} \right) \right. \\
 & \quad \left. - h \left(\sum_{j=1}^m y_{jt} - \sum_{i=1}^n \sum_{j=1}^m q_{ijt} \right) \right] \\
 \text{s.t.} \quad & \sum_{j=1}^m q_{ijt} \leq x_{it}, \quad i \in \mathcal{D}, t = \tau, \tau + 1, \dots, T, \\
 & \sum_{i=1}^n q_{ijt} \leq y_{jt}, \quad 1 \leq j \leq m, \tau \leq t \leq T, \\
 & x_{i,t+1} = \alpha \left(x_{it} - \sum_{j=1}^m q_{ijt} \right) + \lambda_{it}, \quad i \in \mathcal{D}, t = \tau, \dots, T-1, \\
 & y_{j,t+1} = \beta \left(y_{jt} - \sum_{i=1}^n q_{ijt} \right) + \mu_{jt}, \quad 1 \leq j \leq m, \tau \leq t \leq T-1, \\
 & q_{ijt} \geq 0, \quad i \in \mathcal{D}, j \in \mathcal{S}, t = \tau, \dots, T, \\
 & x_{i\tau} = x_i, \quad y_{j\tau} = y_j, \quad i \in \mathcal{D}, j \in \mathcal{S}, \tag{8.3}
 \end{aligned}$$

where $(\mathbf{x}, \mathbf{y}) = (x_1, \dots, x_n, y_1, \dots, y_m)$ is a given initial state at the beginning of period τ .

From the optimal solution to $(\mathbf{P}_\tau^{\mathbf{x}, \mathbf{y}})$, $\{\hat{q}_{ijt}, \hat{x}_{it}, \hat{y}_{jt}\}_{i \in \mathcal{D}, j \in \mathcal{S}, t = \tau, \dots, T}$, we obtain a feasible matching decision $\{\hat{q}_{ij\tau}\}_{i \in \mathcal{D}, j \in \mathcal{S}}$ in period τ for state (\mathbf{x}, \mathbf{y}) , and use it as a heuristic decision. If we start in period 1 with an initial state $(\mathbf{x}_1, \mathbf{y}_1)$, we will solve $(\mathbf{P}_1^{\mathbf{x}_1, \mathbf{y}_1})$ to obtain matching decisions $\{\hat{q}_{ij1}\}_{i \in \mathcal{D}, j \in \mathcal{S}}$ in period 1. Given a realization of demand and supply in period 2 as $\mathbf{D}_2 = \mathbf{d}_2$ and $\mathbf{S}_2 = \mathbf{s}_2$, respectively, the state in period 2 is then $(\mathbf{x}_2, \mathbf{y}_2) = (\alpha \hat{\mathbf{u}} + \mathbf{d}_2, \beta \hat{\mathbf{v}} + \mathbf{s}_2)$. We then solve $(\mathbf{P}_2^{\mathbf{x}_2, \mathbf{y}_2})$ to obtain the heuristic decision for period 2. We proceed until period T to obtain the heuristic decisions along a sample path of demand and supply realization. The above procedure is referred to as the deterministic re-solving heuristic.

Intuitively, one would expect that the uncertainty in demand and supply in the stochastic model would result in a lower expected surplus. This is confirmed by the following proposition.

Proposition 2 (Deterministic upper bound) *The deterministic model provides an upper bound on the optimal total surplus of the stochastic model.*

If we scale up the arrival of demand and supply of all types by a multiplier k (i.e., $D_{it} = \sum_{\ell=1}^k D_{it}^{(\ell)}$ and $S_{jt} = \sum_{\ell=1}^k S_{jt}^{(\ell)}$ for all $i \in \mathcal{D}$ and $j \in \mathcal{S}$, where $D_{it}^{(1)}, \dots, D_{it}^{(k)}$ are independent and identically distributed (i.i.d.) random variables and $S_{jt}^{(1)}, \dots, S_{jt}^{(k)}$ are i.i.d. random variables), the policy suggested by the deterministic re-solving heuristic becomes asymptotically optimal as $k \rightarrow \infty$.

With the scalar k , let $V_t^k(\mathbf{x}, \mathbf{y})$ be the value function and $V_t^{\text{resolve}(k)}(\mathbf{x}, \mathbf{y})$ the value for applying the deterministic re-solving heuristic in system k .

Theorem 7 (Asymptotic optimality of the deterministic heuristic and rate of convergence) *In the stochastic system k , the deterministic re-solving heuristic leads to the relative error $[V_t^k(\mathbf{x}, \mathbf{y}) - V_t^{\text{resolve}(k)}(\mathbf{x}, \mathbf{y})]/V_t^k(\mathbf{x}, \mathbf{y}) = O(1/\sqrt{k})$ as $k \rightarrow \infty$.*

8.3.2.1 Numerical Study

Consider a 10-period dynamic matching problem with 5 supply types and 5 demand types. For each instance of the problem, we generate the time-independent parameters uniformly at random as follows.

Let $r_{ij} \sim \mathcal{U}[50, 150]$ (for all $i \in \mathcal{D}$ and $j \in \mathcal{S}$), $c \sim \mathcal{U}[0, 50]$, $h \sim \mathcal{U}[0, 50]$, $\alpha \sim \mathcal{U}[0, 1]$, $\beta \sim \mathcal{U}[0, 1]$, $\lambda_i = ED_i \sim \mathcal{U}[10, 25]$, $\mu_j = ES_j \sim \mathcal{U}[10, 25]$, $\gamma \sim \mathcal{U}[0.8, 1]$.

In addition, we also randomly generate the initial state $(\mathbf{x}_0, \mathbf{y}_0)$ at the beginning of the first period. We let $x_{i0} \sim \mathcal{U}[0, 30]$ and $y_{j0} \sim \mathcal{U}[0, 30]$ for all $i \in \mathcal{D}$ and $j \in \mathcal{S}$.

We run two sets of numerical experiments described as follows.

- (a) Demand and supply follow a uniform distribution. For given realizations of λ_i and μ_j , we generate $\delta_i^d \sim \mathcal{U}[0, \lambda_i]$ and $\delta_j^s \sim \mathcal{U}[0, \mu_j]$. Then, we let $D_i \sim \mathcal{U}[\lambda_i - \delta_i^d, \lambda_i + \delta_i^d]$ and $S_j \sim \mathcal{U}[\mu_j - \delta_j^s, \mu_j + \delta_j^s]$.

- (b) Demand and supply follow conditional normal distributions (conditioned on nonnegative values). For given realizations of λ_i and μ_j , we generate $\sigma_i^d \sim \mathcal{U}[0, \lambda_i/3]$ and $\sigma_j^s \sim \mathcal{U}[0, \mu_j/3]$. Then, we let D_i (resp. S_j) be the conditional normal distribution with mean λ_i (resp. μ_j) and standard deviation σ_i^d (resp. σ_j^s).

Note that all the parameters are generated independently. For each randomly generated instance, we solve the 10-period deterministic problem (P) and obtain the optimal value V^{det} , which is an upper bound of the optimal value V^{opt} of the stochastic problem. Let \tilde{V} be the optimal value of the expected total discounted reward minus costs when the deterministic heuristic is applied throughout the decision horizon. We calculate \tilde{V} approximately by simulation: For each randomly generated sample path ω , in period t ($t = 1, \dots, T$) with state $(\mathbf{x}_t(\omega), \mathbf{y}_t(\omega))$, apply the optimal decision from solving the $(T - t + 1)$ -period problem with initial state $(\mathbf{x}_t(\omega), \mathbf{y}_t(\omega))$; The total reward minus cost for the sample path ω can be easily calculated; Then we average over 5000 sample paths to obtain the approximate value of \tilde{V} . Since $(V^{\text{opt}} - \tilde{V})/V^{\text{opt}} \leq (V^{\text{det}} - \tilde{V})/V^{\text{det}} \equiv \rho$, the relative error by the deterministic heuristic is even smaller if ρ is small. Thus, we focus on ρ to measure the relative error.

For the set (a) of the experiments, 600 instances are generated. Among the 600 instances, the maximum value of ρ is 21.24%, the mean is 9.84%, and the median is 9.51%.

For the set (b) of the experiments, 820 instances are generated. Among the 820 instances, the maximum value of ρ is 19.24%, the mean is 7.24%, and the median is 6.79%.

8.4 Pricing and Matching with Strategic Suppliers and Customers

In this section, we study the joint pricing and matching decision of an intermediary platform. Consider an intermediary who dynamically matches demand and supply of a single-type product or service. Customers and suppliers have heterogeneous valuations and sequentially arrive at the market over a finite horizon $[0, T]$. The intermediary implements an anonymous posted price mechanism on both demand and supply sides. That is, at each point of time $t \in [0, T]$, the intermediary posts demand-side price π_t^d (i.e., ask price) that she charges customers for each unit of the product and supply-side price π_t^s (i.e., bid price) that she pays to suppliers for each unit of the product.

Customer behavior Over the horizon $[0, T]$, customers with heterogeneous valuations arrive at the intermediary according to a Poisson process with rate λ^d . A customer arriving at time t is endowed with a product valuation $v \in [\underline{v}, \bar{v}]$, as a realization from a willingness-to-pay distribution. We denote by

$$\phi \equiv (t_\phi, v_\phi),$$

the “type” of an arriving customer which specifies her arrival time t_ϕ and valuation v_ϕ . Every customer purchases at most one unit of the product. Customers are forward-looking and can strategically determine when to request to buy the product. Specifically, for every customer ϕ , at each point of time after her arrival, she decides to either continue to monitor price dynamics or stop doing so by either sending a request to the intermediary for buying the product or permanently leaving the market without buying anything. We denote by $\tau_\phi \in [t_\phi, T]$ the time that customer ϕ stops monitoring price dynamics. We denote by $a_\phi \in \{0, 1\}$ the indicator function of whether customer ϕ requests to buy the product at time τ_ϕ . A customer who requests to buy the product keeps on staying in the market until her demand is matched with a unit of supply. As an exception, if a customer’s demand cannot be matched by the end of the horizon, then she withdraws her demand request from the intermediary and leaves the market. We denote by $s_\phi \in [\tau_\phi, T]$ the time when customer ϕ leaves the market. We denote by $m_\phi \in \{0, 1\}$ the indicator function of whether customer ϕ ’s demand request is successfully matched with a supplier at time s_ϕ . Customer ϕ pays $p_\phi = \pi_{\tau_\phi}^d m_\phi$ to the intermediary, i.e., if customer ϕ ’s demand is successfully matched with a supplier, then she pays the demand-side price posted at the time that she requests to buy the product, $\pi_{\tau_\phi}^d$; otherwise, she makes no payment to the intermediary. This business rule is consistent with the practice of the ride-hailing industry. In that setting, riders and drivers submit matching requests and then wait to be matched. If a rider is not matched with any driver for some reason, no payment will be made by the rider, and the rider will also not be compensated by the intermediary for the wait.

We define the tuple

$$y_\phi \equiv (\tau_\phi, a_\phi, s_\phi, m_\phi, p_\phi).$$

Customer ϕ garners utility

$$U^d(\phi, y_\phi) = v_\phi m_\phi - p_\phi - b(s_\phi - t_\phi),$$

where $b \in \mathbb{R}_+$ is customer ϕ ’s per unit of time disutility from staying in the system over $[t_\phi, s_\phi]$, hereafter called as the *customers’ waiting cost parameter*, which is assumed to be common knowledge.

As mentioned, we allow for the heterogeneity of customer arrival times and valuations. A customer’s arrival time and valuation are private information and are independent from each other. We denote the cumulative distribution function (c.d.f.) of customer product valuation by $F^d(\cdot)$ and the corresponding probability density function (p.d.f.) by $f^d(\cdot)$. We denote $\bar{F}^d(\cdot) \equiv 1 - F^d(\cdot)$. In addition to assuming $\bar{F}^d(\cdot)$ has an inverse, denoted by $\bar{F}^{d,-1}(\cdot)$, we make a standard assumption on the valuation distribution:

Assumption 1 (Willingness-to-pay) *The customer virtual value function $V^d(v) \equiv v - \bar{F}^d(v)/f^d(v)$ is increasing in $v \in [\underline{v}, \bar{v}]$.*

Supplier behavior Over the horizon $[0, T]$, suppliers arrive at the intermediary according to a Poisson process with rate λ^s . A supplier arriving at time t has c as a production and delivery cost for a good, or opportunity cost for providing a service, where $c \in [\underline{c}, \bar{c}]$. We assume $\underline{c} \leq \bar{v}$. We denote by

$$\psi \equiv (t_\psi, c_\psi),$$

the “type” of an arriving supplier. Every supplier sells at most one unit of the product. All suppliers are forward-looking and can strategically determine when to request to sell her product. For every supplier ψ , at each point of time after her arrival, she decides to either continue to monitor price dynamics or stop doing so by either sending a request to the intermediary for selling the product or permanently leaving the market without selling anything. We denote by $\tau_\psi \in [t_\psi, T]$ the time that supplier ψ stops monitoring price dynamics. We denote by $a_\psi \in \{0, 1\}$ the indicator function of whether supplier ψ requests to sell her product at time τ_ψ . A supplier who requests to sell her product keeps on staying in the market until she is matched with a customer. Consistent with the demand side, as an exception, if a supplier cannot be matched by the end of the horizon, then she withdraws her supply request from the intermediary and leaves the market. We denote by $s_\psi \in [\tau_\psi, T]$ the time when supplier ψ leaves the market. We denote by $m_\psi \in \{0, 1\}$ the indicator function of whether supplier ψ 's supply is successfully matched with a unit of demand at s_ψ . The intermediary pays the supplier ψ with $p_\psi = \pi_{\tau_\psi}^s m_\psi$, i.e., if supplier ψ 's supply is successfully matched with a demand, then the intermediary pays her the supply-side price posted at the time that she requests to sell the product, $\pi_{\tau_\psi}^s$; otherwise, she receives no payment from the intermediary.

We define the tuple

$$y_\psi \equiv (\tau_\psi, a_\psi, s_\psi, m_\psi, p_\psi).$$

Supplier ψ garners utility

$$U^s(\psi, y_\psi) = p_\psi - c_\psi m_\psi - h(s_\psi - t_\psi),$$

where $h \in \mathbb{R}_+$ is supplier ψ 's per unit of time disutility from staying in the system over $[t_\psi, s_\psi]$, hereafter called as the *suppliers' waiting cost parameter*, which is also assumed to be common knowledge. In addition, like the demand side, we assume there is no heterogeneity in suppliers' waiting cost parameter h .

As mentioned, we allow for the heterogeneity of suppliers' arrival times and supply costs. A supplier's arrival time and supply cost are private information and are independent of each other. We denote the c.d.f. of supplier product producing and delivering cost by $F^s(\cdot)$ and the corresponding p.d.f. by $f^s(\cdot)$. In addition to assuming $F^s(\cdot)$ has an inverse, denoted by $F^{s,-1}(\cdot)$, we make the following assumption on the supply cost distribution:

Assumption 2 (Willingness-to-sell) *The supplier virtual cost function $V^s(c) \equiv c + F^s(c)/f^s(c)$ is increasing in $c \in [\underline{c}, \bar{c}]$.*

Game dynamics and the equilibrium Before the start of the horizon, the intermediary determines a *pricing policy* on both demand and supply sides

$$\pi = \{(\pi_t^d, \pi_t^s) : t \in [0, T]\},$$

and a demand and supply *matching policy*

$$M = \{(s_\phi, m_\phi), (s_\psi, m_\psi) : (\tau_\phi \in [0, T], a_\phi = 1), (\tau_\psi \in [0, T], a_\psi = 1)\}.$$

The pricing policy π and the matching policy M are in general dynamic policies depending on the realized uncertainty up to the decision point. The intermediary commits to implement this pricing policy π and matching policy M over the entire course of the horizon. The intermediary's pricing and matching policies are common knowledge for all customers and suppliers. We denote by $H^t \equiv \{\phi, \psi : t_\phi \leq t, t_\psi \leq t\}$ the set of customer and supplier types that arrive up to time t . Define by $\phi^t \equiv \{(\tau_\phi, a_\phi) : \tau_\phi \leq t, a_\phi = 1\}$ the set of demand-side information that the intermediary collects up to time t . Define by $\phi_M^t \equiv \{(s_\phi, m_\phi) : s_\phi \leq t\}$ the set of matching decisions the intermediary has made on the demand side up to time t . Similarly, define by $\psi^t \equiv \{(\tau_\psi, a_\psi) : \tau_\psi \leq t, a_\psi = 1\}$ the set of supply-side information that the intermediary collects up to time t and by $\psi_M^t \equiv \{(s_\psi, m_\psi) : s_\psi \leq t\}$ the set of matching decisions the intermediary has made on the supply side up to time t . Define by $\pi^{d,t} \equiv \{\pi_{t'}^d : t' \in [0, t]\}$ the historic demand-side prices posted up to time t . Define by $\pi^{s,t} \equiv \{\pi_{t'}^s : t' \in [0, t]\}$ the historic supply-side prices posted up to time t . Define a filtration $\{\mathcal{F}_t : t \geq 0\}$ with $\mathcal{F}_t = \sigma(\phi^{t-}, \phi_M^{t-}, \psi^{t-}, \psi_M^{t-}, \pi^{d,t-}, \pi^{s,t-})$. A feasible pricing policy π requires π_t^d and π_t^s to be \mathcal{F}_t -progressive. Denote by Π the set of all feasible pricing policies. A feasible matching policy M requires $\{s_\phi \leq t\}$ and m_ϕ , and $\{s_\psi \leq t\}$ and m_ψ to be \mathcal{F}_t -progressive, and to satisfy the demand and supply balancing condition that $\sum_{\phi \in H^t} \mathbf{1}\{s_\phi = t, m_\phi = 1\} = \sum_{\psi \in H^t} \mathbf{1}\{s_\psi = t, m_\psi = 1\}$ for all $t \in [0, T]$. Denote by \mathcal{M} the set of all feasible matching policies.

During the horizon, on the demand side, customers are forward-looking and employ (symmetric) stopping and purchasing rules contingent on their types that constitute a symmetric Markov Perfect Equilibrium. The following information structure is mainly motivated by the current practice of ride-hailing apps. Our results would still hold under alternative information structures with an update of the waiting time compensation. For a given customer ϕ , the information available to her at time $t \in [t_\phi, T]$ consists of demand-side price dynamics that she tracks during her stay in the system, $\{\pi_{t'}^d : t' \in [t_\phi, t]\}$, and the number of unmatched supply during her stay in the system, $\{U_{t'-}^s : t' \in [t_\phi, t]\}$, where $U_{t'-}^s = \sum_{\psi \in H^{t'-}} \mathbf{1}\{\tau_\psi < t', a_\psi = 1, s_\psi \geq t'\}$. Therefore, at each point of time $t \in [t_\phi, T]$, the event associated with the stopping decision $\{\tau_\phi \leq t\}$ and the purchasing decision a_ϕ are adapted to $\sigma(\pi_{t'}^d, U_{t'-}^s : t' \in [t_\phi, t])$. Under the intermediary's pricing policy π and matching

policy M , for a given customer type ϕ , the optimal stopping rule $\tau_\phi^{\pi, M}$ and the optimal purchasing rule $a_\phi^{\pi, M}$ are the optimal solution to the following optimization problem

$$\sup_{\tau_\phi \in [t_\phi, T], a_\phi \in \{0,1\}} \mathbb{E} [U^d(\phi, y_\phi) \mid \pi_{t_\phi}^d, U_{t_\phi-}^s, \phi],$$

where the expectation assumes that other customers use symmetric stopping and purchasing rules.

On the supply side, suppliers are forward-looking and employ (symmetric) stopping and selling rules contingent on their types that constitute a symmetric Markov Perfect Equilibrium. For a given supplier ψ , the information available to her at time $t \in [t_\psi, T]$ consists of supply-side price dynamics that she tracks during her stay in the system, $\{\pi_{t'}^s : t' \in [t_\psi, t]\}$. Therefore, at each point of time $t \in [t_\psi, T]$, the event associated with the stopping decision $\{\tau_\psi \leq t\}$ and the selling decision a_ψ are adapted to $\sigma(\pi_{t'}^s : t' \in [t_\psi, t])$. Under the intermediary's pricing policy π and matching policy M , for a given supplier type ψ , the optimal stopping rule $\tau_\psi^{\pi, M}$ and the optimal selling rule $a_\psi^{\pi, M}$ are the optimal solution to the following optimization problem

$$\sup_{\tau_\psi \in [t_\psi, T], a_\psi \in \{0,1\}} \mathbb{E} [U^s(\psi, y_\psi) \mid \pi_{t_\psi}^s, \psi],$$

where the expectation assumes that other suppliers use symmetric stopping and selling rules.

Our goal in this paper is to construct a price process $\pi \in \Pi$ and a matching policy $M \in \mathcal{M}$, and characterize the corresponding customer stopping rule $\tau_\phi^{\pi, M}$ and purchasing rule $a_\phi^{\pi, M}$ and supplier stopping rule $\tau_\psi^{\pi, M}$ and selling rule $a_\psi^{\pi, M}$ to maximize the intermediary's expected profit

$$J^{\pi, M} = \mathbb{E} \left[\sum_{\phi \in H^T} p_\phi - \sum_{\psi \in H^T} p_\psi \right].$$

8.4.1 Upper Bound of the Intermediary's Optimal Profit

In this subsection, we establish an upper bound of the intermediary's optimal profit. As we will see in the subsequent subsections, this upper bound will be use to quantify the performance of our proposed heuristic policy.

Consider the following problem (B) that assumes that the intermediary is clairvoyant that she knows customer and supplier arrival processes H^T over $[0, T]$ at time 0:

$$\begin{aligned}
 & \max_{\{x_{\phi\psi}: \phi, \psi \in H^T\}} && \sum_{\phi, \psi \in H^T} (V^d(v_\phi) - V^s(c_\psi) - b(t_\psi - t_\phi)^+ - \bar{h}(t_\phi - t_\psi)^+) x_{\phi\psi} \\
 & \text{s.t.} && \sum_{\psi \in H^T} x_{\phi\psi} \leq 1, \quad \forall \phi \in H^T, \\
 & && \sum_{\phi \in H^T} x_{\phi\psi} \leq 1, \quad \forall \psi \in H^T, \\
 & && x_{\phi\psi} \in \{0, 1\}, \quad \forall \phi, \psi \in H^T.
 \end{aligned} \tag{B}$$

Note that problem (B) is simply a deterministic assignment problem. It has the following interpretation. If customer ϕ and supplier ψ are matched, then the intermediary collects revenue from customer ϕ with the amount that is equal to her virtual value, $V^d(v_\phi)$, subsidizes supplier ψ with the amount that is equal to her virtual cost, $V^s(c_\psi)$, and suffers from either customer ϕ 's waiting for supplier ψ or vice versa, depending on whoever arrives earlier. We denote by $\bar{J}(H^T)$ the optimal value of problem (B) conditional on customer and supplier arrival processes H^T . We have the following result. All proofs in this section can be found in Chen and Hu (2018).

Lemma 1 *The intermediary's profit under any pricing policy $\pi \in \Pi$ and matching policy $M \in \mathcal{M}$ is upper bounded by the expected optimal value of problem (B),*

$$J^{\pi, M} \leq \mathbb{E}[\bar{J}(H^T)].$$

8.4.2 A Simple Dynamic Policy: Asymptotic Optimality

In this subsection, we begin with characterizing the intermediary's optimal policy in an auxiliary setting wherein all uncertainties are washed away, and all customers and suppliers behave myopically. We then use this policy as a basis to develop another policy for the primary setting that takes into account customers' and suppliers' waiting disutility.

8.4.2.1 Optimal Policy in an Auxiliary Setting

We consider an auxiliary version of the primary stochastic model. In this auxiliary problem, the system is fully deterministic, and customers and suppliers are infinitesimal and myopic. To be precise, in the auxiliary problem, the intermediary solves the following optimization problem:

$$\begin{aligned} \max_{\hat{\pi} \in \Pi} \quad & \int_0^T \lambda^d \hat{\pi}_t^d \bar{F}^d(\hat{\pi}_t^d) dt - \int_0^T \lambda^s \hat{\pi}_t^s F^s(\hat{\pi}_t^s) dt \\ \text{s.t.} \quad & \lambda^d \bar{F}^d(\hat{\pi}_t^d) = \lambda^s F^s(\hat{\pi}_t^s), \quad \forall t \in [0, T], \end{aligned} \quad (\text{D})$$

where $\hat{\pi} = \{(\hat{\pi}_t^d, \hat{\pi}_t^s) : t \in [0, T]\}$ is an arbitrary measurable function from $[0, T]$ to \mathbb{R}_+^2 .

In problem (D), the intermediary determines the demand-side price trajectory $\{\hat{\pi}_t^d : t \in [0, T]\}$ and the supply-side price trajectory $\{\hat{\pi}_t^s : t \in [0, T]\}$ at time 0. The intermediary's pricing policy $\hat{\pi}$ is feasible if it clears the market at each point of time t , $\lambda^d \bar{F}^d(\hat{\pi}_t^d) = \lambda^s F^s(\hat{\pi}_t^s)$. Under the pricing policy $\hat{\pi}$, over the entire season $[0, T]$, the total revenue that the intermediary collects from customers is $\int_0^T \lambda^d \hat{\pi}_t^d \bar{F}^d(\hat{\pi}_t^d) dt$, and the total cost that the intermediary incurs from compensating suppliers is $\int_0^T \lambda^s \hat{\pi}_t^s F^s(\hat{\pi}_t^s) dt$. The intermediary aims at maximizing her net profit over the entire season.

Now, we characterize the intermediary's optimal policy and profit in this auxiliary setting.

Proposition 3 (Optimal solution to the deterministic problem) *The optimal solution to problem (D) is that the intermediary simply posts fixed prices p^* and w^* for customers and suppliers, respectively, throughout the horizon, where prices p^* and w^* always exist and are determined by the following conditions:*

(i) *(Demand-supply balancing condition)*

$$\lambda^d T \bar{F}^d(p^*) = \lambda^s T F^s(w^*) \equiv \mu^*; \quad (8.4)$$

(ii) *(Virtual value-cost balancing condition)*

$$\mu^* = \max\{\mu \in [0, \min\{\lambda^d T, \lambda^s T\}] : V(\mu) \geq 0\}, \quad (8.5)$$

$$\text{where } V(\mu) \equiv V^d(\bar{F}^{d,-1}(\mu/(\lambda^d T))) - V^s(F^{s,-1}(\mu/(\lambda^s T))).$$

Moreover, $p^* \geq w^*$. The optimal value of program (D) is

$$\bar{J}^* = (p^* - w^*)\mu^*. \quad (8.6)$$

We observe that the optimal price pair (p^*, w^*) is determined by Eqs. (8.4) and (8.5). Equation (8.4) is the market-clearing condition. Under this condition, the number of customers who purchase the product is equal to the number of suppliers who sell the product. Equation (8.5) entails that either the intermediary has matched the most number of customers and suppliers under the optimal price pair (p^*, w^*) and it is infeasible to match an additional pair of customer and supplier, $\mu^* = \min\{\lambda^d T, \lambda^s T\}$, or although it is feasible to match more pairs of customers and suppliers, $\mu > \mu^*$, by adjusting the price pair (p^*, w^*) , the marginal revenue that the intermediary collects from enabling one additional customer to get the

product is less than the marginal cost that the intermediary incurs from enabling one additional supplier to sell the product, i.e., the intermediary's marginal net profit from matching one additional customer with one additional supplier is negative, $V(\mu) < 0$.

The intermediary's optimal profit in this auxiliary setting, \bar{J}^* , has the following property.

Lemma 2 *We have*

$$E[\bar{J}(H^T)] \leq \bar{J}^*.$$

8.4.2.2 Waiting Adjusted Fixed Pricing Policy

We use the optimal policy in the auxiliary setting characterized above to propose a simple dynamic pricing and matching policy and show that this policy is nearly optimal. We begin by presenting the matching part of our policy, the *greedy* matching policy, denoted by M^g . Under this policy, the intermediary matches each demand (resp. supply) request as soon as a unit of supply (resp. demand) is available on a first-come-first-served basis. Therefore, the intermediary minimizes demand and supply mismatch at each point of time. Define

$$I_t \equiv \sum_{\psi \in H^s} \mathbf{1}\{\tau_\psi \leq t, a_\psi = 1\} - \sum_{\phi \in H^d} \mathbf{1}\{\tau_\phi \leq t, a_\phi = 1\}.$$

Therefore, at each point of time t , the number of unmatched supply is $(I_t)^+$ and the number of unmatched demand is $(I_t)^-$. This matching policy is natural, practical and fair.

Along with the simple greedy matching part of our policy, we next present the pricing part, the *waiting adjusted fixed pricing (FP) policy*, denoted by $\pi^{\text{WFP}} = \{\pi_t^{\text{WFP},d}, \pi_t^{\text{WFP},s} : t \in [0, T]\}$. This policy is constructed in the following way. Recall from the previous subsection that the intermediary's optimal pricing policy in the auxiliary deterministic myopic customer and supplier model is to post fixed prices p^* and w^* on the demand and supply sides, respectively. However, in our original stochastic system, although this policy is easy to implement, it does not lead to simple customers' and suppliers' behavior. The presence of customers' and suppliers' arrival uncertainty in the original model can cause them to wait to be matched with waiting disutilities, even though customers and suppliers do not strategize their entry to the matching pool given fixed prices. Therefore, if the intermediary posts prices p^* and w^* on the demand and supply sides, respectively, then a customer (resp. supplier) cannot make the purchasing (resp. selling) decision by merely comparing her valuation (resp. cost) with p^* (resp. w^*). She has to take into account the joint effects of the price p^* (resp. w^*) and the waiting disutility. To alleviate customers' and suppliers' computational burdens and ease their decisions, we require the intermediary to adjust the fixed prices p^* and w^* by

taking into account the waiting disutility that customers and suppliers incur, such that a customer (resp. supplier) can make an easy decision by merely comparing her valuation (resp. cost) with p^* (resp. w^*).

Formally, under policy π^{WFP} , the prices posted at each point of time $t \in [0, T]$ on demand and supply sides, respectively, are given by

$$\pi_t^{\text{WFP},d} = p^* - b \frac{E_{I_t^-}[s_\phi - t_\phi | t_\phi = t, v_\phi \geq p^*, (I_t^-)^+]}{\Pr(m_\phi = 1 | t_\phi = t, v_\phi \geq p^*, (I_t^-)^+)},$$

$$\pi_t^{\text{WFP},s} = w^* + h \frac{E[s_\psi - t_\psi | t_\psi = t, c_\psi \leq w^*]}{\Pr(m_\psi = 1 | t_\psi = t, c_\psi \leq w^*)},$$

where p^* and w^* are determined by Eqs. (8.4) and (8.5), respectively, and the expectations and the supply-demand mismatch quantity I_t are computed under the assumption that all customers (resp. suppliers) behave myopically, i.e., every customer ϕ (resp. supplier ψ) makes her purchasing (resp. selling) decision at her arrival time, $\tau_\phi = t_\phi$ (resp. $\tau_\psi = t_\psi$), and decides to purchase (resp. sell) if and only if her valuation (resp. cost) is no less (resp. more) than p^* (resp. w^*).

The waiting compensation terms have the following properties. First, the probability of being matched for a customer (resp. supplier) is in the denominator of the waiting compensation, because the monetary funds exchange hands only when a match is realized. Second, due to customers' and suppliers' different information structures, on the demand side, each customer ϕ 's expected time of staying in the system and the probability of being matched are conditional on the number of unmatched supply, $(I_{t_\phi^-})^+$. Since $(I_{t_\phi^-})^+$ is a random variable, the demand-side compensation term is random. As a result, the demand-side pricing policy $\pi^{\text{WFP},d}$ is a *contingent* policy. In contrast, on the supply side, because a supplier does not have the information of the number of unmatched supply or demand at any point of time, the supplier's expected time of staying in the system and the probability of being matched are not conditional on the number of unmatched supply or demand. Therefore, the supply-side compensation for each point of time is deterministic. As a result, the supply-side pricing policy $\pi^{\text{WFP},s}$ is a *deterministic* policy.

Now, we establish an equilibrium stopping and requesting rules for customers and suppliers when the intermediary follows the waiting adjusted FP and the greedy matching policy.

Theorem 8 (Strategic myopia) *Assume that the intermediary adopts the waiting adjusted FP policy π^{WFP} and the greedy matching policy M^g . Then in an equilibrium, all forward-looking buyers and sellers behave myopically, i.e., they will submit a request for matching upon arrival without delay if the buyer's valuation is no less than p^* or the seller's cost is no more than w^* .*

In our model, the stochastic nature of customers' and suppliers' arrival processes makes them to wait to be matched and hence to incur waiting disutilities. Therefore, the waiting adjustment terms play a vital role in compensating their losses and then induce them to behave myopically.

So far, we have shown with waiting compensation, the two-sided pricing policy π^{WFP} is dynamic. The following proposition characterizes the variability of the pricing policy in an asymptotic sense.

Proposition 4 (Waiting compensation) *For any $t \in [0, T]$ and any $k \in [0, 1]$, we have*

$$\begin{aligned} E_{I_t^-} [p^* - \pi_t^{\text{WFP},d}] &\leq b \frac{\min\{kT, T-t\} + T \min\{K, 1\}}{1 - \min\{K, 1\}}, \\ \pi_t^{\text{WFP},s} - w^* &\leq h \frac{\min\{kT, T-t\} + T \min\{K, 1\}}{1 - \min\{K, 1\}}, \end{aligned}$$

where $K \equiv 20/(\mu^* \min\{k^2, (1-t/T)^2\})$. In addition, we have the following results:

- (i) *If customers are fully patient, i.e., $b = 0$, then $\pi_t^{\text{WFP},d} = p^*$. If suppliers are fully patient, i.e., $h = 0$, then $\pi_t^{\text{WFP},s} = w^*$.*
- (ii) *Consider a sequence of systems. In the n -th system, $\lambda^{d,(n)} = n^{\alpha_d} \lambda^d$ with $\alpha_d > 0$ and $\lambda^{s,(n)} = n^{\alpha_s} \lambda^s$ with $\alpha_s > 0$. Denote $\underline{\alpha} \equiv \min\{\alpha_d, \alpha_s\}$. For any $t \in [0, T]$, we have*

$$\begin{aligned} \limsup_{n \rightarrow \infty} E_{I_t^-} [p^{*,(n)} - \pi_t^{\text{WFP},d,(n)}] &\leq O\left(\frac{1}{n^{\underline{\alpha}/3}\right), \\ \limsup_{n \rightarrow \infty} \pi_t^{\text{WFP},s,(n)} - w^{*,(n)} &\leq O\left(\frac{1}{n^{\underline{\alpha}/3}\right). \end{aligned}$$

We make the following observations from this proposition. First, if customers (resp. suppliers) are fully patient, i.e., $b = 0$ (resp. $h = 0$), they do not incur any waiting disutility, although they may spend time in waiting to be matched. Therefore, π^{WFP} does not need to be adjusted from the base fixed prices (p^*, w^*) . Second, in the high-volume regime in which customers' and suppliers' arrival rates grow large (scaled by n), regardless of whether they grow at the same or different speeds (measured by α_d and α_s), the waiting adjusted terms on both demand and supply sides in the policy π^{WFP} diminish to zero, i.e., the policy π^{WFP} tends to be the fixed pricing policy (p^*, w^*) . In addition, as n grows large, the variability of π^{WFP} decays to zero at a speed that is no slower than $1/n^{\underline{\alpha}/3}$.

Next, we analyze the performance of the waiting adjusted FP and the greedy matching policy M^g .

Theorem 9 (Performance guarantee) *Under the waiting adjusted FP policy π^{WFP} and the greedy matching policy M^g ,*

$$\frac{J\pi^{\text{WFP},M^g}}{E[\bar{J}(HT)]} \geq \frac{J\pi^{\text{WFP},M^g}}{\bar{J}^*} \geq 1 - \left(1 + \frac{2(b+h)T}{3(p^* - w^*)}\right) \frac{1}{\sqrt{\mu^*}}.$$

In addition, consider a sequence of systems. In the n -th system, $\lambda^{d,(n)} = n^{\alpha_d} \lambda^d$ with $\alpha_d > 0$ and $\lambda^{s,(n)} = n^{\alpha_s} \lambda^s$ with $\alpha_s > 0$. Denote $\underline{\alpha} \equiv \min \{\alpha_d, \alpha_s\}$. Therefore,

$$\frac{J\pi^{\text{WFP}, M^g, (n)}}{E[\bar{J}^{(n)}(HT)]} \geq \frac{J\pi^{\text{WFP}, M^g, (n)}}{\bar{J}^{*,(n)}} \geq 1 - O\left(\frac{1}{\sqrt{n^{\underline{\alpha}}}}\right).$$

Theorem 9 has the following implications. First, as both customers' and suppliers' arrival rates grow large (scaled by n), regardless of whether they grow up at the same or different speeds (measured by α_d and α_s), the simple, waiting adjusted FP policy π^{WFP} and greedy matching policy M^g , are asymptotically optimal. Second, as n grows large, the relative profit loss of implementing the simple heuristic policy, compared to the optimal mechanism, converges to zero at a speed that is no slower than $1/\sqrt{n^{\underline{\alpha}}}$, where $\underline{\alpha} = \min \{\alpha_d, \alpha_s\}$. Put differently, the relative additional benefit of implementing any more sophisticated policy, than our simple heuristic, is no more than a magnitude of $1/\sqrt{n^{\underline{\alpha}}}$.

8.5 Conclusion

Operations management is about matching supply with demand at the operational level. We study pricing and matching problems for a sharing economy platform to coordinate demand and supply. First, we show that the commonly used fixed commission contract by a platform achieves a guaranteed portion of the optimal expected profit under full flexibility of optimally choosing both wage and price for every possible market condition. Second, we formulate a stochastic dynamic programming model to study the problem of matching heterogeneous types of demand and supply. For that model, we propose the modified Monge condition to establish the priority structure of the optimal matching policy, as well as a deterministic re-solving heuristic for computing the optimal matching decisions. Lastly, we study the joint pricing and matching decision by a platform for a single service and take into account suppliers' and customers' forward-looking behavior. We propose a simple pricing and matching policy under which suppliers and customers behave myopically. We use the mechanism design approach to show that our policy is asymptotically optimal when the market sizes of both sides become sufficiently large.

References

- Chen Y, Hu M (2018, forthcoming) Pricing and matching with forward-looking buyers and sellers. *Manuf Serv Oper Manag*. <http://ssrn.com/abstract=2859864>
- Chen L, Mislove A, Wilson C (2015) Peeking beneath the hood of uber. In: Proceedings of the 2015 ACM conference on Internet measurement conference. ACM, New York, pp 495–508

- Hu M, Zhou Y (2018a) Dynamic type matching, available at SSRN: <http://ssrn.com/abstract=2592622>
- Hu M, Zhou Y (2018b) Price, wage and fixed commission in on-demand matching, available at SSRN: <http://ssrn.com/abstract=2949513>
- Xu P, Allgor R, Graves S (2009) Benefits of reevaluating real-time order fulfillment decisions. *Manuf Serv Oper Manag* 11(2):340–355
- Yu Y, Chen X, Zhang F (2015) Dynamic capacity management with general upgrading. *Oper Res* 63(6):1372–1389

Chapter 9

Large-Scale Service Marketplaces: The Role of the Moderating Firm



Gad Allon, Achal Bassamboo, and Eren B. Çil

Abstract Recently, large-scale, web-based service marketplaces, where many small service providers compete among themselves in catering to customers with diverse needs, have emerged. Customers who frequent these marketplaces seek quick resolutions and thus are usually willing to trade prices with waiting times. The main goal of the paper is to discuss the role of the moderating firm in facilitating information gathering, operational efficiency, and communication among agents in service marketplaces. Surprisingly, we show that operational efficiency may be detrimental to the overall efficiency of the marketplace. Further, we establish that to reap the “expected” gains of operational efficiency, the moderating firm may need to complement the operational efficiency by enabling communication among its agents. The study emphasizes the scale of such marketplaces and the impact it has on the outcomes (This chapter is based on our published paper (see Allon et al., *Manag Sci* 58:1854–1872, 2012).).

9.1 Introduction

Recently, large-scale, web-based service marketplaces, where many small service providers (agents) compete among themselves in catering to customers with diverse needs, have emerged. Customers who frequent these marketplaces seek quick resolutions for their temporary problems and thus are usually willing to trade prices with waiting times. These marketplaces are typically operated by an independent

G. Allon (✉)
University of Pennsylvania, Philadelphia, PA, USA
e-mail: gadallon@wharton.upenn.edu

A. Bassamboo
Northwestern University, Evanston, IL, USA
e-mail: a-bassamboo@kellogg.northwestern.edu

E. B. Çil
University of Oregon, Eugene, OR, USA
e-mail: erencil@uoregon.edu

firm, which we shall refer to as the *moderating firm*. The moderating firm establishes the infrastructure for the interaction between customers and agents. In particular, it provides the customers and the agents with the information required to make their decisions. These moderating firms vary with respect to their involvement in the marketplace. They can introduce operational tools that specify how the customers and the agents are matched together. For instance, while some of the moderating firms allow customers to choose a specific service provider directly, others allow customers to post their needs and let service providers apply, postponing the service provider selection decision of the customers until they obtain enough information about agents' availability. Moreover, moderating firms can introduce strategic tools that allow communication and collaboration among the agents themselves. These different involvements result in different economic and operational systems, and thus vary in their level of efficiency and the outcomes for both customers and service providers.

A typical example of such a marketplace is UpWork.com, where over 10,000,000 registered freelancers compete to provide online solutions. UpWork.com allows for two types of interaction between customers and service providers. Customers can go directly to a programmer and ask him to provide the service. The customers are then queued for this specific agent. In this type of interaction, most of the time is spent waiting for the agent to complete his previous jobs (36% of the waiting time is spent from the moment the customer chooses the agent until the agent begins working.¹). On the other hand, UpWork.com also allows customers to post jobs and wait while agents apply for the job. In this type of interaction, a negligible amount of time passes until more than 10 agents apply, leaving the decision at the hands of the customer. Another large-scale, online service marketplace is ServiceLive.com, which is a start-up owned by Sears Holding Company. ServiceLive.com (with the slogan of "your price, your time") caters to time and price-conscious customers and service providers in the home repair and improvement arena. ServiceLive.com allows its customers to choose among multiple agents once they describe their projects. This type of interaction between customers and service providers is equivalent to the second one described for UpWork.com. Both UpWork.com and ServiceLive.com receive 10% of the revenue obtained by the providers at service completion. In both marketplaces, the moderating firms allow customers to browse among tens of thousands of agents and communicate with different providers.

Both UpWork.com and ServiceLive.com are part of a growing industry of online service marketplaces. According to a survey conducted by Upwork.com (see Upwork Press Release 2016), freelancers contributed over \$1 trillion in freelance earnings to the economy. The same surveys reports that more than three-quarters of freelancers view their jobs as more appealing than a traditional job. More importantly, nearly half of full-time freelancers surveyed by upwork.com mentioned that they raised their rates in the past year, and more than half plan indicate their plans to raise their rates next year.

¹This is based on data obtained from UpWork.com for about 10,000 randomly chosen transactions.

Motivated by these online service marketplaces, we aim to study the moderating firm's role in the service marketplace where the objective of the individual players, customers as well as service providers, is to maximize their own utility. We distinguish between three degrees of moderating firms' involvement in such markets:

1. *No-intervention*: the moderating firm restricts its involvement to providing the facility for agents to advertise their services and set their prices, and for customers to compare the different agents.
2. *Operational efficiency*: the moderating firm provides additional mechanisms that facilitate efficient matching between customers and service providers. These mechanisms aim at reducing the inefficiency associated with having the right agent with the right capability idle while a customer with similar needs is waiting in line for another agent. As we will discuss, a system in which customers post their needs and wait for agents' applications is an example of such a mechanism.
3. *Enabling Communication*: the moderating firm may allow providers to communicate among themselves and exchange information on prices and job requirements.

To study the different configurations possible in such marketplaces we consider a sequence of related games where the set of possible strategies and the solution concepts vary to reflect the different modes of interaction available in the marketplace, either between the customer and the agents or between the agents themselves. Specifically, we study the following three games:

No-intervention Model In this game, each agent chooses his price and operates as a single-server queue. Customers then choose agents based on prices and waiting times. We characterize the Subgame Perfect Nash equilibrium in this game.

Operational Efficiency Model In this game, the mechanism introduced by the moderating firm efficiently matches customers interested in purchasing the service at a particular price with the available agents charging that or a lower price. This mechanism achieves the desired level of efficiency by virtually grouping all agents charging the same price. In contrast to the no-intervention model, customers do not need to commit to a specific agent upon their arrival.

Communication Enabled Model In this game, agents can exchange information in a non-committal, costless manner. As in the model with operational efficiency, all the agents charging the same price are virtually grouped, and customers choose the price/sub-pool. We would be interested in allowing limited pre-play communication among the agents within a noncooperative structure; i.e., the agents are free to discuss their pricing strategies but not allowed to make binding commitments. Ray (1996) claims that the possibility of pre-play communication have motivated the notion of strong Nash equilibrium, see Aumann (1959), which requires stability against deviations by every conceivable coalition. Following this idea, we use a refinement of the Subgame Perfect Nash Equilibrium concept that requires the equilibrium to be (limited size) coalition proof.

We next state our key findings along with the contributions of the paper:

1. We appear to be the first to distinguish between tools aimed at increasing the operational efficiency and tools aimed at changing the nature of the strategic interaction by enabling communication. We show these tools have a non-trivial impact on the outcomes for all involved parties.
2. In analyzing a market with operational efficiency, we first show that only the prices in a small neighborhood of the operating cost of agents are sustained as equilibrium outcomes when supply exceeds demand. Further, when demand exceeds supply, we are able to show that operational efficiency leads to multiple equilibria in markets with a sufficiently large number of agents. In many of these equilibria, the emerging prices are lower than those arising in the market with no-intervention.
3. We show that to overcome the possible deterioration of the profits discussed above, the moderating firm can allow for communication among the agents, even if done through a non-binding mechanism. The main contribution of this result is in showing that the operational efficiency needs to be complemented with the ability to communicate in order to obtain desirable outcomes for the involved parties. These desirable outcomes are only achievable in a marketplace where demand exceeds supply. Therefore, the contribution is also in highlighting the fact that it is crucial to understand the specific market conditions in terms of the ratio between demand and supply.

9.2 Literature Review

The previous work related with this chapter can be divided into two categories. The first category consists of research that studies the applications of queuing theory in service systems. The second one consists of research focused on developing approximations to analyze complex service systems.

Service systems with customers, who are both price and time sensitive, have attracted the attention of researchers for many years. The analysis of such systems dates back to Naor's seminal work (see Naor 1969), which analyzes customer behavior in a single-server queuing system. Motivated by his work, many researchers study the service systems facing price- and delay-sensitive customers in various settings. We refer the reader to Hassin and Haviv (2003) for an extensive summary of the early attempts in this line of research. More recently, Cachon and Harker (2002) and Allon and Federguen (2007) studies the competition between multiple firms offering substitute but differentiated services by modeling the customer behavior implicitly via an exogenously given demand function. An alternative approach is followed in Chen and Wan (2003), where authors examine the customers' choice problem explicitly by embedding it into the firms' pricing problem. Other notable examples focusing on the customers' demand decision in competition models are Ha et al. (2003), and Cachon and Zhang (2007).

The pricing and the capacity planning problem of the service systems can easily become analytically intractable when trying to study more complex models, such as a multi-server queueing systems. Recognizing this difficulty, many researchers seek robust and accurate approximations to analyze multi-server queues. Halfin and Whitt (1981) is the first paper that proposes and analyzes a multi-server framework. This framework is aimed at developing approximations, which are asymptotically correct, for multi-server systems. It has been applied by many researchers to study the pricing and service design problem of a monopoly in more realistic and detailed settings. Armony and Maglaras (2004) and Maglaras and Zeevi (2005) are examples of recent work using the asymptotic analysis to tackle complexity of these problems. Furthermore, Garnett et al. (2002), Ward and Glynn (2003), and Zeltyn and Mandelbaum (2005) extend the asymptotic analysis of Markovian queueing system by considering customer abandonments.

The idea of using approximation methods can also be applied to characterize the equilibrium behavior of the firms in a competitive environment. To our knowledge, Allon and Gurvich (2010) is the first paper studying competition among complex queueing systems by using asymptotic analysis to approximate the queueing dynamics. Another recent paper studying the equilibrium characterization of a competitive marketplace using asymptotic analysis is Chen et al. (2008). They consider a marketplace with multiple suppliers competing with each other over their prices and target lead times. There are two main differences between these two papers and our work. First, both of them study a service environment with a fixed number of decision makers (firms) while the number of decision makers in our marketplace (agents) is growing. Second, they only consider a competitive environment where the firms behave individually. In contrast, we study the non-cooperative case as well as the case where the agents have a limited level of collaboration.

In the field of operations management (OM), the majority of the papers employing game-theoretic foundations study non-cooperative settings. For an excellent survey, we refer to Cachon and Netessine (2004). There is also a growing literature that studies the OM problems in the context of cooperative game theory. Nagarajan and Sošić (2008) provide an extensive summary of the applications of cooperative game theory in supply chain management. Notable examples are the formation of coalitions among retailers to share their inventories, suppliers, and marketing powers (see Granot and Sošić 2005; Sošić 2006; Nagarajan and Sošić 2007). This body of research is related with our work, where we look for the limited collaboration among agents.

Our work may also be viewed as related to the literature on labor markets that studies the wage dynamics (see Burdett and Mortensen 1998; Manning 2003, 2004; Michaelides 2010). In both our model and labor economics literature, people or firms with service needs seek an employee or an agent to perform the job they requested. In our model, service seekers trade-off time they need to wait until their job starts and cost, the phenomenon generally disregarded in labor economics literature. Further, our focus is on a market for temporary help, which means that the engagement between sides ends upon the service completion. This stands in

contrast to the labor economics literature in which the engagement is assumed to be permanent. It is also important to note the difference between interventions studied in our model and the ones in the labor economics literature. Unlike the interventions we studied, which focus on improving operational efficiency, the interventions discussed in labor economics are usually aimed at regulating wages directly. The model we consider in this chapter also differs from the literature on market microstructure. This body of literature studies market makers who can set prices and hold inventories of assets in order to stabilize markets (see Garman 1976; Amihud and Mendelson 1980; Ho and Stoll 1983) and a comprehensive survey by Biais et al. (2005). However, the moderating firm considered in this chapter has no direct price-setting power and cannot respond to customers' service requests. Furthermore, papers studying market microstructure disregard the operational details such as waiting and idleness.

9.3 Model Formulation

Consider a service marketplace where agents and customers make their decisions in order to maximize their individual utilities. Customers' need for the service is generated according to a Poisson process with rate Λ . This forms the "potential demand" for the marketplace. A customer decides whether to join the marketplace or not: If she decides not to join the system, her utility is zero. If she joins the system, she decides who would process her job. The customers who join the marketplace form the "effective demand" for the marketplace. The exact nature of this decision depends on the specific structure of the marketplace, decided upfront by the moderating firm. We shall elaborate on the choices of customers in Sects. 9.4, 9.5, and 9.6. We assume that the service time required to satisfy the requests of a given customer is exponentially distributed with rate μ . Without loss of generality, we let $\mu = 1$. When the service of a customer is successfully completed, she pays the price of the service, earns a reward of R , and incurs a waiting cost of c per unit time until her service commences.² As the customers visiting the marketplace seek temporary help, a customer joining the system may become impatient while waiting for her service to start and abandon. In this case, the abandoning customer does not pay any price or earn any reward, but she does incur a waiting cost for the time she spends in the system. We assume that customers' abandonment times are independent of all other stochastic components and are exponentially distributed with mean m_a . Customers decide whether to request service or not and by whom to be served according to their expected utility. The expected utility of a customer is based on the reward, the price and the anticipated waiting time.

²Our model can also be used to study a setting where customers incurs waiting cost also during their service. One can incorporate that by modifying the customer reward from R to $R - c/\mu$.

The above summarizes the demand arriving to the marketplace. Next, we discuss the service provision in a marketplace with k ex-ante identical agents.³ The only decision of an agent is to choose a price for his service; each agent makes this decision independently. Let (p_1, \dots, p_k) denote the resulting price vector with p_n being the price chosen by the n th agent. We normalized the operating cost of the agents to zero for notational convenience. The expected revenue of an agent depends on the price he chooses and his demand volume.

We refer to the ratio $\Lambda/(\mu k)$ as the demand-supply ratio of the system and denote it by ρ . The demand-supply ratio is a first order measure for the mismatch between aggregate demand and the total processing capacity. Marketplaces vary with respect to their demand-supply ratio, ρ , and, as we shall discuss, ρ has a significant impact on the market outcome. We broadly categorize marketplaces into two: Buyer's market where $\rho \leq 1$, and seller's market where $\rho > 1$.

9.4 No-Intervention Model

The essential role of the moderating firm in a large scale marketplace is to set up the infrastructure for the interaction between players. This is crucial because all players have to be equipped with the necessary information, such as prices to make their decisions, yet individual players cannot gather this information on their own. When the moderating firm provides only the required information, it has no impact on the strategic interaction taking place in the marketplace. We thus refer to such a setting as the no-intervention model. We analyze the dynamics of a large-scale marketplace in the no-intervention model not only to derive insights about the behavior of the self-interested and competing players in such a system, but also to build a benchmark for the cases in which the moderating firm introduces additional features which change the nature of the marketplace. Therefore, in this section, we study the behavior of a marketplace where the moderating firm confines itself to aggregating and providing information.

We model the strategic interaction between the agents and the customers as a sequential move game. Given the setup of Sect. 9.3, along with the above mentioned role of the moderating firm, the agents first announce their prices. Each arriving customer observes these prices and decides whether to request service or not. Further, if a customer decides to join the system, she also chooses the agent who processes her service request. The service of a customer starts immediately if the agent she chooses is available. Otherwise, she joins the queue in front of the agent and waits for her service to commence. We denote the fraction of customers choosing agent- n by D_n . Then, ΛD_n is the demand volume for agent- n .

³We will discuss a model with heterogeneous agents in Sect. 9.7.

More specifically, each agent's operations can be modeled as an $M/M/1 + M$ queueing system⁴ where the arrival rate of customers depends on the strategies of customers and agents.⁵ If the rate of customers who request service from an agent charging price p is λ , the utility of a customer requesting service from this agent is $U(\lambda, p) = (R - p)[1 - \beta(\lambda)] - W(\lambda)c$, where $\beta(\lambda)$, which will be referred to as the abandonment function, is the probability of abandonment, and $W(\lambda)$ is the expected waiting time, in an $M/M/1 + M$ system with arrival rate λ , service rate 1, and abandonment rate $1/m_a$. Using queueing theory, the utility of customers can be rewritten as $U(\lambda, p) = (R - p + cm_a)[1 - \beta(\lambda)] - cm_a$. Similarly, the revenue of that agent is $V(\lambda, p) = p\lambda[1 - \beta(\lambda)]$. It is important to note that $V(\lambda, p)$ is the revenue rate of an agent, but throughout the paper we will refer to it as the revenue for ease of exposition.

As we consider a sequential move game, we are interested in the Subgame Perfect Nash Equilibrium (SPNE) of the game. We begin by characterizing the equilibrium in the second stage game where customers make their service requests given the agents' pricing decisions. Then, based on the second stage equilibrium, we derive the equilibrium of the first stage in which only agents make pricing decisions.

Fixing the agents' strategies $(p_n)_{n=1}^k$, an arriving customer observes the agents' prices and chooses the agent who maximizes her utility, anticipating the behavior of all other customers. Therefore, in equilibrium a customer chooses an agent only if the utility she obtains from him (weakly) dominates her utility from any other agent. This is also known as "Nash Flow Equilibrium" (see Roughgarden 2005) in the congestion games literature. We formally define the Customer Equilibrium as follows:

Definition 1 (Customers Equilibrium) Given $(p_n)_{n=1}^k$, we say that $(D_n)_{n=1}^k$ is a Customers Equilibrium if the following conditions are satisfied:

1. For any n with $D_n > 0$, we have that $U(\Delta D_n, p_n) \geq U(\Delta D_m, p_m) \geq 0$, for all $m \leq k$.
2. If $U(\Delta D_n, p_n) > 0$ for some $n \leq k$, then $\sum_{n=1}^k D_n = 1$.

The first condition of the Customer Equilibrium requires that customers request service from an agent in equilibrium only if that agent is one of their best alternatives. Moreover, the second condition ensures that all customers join the system if it is possible to earn strictly positive utility by requesting service from an agent. Customer Equilibrium exists by the continuity of the utility functions and Rath (1992). In the following proposition, we show that for any given price vector, the second stage game has a unique equilibrium.

⁴ $+M$ notation denotes the exponential abandonment times.

⁵Note that an agent can process more than one jobs at the same time in certain settings. In such settings, a processor sharing model will be a more appropriate queueing model, yet these models are known to be significantly more complex than our queueing model. Our model can be viewed as an approximation of such settings.

Proposition 1 *Given a price vector $(p_n)_{n=1}^k$, there is a unique Customer Equilibrium.*

Since the Customer Equilibrium is unique for any given price vector, we denote the fraction of customers requesting service from agent- n in equilibrium by $D_n^{CE}(p_1, \dots, p_k)$ when (p_1, \dots, p_k) are the prices announced by agents. $D_n^{CE}(p_1, \dots, p_k)$ is well defined in the light of Proposition 1.

We can now move to the first stage game which is played only among the agents. An equilibrium in this stage requires that none of the agents can improve his revenues by deviating unilaterally while taking the customers' response into account. We formalize this in the following definition:

Definition 2 (Subgame Perfect Nash Equilibrium) Let $(D_n, p_n)_{n=1}^k$ summarize the strategy of all players in the market for all $n = 1, \dots, k$. Then, $(D_n, p_n)_{n=1}^k$ is a SPNE if the following conditions are satisfied:

1. $D_n = D_n^{CE}(p_1, \dots, p_k)$ for all $n \leq k$.
2. For any $\ell \leq k$, we have

$$V(\Delta D_\ell, p_\ell) = \max_{p'} V(\Delta D_\ell^{CE}(p_1, \dots, p_{\ell-1}, p', p_{\ell+1}, \dots, p_k), p').$$

The first condition requires that $(D_n)_{n=1}^k$ arises in equilibrium in the second stage game. The second condition states that none of the agents has incentive to change his price. Note that agents take into account the impact price changes have on the Customer Equilibrium, and thus on demand.

9.4.1 Characterization of SPNE

In this section, we restrict attention to symmetric SPNE where all agents charge the same price p in the first stage. This is a natural choice since all agents are identical. We will discuss non-symmetric equilibria in Sect. 9.7.

A price p emerges in equilibrium in the first stage if a single agent chooses to charge p to maximize his revenues given that all other agents announce p . When all other $k - 1$ agents announce p , a generic agent, say agent- ℓ , solves the following maximization problem to determine his best-response:

$$\max_{p_\ell \geq 0} p_\ell \Delta D_\ell^{CE}(p, \dots, p, p_\ell, p, \dots, p) [1 - \beta(\Delta D_\ell^{CE}(p, \dots, p, p_\ell, p, \dots, p))] \quad (9.1)$$

In this problem, the objective function is the revenue of agent- ℓ when he charges p_ℓ and the remaining agents charge p . Thus, p is a symmetric equilibrium in the first stage game if it is a solution to the above problem. We denote the symmetric SPNE by (D^*, p^*) where all agents charge p^* and each agent has a demand of ΔD^* , i.e. $D_n^{CE}(p, \dots, p) = D^*$ for any $n \leq k$. We characterize the symmetric SPNE in the following theorem:

Theorem 1 *If $\beta(\lambda)$ is concave, then there exists a symmetric SPNE. Furthermore, the symmetric SPNE is characterized as follows:*

1. *If $\Lambda \geq k\lambda^0$, then the symmetric SPNE is*

$$(D^*, p^*) = \left(\frac{\min\{\lambda^{\text{mon}}, \rho\}}{\Lambda}, R + cm_a - \frac{cm_a}{1 - \beta(\min\{\lambda^{\text{mon}}, \rho\})} \right).$$

2. *If $\Lambda \leq k\lambda^0$, then the symmetric SPNE is*

$$(D^*, p^*) = \left(\frac{1}{k}, (R + cm_a) - \frac{(R + cm_a)(k - 1)}{k/(1 - v(\rho)) - 1} \right).$$

Here λ^{mon} is the unique solution to $1 - \beta(\lambda) - \lambda\beta'(\lambda) = cm_a/(R + cm_a)$, λ^0 is the unique solution to

$$(R + cm_a)(k - 1) - \frac{cm_a}{1 - \beta(\lambda)} \left(\frac{k}{1 - v(\lambda)} - 1 \right) = 0,$$

and $v(\lambda) = \lambda\beta'(\lambda)/(1 - \beta(\lambda))$.

Similar to Theorems 1–3 in Chen and Wan (2003), the above result suggests that agents behave as local monopolists and charge their monopoly prices when the arrival rate is sufficiently high. Moreover, in this case, agents may choose not to cover the market completely. However, once the arrival rate becomes less than λ^0 , the equilibrium price will be pushed down as the agents are engaged in a cut-throat competition, where intensity of competition can be quantified by the strictly positive utility left for customers in the equilibrium. It is also worth noting that utility of customers in the equilibrium increases as the arrival rate decreases.

Remark 1 Concavity of the abandonment function, $\beta(\lambda)$, is a sufficient condition for the existence of symmetric equilibrium. In Lemma 1 in Allon et al. (2012), we show that $\beta(\lambda)$ is concave when $m_a \leq 1$, i.e. abandonment rate is higher than service rate. Furthermore, conducting a numerical study, we observe that $\beta(\lambda)$ is concave even for $1 \leq m_a \leq 2$. However, for higher values of m_a , the function $\beta(\lambda)$ is not concave in λ . This is not surprising given the complicated structure of queueing systems with impatient customers. For instance, Armony et al. (2009) shows the difficulty of proving the convexity of the expected head-count in the steady state of a system with customer abandonments. Even though $\beta(\lambda)$ is not concave, there can be a symmetric SPNE, and the above theorem characterizes this symmetric equilibrium. Numerically, we see that the equilibrium candidate characterized above still emerges as the symmetric SPNE when $\beta(\lambda)$ is not concave. In this numerical study, we consider a marketplace where $R = 1$, $c \in \{0.05, 0.06, \dots, 0.2\}$, and $k = 50$. Then, we study five scenarios that differ in the average abandonment time m_a and lead to non-concave $\beta(\lambda)$. We assume $m_a \in \{5, 6, \dots, 10\}$. For each of these scenarios, we show that the price proposed as equilibrium price in Theorem 1 is equilibrium by varying the arrival rate Λ on a grid from 10 to 50 with a step size of 1.

9.5 Operational Efficiency Model

In the previous section, we characterized the market outcome in the absence of any intervention on the part of the moderating firm. We now turn to discuss the impact of different mechanisms used by the moderating firm. As we discussed in the introduction, the moderating firm may provide a mechanism that improves the operational efficiency of the whole system by efficiently matching customers and agents. This mechanism aims at reducing inefficiency due to the possibility of having a customer waiting in line for a busy agent while an agent who can serve her is idle. This efficiency improvement is equivalent to virtually grouping the agents charging the same price. For instance, UpWork.com achieves this goal by allowing customers to post their needs and allowing service providers to apply to these postings. When a customer posts a job at UpWork.com, agents that are willing to serve this customer apply to the posting. Among the applicants charging less than what the customer wants to pay, the customer will favor agents based on their immediate availability. The main driver of the operational efficiency in this setting is the fact that customers no longer need to specify an agent upon their arrival because the job posting mechanism allows customers to postpone their service request decisions until they have enough information about the availability of the providers.

In this section, we modify the service marketplace considered in Sect. 9.4 by assuming that the mechanism introduced by the moderating firm ensures that customers do not stay in line when there is an idle agent willing to serve them by charging the price they want to pay or less. This can be modeled as a queuing network where the agents announcing the same price are virtually grouped together. Once each agent announces a price per customer to be served, we can construct a resulting price vector $(p_n)_{n=1}^N$ where $N \leq k$ is the number of different prices announced by the agents. We refer to the agents announcing the price p_n as sub-pool- n and denote the number of agents in the sub-pool- n by y_n . Hence, $(p_n, y_n)_{n=1}^N$ summarizes the strategy of all agents.

Under this mechanism, we model the customer decision making and experience as follows: If there are different prices announced by the agents, i.e., $N > 1$, the customer chooses a sub-pool from which she requests the service. We refer to the price charged by this sub-pool as the “preferred price”. Each customer who decides to join the system enters the service immediately if there is an available agent either in the sub-pool she chooses or in any sub-pool announcing a price less than her preferred price. Moreover, the customer is served by the sub-pool offering the lowest price among all available sub-pools. Otherwise, she waits in a queue until an agent, who charges a price less than or equal to her preferred price, becomes available. We denote the fraction of customers requesting service from sub-pool- n by D_n . In this model of customer experience, there are two crucial features: (1) The service of an arriving customer commences immediately when there are available agents charging less than or equal to her preferred price, (2) If they have to wait, customers no longer wait for a specific agent rather for an available agent.

As we model the marketplace as a queuing network, the operations of each sub-pool depend on the operations of the other sub-pools. For instance, each sub-pool may handle customers from the other sub-pools (giving priority to its “own” customers) while some of the other sub-pools are serving its customers. Therefore, given the strategies of agents, $(p_n, y_n)_{n=1}^N$, and the service decisions of customers, $(D_n)_{n=1}^N$, the expected utility of a customer choosing the sub-pool- ℓ depends on all of these decisions, and can be written as:

$$U_\ell(D_1, \dots, D_N; p_1, \dots, p_N; y_1, \dots, y_N) = \text{PServ}_{\ell\ell} [(R - p_\ell + cm_a)(1 - \beta_\ell) - cm_a] \\ + \sum_{m \neq \ell} \text{PServ}_{\ell m} (R - p_m),$$

where $\beta_\ell(D_1, \dots, D_N; p_1, \dots, p_N; y_1, \dots, y_N)$ denotes the probability of abandonment in the sub-pool- ℓ , and $\text{PServ}_{\ell m}(D_1, \dots, D_N; p_1, \dots, p_N; y_1, \dots, y_N)$ denotes the probability that a customer choosing the sub-pool- ℓ is served by the sub-pool- m when ΛD_n is the rate of customer arrival to the sub-pool- n for $n = 1, \dots, N$. We want to note that for any sub-pool- ℓ , $\text{PServ}_{\ell m} = 0$ for any m such that $p_m > p_\ell$ since customer choosing sub-pool- ℓ cannot be served by a sub-pool charging more than p_ℓ . Furthermore, the revenue of an agent in sub-pool- ℓ is: $V_\ell(D_1, \dots, D_N; p_1, \dots, p_N; y_1, \dots, y_N) = p_\ell \sigma_\ell(D_1, \dots, D_N; p_1, \dots, p_N; y_1, \dots, y_N)$, where $\sigma_\ell(\dots; \dots; \dots)$ is utilization of agents in sub-pool- ℓ when ΛD_n is the rate of customer arrival to the sub-pool- n for $n = 1, \dots, N$. Here, we assume that a customer choosing the sub-pool- ℓ pays p_m when she is served by sub-pool- m for $m \neq \ell$.

It is also worth noting that a marketplace operates as an $M/M/k + M$ system when all agents charge the same price. This allows us to employ the well-known limiting behavior of the multi-server systems to characterize the market outcome. Furthermore, in the case, where the agents announce different prices, we will show that the interdependency between the sub-pools announcing different prices diminishes as the market grows. In fact, large-scale marketplaces operate “almost like” the combination of independent multi-server systems.

The strategic interaction between the agents and the customers is modeled, as before, as a sequential move game. However, we use a slightly different second stage equilibrium than the one in Definition 1 since the customers decision and utility is changed by the new mechanism. The new customer equilibrium, which we refer to as Market Customer Equilibrium, uses the concept of Nash Flow Equilibrium with the requirement that customers only care for the prices announced by the sub-pools instead of individual prices.

Definition 3 (Market Customers Equilibrium) Given $(p_n, y_n)_{n=1}^N$, we say that $(D_n)_{n=1}^N$ is a Market Customers Equilibrium (MCE) if the following conditions are satisfied:

1. For any ℓ with $D_\ell > 0$, we have that $U_\ell(D_1, \dots, D_N; p_1, \dots, p_N; y_1, \dots, y_N) \geq U_m(D_1, \dots, D_N; p_1, \dots, p_N; y_1, \dots, y_N)$, for all $m \leq N$.
2. If $U_\ell(D_1, \dots, D_N; p_1, \dots, p_N; y_1, \dots, y_N) > 0$ for some $\ell \leq N$, then $\sum_{n=1}^N D_n = 1$.

While *MCE* always exists by the continuity of the utility functions and Rath (1992), its uniqueness cannot be guaranteed. For notational convenience, we shall assume that the best outcome from the customer perspective arises when there are multiple *MCE* (In fact, it can be shown that the limit of all *MCEs* is unique as the number of agents in the market grows). As the outcome is assumed to be unique, we denote the fraction of customers requesting service from sub-pool- n in a Market Customer Equilibrium by $D_n^{MCE}(p_1, \dots, p_N; y_1, \dots, y_N)$ when $(p_n, y_n)_{n=1}^N$ is a tuple of two vectors whose components are the prices and the number of agents announcing them.

Agents make pricing decisions in the first stage of the game. Unlike the no-intervention model, we need to account for two types of unilateral deviation of agents: an agent can either choose to deviate by joining an existing sub-pool or announce a new price. Therefore, an equilibrium in the first stage should be immune to any of these two deviations. One can show that, as the market grows, there exists a profitable unilateral deviation from any price in a buyer's market. In analyzing such markets, we would like to highlight the following two observations: (1) The arising system dynamic is too complex for exact analysis yet amenable to asymptotic analysis. (2) While a single agent, indeed, may have profitable deviations from every price in a buyer's market, the gains from deviations are small and diminish as the market grows. Thus, following Dixon (1987) and recently Allon and Gurvich (2010), we study a somewhat weaker notion of equilibrium, which allows us to characterize the market outcome (if one exists), as the market grows even when Nash equilibrium does not exist. To this end, we consider a sequence of marketplaces indexed by the number of agents, i.e., there are k agents in the k th marketplace. The arrival rate in the k th marketplace is assumed to be $\Lambda^k = \rho k$. This ensures that the demand-supply ratio is constant along the sequence of marketplaces. Then, in each market, we focus on an equilibrium concept, which requires immunity against only deviations that improve the revenue of an agent by at least $\epsilon \geq 0$ as formally stated in Definition 4 (See below). We refer to ϵ as the level of equilibrium approximation. We denote the level of equilibrium approximation in the k th market by ϵ^k , and we assume that $\epsilon^k \rightarrow 0$ and $\epsilon^k \sqrt{k} \rightarrow \infty$ as $k \rightarrow \infty$. We study the behavior of the equilibrium along the sequence of marketplaces we described above in order to derive the equilibrium in a marketplace with large number of agents.

Definition 4 (ϵ -Market Equilibrium) Let $(D_n^k, p_n^k, y_n^k)_{n=1}^N$ summarize the strategy of all players in the k th market with $y_n^k > 0$ for all $n = 1, \dots, N$. Then, $(D_n^k, p_n^k, y_n^k)_{n=1}^N$ is an ϵ -Market Equilibrium if the following conditions are satisfied:

1. $D_n^k = D_n^{MCE}(p_1^k, \dots, p_N^k; y_1^k, \dots, y_N^k)$ for all $n \leq N$.
2. For any $\ell \leq N$ and $m \leq N$, we have that

$$\begin{aligned} V_\ell(D_1^k, \dots, D_N^k; p_1^k, \dots, p_N^k; y_1^k, \dots, y_N^k) \\ \geq V_\ell(\hat{D}_1^k, \dots, \hat{D}_N^k; p_1^k, \dots, p_N^k; \hat{y}_1^k, \dots, \hat{y}_N^k) - \epsilon^k, \end{aligned}$$

where $\hat{y}_n^k = y_n^k - 1$ if $n = \ell$, $\hat{y}_n^k = y_n^k + 1$ if $n = m$, $\hat{y}_n^k = y_n^k$ otherwise, and $\hat{D}_n^k = D_n^{MCE}(p_1^k, \dots, p_N^k; \hat{y}_1^k, \dots, \hat{y}_N^k)$ for all $n \leq N$.

3. For any $\ell \leq N$ and $p' \neq p_n^k$ for all $n = 1, \dots, N$, we have that

$$\begin{aligned} V_\ell(D_1^k, \dots, D_N^k; p_1^k, \dots, p_N^k; y_1^k, \dots, y_N^k) \\ \geq V_{N+1}(\hat{D}_1^k, \dots, \hat{D}_{N+1}^k; p_1^k, \dots, p_N^k, p'; \hat{y}_1^k, \dots, \hat{y}_{N+1}^k) - \epsilon^k, \end{aligned}$$

where $\hat{y}_n^k = y_n^k - 1$ if $n = \ell$, $\hat{y}_n^k = 1$ if $n = N + 1$, $\hat{y}_n^k = y_n^k$ otherwise, and $\hat{D}_n^k = D_n^{MCE}(p_1^k, \dots, p_N^k, p'; \hat{y}_1^k, \dots, \hat{y}_{N+1}^k)$ for all $n \leq N + 1$.

The first condition in the above definition requires that the vector $(D_n^k)_{n=1}^N$ forms an equilibrium among the customers if the agents choose the strategy $(p_n^k, y_n^k)_{n=1}^N$. The second and third conditions characterize the equilibrium in the first stage game: The second condition states that an agent cannot improve his revenue by more than ϵ^k when he joins an existing sub-pool, while the third condition states that an agent cannot improve his revenue by more than ϵ^k when he introduces a new sub-pool. We next turn to characterize the equilibrium in the k th marketplace. Note that if $\epsilon^k \equiv 0$ for all k , then the above definition reduces to that of the Nash Equilibrium.

9.5.1 Characterization of the Market Equilibrium

In this subsection, we study the symmetric equilibrium for the sequence of marketplaces we constructed above. As a first step towards characterizing the symmetric equilibrium, we derive the revenues of agents when they announce the same price in the k th marketplace. As we noted before, such a marketplace operates as an $M/M/k + M$ system with arrival rate $\Lambda^k D_1^{MCE}(p^k; k)$, service rate 1, and abandonment rate $1/m_a$, where $D_1^{MCE}(p^k; k)$ is the Market Customer Equilibrium when all k agents charge p^k . Therefore, the revenue of an agent in this case is given by

$$V_1(D_1^{MCE}(p^k; k); p^k; k) = p\rho D_1^{MCE}(p^k; k)[1 - \beta^M(\Lambda^k D_1^{MCE}(p^k; k); k)], \quad (9.2)$$

where $\beta^M(\lambda; k)$ is probability of abandonment in $M/M/k + M$ system with arrival rate λ , service rate 1, and abandonment rate $1/m_a$.

In order to characterize an ϵ^k -symmetric Market Equilibrium, we need to verify that a single agent does not have any incentive to deviate to a price other than p^k in the k th marketplace. Recall that if an agent chooses $p' \neq p^k$, this amounts to creating his own sub-pool, and his revenue is given by $V_2(D_1^{MCE}(p^k, p'; k-1, 1), D_2^{MCE}(p^k, p'; k-1, 1); p^k, p'; k-1, 1)$, where $(D_n^{MCE}(p^k, p'; k-1, 1))_{n=1}^2$ is the Market Customer Equilibrium given that $k-1$ agents charge p^k and one agent charges p' . We then say that a price p^k emerges as the symmetric ϵ^k -Market Equilibrium if

$$\begin{aligned} & V_1(D_1^{MCE}(p^k; k), p^k, k) \\ & \geq \max_{0 \leq p' \leq R} V_2(D_1^{MCE}(p^k, p'; k-1, 1), D_2^{MCE}(p^k, p'; k-1, 1); p^k, p'; k-1, 1) \\ & \quad - \epsilon^k, \end{aligned} \tag{9.3}$$

where the left-hand side is the revenues of agents when all agents charge p^k , and the right-hand side is the maximum revenue that a single agent can obtain by deviating from p^k .

To understand the behavior of the market outcome in large markets, we shall first study the left-hand side of (9.3) along the trajectory of marketplaces in which all k agents charge p^k and $p^k \rightarrow p$ as $k \rightarrow \infty$. In a buyer's market, we show that all customers join the system in equilibrium as long as $p < R$ since they experience negligible waiting times and obtain approximately the utility of $R - p$ by joining in a marketplace with a large number of agents. Therefore, the revenue of each agent is approximated by $p\rho$ in a buyer's market when $p < R$. In a seller's market, some of the customers leave the market immediately due to the high congestion level even if $p < R$, but the rate of customers requesting service should, in equilibrium, be higher than the processing capacity when $p < R$. Therefore, agents are always "over-utilized" in a seller's market and the revenue of each agent is approximately p when $p < R$. When $p = R$, the rate of customers requesting service depends on the convergence rate of p^k both in a buyer's and a seller's market. Thus, $p \min\{\rho, 1\}$ constitutes an upper bound for the revenue of each agent if $p = R$. The following proposition presents these results formally.

Proposition 2 *Let $D_1^{MCE}(p^k; k)$ be the Market Customer Equilibrium when all agents charge p^k in the k th marketplace such that $\lim_{k \rightarrow \infty} p^k = p$. When $p < R$, we have that $\lim_{k \rightarrow \infty} D_1^{MCE}(p^k; k) = \min\{1, (R - p + cm_a)/(\rho cm_a)\}$ and*

$$\lim_{k \rightarrow \infty} V_1(D_1^{MCE}(p^k; k); p^k; k) = \begin{cases} p\rho & \text{if } \rho \leq 1 \\ p & \text{if } \rho > 1 \end{cases}.$$

When $p = R$, we have that $\limsup_{k \rightarrow \infty} D_1^{MCE}(p^k; k) \leq \min\{1, 1/\rho\}$, and

$$\lim_{k \rightarrow \infty} V_1(D_1^{MCE}(p^k; k); p^k; k) \leq \begin{cases} p\rho & \text{if } \rho \leq 1 \\ p & \text{if } \rho > 1 \end{cases}.$$

After approximating the revenue of the agents when they charge the same price, we now focus on the maximum revenue that an agent can obtain by creating his own sub-pool. As we did above, we again distinguish between buyer's and seller's markets.

9.5.1.1 Buyer's Market

When all agents charge the same price p^k in a buyer's market, we next show that a single agent can improve his revenue when he decreases his price. Such a cut will allow a single agent to serve not only his own customers but also the customers choosing the price p^k . In fact, his revenue can be arbitrarily close to p^k following a small price cut as long as the rate of customers requesting service is bounded away from zero when all agents charge p^k , i.e., $\lim_{k \rightarrow \infty} D_1^{MCE}(p^k; k) > 0$. The following proposition proves this observation formally.

Proposition 3 *Let*

$$V'(p^k; k) = \max_{0 \leq p' < p^k} V_2 \left(D_1^{MCE}(p^k, p'; k-1, 1), \right. \\ \left. D_2^{MCE}(p^k, p'; k-1, 1); p^k, p'; k-1, 1 \right)$$

for any sequence of p^k with $\lim_{k \rightarrow \infty} p^k = p$. Then, we have that $\liminf_{k \rightarrow \infty} V'_k(p^k; k) > 0$ when $p > 0$. Furthermore, when $\lim_{k \rightarrow \infty} D_1^{MCE}(p^k; k) > 0$, we have the relation $\lim_{k \rightarrow \infty} V'(p^k; k) = p$.

As we established in Proposition 2, the revenue of an agent when all agents charge the same price p^k can be bounded from above by $p^k \rho$ in large marketplaces. Then, Proposition 3 implies that any p^k satisfying $\lim_{k \rightarrow \infty} p^k = p > \epsilon^k / (1 - \rho)$ cannot emerge as the equilibrium price of a symmetric ϵ^k -Market Equilibrium for large k . Thus, as $\lim_{k \rightarrow \infty} \epsilon^k = 0$, we obtain that any sequence of prices except the ones converging to zero cannot be sustained as the equilibrium price of a symmetric ϵ^k -Market Equilibrium along the trajectory of marketplaces. Note that we do not need to analyze the revenue of an agent after a price increase because it is sufficient to demonstrate the existence of one profitable deviation in order to show that a given price cannot be an equilibrium outcome. We formalize these observations in the following theorem.

Theorem 2 *In a buyer's market with $\rho < 1$,*

1. *Let p_{EQ}^k be a price emerging as the equilibrium price of a symmetric ϵ^k -Market Equilibrium in the k th marketplace. Then, for any $\xi > 0$, there exists a K such that $p_{EQ}^k < \xi$ for all $k > K$.*
2. *There exists a K such that zero is an equilibrium price of a symmetric ϵ^k -Market Equilibrium in the k th marketplace for all $k > K$.*

3. Let Π_{OE}^k and Π_{NI}^k be the total revenue in the k th marketplace with and without operational efficiency, respectively. Then, for any $\xi > 0$, there exists a K such that $\Pi_{OE}^k/\Pi_{NI}^k < \xi$ for all $k > K$.

The above theorem states that if a moderating firm provides efficient matching in a buyer's market, the equilibrium outcome of the marketplace will converge to zero. As the profit of the firm is the share of the revenue generated in the marketplace, providing efficient matching deteriorates the profit of the firm compared to the no-intervention case as well as the revenue of the agents. In fact, we show that the ratio between the total revenue generated in a marketplace under operational efficiency and under the no-intervention converges to zero. We also establish that zero can emerge as the equilibrium price in large marketplaces. In Sect. 9.7, we discuss the extension of the above theorem, which is based on showing that the revenues of agents converges to zero even in a non-symmetric equilibrium.

9.5.1.2 Seller's Market

After discussing the impact of providing efficient matching in a buyer's market, we now focus on a seller's market. Unlike in a buyer's market, a single agent cannot improve his revenue after a price cut since it does not improve his utilization significantly. Note that agents are already "over-utilized," and earning a revenue of p^k while they are charging the same price p^k in a seller's market. Therefore, in a seller's market, the only possible profitable deviation for a single agent is to increase his price in large enough marketplaces. In such a deviation, a single agent loses some of his customers because of his high price, and he also loses the benefits of efficient matching since he becomes an individual provider. Both of these factors will limit his ability to make higher profit. In fact, the following proposition establishes an upper bound on the asymptotic revenue which a single agent can generate by increasing his price.

Proposition 4 *Let*

$$V'(p^k; k) = \max_{p^k \leq p' \leq R} V_1 \left(D_1^{MCE}(p', p^k; 1, k-1), \right. \\ \left. D_2^{MCE}(p', p^k; 1, k-1); p', p^k; 1, k-1 \right)$$

for any given sequence of prices p^k such that $\lim_{k \rightarrow \infty} p^k = p$. When $p < R$ in a seller's market ($\rho > 1$), we have that

$$\limsup_{k \rightarrow \infty} V'(p^k; k) \leq (R + cm_a) \lambda^\Delta(p; R) [1 - \beta(\lambda^\Delta(p; R))] \\ - \lambda^\Delta(p; R) (\Delta(p; R) + cm_a),$$

where $\Delta(p; R) = \max\{0, (R - p + cm_a)/\rho - cm_a\}$, and $\lambda^\Delta(p; R)$ is the unique solution to $1 - \beta(\lambda) - \lambda\beta'(\lambda) = (\Delta(p; R) + cm_a)/(R + cm_a)$.

When a single provider increases his price, we show that the demand for agents, who do not change their prices, is almost the same as their original demand before deviation. Hence, the utility of customers choosing the sub-pool consisting of $k - 1$ agents is $\Delta(p; R)$, which is the utility that the customers obtain in the Market Customer Equilibrium in a large marketplace when all agents charge p^k . Then, to approximate the maximum post-deviation revenue, one can treat the deviating agent as a monopoly whose customers have an outside option with the value of $\Delta(p; R)$. In fact, the above proposition shows that this approximation constitutes an upper bound on the agent's post-deviation revenue. A monopoly always makes sure that the utility of customers is exactly equal to their outside option, by setting the price to $R + cm_a - (\Delta(p; R) + cm_a)/(1 - \beta(\lambda))$ for any given target of demand rate λ . He then picks λ , maximizing his revenue and sets his price accordingly. We refer the reader to the proof of Proposition 4 for a more detailed discussion on the revenue maximization problem of a monopoly.

Combining the two observations above, it is clear that in a large marketplace, a price p^k emerges as the symmetric ϵ^k -Market Equilibrium outcome if p^k is greater than the profit of a monopoly serving customers with outside option $\Delta(p; R)$. We state this result in the following theorem.

Theorem 3 *In a seller's market ($\rho > 1$), let*

$$p^* \in \mathcal{P}(\rho; R) \equiv \left\{ p : p > (R + cm_a)\lambda^\Delta(p; R)[1 - \beta(\lambda^\Delta(p; R))] \right. \\ \left. - \lambda^\Delta(p; R)(\Delta(p; R) + cm_a), 0 \leq p < R \right\},$$

where $\Delta(p; R)$, and $\lambda^\Delta(p; R)$ are defined as in Proposition 4. Then, for any given sequence of prices p^{*k} that converges to p^* as $k \rightarrow \infty$, there exists a K such that p^{*k} emerges as the equilibrium price of a symmetric ϵ^k -Market Equilibrium in the k th marketplace for all $k > K$. Furthermore, for any $\rho_1 > \rho_2$, we have that $\mathcal{P}(\rho_1; R) \subseteq \mathcal{P}(\rho_2; R)$.

The above theorem characterizes the set of symmetric ϵ^k -Market Equilibria for large marketplaces. The theorem does not guarantee the uniqueness of such an equilibrium, i.e. $\mathcal{P}(\rho; R)$ may not be a singleton. In fact, $\mathcal{P}(\rho; R)$ may consist of uncountably many prices. Furthermore, we show that $\mathcal{P}(\rho; R)$ shrinks as ρ increases. As the demand-supply ratio increases, customers experience significant waiting times even if they are served by a price-generated pool. Therefore, the level of customer surplus that a deviating agent has to forego declines as ρ rises. As a result of this, a single agent has more room to deviate and improve his revenue when demand is high. It is also worth highlighting that a single agent has such a profitable deviation opportunity even though the number of agents grows to infinity.

Characterizing the set of symmetric equilibria, $\mathcal{P}(\rho; R)$, is difficult in general. For illustrative purposes, we consider the case where the abandonment rate is equal to the service rate. We show that a similar structure holds for the settings when $\mu \neq m_a$ using a numerical study (see Allon et al. 2012). The next corollary

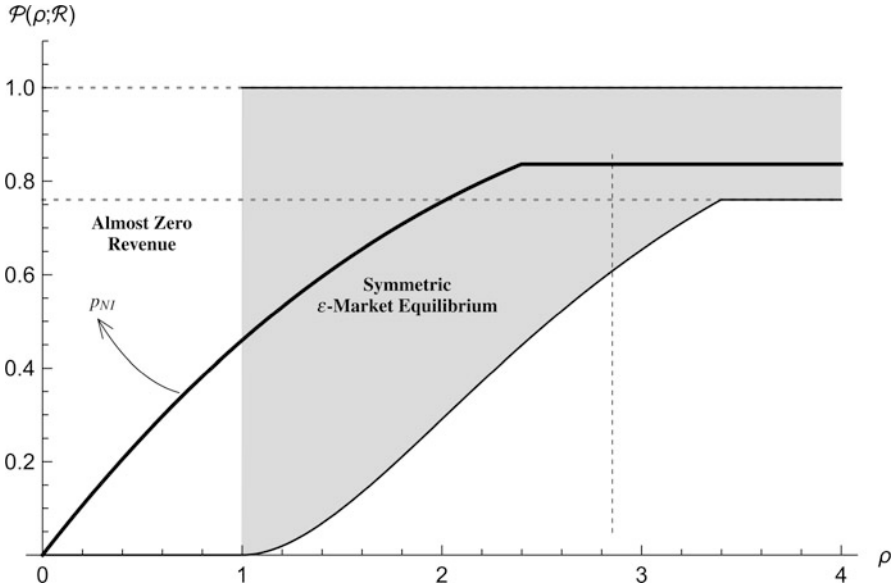


Fig. 9.1 The prices that form a symmetric market equilibrium as a function of the demand-supply mismatch (ρ). The service rates and abandonment rates are assumed to be one

characterizes the correspondence $\mathcal{P}(\rho; R)$ as well as the asymptotic behavior of the unique equilibrium price under the no-intervention model.

Corollary 1 *Suppose the abandonment rate is equal to the service rate. Then, we have that*

1. $\lambda^{\Delta}(p; R) = \log[(R + c)/(\Delta(p; R) + c)]$ where $\Delta(p; R)$ is defined as in Proposition 4. Furthermore, the correspondence $\mathcal{P}(\rho; R)$ defined in Theorem 3 can be expressed as

$$\mathcal{P}(\rho; R) = \left\{ p : p > R + c - \left(1 + \log \left(\frac{R + c}{\Delta(p; R) + c} \right) \right) [\Delta(p; R) + c], 0 \leq p < R \right\}.$$

2. $\lim_{k \rightarrow \infty} p_{NI}^k = p_{NI} \equiv (R + c) \min\{1 - \rho/(e^{\rho} - 1), 1 - (c/R) \log((R + c)/c)\}$, where p_{NI}^k is the unique equilibrium price under no-intervention setting in the k th marketplace.

Figure 9.1 displays the correspondence $\mathcal{P}(\rho; R)$ and the limit p_{NI} . More specifically, the gray area represents the prices that can emerge as the equilibrium price of a symmetric equilibrium in a large marketplace and the bold curve depicts p_{NI} . We observe that for all $\rho > 1$, the set $\mathcal{P}(\rho; R)$ is not a singleton. In fact, we have a wide range of prices that can form an equilibrium. Furthermore, many of the possible equilibrium prices in $\mathcal{P}(\rho; R)$ are lower than p_{NI} . The intuition

behind this result is the following: In a marketplace where the moderating firm efficiently matches customers and agents, a single agent, who deviates by increasing his price, loses benefits of efficient matching, and thus cannot sustain the same quality of service (in terms of waiting times) as his “original” pool. It turns out that the deviating agent cannot improve his “original” revenue by decreasing his price either. Thus, in a seller’s market, the price-generated pool serves as a deterrent against single agent deviations even if prices are unappealing from a system point of view. It is also important to note that such lower prices lead to loss in total revenue for the marketplace compared to the no-intervention setting. While one may expect operational efficiency tools to be a leverage for higher revenues in the market, it is surprising to see that reducing the unnecessary waiting and idleness present in a system with no-intervention may deteriorate the revenues.

Myerson (1991) argues that the question of which equilibrium would emerge as the outcome of a game with multiple equilibria can be answered with the focal-point effect phenomenon.⁶ Our goal in this chapter is not to conclude that only the low prices can be a focal-point. In fact, when comparing the equilibrium prices in a market with and without operational efficiency, one should also observe that operational efficiency does not only serve as a deterrent for deviations from low prices but also prevents deviations from high prices for any level of demand-supply ratio. Moreover, when the aggregate demand is sufficiently high, efficient matching always leads to higher profits, although the equilibrium prices under operational efficiency may be slightly lower than the unique equilibrium in a market without operational efficiency.

Our analysis in this section assumes that a customer with a preferred price p_ℓ pays p_m when she is served by sub-pool- m for any $m \neq \ell$. In real marketplaces such as UpWork.com, customers may end up paying a price between their preferred price and the prices asked by the providers when these prices are different. To account for that, one may envision an extension of our model, in which a customer choosing sub-pool- ℓ pays $\phi p_\ell + (1 - \phi)p_m$, where $\phi \in (0, 1)$, when an agent from sub-pool- m , with $m \neq \ell$ serves her. In such a model, our key findings, which are equilibrium prices are close to zero in a buyer’s market, and some of the equilibrium outcomes may lead to profit loss in a seller’s market, would continue to hold.

In this section, we study a specific mechanism that the moderating firm uses to achieve operational efficiency. There are other mechanisms, such as providing real-time congestion information, that may be used by the moderating firm. When customers are able to obtain real-time congestion information, Allon et al. (2012) shows that our analytical results for a buyers market continue to hold, while a simulation experiment demonstrates that multiple, possibly harmful equilibria also exist in a sellers market.

⁶Focal-point effects are any psychological or cultural norms that tends to focus players’ attention on one equilibrium.

9.6 Communication Enabled Model

In this section, we continue to study the impact of different mechanisms used by the moderating firm. As we mentioned in the introduction, the moderating firm may complement its operational tool discussed in the previous section with a strategic tool, which changes the nature of the interaction among agents. In a marketplace such as UpWork.com, service providers are offered discussion boards in which they are allowed to exchange information. Moreover, the market supports the creation of affiliation groups, which are self-enforcing entities. We will thus focus on the impact of enabling communication among agents on the market outcome.

The economics literature suggests that, when the players have the opportunity to perform non-binding pre-play communication among themselves, the stability of an outcome can be threatened by potential deviations formed by coalitions, even in noncooperative games. Following this idea, the well-know notion of Strong Nash Equilibrium (SNE) requires stability against deviations formed by any conceivable coalitions (see Aumann 1959). The main drawback of SNE is that many of the games do not have any SNE.

In this section, we modify the marketplace we study in the previous section by assuming that agents have opportunities to make non-binding communication prior to making their decisions, so that they can try to self-coordinate their actions in a mutually beneficial way, despite the fact that each agent selfishly maximizes his own utility.

Echoing the ideas in the economics literature, allowing communication among agents changes the equilibrium concept we use to characterize the outcome in the marketplace. We model this by proposing a new equilibrium concept that allows several agents to deviate together. More specifically, the new concept requires that a strategy of agents should be immune to any coalitions. Since a marketplace tends to be large, e.g., there are hundreds of thousands of agents in UpWork.com, one has to restrict the possible size of a coalition. We denote the largest fraction of agents that is allowed to deviate together by $\delta \in (1/k, 1]$. As in Sect. 9.5, we focus on the deviations that improve the revenues of agents at least by $\epsilon \geq 0$. Furthermore, we again study the behavior of the equilibrium along the sequence of marketplaces we described in Sect. 9.5. Recall that there are k agents, the arrival rate is $\Lambda^k = \rho k$, and the level of equilibrium approximation is ϵ^k , with the same asymptotic properties as in Sect. 9.5, in the k th marketplace. We let δ^k be the largest fraction of agents that is allowed to deviate together in the k th marketplace. We assume that $\delta^k k \rightarrow \infty$ as $k \rightarrow \infty$. This condition states that the number of agents allowed to deviate increases without bound as the market size increases. We refer to our new equilibrium concept as (δ, ϵ) -Market Equilibrium which is defined as follows:

Definition 5 ((δ, ϵ)-Market Equilibrium) Let $(D_n^k, p_n^k, y_n^k)_{n=1}^N$ summarize the strategy of all players in the k th market with $y_n^k > 0$ for all $n = 1, \dots, N$. Then, $(D_n^k, p_n^k, y_n^k)_{n=1}^N$ is a (δ^k, ϵ^k) -Market Equilibrium if the following conditions are satisfied:

1. $D_n^k = D_n^{MCE}(p_1^k, \dots, p_N^k; y_1^k, \dots, y_N^k)$ for all $n \leq N$.

2. For any $\ell \leq N$, $m \leq N$, and $0 < d \leq \min\{y_\ell^k, \lfloor \delta^k k \rfloor\}$, we have that

$$\begin{aligned} & V_\ell(D_1^k, \dots, D_N^k; p_1^k, \dots, p_N^k; y_1^k, \dots, y_N^k) \\ & \geq V_\ell(\hat{D}_1^k, \dots, \hat{D}_N^k; p_1^k, \dots, p_N^k; \hat{y}_1^k, \dots, \hat{y}_N^k) - \epsilon^k, \end{aligned}$$

where $\hat{y}_n^k = y_n^k - d$ if $n = \ell$, $\hat{y}_n^k = y_n^k + d$ if $n = m$, $\hat{y}_n^k = y_n^k$ otherwise, and $\hat{D}_n^k = D_n^{MCE}(p_1^k, \dots, p_N^k; \hat{y}_1^k, \dots, \hat{y}_N^k)$ for all $n \leq N$.

3. For any $\ell \leq N$, $0 < d \leq \min\{y_\ell^k, \lfloor \delta^k k \rfloor\}$, and $p' \neq p_n$ for all $n = 1, \dots, N$, we have that

$$\begin{aligned} & V_\ell(D_1^k, \dots, D_N^k; p_1^k, \dots, p_N^k; y_1^k, \dots, y_N^k) \\ & \geq V_{N+1}(\hat{D}_1^k, \dots, \hat{D}_{N+1}^k; p_1^k, \dots, p_N^k, p'; \hat{y}_1^k, \dots, \hat{y}_{N+1}^k) - \epsilon^k, \end{aligned}$$

where $\hat{y}_n^k = y_n^k - d$ if $n = \ell$, $\hat{y}_n^k = d$ if $n = N + 1$, $\hat{y}_n^k = y_n^k$ otherwise, and $\hat{D}_n^k = D_n^{MCE}(p_1^k, \dots, p_N^k, p'; \hat{y}_1^k, \dots, \hat{y}_{N+1}^k)$ for all $n \leq N + 1$.

The above definition is closely related to the definition of ϵ -Market Equilibrium in Sect. 9.5. The key difference between these two equilibrium definitions is that (δ, ϵ) -Market Equilibrium allows a group of agents to deviate by either forming a new sub-pool or joining an existing one. In fact, our new equilibrium concept is a refinement of the ϵ -Market Equilibrium. Therefore, any (δ, ϵ) -Market Equilibrium is also a ϵ -Market Equilibrium. Employing the (δ, ϵ) -Market Equilibrium concept, we expect that the set of prices that can be sustained as a ϵ -Market Equilibrium will shrink since (δ, ϵ) -Market Equilibrium is more restrictive. Kalai (2004) and Gradwohl and Reingold (2008) study large games and shows that all Nash Equilibria of certain large games are resilient to deviations by coalitions. Such a phenomena does not exist in our model.⁷

9.6.1 Characterization of the (δ, ϵ) -Market Equilibrium

Similar to Sect. 9.5, we focus on the symmetric (δ, ϵ) -ME where all agents charge the same price. The revenue of an agent when all agents charge the same price p^k is the same as in (9.2), and thus Proposition 2 establishes its asymptotic behavior.

In a buyer's market with $\rho < 1$, we showed that only the prices in a small neighborhood of zero can emerge as a symmetric ϵ -Market Equilibrium in large

⁷According to the definition in Gradwohl and Reingold (2008), a Nash Equilibrium is resilient to coalitions if players cannot improve their revenues “too much” even after a coordinated deviation. In our setting, “too much” has to be almost as much as the customer reward, R , in order to apply their results to our game. Clearly, this makes the definition of resilience vacuous because none of the agents can increase his revenue by more than R .

marketplaces. As a direct implication of the fact that (δ, ϵ) -Market Equilibrium is a refinement of the ϵ -Market Equilibrium, any sequence of prices that emerge as symmetric (δ, ϵ) -Market Equilibrium converges to zero as the market size grows. Furthermore, we show that $p = 0$ can emerge as the equilibrium price in large marketplaces.

Theorem 4 *Let p_{EQ}^k be a price emerging as a symmetric (δ^k, ϵ^k) -Market Equilibrium in the k th marketplace where $\rho < 1$. Then, for any $\xi > 0$, there exists a K such that $p_{EQ}^k < \xi$ for all $k > K$. Furthermore, when $\lim_{k \rightarrow \infty} \delta^k = 0$, there exists a K such that zero is an equilibrium price of a symmetric (δ^k, ϵ^k) -Market Equilibrium in the k th marketplace for all $k > K$.*

In a seller's market, Proposition 2 shows that the rate of customers requesting service will exceed the processing capacity of agents when all agents charge a price lower than R . Therefore, customers experience significant waiting times, and not only pay the price of the service but also incur a strictly positive waiting cost. Then, we show that a small group of agents can use the fact that customers pay an extra cost to increase their prices while ensuring that they are still "over-utilized" after the price increase. Since this small group of agents increases their prices without hurting their utilization, this deviation clearly improves their revenues (This is in contrast to the setting in Sect. 9.5 where the utilization of a single agent does drop after a price decrease). Thus, in a seller's market, only the prices, which are very close to R , can emerge as the equilibrium price of a symmetric (δ, ϵ) -Market Equilibrium in large marketplaces. To contrast this result with the result in Theorem 3, it is worth noting that a single agent has only a limited opportunity to improve his revenue by increasing his price as in most cases, the revenue improvement due to the price increase is overcome by the drop in utilization. Therefore, without the communication opportunity, it was possible to observe low prices as the market outcome even though demand exceeds supply.

Theorem 5 *Let p_{EQ}^k be a price emerging as a symmetric (δ^k, ϵ^k) -Market Equilibrium in the k th marketplace where $\rho > 1$. Then, for any $\xi > 0$, there exists a K such that $p_{EQ}^k > R - \xi$ and $D_1^{MCE}(p_{EQ}^k; k) > 1/\rho - \xi$ for all $k > K$. Furthermore, there exist a sequence p^{*k} and a K such that p^{*k} forms a symmetric (δ^k, ϵ^k) -Market Equilibrium in the k th marketplace, for all $k > K$.*

The above result shows that agents can sustain a price, which extracts all of the customer surplus, as the equilibrium outcome in a seller's market. Moreover, it also implies that the marketplace cannot be congested in the equilibrium even in a seller's market since any level of congestion can be capitalized by agents through a price increase.

Theorem 5 characterizes the unique limit of symmetric (δ, ϵ) -Market Equilibrium, but this result can be extended by showing that R is indeed the unique limit of all possible (δ, ϵ) -Market Equilibria as discussed in Sect. 9.7. Furthermore, Allon et al. (2012) shows that the ability to communicate leads to high equilibrium prices

when the moderating firm provides real-time queue information in order to reduce the mismatch between customers and agents as long as the largest fraction of agents that is allowed to deviate together is close to 1.

9.7 A Marketplace with Non-identical Agents

In Sect. 9.3, we introduce a model where all of the agents in the marketplace are a priori identical. However, it is natural to imagine that large service marketplaces attract service providers with different skill sets, which provide their customers different values for the service. Thus, we explore the robustness of the conclusions of the previous sections to the heterogeneity among providers.

To this end, we consider a marketplace where agents provide the same service but in different quality levels, say low and high. We assume that customers value the service with respect to its quality. Particularly, customers earn a reward of R_H and R_L when they are served by a high-quality and a low-quality agent, respectively. Without loss of generality, we assume $R_L \leq R_H$. The model set-up is the same as in Sect. 9.3, and we use a similar mode of analysis as in Sects. 9.4, 9.5, and 9.6.

The behavior of the marketplace when the moderating firm confines itself to setting up the necessary infrastructure is very similar to the equilibrium in Theorem 1: Agents may behave as local monopolists when the arrival rate is sufficiently high. Furthermore, once the arrival rate is less than a certain threshold, customers observe lower prices, which allow them to earn strictly positive utility, due to the intensified competition. However, we also encounter new results when we allow for heterogeneous agents. First, unlike the identical agent model, we observe that the main driver of equilibrium outcomes for certain parameters is not only the competition between providers but also the fact that agents offer different quality of service. For instance, when the demand rate is in a certain range, high-quality agents charge a low price and forego a significant customer surplus both because of the low demand and the fact that they want to keep the low-quality agents out of the marketplace. We also show that it is possible to have a continuum of symmetric equilibria, whereas we always have a unique symmetric equilibrium with the identical agents.

The impact of improving the operational efficiency in a marketplace with non-identical agents is also similar to our findings in the identical agents model: When demand is sufficiently low in a buyer's market, the revenues of agents are always in a small neighborhood of zero in large marketplaces. In a seller's market, there are multiple equilibria, which may lead to profit loss for the firm compared to the no-intervention model. Unlike the identical agents model, we show that there may be multiple equilibria even in a buyer's market as long as demand exceeds the total capacity of high-quality agents. However, most of these equilibrium prices may be very low compare to the equilibrium outcome in the no-intervention model. Thus, providing tools to improve the operational efficiency may still deteriorate the moderating firm's profit.

Finally, we explore the impact of enabling communication among agents in a market with non-identical agents. As in Sect. 9.6, we establish that pre-play communication helps agents sustain the profit maximizing one among the multiple equilibria arising due to providing operational efficiency.

Our results in the non-identical agents model also provide insights about the non-symmetric equilibrium outcomes in the identical agents model. In particular, our model with non-identical agents helps us to prove that the non-symmetric equilibrium may exist only for a small range of demand- supply ratio ρ in the no-intervention model with identical agents, and this range becomes negligible as the number of agents grow. Furthermore, using our results in this section, we show that in the operational efficiency model with identical agents, the revenues of all agents in any non-symmetric equilibrium (if exists) should be in a small neighborhood of zero in a buyer's market. We also show that, in the communication model with identical agents, even if there are any non-symmetric equilibria in a seller's market, the revenue of each agent in equilibrium should converge to R as well as the price they charge.

We refer the reader to Allon et al. (2012) for a detailed discussion of our findings in this section.

9.8 Conclusion

In this chapter, we study a marketplace in which many small service providers compete with each other in providing service to self-interested customers looking for temporary help. The main focus of the paper is on the role of the moderating firm, which sets up the marketplace and creates the infrastructure where agents and customers interact. To this end, we explore the impact of different strategies employed by the moderating firm by considering three market models.

We characterize the market outcomes in each of these models. We observe that outcomes critically depend on the moderating firm's involvement and market conditions, i.e., whether it is a buyer's or a seller's market. Since different types of involvement of the moderating firm result in different equilibrium prices and customer demand, the moderating firm aims to intervene in the marketplace in order to make sure that the "right" prices and customer demand emerge in equilibrium. Specifically, the moderating firm tries to maximize the revenues of agents since its profit is a share of the agents' revenues.

We show that when the firm ensures efficient operational matching and enables agent communication in a seller's market, the natural upper-bound on the revenue generated in a marketplace⁸ is asymptotically achievable, and thus, using these two

⁸In a given marketplace, the total revenues of the agents cannot exceed $\min\{\Lambda, k\}R$ since they cannot charge more than R , and their effective demand is the minimum of their processing capacity and the aggregate demand.

tools together dominates any other strategy from the moderating firm's perspective in a seller's market. We also show that efficient operational matching in a buyer's market leads to arbitrarily small total marketplace revenue compared to the total revenue under the no-intervention model. Hence, using the matching mechanism we discuss in this chapter is not advisable in a buyer's market despite the fact that it reduces the mismatch between demand and supply. This result is somewhat counter-intuitive, because the efficiency improvement due to better matching is not necessarily translated into additional profits. It seems other tools aimed at improving the operational efficiency, such as providing real-time queue information, will have a similar impact on the moderating firm's profit in a buyer's market.

Both UpWork.com and ServiceLive.com are currently in their growth stage and have not achieved their full potential in terms of demand for their services. However, both firms can and should project the "mature" market conditions and decide on their appropriate measures to adopt. Given the moderate level of congestion in UpWork.com, one may infer that the marketplace can be identified as a seller's market. Following the discussion before, UpWork.com's decision to offer operational tools complemented with strategic tools is well justified.

In this chapter, we focus on the operational and strategic tools that the firm can use to be involved in the marketplace. There are also other possible ways for a moderating firm to intervene in the marketplace including introducing a skill screening mechanism. In general, these mechanisms take the form of skill tests and/or certification programs that are run by moderating firms. For instance, UpWork.com offers various exams to test the ability of the candidate providers. If necessary, UpWork.com can use these exams to disqualify some of the agents, and thus control the portfolio of different agent types (e.g., flexible, dedicated) and the service capacity in the marketplace. It is in the best interest of the moderating firm to use its skill tests in order to make sure that the "right" prices and customer demand emerge in the marketplace. To gain insights about the effectiveness of skill screening as a revenue maximization tool, in Allon et al. (2017), we analyze how much benefit the firm obtains after each additional skill test. Our findings in this paper suggest that the firm does not need to regulate the marketplace via skill screening when agents are endowed with highly compatible skills. As the compatibility of agent skills weakens, we show that the firm starts to experience substantial revenue benefits from skill screening. We also show that the skill screening becomes more effective as the different classes of customers start to vary in terms of their processing time needs. When intervention is needed, we establish that the marginal benefit from skill testing may decline sharply. In fact, under certain demand structures, firm does not gain any revenue benefits from an additional skill test.

It is also possible that a moderating firm can be involved in the marketplace via contracting with agents or providing a suggested price. Particularly, the setting in which the firm provides a suggested price can be viewed as pre-play communication and will indeed shrink the set of equilibria. However, these type of interactions between the moderating firm and agents are outside the scope of this chapter as these settings are not a market per-se anymore. In such environments, the firm would decide on prices as well as the allocation of agents to customers.

While modeling operational efficiency, we assume that agents give priority to their own customers. One may consider an extension of our model in which agents are allowed to choose both priority and prices, simultaneously. The equilibria that arise in our model with fixed priority rule would still be sustained in such an extended game. Hence, the main spirit of our findings, namely, the fact that providing operational efficiency may lead to profit loss, would not change. Additional equilibria would be possible in the extended model only when demand exceeds supply.

References

- Allon G, Federgruen A (2007) Competition in service industries. *Oper Res* 55:37–55
- Allon G, Gurvich I (2010) Pricing and dimensioning competing large-scale service providers. *Manuf Serv Oper Manag* 12:449–469
- Allon G, Bassamboo A, Çil E (2012) Large-scale service marketplaces: the role of the moderating firm. *Manag Sci* 58:1854–1872
- Allon G, Bassamboo A, Çil E (2017) Skill management in large-scale service marketplaces. *Prod Oper Manag* 26(11):2050–2070
- Amihud Y, Mendelson H (1980) Dealership market: market-making with inventory. *J Financ Econ* 8:31–53
- Armony M, Maglaras C (2004) On customer contact centers with a call-back option: customer decisions, routing rules and system design. *Oper Res* 52:271–292
- Armony M, Plambeck E, Seshadri S (2009) Sensitivity of optimal capacity to customer impatience in an unobservable M/M/S queue (Why you shouldn't shout at the DMV). *Manuf Serv Oper Manag* 11:19–32
- Aumann RJ (1959) Acceptable points in general cooperative n -person games. In: Kuhn HW, Tucker AW (eds) *Contributions to the theory of games IV*. Princeton University Press, Princeton
- Biais B, Glosten L, Spatt C (2005) Market microstructure: a survey of microfoundations, empirical results, and policy implications. *J Financ Markets* 8:217–264
- Burdett K, Mortensen D (1998) Wage differentials, employer size, and unemployment. *Int Econ Rev* 39:257–273
- Cachon G, Harker PT (2002) Competition and outsourcing with scale economies. *Manag Sci* 48:1314–1333
- Cachon G, Netessine S (2004) Game theory in supply chain analysis. In: Simchi-Levi D, Wu SD, Shen ZJ(M) (eds) *Handbook of quantitative supply chain analysis: modeling in the E-business era*. Springer, Boston
- Cachon G, Zhang F (2007) Obtaining fast service in a queuing system via performance-based allocation of demand. *Manag Sci* 53:408–420
- Chen H, Wan Y (2003) Price competition of make-to-order firms. *IIE Trans* 35:817–832
- Chen Y, Maglaras C, Vulcano G (2008) Design of an aggregated marketplace under congestion effects: asymptotic analysis and equilibrium characterization. Working paper
- Dixon H (1987) Approximate Bertrand equilibria in a replicated industry. *Rev Econ Stud* 54(1):47–62
- Garman M (1976) Market microstructure. *J Financ Econ* 3:257–275
- Garnett O, Mandelbaum A, Reiman M (2002) Designing a call center with impatient customers. *Manuf Serv Oper Manag* 4:208–227
- Gradwohl R, Reingold O (2008) Fault tolerance in large games. In: *Proceedings of the 9th ACM conference on electronic commerce*. ACM, New York, pp 274–283

- Granot D, Sošić G (2005) Formation of alliance in Internet-based supply exchanges. *Manag Sci* 51:92–105
- Ha AY, Li L, Ng S (2003) Price and delivery logistics competition in a supply chain. *Manag Sci* 49:1139–1153
- Halfin S, Whitt W (1981) Heavy-traffic limits for queues with many exponential servers. *Oper Res* 29:567–588
- Hassin R, Haviv M (2003) *To queue or not to queue: equilibrium behavior in queueing systems*. Kluwer Academic Publishers, Boston
- Ho T, Stoll H (1983) The dynamics of dealer markets under competition. *J Finance* 38:1053–1074
- Kalai E (2004) Large robust games. *Econometrica* 72:1631–1665
- Maglaras C, Zeevi A (2005) Pricing and design of differentiated services: approximate analysis and structural insights. *Oper Res* 53:242–262
- Manning A (2003) The real thin theory: monopsony in modern labour markets. *Labour Econ* 10:105–131
- Manning A (2004) Monopsony and the efficiency of labour market interventions. *Labour Econ* 11:145–163
- Michaelides M (2010) Labour market oligopsonistic competition: the effect of worker immobility on wages. *Labour Econ* 17:230–239
- Myerson RB (1991) *Game theory: analysis of conflict*. Harvard University Press, Cambridge
- Nagarajan M, Sošić G (2007) Stable farsighted coalitions in competitive markets. *Manag Sci* 53:29–45
- Nagarajan M, Sošić G (2008) Game-theoretic analysis of cooperation among supply chain agents: review and extensions. *Eur J Oper Res* 187:719–745
- Naor P (1969) The regulation of queue size by levying tolls. *Econometrica* 37:15–24
- Rath KP (1992) A direct proof of the existence of pure strategy equilibria in games with a continuum of players. *Econ Theory* 2:427–433
- Ray I (1996) Coalition-proof correlated equilibrium: a definition. *Games Econ Behav* 17:56–79
- Roughgarden T (2005) *Selfish routing and the price of anarchy*. The MIT Press, Cambridge/London
- Sošić G (2006) Transshipment of inventories among retailers: Myopic vs. farsighted stability. *Manag Sci* 52:1491–1508
- Upwork Press Release (2016) New study finds freelance economy grew to 55 million Americans this year, 35% of total U.S. workforce <https://www.upwork.com/press/2016/10/06/freelancing-in-america-2016/> (Oct 6)
- Ward AR, Glynn PW (2003) Diffusion approximation for a Markovian queue with reneging. *Queueing Syst* 43:103–128
- Zeltyn S, Mandelbaum A (2005) Call centers with impatient customers: many-server asymptotics of the $M/M/n + G$ queues. *Queueing Syst* 51:361–402

Chapter 10

Inducing Exploration in Service Platforms



Kostas Bimpikis and Yiangos Papanastasiou

Abstract Crowd-sourced content in the form of online product reviews or recommendations is an integral feature of most Internet-based service platforms and marketplaces, including Yelp, TripAdvisor, Netflix, and Amazon. Customers may find such information useful when deciding between potential alternatives; at the same time, the process of generating such content is mainly driven by the customers' decisions themselves. In other words, the service platform or marketplace “explores” the set of available options through its customers' decisions, while they “exploit” the information they obtain from the platform about past experiences to determine whether and what to purchase. Unlike the extensive work on the trade-off between exploration and exploitation in the context of multi-armed bandits, the canonical framework we discuss in this chapter involves a principal that explores a set of options through the actions of self-interested agents. In this framework, the incentives of the principal and the agents towards exploration are misaligned, but the former can potentially incentivize the actions of the latter by appropriately designing a payment scheme or an information provision policy.

10.1 Introduction

An important function of most Internet-based platforms that act as intermediaries between customers and service providers is the provision of information regarding the quality of the potential alternatives faced by the consumers. As the service platform landscape continues to evolve, the dominant form of generating such information is through *crowdsourcing*: after transacting with a service provider, a

K. Bimpikis (✉)
Graduate School of Business, Stanford University, Stanford, CA, USA
e-mail: kostasb@stanford.edu

Y. Papanastasiou
Haas School of Business, University of California, Berkeley, CA, USA
e-mail: yiangos@haas.berkeley.edu

customer may provide feedback on the provider's performance; this feedback is recorded by the platform and may become available to subsequent customers and assist them with their decision-making.

While soliciting feedback from customers is both straightforward and cost-effective, the crowdsourcing process through which information about the quality of the providers is generated is inherently inefficient from a system perspective, since it relies on the customers' self-interested choices. For an illustration of this inefficiency, consider the following example: a customer arrives at the platform and is presented with a choice between two providers, *A* and *B*. Provider *A* has eight "good" reviews and two "bad" reviews; Provider *B* has one of each. Given the available information (we assume that the customer is risk-neutral), provider *A* appears to be the better option; thus, the customer chooses *A*, and subsequently provides feedback on her choice. In fact, as long as provider *A* maintains a higher number of "good" reviews than "bad," he will always be preferred to provider *B*. However, this may not be the optimal outcome from a system perspective, which here refers to the outcome that maximizes the expected utility of the entire population of customers, because the customers' self-interested choices do not generate sufficient information on provider *B* to determine that he is, in fact, the inferior option.

The above example describes a phenomenon known in the experimentation literature as "under-exploration," as the self-interested individuals tend to take actions that "over-exploit" the information available to them. This chapter takes the perspective of a principal (e.g., the platform designer) who is interested in the efficient generation of information in such a system, where efficiency entails balancing exploration against exploitation with the goal of maximizing a long-run objective. Because the principal cannot dictate to the agents which action to take, she must find ways of incentivizing them to take system-optimal actions. Although we discuss a number of ways of achieving this, our main focus is on the active use of *information disclosure*, and in particular on the design of informational mechanisms that incentivize exploration in decentralized learning settings.

10.2 Related Literature

Studying the tradeoff between exploration and exploitation has a long research tradition in the context of the *multi-armed bandit* problem. In its classic version, a forward-looking decision maker makes a choice sequentially among a set of alternative arms, each of which generates rewards according to an ex ante unknown distribution. Every time an arm is chosen, the decision maker receives a reward, which, apart from its intrinsic value, is used to learn about the arm's underlying reward distribution. When deciding which arm to play, the decision maker faces the tradeoff between the arm that she currently believes to be superior (exploitation) given the information she has at her disposal, or an alternative arm with the goal of acquiring knowledge that can be used to make better-informed decisions in the

future (exploration). Since its inception, the multi-armed bandit framework has found numerous applications in various real-world settings (e.g., Caro and Gallien (2007), study dynamic assortment of seasonal goods in the presence of demand learning, while Bertsimas and Mersereau (2007), consider learning in the context of developing marketing strategies).

In most existing applications of the multi-armed bandit framework, a single decision maker dynamically decides on the actions to be taken while observing the outcomes of her past actions. As such, the decision maker fully internalizes the benefits of exploration when taking actions that may not be optimal as far as maximizing her present payoff is concerned. In contrast, this chapter focuses on settings that can be essentially cast as decentralized multi-armed bandits problems: there is a forward-looking principal (the designer) who seeks to maximize a long-term objective, while actions are taken by a series of (short-lived) agents. In particular, we discuss recent work along this direction that is mostly motivated by the growing popularity of online recommendation platforms. In a nice contribution, Kremer et al. (2014) focus on eliciting experimentation in an environment where outcomes are deterministic, while Papanastasiou et al. (2017) consider a stochastic environment, in which the designer is effectively tasked with managing a dynamic exploration-exploitation trade-off. Furthermore, Che and Horner (2017) consider a single-product setting where a designer at any time optimally “spams” a fraction of consumers to learn about the product’s quality. Frazier et al. (2014) aim to investigate how the principal can incentivize the agents to take her desired actions by offering direct monetary payments, i.e., their focus is not on the role of information disclosure policies (there is no *ex ante* or *ex post* asymmetry of information between the designer and the agents). Finally, Hörner and Skrzypacz (2016) also survey recent related work that combines ideas from experimentation, learning, and strategic interactions, with a particular emphasis on understanding how information but also delegation can be employed to deal with agents’ incentives.¹

Given the emphasis on the role of information the principal shares with the agents, the work we discuss here is related to, but quite distinct from, the well-developed literature on “cheap talk” (e.g., Crawford and Sobel 1982; Allon et al. 2011). In cheap-talk games, the principal privately observes the realization of an informative signal, after which she (costlessly) communicates any message she wants to the agent. In this work, there is emphasis on how the message received by the agent is interpreted, and whether any information can be credibly transmitted by the principal. In contrast, the principal in the settings we consider commits *ex ante* to an information-provision policy, which maps realizations of the informative signal to messages. Once this policy has been decided and implemented, the principal cannot manipulate the information she discloses (e.g., by misrepresenting the signal realization). In this case, there is no issue of how the agents will interpret the messages; rather, the focus is on how the principal should structure credible messages in a manner that internalizes the misalignment between her and the consumers’ objectives.

¹Kleinberg and Slivkins (2017) also presented recently a comprehensive tutorial related to these issues.

As such, this chapter discusses work that is more in the spirit of the recent stream of literature that examines how a principal can design/re-structure informative signals in ways that render agents *ex ante* more likely to take desirable actions. Bimpikis and Drakopoulos (2016) find that in order to overcome the adverse effects of free-riding, teams of agents working separately towards the same goal should initially not be allowed to share their progress for some pre-determined amount of time. Bimpikis et al. (2018) investigate innovation contests and demonstrate how award structures should be designed so as to implicitly enforce information-sharing mechanisms that incentivize participants to remain active in the contest. Kamenica and Gentzkow (2011) and Rayo and Segal (2010) illustrate an explicit technique for structuring informative signals—referred to as “Bayesian persuasion”—in static (i.e., one-shot) settings.

Furthermore, the discussion here connects to the work on social learning. The basic setup involves agents (e.g., consumers) that are initially endowed with private information regarding some unobservable state of the world (e.g., product quality). When actions (e.g., purchase decisions) are taken sequentially and are commonly observable, the seminal papers by Banerjee (1992) and Bikhchandani et al. (1992) demonstrate that herds may be triggered, whereby agents rationally disregard their private information and simply mimic the action of their predecessor. This classic paradigm has since been extended in multiple directions (e.g., representative references along this direction include Acemoglu et al. 2011, 2014; Lobel and Sadler 2015; Besbes and Scarsini 2017).

While the above papers focus on studying features of the learning process itself, another stream of literature investigates how firms can use their operational levers to steer the social-learning process to their advantage. Bose et al. (2006) and Crapis et al. (2017) investigate dynamic pricing in the presence of social learning that occurs on the basis of actions (i.e., purchase decisions) and outcomes (i.e., product reviews), respectively. Veeraraghavan and Debo (2009) and Debo et al. (2012) consider how customers’ queue-joining behavior depends on observable queue-length, and how service-rate decisions may be used to influence this behavior. Papanastasiou and Savva (2017) and Feldman et al. (2018) highlight how pricing and product-design policies are affected by the interaction between product reviews and strategic consumer behavior (see also Swinney (2011) for additional related work), while Allon and Zhang (2017) explore service-level differentiation for service organizations whose customers engage in communication through their social networks. Complementing this literature, the present chapter explores how the firm (platform) can influence consumer decisions and learning through its information-provision policy, a lever, which may also be used in conjunction with other operational levers (e.g., pricing, inventory).

Finally, the chapter is also broadly related to a recent line of work that studies operational decisions in the context of Internet-enabled business models. Among others, Marinesi et al. (2017) and Hu et al. (2013) study group-buying platforms; Balseiro et al. (2014, 2015) consider the design and operations of ad-exchanges;

Kanoria and Saban (2017) address search in two-sided platforms; and Taylor (2018), Cachon et al. (2017), and Bimpikis et al. (2017a) explore optimal pricing and compensation policies in on-demand service platforms.²

10.3 Illustrative Example

The following example, which is taken from Kremer et al. (2014), provides a nice illustration of the setting and the questions we explore in this chapter.

Example 1 Agents choose sequentially between products A and B. Agent i makes her decision based on her prior on the quality of the two products and the information she obtains from the principal. In turn, the principal observes the choices and resulting payoffs of agents $1, \dots, i - 1$, and makes a recommendation to agent i , i.e., whether to purchase product A or B. The principal commits ex ante to the mechanism that generates the recommendation for agent i , i.e., the function that maps the actions and payoffs of agents $1, \dots, i - 1$ to a binary recommendation. Furthermore, agents know the mechanism set by the principal for generating recommendations and take it into account when they form their (posterior) beliefs about the quality of the two products.

Assume that the agents' common prior is that the quality of product A is uniformly distributed in $[-1, 5]$ whereas the quality of product B is uniformly distributed in $[-5, 5]$. Also, assume that when an agent buys a product, her (realized) payoff is equal to the quality of the product, i.e., one purchase is enough to reveal a product's (true) quality. Finally, suppose that the principal aims to explore both alternatives as soon as possible (so that she recommends the best one to future agents).

If information about past choices and outcomes were observable by the agents, the second agent would choose to take action B only if the payoff of the first agent (that would optimally take action A) was negative. Otherwise, i.e., if product A has positive payoff, the second agent (and subsequently all future agents) would choose product A and no agent would find it optimal to explore product B (which, nevertheless, could have been the optimal choice).

On the other hand, if agents do not directly observe prior choices and outcomes, the principal could induce more exploration by recommending action B to the second agent whenever the payoff associated with product A is less than one. In other words, the principal could send a binary message to the second agent: choose A if the first agent's payoff was higher than one and choose B, otherwise. Similarly, the principal can employ the following policy for the third agent: recommend choosing product B if (i) the second agent was recommended to choose product B and it turned out that B's payoff is higher than A's; or (ii) both the first and

²There is also recent empirical work exploring operational issues on online marketplaces, e.g., Moon et al. (2017), Li and Netessine (2017), and Bimpikis et al. (2017b).

the second agent chose product A but its payoff is between 1 and 3.23. It is straightforward to show that following this policy guarantees that agents would have explored both options by the third time period unless the payoff for product A is higher than 3.23 (one can similarly extend the policy for the fourth agent to ensure that by the fourth time period both options are explored with certainty).

In sum, agents find it optimal to follow the principal's recommendations, which, in turn, leads to more exploration and better outcomes on aggregate (assuming that the population of agents is large enough). The simple takeaway message that one can draw from this example is that by coarsening the information that the principal shares with the agents, she is able to mitigate their misalignment of interests.

10.4 Benchmark Model

Building on the discussion above, we consider a setting, where a series of agents interact with a principal who manages the disclosure of information regarding the experiences of their predecessors. For concreteness and to be in line with Sect. 10.3, we anchor our exposition in the example of an online platform which is operated by a designer and is used by customers to assist with their choice of a service provider. We assume that the marketplace features two providers, A and B ; let $S = \{A, B\}$.³ Each provider $i \in S$ is fully characterized by a probability p_i , which represents the provider's service quality. Upon using provider i , a customer receives reward equal to one with probability p_i , and equal to zero otherwise; that is, service outcomes constitute independent draws from a Bernoulli distribution with success probability p_i . Initially, p_i is known to the designer and the customers only to the extent of a common prior belief, which is expressed in our model through a Beta random variable with shape parameters $\{s_1^i, f_1^i\}$, with $s_1^i, f_1^i \in \mathbb{Z}_+$.^{4,5}

At the beginning of each time period $t \in T$, $T = \{1, 2, \dots\}$, a single customer visits the platform, observes information pertaining to the experiences of past customers, and chooses a provider. We assume that upon completion of service, and before the end of period t , the customer reports to the platform whether her experience was positive or negative (i.e., the realization of the Bernoulli random variable associated with her service experience). At any time t , the knowledge accumulated by the platform is summarized by the *information state* (henceforth

³Our analysis can be readily extended to the case of more than two providers.

⁴The probability density function of a $Beta(s, f)$ random variable is given by

$$g(x; s, f) = \frac{x^{s-1}(1-x)^{f-1}}{B(s, f)}, \text{ for } x \in [0, 1].$$

⁵The platform and the customers hold the same prior belief, so that platform actions (e.g., the choice of an information-provision policy) do not convey any additional information on provider quality to the customers (e.g., Bergemann and Välimäki 1997; Bose et al. 2006; Papanastasiou and Savva 2017).

“state”) $x_t = \{x_t^A, x_t^B\}$, where $x_t^i = \{s_t^i, f_t^i\}$ and s_t^i (f_t^i) is the accumulated number of successful (failed) service outcomes for provider i up to period t (this includes the initial successes and failures, s_1^i and f_1^i , specified in the prior belief). When the system state is x_t , the Bayesian posterior belief over the quality p_i is $Beta(s_t^i, f_t^i)$, and the expected utility for the next customer if she uses i is $r(x_t, i) = s_t^i / (s_t^i + f_t^i)$.

In general, the history of service outcomes (i.e., the system state x_t) is not directly observable to the customers. Instead, there is a platform designer who *commits* upfront to a “messaging policy” that acts as an instrument of information-provision to the customers.⁶ This policy specifies the *message* that is displayed on the platform, given any underlying system state. In addition, the platform may accompany messages with monetary payments to customers as a further incentive to induce them to take certain actions (in fact, Frazier et al. (2014), exclusively explores the case where all generated information is observable to customers and the platform has the discretion to incentivize their actions through monetary transfers in the form of “coupons”).⁷ The designer’s objective in choosing her messaging policy is to maximize the expected sum of customers’ discounted rewards over an infinite horizon (i.e., customer surplus), applying a discount factor of $\delta \in [0, 1)$.⁸ Customers are modeled as homogeneous, short-lived, rational agents. In our main analysis, we assume that customers know the period of their arrival (however, the qualitative insights we obtain are robust to relaxing this assumption). Upon visiting the platform, each customer observes a message generated by the designer’s policy and chooses a service provider with the goal of maximizing her individual expected reward.

The designer’s choice of messaging policy (and potential monetary transfers to customers), along with the customers’ choices of service provider in response to this policy, simultaneously govern the dynamics of both the learning process and the customers’ reward stream.

10.5 Inducing Exploration

The section explores whether the designer can incentivize customers to take actions that contribute to her long-run objective of generating information about the available service providers using mainly the platform’s messaging policy. At the

⁶Commitment is a reasonable assumption in the context of online platforms, where information provision occurs on the basis of pre-decided algorithms and the large volume of products/services hosted renders ad-hoc adjustments of the automatically-generated content prohibitively costly.

⁷The generic term “message” refers to a specific configuration of information that is observed by the customer; examples of messages include detailed outcome histories (i.e., distributions of customer reviews), relative rankings of providers, or recommendations for a specific product.

⁸More generally, our analysis is relevant for cases where the platform has a different (e.g., longer-run) objective than its users.

end of the section, we also report on work that has studied the use of monetary transfers in a similar setting.

Equilibrium and Model Dynamics We begin our analysis by formalizing the strategic interaction between the designer and the customers. There are two main features of this interaction. First, the designer's *messaging policy*, which takes the platform state as an input and generates a message to be displayed by the platform to the next incoming customer. Second, the customers' *choice strategy*, which takes the platform's message in any given period as an input and determines the customer's action (choice of provider).

Let $X \subseteq \mathbb{Z}_+^4$ denote the set of possible states of the platform such that $x_t \in X$ for all $t \in T$, and define the discrete set M of feasible messages that the platform can display to an incoming customer in period t (see Footnote 7).

A messaging policy $g(\cdot)$ is a (possibly stochastic) mapping from the set of states X to the set of messages M ; that is, a messaging policy g associates with each state $x_t \in X$ a probability $P(g(x_t) = m)$ that message $m \in M$ is displayed on the platform. Let \mathcal{G} be the set of possible messaging policies.

In each period t , a single customer enters the system, observes the platform's message and chooses a service provider from the set S . The period- t customer's choice strategy, denoted by $c_t(\cdot)$, is a mapping from the set of messages M to the set of service providers S . Let \mathcal{C}_t be the set of possible choice strategies for the period- t customer, and define $c(\cdot) := [c_1(\cdot), c_2(\cdot), \dots]$.

The designer's messaging policy g along with the customers' choice strategy c generate a *controlled Markov chain* characterized by the stochastic state-action pairs $\{(x_t, y_t); t \in T\}$, where the actions y_t that accompany the states x_t are determined by the designer's policy and the customers' strategy via $y_t = c_t(g(x_t))$. When the state of the system is x_t , the expected reward of a customer that uses provider i is $r(x_t, i) = s_t^i / (s_t^i + f_t^i)$. Transitions between system states occur as follows. The initial state x_1 is determined by the prior belief over the two providers; when the state of the system is x_t and action y_t is chosen by the period- t customer, the state in period $t + 1$, $x_{t+1} = \{x_{t+1}^A, x_{t+1}^B\}$, is determined as follows:

$$x_{t+1}^i = x_t^i \text{ for } i \neq y_t, \quad x_{t+1}^i = \begin{cases} \{s_t^i + 1, f_t^i\} & \text{w.p. } r(x_t, i) \\ \{s_t^i, f_t^i + 1\} & \text{w.p. } 1 - r(x_t, i) \end{cases} \quad \text{for } i = y_t.$$

The above transition probabilities reflect the learning dynamics of the system: new information regarding the quality of provider i is generated in period t only if the provider is chosen by the period- t customer.⁹

The sequence of events in our model is described in reverse chronological order as follows. Each customer observes the designer's messaging policy and chooses

⁹Note that for the case of a Bernoulli reward process the current probability of success (i.e., the Bayesian probability of the next trial being a success given the current state of the system) is equal to the immediate expected reward, $r(x_t, i)$ (e.g., Gittins et al. 2011).

a choice strategy c_t to maximize her individual expected reward. In particular, the period- t customer's response to message m , $c_t^*(m)$ maximizes:

$$E_{x_t}[r(x_t, c_t) \mid g(x_t) = m].^{10}$$

At the beginning of the time horizon, the designer (taking into account the customers' response to any messaging policy), commits to a policy that maximizes the expected sum of customers' discounted rewards. In particular, the designer's messaging policy $g^*(x_t)$ maximizes

$$E \left[\sum_{t \in T} \delta^{t-1} r(x_t, y_t) \right], \quad \text{for } y_t = c_t^*(g(x_t)).$$

Incentive-Compatible Recommendation Policies In general, multiple equilibria exist that result in the same payoff for the designer and the customers, and the same dynamics in the learning process, not least because the same information can be conveyed from the designer to the customers through a multitude of interchangeable messages contained in M . We follow Allon et al. (2011) in referring to such equilibria as being “dynamics-and-outcome equivalent”. In our analysis, we will employ the result of Lemma 1 below to simplify the exposition and focus attention on the informational content of equilibria, rather than on the alternative ways in which these equilibria can be implemented. Before stating the lemma, we define a subclass of messaging policies, which we refer to as “incentive-compatible recommendation policies.”

Definition 1 (ICRP: Incentive-Compatible Recommendation Policy) A recommendation policy is a messaging policy defined as

$$g(x_t) = \begin{cases} A & \text{w.p. } q_{x_t} \\ B & \text{w.p. } 1 - q_{x_t}, \end{cases} \quad (10.1)$$

where $q_{x_t} \in [0, 1]$ for all $x_t \in X$. A recommendation policy is said to be incentive-compatible if for all $x_t \in X$, $t \in T$, we have $c_t^*(g(x_t)) = g(x_t)$.

Put simply, under an ICRP the platform recommends either provider A or provider B to the period- t customer, and the customer finds it Bayes-rational to follow this recommendation. We may now state the following result, which is analogous to the revelation principle in the mechanism-design literature, and suggests that any feasible platform payoff can be achieved through some ICRP.

¹⁰This expectation can be computed by the period- t customer, since the ex ante probability that the state in period t is x_t (i.e., unconditional on the message $g(x_t)$) is known to the customer through her knowledge of the designer's policy in previous periods and the preceding customers' best response to this policy.

Lemma 1 *For any arbitrary messaging policy g , there exists an ICRP g' which induces a dynamics-and-outcome equivalent equilibrium in the game between the designer and the customers.*

A proof for Lemma 1 can be found in Papanastasiou et al. (2017).

First Best As a primer to our main analysis, we consider how the designer would direct individual customers to the two providers, had the customers' actions been under her *full control*. The solution to the designer's full-control problem is due to Gittins and Jones (1974) and consists of directing customers in each period to the provider with the highest *Gittins Index*. The Gittins index for service i when in state z^i is denoted by $G_i(z^i)$ and given by:

$$G_i(z^i) = \sup_{\tau > 0} \frac{E \left[\sum_{t=0}^{\tau-1} \delta^t r(x_t^i, i) \mid x_0^i = z^i \right]}{E \left[\sum_{t=0}^{\tau-1} \delta^t \mid x_0^i = z^i \right]}, \quad (10.2)$$

where τ is a past-measurable stopping time (i.e., measurable with respect to the information obtained up to time τ) and $r(x_t^i, i)$ is the instantaneous expected reward of provider i in state x_t^i .

In the decentralized system, the designer's ability to direct customers to her desired provider will be limited by the customers' self-interested behavior. Each customer knows (i) the prior belief summarized by the initial state, x_1 ; (ii) the time period, t ; and (iii) the designer's messaging policy, g . Upon visiting the platform, the customer observes a message m , updates her belief over the current system state, x_t , and selects the provider which maximizes her individual expected reward. As a consequence, the designer will be able to achieve first-best only if she can design a messaging policy which induces customers to make Gittins-optimal decisions in all periods and in all system states—a sufficient condition for at least one such messaging policy to exist is the existence of an ICRP which always recommends the provider of highest Gittins index.

Throughout the following analysis we will refer to provider choices that are desirable from the platform's perspective as being "system-optimal."

10.5.1 Strategic Information Disclosure

Typically, the provider will not be able to achieve the first best given that Gittins-based recommendations (system optimal provider choice) are not incentive compatible in general. This section provides a characterization of the designer's optimal policy in the presence of incentive constraints resulting from the customers' decision making.

By Lemma 1, the designer in our model seeks to find the best possible ICRP, that is, to choose optimally the probabilities q_{x_t} that define the recommendations received by the period- t customer in each possible system state:

$$g(x_t) = \begin{cases} A & \text{w.p. } q_{x_t} \\ B & \text{w.p. } 1 - q_{x_t}, \end{cases}$$

while at the same time ensuring that any recommendation received by the period- t customer is incentive compatible. The designer's general problem may be framed as the following *Constrained* Markov Decision Process (CMDP; see Altman 1999),

$$\begin{aligned} \max_{g(x_t)} E \left[\sum_{t \in T} \delta^{t-1} r(x_t, g(x_t)) \right] \\ \text{s.t. } E_{x_t}[r(x_t, A) \mid g(x_t) = A] \geq E_{x_t}[r(x_t, B) \mid g(x_t) = A], \forall t \in T, \\ E_{x_t}[r(x_t, B) \mid g(x_t) = B] \geq E_{x_t}[r(x_t, A) \mid g(x_t) = B], \forall t \in T, \end{aligned} \quad (10.3)$$

where the constraints state that any recommendation that is generated by policy g in period t is found to be incentive compatible (and is therefore followed) by the period- t customer.

The presence of the IC constraints introduces both direct and indirect complications. The direct complication is that recommendations generated by the designer's policy in all states that *could* occur in period t must now be viewed jointly, since such recommendations are coupled by the need to satisfy the period- t customer's IC constraints. The indirect complication is that the designer's choice of policy up to period t affects the beliefs of customers that visit the platform in periods $t + 1$ onwards, and therefore (through the IC constraints) also affects the feasible region of recommendations in future periods.

To facilitate exposition of the result that follows, we introduce the following additional notation. Let X_t be the set of states that are reachable from the initial state x_1 (under some policy) in period t , so that the total state space is $X = \bigcup_{t \in T} X_t$. Denote by \mathcal{P}_{kiz} the transition probability from state k to state z when provider i is used (note that these probabilities have been specified in Sect. 10.5), and let Δ_a denote the Dirac delta function concentrated at a .¹¹

Proposition 1 *The optimal ICRP is given by*

$$q_k^* = \frac{\rho(k, A)}{\sum_{i \in S} \rho(k, i)},$$

¹¹The result of Proposition 1 extends readily to the case of $|S| = n$ providers (in this case, an ICRP consists of n possible recommendations, and each recommendation must satisfy $n - 1$ IC constraints per period), as well as to alternative platform objective functions (by replacing $r(k, i)$ with suitable reward functions).

where $\rho(k, i)$ solve

$$\begin{aligned}
& \max_{\rho} \sum_{k \in X} \sum_{i \in S} \rho(k, i) r(k, i) \\
& \text{s.t.} \sum_{k \in X_t} \rho(k, B) [r(k, B) - r(k, A)] \geq 0, & \forall t \in T, \\
& \sum_{k \in X} \sum_{i \in S} \rho(k, i) (\Delta_z(k) - \delta \mathcal{P}_{kiz}) = \Delta_{x_1}(z), & \forall z \in X, \\
& \rho(k, i) \geq 0, & \forall k \in X, i \in S. \quad (10.4)
\end{aligned}$$

A few comments on the solution technique of Proposition 1 are warranted. To solve the designer's problem, the objective and constraints of the CMDP (10.3) are first expressed as sums of the immediate expected reward in each state-action pair, $r(k, i)$, multiplied by the time-discounted "occupancy" of the pair, $\rho(k, i)$. Then, the LP (10.4) optimizes over the admissible set of occupancy measures, which is described by the LP's constraints. In particular, in the context of our problem, any admissible occupancy measure must be consistent with (i) the customers' incentives (this is captured by the period-specific inequality constraints, which ensure that each period- t customer finds the recommendation she receives IC), and (ii) the system's dynamics (this is captured by the state-specific equality constraints, which ensure that the occupancy of each state is consistent with the system's state-transition probabilities).¹² Finally, once the optimal occupancy measure has been identified, the probabilities q_k^* are chosen in a manner that induces this measure.

To gain insight into the structure of optimal policies, it is instructive to consider a finite-horizon version of the problem, consisting of T_F time periods. In this case, applying Theorem 3.8 of Altman (1999) reveals that the optimal ICRP uses randomized recommendations in at most T_F states. As the horizon length T_F increases, the state space grows exponentially, but the number of states in which randomization occurs grows only linearly (for instance, the number of possible states for $T_F = 20$ is of the order 10^{12} , but randomization occurs in at most 20 states). This suggests that optimal policies consist mainly of deterministic recommendations, relying extensively on the merging different information states that could correspond to different optimal actions for the customers to "persuade" them to explore.

¹²Note that the solution to LP (10.4) can also be used to retrieve the period- t customer's belief over the system state upon entry to the platform; specifically, this belief is given by $P(x_t = z) = \sum_{i \in S} \rho(z, i) / (\sum_{k \in X_t} \sum_{i \in S} \rho(k, i))$.

10.5.2 The Value of Information Obfuscation

The “curse of dimensionality” renders the optimal solution to the designer’s general problem computationally intractable. However, by combining the structural insights yielded by our analysis (i.e., state-merging, limited randomizations, sufficiency of two-message policies), it is possible to generate tractable and effective heuristic solutions. In this section, we consider one such heuristic and use it to establish that the value of information obfuscation is significant, even if this is implemented in a simple and intuitive manner (we note that the payoff under any heuristic serves as a lower bound on the payoff of the optimal policy described in Proposition 1).

Consider the following Gittins-based heuristic, which combines our preceding analysis with the centralized solution to the designer’s problem to deliver IC recommendations. Let p_{x_t} denote the probability that the state in period t is x_t . The heuristic is initialized by choosing the starting state x_1 and proceeds by repeating two steps. First, it solves the period- t linear program:

$$\begin{aligned} \max_{0 \leq q_{x_t} \leq 1} \quad & \sum_{x_t \in X} p_{x_t} q_{x_t} [G_A(x_t) - G_B(x_t)] \\ \text{s.t.} \quad & \sum_{x_t \in X} p_{x_t} (1 - q_{x_t}) [r(x_t, B) - r(x_t, A)] \geq 0, \end{aligned} \quad (10.5)$$

and stores the solution q_{x_t} (this is the designer’s recommendation policy for period t); second, the period- t solution is used along with the probabilities p_{x_t} to calculate the probabilities $p_{x_{t+1}}$. The two steps are repeated until a pre-specified period $t = K$ is reached, after which a full-information policy is employed (or, equivalently, an ICRP which always recommends the provider of highest expected reward). Essentially, in each of the first K periods of the horizon, the heuristic employs state-merging to deliver recommendations that maximize the expected Gittins index, subject to the recommendations being IC.

To evaluate the benefits of information obfuscation (in the sense of the Gittins-based heuristic), we conduct the numerical experiments presented in Table 10.1. The table focuses on the added “learning value” of obfuscation in comparison to that of a FI policy. Specifically, we first calculate the difference $(\pi^* - \pi^{NI})$, i.e., the difference between the platform’s payoff when no social learning takes place (π^{NI}) and when social learning takes place optimally (π^*). This difference is an upper bound on the learning value that can be achieved by the designer in the decentralized system through information-provision. We then calculate the percentage of this value achieved under FI ($\Delta\pi^{FI}$) and under the Gittins-based heuristic ($\Delta\pi(\hat{g})$).

The upper half of the table pertains to initial states which are “unfavorable” for the designer, in the sense that there is an ex ante misalignment between the provider of highest expected reward and the provider of highest Gittins index; by contrast, the lower part of the table pertains to “favorable” initial states. Across all instances we consider, the heuristic performs significantly better than full

Table 10.1 Proportion of first-best learning value captured in the decentralized system by FI , defined as $\Delta\pi^{FI} = (\pi^{FI} - \pi^{NI})/(\pi^* - \pi^{NI})$, and by the Gittins-based heuristic \hat{g} with $K = 50$, defined as $\Delta\pi(\hat{g}) = (\pi(\hat{g}) - \pi^{NI})/(\pi^* - \pi^{NI})$ (where π^* , π^{FI} , π^{NI} and $\pi(\hat{g})$ denote expected platform payoff under first best, FI , NI and the Gittins-based heuristic, respectively). $r(x_1, i)$ and $\text{std}(x_1, i)$ denote, respectively, the expectation and standard deviation of the reward of provider $i \in \{A, B\}$ at the initial state x_1 . Parameter values: $\delta = 0.99$

$x_1 = \{(a_1^A, b_1^A), (a_1^B, b_1^B)\}$	$r(x_1, A)$	$\text{std}(x_1, A)$	$r(x_1, B)$	$\text{std}(x_1, B)$	$\Delta\pi^{FI}(\%)$	$\Delta\pi(\hat{g})(\%)$
$\{(6, 3), (1, 1)\}$	0.67	0.15	0.5	0.29	47.2	96.3
$\{(12, 6), (1, 1)\}$	0.67	0.11	0.5	0.29	18.6	85.0
$\{(18, 9), (1, 1)\}$	0.67	0.09	0.5	0.29	6.0	83.7
$\{(15, 6), (2, 1)\}$	0.71	0.10	0.67	0.24	58.1	97.8
$\{(15, 6), (4, 2)\}$	0.71	0.10	0.67	0.18	66.0	90.7
$\{(15, 6), (6, 3)\}$	0.71	0.10	0.67	0.15	71.7	93.0
$\{(1, 1), (3, 6)\}$	0.5	0.29	0.33	0.15	87.6	100
$\{(1, 1), (6, 12)\}$	0.5	0.29	0.33	0.11	81.0	95.9
$\{(1, 1), (9, 18)\}$	0.5	0.29	0.33	0.09	80.0	100
$\{(1, 1), (3, 6)\}$	0.5	0.29	0.33	0.15	85.4	94.6
$\{(3, 3), (3, 6)\}$	0.5	0.19	0.33	0.15	85.9	94.6
$\{(6, 6), (3, 6)\}$	0.5	0.14	0.33	0.15	51.1	96.2

information. Furthermore, we observe that the benefit is highest when the initial state is unfavorable: in such cases, under full information the customers tend to stick with the ex ante preferable provider and only rarely engage in experimentation with the alternative option. Next, notice that in each of the four subgroups of initial states, the ex ante expected reward of the two providers is maintained constant, but the variance of one of the two changes; this allows us to capture different environments in terms of the potential benefits of exploration. Here, intuitively, we observe that the benefits of information obfuscation are especially pronounced when the quality of the ex ante preferable provider is relatively certain while the quality of the alternative provider is relatively uncertain.

10.5.3 Minimizing Regret

In the setting we have considered so far, the designer’s objective was to maximize the expected discounted sum of the customers’ rewards over an infinite horizon. A related objective that has been studied in the literature is that of minimizing the designer’s long-run *regret*. Typically, focusing on regret as the designer’s objective simplifies the analysis (at least to some extent) and, thus, allows for a different set of results that mainly provide reasonable guarantees of performance for relatively simple strategies.

The discussion in this section follows Kremer et al. (2014) and Mansour et al. (2015). In the context of minimizing regret, the performance of a suggested policy

is compared to the case that the designer knows the stochastic process generating the rewards for its customers, i.e., knows p_i for $i \in \{A, B\}$, and, thus, would always direct them towards the better of the two service providers.

Kremer et al. (2014) consider a horizon of T time periods (thus, T incoming customers) and propose the following policy for generating messages (recommendations) to them:

1. Customers are partitioned into $\lceil T/m \rceil$ blocks of m customers each. Customers belonging to the same block receive the same recommendation (and, thus, end up using the same service provider at the induced equilibrium). Customers in the first block are recommended to visit the provider that is ex ante more likely to generate higher expected rewards (based on the common prior), say provider A .
2. The designer observes the realizations of the rewards for the first m agents and computes the average empirical mean reward, $\hat{\mu}_A$, for provider A , i.e., the provider that is ex ante more likely to be a better choice for the customers.
3. Keeping $\hat{\mu}_A$ fixed throughout the horizon, the designer recommends provider B to the i -th block of customers, if $\hat{\mu}_A \in (\theta_{i-1}, \theta_i]$, where $\{\theta_i\}_{i=1}^{\lceil T/m \rceil}$ is a set of thresholds that the designer determines so that the recommendations she makes to customers are incentive compatible (essentially θ_i is such that customers would be indifferent between following the designer's recommendation and choosing the provider that is ex ante more likely to generate higher rewards, i.e., provider A , if $\hat{\mu}_A$ was exactly equal to the threshold).¹³ The first time customers are recommended to use provider B (and, as a result, end up using B), the designer computes $\hat{\mu}_B$ based on their realized rewards.
4. After the designer recommends provider B for the first time and computes $\hat{\mu}_B$, the messaging policy takes the following form:
 - If $\hat{\mu}_A \leq \theta_{i-1}$, the designer recommends the provider that corresponds to the highest empirical mean, i.e., she recommends provider A if $\hat{\mu}_A \geq \hat{\mu}_B$, and provider B otherwise.
 - $\hat{\mu}_A > \theta_i$, the designer recommends provider A .

Kremer et al. (2014) provide the following theorem for the performance of the messaging policy described above:

Theorem 1 *Setting the size of each block to $T^{2/3} \ln T$, i.e., $m = T^{2/3} \ln T$, guarantees that the average regret per customer, i.e., the expected difference between taking the best possible action and following the designer's recommendation, is bounded above by:*

$$C \frac{\ln T}{T^{1/3}},$$

where C is a constant that depends only on the priors.

¹³This is a natural generalization of the computation in the example of Sect. 10.3.

In other words, the theorem above implies that as the horizon gets longer (equivalently, the population of customers getting recommendations from the platform increases), the average regret per customer becomes negligible. Thus, this simple policy that appropriately partitions customers into different blocks, achieves a reasonably good asymptotic performance compared to always choosing the best available service provider.

Mansour et al. (2015) extend Kremer et al. (2014) by considering a setting where in each of T time periods n new customers interact with the platform and simultaneously take an action (out of $k \geq 2$ potential alternatives). The payoff of each customer is determined by an underlying state (that captures the quality of each of the alternatives) and the actions of the rest of the customers in her cohort. As in Kremer et al. (2014) and Mansour et al. (2015) present results on (simple) policies that are near-optimal (in a regret minimization sense) compared to the best-in-hindsight policy.¹⁴

10.5.4 *Incentivizing Customers Using Payments*

Although we mainly focus on the role of the designer's information disclosure policy to induce system-optimal actions, it is reasonable to consider a setting where the designer may use monetary transfers as a way to promote exploration among the platform's customers. In particular, a very interesting direction for future work would be to extend the modeling framework in Sect. 10.4 to a setting where the designer can combine her messaging policy with monetary transfers with the objective of maximizing the discounted sum of the customers' expected rewards minus the corresponding transfers.

Relatedly, Frazier et al. (2014) explore the use of monetary transfers as a way to incentivize exploration when the information generated by the customers' past interactions with the service providers is available to both the designer but also to future customers, i.e., there is no ex post information asymmetry between the designer and the customers. For example, two extreme policies that one could consider (and Frazier et al. (2014) discuss briefly) would be the following:

1. The designer never compensates customers for taking an action. Then, each customer chooses the action that maximizes her one time period payoff based on the information generated by past actions. In other words, customers never "explore" and always "exploit" (based on the history of actions and payoffs that they observe). Such myopic behavior is typically suboptimal given that the rate of exploration is inefficient from the designer's perspective.

¹⁴Che and Horner (2017) also consider the problem of optimally designing recommendation policies in a setting where information about the quality of two potential alternatives arrives continuously over time—their setting uses the exponential bandit framework of Keller et al. (2005) as a building block.

2. On the other extreme, the designer compensates customers sufficiently to induce customers to take the system-optimal action at every time period, i.e., the designer offers a payment to the customer taking action at time t , which is equal to the difference between the expected reward corresponding to the service provider with the highest Gittins index at t and the service provider that would be myopically optimal to choose given the available information. Obviously, this policy generates the system-optimal rate of exploration, but it may lead to large cumulative payments from the designer.

Frazier et al. (2014) provide a characterization of the extent to which payments from the designer to the customers can mitigate their incentive constraints and recover the optimal reward on aggregate. In particular, letting OPT denote the first best cumulative rewards, i.e., the discounted sum of the expected rewards corresponding to always choosing the service provider with the highest Gittins index, they call a point $(\alpha, \beta) \in [0, 1]^2$ *achievable at discount rate δ* if there exists a payment policy for the designer, i.e., a mapping from the history of observations to payments to customers, that satisfies the following two conditions:

1. The discounted sum of the expected rewards corresponding to the policy is at least as high as $\alpha \cdot \text{OPT}$.
2. The discounted sum of the expected payments corresponding to the policy is at most as high as $\beta \cdot \text{OPT}$.

The main result in Frazier et al. (2014) is the following theorem that provides a remarkably clean characterization of the set of points that can be achieved by the designer.

Theorem 2 *Let (α, β) be a point in $[0, 1]^2$. Then, (α, β) is achievable at discount rate δ if and only if the following condition holds:*

$$\sqrt{\beta} + \sqrt{1 - \alpha} > \sqrt{\delta}.$$

Theorem 2 provides some insight on what the designer can (and cannot) do using payments (“coupons”) to incentivize the platform’s customers. A basic ingredient of the proof is a set of policies that involve mixing between the two extremes described above, i.e., the policy that involves no payments and the one where payments are large enough to induce customers to take the action with the highest Gittins index (thus, giving some idea on what type of policies may lead to good performance for the designer). Note that this is a worst-case result, i.e., it bounds the designer’s performance against any distribution of rewards. Thus, the designer could potentially achieve a higher discounted sum of expected rewards in environments where the uncertainty in payoffs takes a more specific form, like the one specified in Sect. 10.4.¹⁵ Importantly though, Theorem 2 assumes that customers can observe the entire history of actions and their corresponding rewards; thus, it does not offer any insight on how the designer might appropriately disclose information to increase the set of achievable points.

¹⁵However, the analysis may, in general, be quite challenging.

10.6 Promising Directions

So far, we have mainly used the example of an online recommendation platform to describe the main questions that motivate this chapter and illustrate a number of key findings. However, the idea that a platform (principal) can appropriately design the information flow to its users (agents) as a way to incentivize them to take actions that may not be myopically optimal is more widely applicable and, to a large extent, still unexplored. In this section, we briefly describe two other application settings that may provide interesting starting points for future work in the area.

10.6.1 *Learning in Dynamic Contests*

Innovation contests are gaining in popularity as a tool that firms and institutions employ to outsource their R&D and innovation efforts to the crowd. An open call is placed for a project that participants compete to finish and the winners, if any, are awarded a prize. Recent successful examples include The NetFlix Prize and the Heritage Prize,¹⁶ and a growing number of ventures like Innocentive, TopCoder, and Kaggle provide online platforms to connect innovation seekers with potential innovators.

The objective of the contest designer is to maximize the probability of reaching the innovation goal while minimizing the time it takes to complete the project. Obviously, the success of a contest depends crucially on the pool of participants and the amount of effort they decide to exert. Typically, innovation projects have the following three key features. First, progress towards the end goal takes the form of a series of breakthroughs interspersed between long intervals of seeming stagnation. Second, and quite importantly, it is not clear at the onset whether the end goal is attainable, even if it is clearly specified, or which of potentially many alternative approaches would be the best one to use. Finally, a third feature that distinguishes innovation contests from other settings involving competition among agents, is that agents can learn from one another: an agent's (partial) progress towards the goal provides useful information to the rest about the feasibility of the project and/or the best approach to follow.¹⁷

¹⁶The NetFlix Prize offered a million dollars to anyone who succeeded in improving the company's recommendation algorithm by a certain margin and was concluded in 2009. The Heritage Prize was a multi-year contest whose goal was to provide an algorithm that predicts patient readmissions to hospitals. A successful breakthrough was obtained in 2013.

¹⁷In addition to the work that we discuss here, which mainly focuses on the dynamics of learning and competition in contests, there is also an extensive body of work that explore a number of questions in a static framework, e.g., Terwiesch and Xu (2008), Ales et al. (2017), and Körpeoğlu and Cho (2017).

These three features imply that news about a participant's progress has the following interesting dual role: it makes agents more optimistic about the state of the world, as the goal is more likely to be attainable; thus, agents have a higher incentive to exert costly effort. We call this the *encouragement effect*.¹⁸ At the same time, such information implies that one of the participants has a lead, which might negatively affect effort provision from the remaining agents as the likelihood of them beating the leader and winning the prize becomes slimmer. We refer to this as the *competition effect*. These two effects interact with one other in subtle ways over the duration of the contest, and understanding this interaction is of first-order importance for contest design.

Thus, the contest's information disclosure policy (e.g., through intermediate milestone awards) may have a large effect on the agents' participation and effort provision and, consequently, on the likelihood that the contest will be successful. In recent work, Bimpikis et al. (2018) consider the question of *whether* and *when* should the contest designer disclose information regarding the competitors' (partial) progress with the goal of maximizing her expected payoff. Interestingly, they illustrate the benefits of non-trivial information disclosure policies, where the designer withholds information from the agents and only releases it after a certain amount of time has elapsed. Such designs further highlight the active role that information may play in incentivizing agents to participate in the contest.

Second, they identify the role of intermediate awards as a way for the designer to implement the desired information disclosure policy—the policy that maximizes the effort provision of the agents and consequently the chances of innovation taking place. Intermediate awards are very common in innovation contests (the aforementioned Netflix and Heritage prizes are examples of contests that have employed intermediate awards) and the discussion here provides a potential reason for their ubiquity.

A simple illustration of the main ideas in Bimpikis et al. (2018) is the following. Consider an innovation contest that consists of well-defined milestones. For example, the goal of the Netflix prize was to achieve an improvement of 10% over the company's proprietary algorithm, with a first progress prize set at 1% improvement. In this example, reaching the milestone of 1% improvement constitutes *partial progress* towards the goal, and we assume that the agents and the designer are able to verifiably communicate this. Assume for now that the innovation is attainable with certainty given enough effort, and that agents are fully aware of that. The lack of progress towards the goal is then solely a result of the stochastic return on effort. When no information is disclosed about the agents' progress, they become progressively more pessimistic about the prospect of them winning, as they believe that someone must have made progress and that they are now lagging behind in

¹⁸The term “encouragement” originates from the literature on strategic experimentation (e.g., Bolton and Harris 1999; Keller et al. 2005).

the race towards the end goal. This might lead them to abandon the contest, thus decreasing the aggregate level of effort and consequently increasing the time to complete the innovation project.

In contrast, when there is uncertainty about the feasibility of the end goal, agents that have made little or no progress towards the goal become pessimistic about whether it is even possible to complete the contest. If this persists, an agent may drop out of the competition as she believes that it is not worth putting the effort for what is likely an unattainable goal, reducing aggregate experimentation in the process and decreasing the chances of reaching a feasible innovation.

This highlights the complex role that information about the agents' progress play in this environment. In the first scenario, when the competition effect is dominant (since there is no uncertainty about the attainability of the end goal), disclosing that one of the participants is ahead may deter future effort provision as it implies that the probability of winning is lower for the laggards. In the second case, when the encouragement effect dominates, an agent's progress can be perceived as good news, since it reduces the uncertainty regarding the feasibility of the end goal.¹⁹

Several directions may be worth pursuing following the main ideas presented in this section. Bimpikis et al. (2018) and Halac et al. (2017) consider contests with well-defined goals that may be unattainable. Alternatively, one could consider a setting where the contest designer's goal is to obtain the best solution or implementation possible to a given project. There can exist multiple approaches that one could employ and observing each other's progress reveals information about their relative merits.²⁰ What is then the optimal way to disclose information as way to balance the tradeoff between exploring the space of alternative approaches and driving effort towards those that look most promising?

An agent's progress provides information not only about the feasibility of the innovation project or the quality of the approach she is employing but also her skill level. Thus, in a setting where there is uncertainty about how good the competition is, appropriately designing an information disclosure policy may play an important role in keeping agents engaged and willing to exert effort.²¹

10.6.2 *Dealing with Misinformation*

The 2016 US presidential elections and the associated "fake news" phenomenon highlighted the importance of incentivizing a new form of "exploration" in the

¹⁹Bimpikis and Drakopoulos (2016) and Halac et al. (2017) also consider a strategic experimentation framework to study the interplay between a principal and the agents' incentives and how appropriately designing information disclosure mechanisms may increase welfare.

²⁰Girotra et al. (2010), Kornish and Ulrich (2011), Huang et al. (2014), and Jiang et al. (2016) are recent empirical studies that consider the role of learning and feedback in crowdsourcing contests and, more broadly, in the innovation process.

²¹There are also a number of notable recent papers that consider different aspects of contest design and its applications, e.g., Seel and Strack (2016), Hu and Wang (2017), and Strack (2016).

online space. Rather than exploration with the goal of identifying the quality of a product, the term here refers to exploration with the goal of evaluating the quality of an information source; for instance, if the information source in question is a news article circulating in a social media platform, exploration refers to the process of “fact-checking” the article’s content to determine its validity.

With this context in mind, Papanastasiou (2017) develops a sequential model of news propagation with endogenous fact-checking, and identifies pathological outcomes whereby fabricated news articles fail to be detected and subsequently spread throughout the society, manipulating the beliefs of the agents in the process. In this setting, there is “under-exploration” from a system perspective in the sense that individual agents do not internalize the impact of their fact-checking decisions on the actions and beliefs of their downstream peers, which in turn may result in inadequate levels of fact-checking. The study proceeds by analyzing a first-order defense against the propagation of fake news involving a social media platform that decides whether and when to intervene with the sharing of a news article by conducting its own fact-check (an approach recently adopted by Facebook).

An alternative to the platform conducting its own fact-checks is that of incentivizing the agents to do so through direct monetary payments, in a manner analogous to the setup of Frazier et al. (2014). Perhaps a more interesting avenue, however, is the design of appropriate information-disclosure mechanisms that may be able to achieve the same effect in a cost-effective way. One insight from Papanastasiou (2017) is that the fact-checking decisions of agents are influenced to a large extent by the number of times the information in question has been shared between their peers. Thus, one might expect that the concealment of such information by the platform may increase the amount of scrutiny an article undergoes, thereby reducing the propagation of fabricated information. At the same time, too much fact-checking is also an inefficient outcome: every fact-check incurs a cost to the agents which may be unnecessary. It follows that, as in Papanastasiou et al. (2017), the optimal information policy must strike a balance between allowing the agents to exploit the information generated by their peers, while also motivating them to explore (fact-check) at a system-optimal level.

In a somewhat related direction, Candogan and Drakopoulos (2017) study the tradeoff between user engagement and misinformation in the context of an online social networking platform. The content available on the platform may contain inaccuracies and false claims. The platform, which knows the quality of its content, may use a signaling device, e.g., recommend whether agents engage or not with the content, so as to induce a desired engagement behavior. A main emphasis in this line of work is the interplay between the platform’s (signaling) policy and the structure of the agents’ social network.

10.7 Concluding Remarks

This chapter showcases that choosing whether, when, and what information to disclose to agents may have a first-order impact on the payoff of a principal. Most of the exposition centers around the example of an online recommendation platform

(e.g., Yelp or Tripadvisor) but as we highlight in Sect. 10.6 these ideas may apply to many more real-world settings. Our hope is that the discussion provided here makes clear that information disclosure policies may be effective operational levers especially in the context of online platforms that rely on their users for ensuring a high quality of service. We believe that the role of information flows in mitigating the potential misalignment of interests between a principal and an agent/set of agents is quite important and relatively unexplored, and may thus provide a fruitful avenue for future research.²² Although the scope of the ideas presented here is quite broad, we expect that they will be particularly relevant in the design and operations of online platforms and marketplaces.

References

- Acemoglu D, Dahleh MA, Lobel I, Ozdaglar A (2011) Bayesian learning in social networks. *Rev Econ Stud* 78(4):1201–1236
- Acemoglu D, Bimpikis K, Ozdaglar A (2014) Dynamics of information exchange in endogenous social networks. *Theor Econ* 9(1):41–97
- Ales L, Cho SH, Körpeoğlu E (2017, Forthcoming) Optimal award scheme in innovation tournaments. *Oper Res* 65(3):693–702
- Allon G, Zhang DJ (2017) Managing service systems in the presence of social networks. Working paper
- Allon G, Bassamboo A, Gurvich I (2011) “We will be right with you”: managing customer expectations with vague promises and cheap talk. *Oper Res* 59(6):1382–1394
- Altman E (1999) *Constrained Markov decision processes*. CRC Press, Boca Raton
- Balseiro SR, Feldman J, Mirrokni V, Muthukrishnan S (2014) Yield optimization of display advertising with ad exchange. *Manag Sci* 60(12):2886–2907
- Balseiro SR, Besbes O, Weintraub GY (2015) Repeated auctions with budgets in ad exchanges: approximations and design. *Manag Sci* 61(4):864–884
- Banerjee A (1992) A simple model of herd behavior. *Q J Econ* 107(3):797–817
- Bergemann D, Välimäki J (1997) Market diffusion with two-sided learning. *RAND J Econ* 28(4):773–795
- Bertsimas D, Mersereau A (2007) A learning approach for interactive marketing to a customer segment. *Oper Res* 55(6):1120–1135
- Besbes O, Scarsini M (2017, Forthcoming) On information distortions in online ratings. *Oper Res* 66(3):597–610
- Bikhchandani S, Hirshleifer D, Welch I (1992) A theory of fads, fashion, custom, and cultural change as informational cascades. *J Polit Econ* 100(5):992–1026
- Bimpikis K, Drakopoulos K (2016) Disclosing information in strategic experimentation. Working paper
- Bimpikis K, Candogan O, Saban D (2017a) Spatial pricing in ride-sharing networks. Working paper
- Bimpikis K, Elmaghraby WJ, Moon K, Zhang W (2017b) Managing market thickness in online B2B markets. Working paper

²²There is currently significant interest in the role of information in mitigating the potential misalignment of interests between a principal and an agent/set of agents, e.g., Renault et al. (2017), Ely (2017), and Orlov et al. (2017).

- Bimpikis K, Ehsani S, Mostagir M (2018, Forthcoming) Designing dynamic contests. *Oper Res*
- Bolton P, Harris C (1999) Strategic experimentation. *Econometrica* 67(2):349–374
- Bose S, Orosel G, Ottaviani M, Vesterlund L (2006) Dynamic monopoly pricing and herding. *RAND J Econ* 37(4):910–928
- Cachon GP, Daniels KM, Lobel R (2017) The role of surge pricing on a service platform with self-scheduling capacity. *Manuf Serv Oper Manag* 19(3):368–384
- Candogan O, Drakopoulos K (2017) Optimal signaling of content accuracy: engagement vs. misinformation. Working paper
- Caro F, Gallien J (2007) Dynamic assortment with demand learning for seasonal consumer goods. *Manag Sci* 53(2):276–292
- Che YK, Horner J (2017) Recommender systems as incentives for social learning. Working paper
- Crapis D, Ifrach B, Maglaras C, Scarsini M (2017) Monopoly pricing in the presence of social learning. *Manag Sci* 63(11):3586–3608
- Crawford V, Sobel J (1982) Strategic information transmission. *Econometrica* 50(6):1431–1451
- Debo L, Parlour C, Rajan U (2012) Signaling quality via queues. *Manag Sci* 58(5):876–891
- Ely JC (2017) Beeps. *Am Econ Rev* 107(1):31–53
- Feldman P, Papanastasiou Y, Segev E (2018, Forthcoming) Social learning and the design of new experience goods. *Manag Sci*
- Frazier P, Kempe D, Kleinberg J, Kleinberg R (2014) Incentivizing exploration. In: Proceedings of the 15th ACM conference on economics and computation, Palo Alto. ACM, pp 5–22
- Girotra K, Terwiesch C, Ulrich KT (2010) Idea generation and the quality of the best idea. *Manag Sci* 56(4):591–605
- Gittins J, Jones D (1974) A dynamic allocation index for the sequential design of experiments. *Progress in statistics*, pp 241–266. Read at the 1972 European Meeting of Statisticians, Budapest
- Gittins J, Glazebrook K, Weber R (2011) Multi-armed bandit allocation indices. Wiley, Chichester
- Halac M, Kartik N, Liu Q (2017) Contests for experimentation. *J Polit Econ* 125(5):1523–1569
- Hörner J, Skrzypacz A (2016) Learning, experimentation and information design. Survey Prepared for the 2015 econometric summer meetings in Montreal
- Hu M, Wang L (2017) Joint vs. separate crowdsourcing contests. Working paper
- Hu M, Shi M, Wu J (2013) Simultaneous vs. sequential group-buying mechanisms. *Manag Sci* 59(12):2805–2822
- Huang Y, Vir Singh P, Srinivasan K (2014) Crowdsourcing new product ideas under consumer learning. *Manag Sci* 60(9):2138–2159
- Jiang ZZ, Huang Y, Beil DR (2016) The role of feedback in dynamic crowdsourcing contests: a structural empirical analysis. Working paper
- Kamenica E, Gentzkow M (2011) Bayesian persuasion. *Am Econ Rev* 101(6):2590–2615
- Kanoria Y, Saban D (2017) Facilitating the search for partners on matching platforms: restricting agent actions. Working paper
- Keller G, Rady S, Cripps M (2005) Strategic experimentation with exponential bandits. *Econometrica* 73(1):39–68
- Kleinberg RD, Slivkins A (2017) Tutorial: incentivizing and coordinating exploration. In: Proceedings of the 18th ACM conference on economics and computation, Cambridge
- Kornish LJ, Ulrich KT (2011) Opportunity spaces in innovation: empirical analysis of large samples of ideas. *Manag Sci* 57(1):107–128
- Körpeoğlu E, Cho SH (2017, Forthcoming) Incentives in contests with heterogeneous solvers. *Manag Sci* 64:2709–2715
- Kremer I, Mansour Y, Perry M (2014) Implementing the “wisdom of the crowd.” *J Polit Econ* 122(5):988–1012
- Li J, Netessine S (2017) Market thickness and matching (in) efficiency: evidence from a quasi-experiment. Working paper
- Lobel I, Sadler E (2015) Preferences, homophily, and social learning. *Oper Res* 64(3):564–584

- Mansour Y, Slivkins A, Syrgkanis V, Wu ZSW (2015) Bayesian exploration: incentivizing exploration in Bayesian games. In: Proceedings of the 16th ACM conference on economics and computation, Portland. ACM, pp 565–582
- Marinesi S, Girotra K, Netessine S (2017, Forthcoming) The operational advantages of threshold discounting offers. *Manag Sci* 64:2690–2708
- Moon K, Bimpikis K, Mendelson H (2017, Forthcoming) Randomized markdowns and online monitoring. *Manag Sci* 64:1271–1290
- Orlov D, Skrzypacz A, Zryumov P (2017) Persuading the principal to wait. Working paper
- Papanastasiou Y (2017) Fake news propagation and detection: a sequential model. Working paper
- Papanastasiou Y, Savva N (2017) Dynamic pricing in the presence of social learning and strategic consumers. *Manag Sci* 63(4):919–939
- Papanastasiou Y, Bimpikis K, Savva N (2017, Forthcoming) Crowdsourcing exploration. *Manag Sci* 64:1727–1746
- Rayo L, Segal I (2010) Optimal information disclosure. *J Polit Econ* 118(5):949–987
- Renault J, Solan E, Vieille N (2017) Optimal dynamic information provision. *Games Econ Behav* 104:329–349
- Seel C, Strack P (2016) Continuous time contests with private information. *Math Oper Res* 41(3):1093–1107
- Strack P (2016) Risk-taking in contests: the impact of fund-manager compensation on investor welfare. Working paper
- Swinney R (2011) Selling to strategic consumers when product value is uncertain: the value of matching supply and demand. *Manag Sci* 57(10):1737–1751
- Taylor T (2018) On-demand service platforms. *Manuf Serv Oper Manag* 20(4):704–720
- Terwiesch C, Xu Y (2008) Innovation contests, open innovation, and multiagent problem solving. *Manag Sci* 54(9):1529–1543
- Veeraraghavan S, Debo L (2009) Joining longer queues: information externalities in queue choice. *Manuf Serv Oper Manag* 11(4):543–562

Chapter 11

Design of an Aggregated Marketplace Under Congestion Effects: Asymptotic Analysis and Equilibrium Characterization



Ying-Ju Chen, Costis Maglaras, and Gustavo Vulcano

Abstract We study an aggregated marketplace where potential buyers arrive and submit requests-for-quotes (RFQs). There are n independent suppliers modeled as $M/GI/1$ queues that compete for these requests. Each supplier submits a bid that comprises of a fixed price and a dynamic target leadtime, and the cheapest supplier wins the order as long as the quote meets the buyer's willingness to pay. We characterize the asymptotic performance of this system as the demand and the supplier capacities grow large, and subsequently extract insights about the equilibrium behavior of the suppliers. We show that supplier competition results in a mixed-strategy equilibrium phenomenon that is significantly different from the centralized solution. In order to overcome the efficiency loss, we propose a compensation-while-idling mechanism that coordinates the system: each supplier gets monetary compensation from other suppliers during his idle periods. This mechanism alters suppliers' objectives and implements the centralized solution at their own will.

Y.-J. Chen (✉)

School of Business and Management & School of Engineering, The Hong Kong University of Science and Technology, Kowloon, Hong Kong
e-mail: imchen@ust.hk

C. Maglaras

Columbia Business School, New York, NY, USA
e-mail: c.maglaras@gsb.columbia.edu

G. Vulcano

School of Business, Universidad Torcuato di Tella, Buenos Aires, Argentina
e-mail: gvulcano@utdt.edu

© Springer Nature Switzerland AG 2019

M. Hu (ed.), *Sharing Economy*, Springer Series in Supply Chain Management 6,
https://doi.org/10.1007/978-3-030-01863-4_11

217

11.1 Introduction

11.1.1 Background and Motivation

Electronic markets (e-markets) have proliferated in the last two decades or so as means to efficiently aggregate supply and demand for services or goods, in an effort to reduce search and transaction costs, improve market outcomes, and benefit both participants that supply and demand services. We study a mathematical model motivated by such marketplaces for services or goods that are subject to congestion effects, manifested in terms of delays until the product is delivered. The focus is on analyzing the market dynamics and gaining insight regarding the competition across suppliers when services are produced in a make-to-order manner (that is modeled via a queueing facility) and where congestion signals are state-dependent.

As part of the dynamics of the e-market evolution, several large-scale, marketplaces have emerged. Examples include global freelancing platforms like Upwork, Guru and Freelancer, where business entities and independent workers connect and collaborate remotely. To illustrate the magnitude of the phenomenon, by 2017 Upwork has connected 12M registered freelancers with 5M registered clients to execute 3M jobs posted annually, worth a total of \$1 billion USD.

In these platforms, many small service providers (e.g., individual computer programmers) seek work orders. Customers (e.g., employers) post job descriptions in the form of RFQs, and have service providers bid on the work. Customers then look at previous ratings and work history of the different candidates before settling on either a contract rate, or a pay-per hour agreement. Generally, money is escrowed by each of the websites (intermediaries), which release the payment to the service provider when the work is completed, while skimming a commission – typically 5–20% of money that changes hands. In addition, sometimes intermediaries also charge a membership fee to the parties involved to enjoy more benefits (e.g., Freelancer offers different monthly and annual plans).

An important portion of the projects auctioned out via these markets are complex and could be better addressed via a team as opposed to an individual. To better serve and bid in such cases, many freelance workers are represented via agents that pool capacity as traditional agents would do, and also provide project management in executing the complex projects so as to best use the pooled resources. These agents aggregate the capacity of many individual freelance providers.

In these settings, customers usually require specific skill sets, quality, and timeliness from their providers, and account for these needs as well as for cost in their utility function. This multidimensional assessment can be captured by a scoring index that a customer assigns to each potential provider. The bids are ranked, and the order is then awarded to the most “desirable” service provider. In this way, the final allocation for each work order is decided based on a reverse auction. Even though multiple service attributes can be subsumed in the scoring function, two of the most relevant ones are expected delay and price. Customers visiting these marketplaces usually seek quick solutions and are willing to trade prices

with waiting times. In this regard, aggregated marketplaces like the aforementioned ones raise several interesting practical and theoretical questions. Sometimes the role of the intermediary is passive, neutral, and limits to pooling supply and demand in an exchange platform. In this case: How does the system dynamics evolve as service providers (suppliers) compete for each potential order by posting a price and a state-dependent delay estimate? How should these suppliers determine their bidding strategies? Is the market efficient under competitive behavior? What is the typical distribution of completion delays that emerges in equilibrium (e.g., what set of different delays would a randomly arriving RFQ encounter)?

When the market is inefficient, the intermediary may take a more active role in the sense of proposing a coordination scheme to align the suppliers' incentives. Is it possible to achieve this centralized optimal solution? If so, how can it be implemented in practical terms?¹ For instance, nowadays AirBnb helps the hosts to make pricing decisions (to maximize the total pie of collected revenue) and eBay helps sellers to set prices by suggesting the trending price for an item.

We make some initial progress in addressing these questions. Specifically, we introduce a stylized mathematical model to study the aggregated marketplace in settings characterized by high volume of transactions. The first goal of our study is to understand the dynamics and performance of the system, and gain insight on the pricing and capacity game among suppliers. Second, to gain insight on the dynamics of such a marketplace, specifically understand how congestion delays reduce overall demand, and how this reduction, in turn, reduces the congestion offered by the suppliers. This moderating effect of congestion stems from the fact that customers are delay sensitive, and, furthermore, that their choice of the "cheaper" supplier to send their order will reduce the workload directed to more congested suppliers. Moreover, the strategic choice made by the arriving customers will couple the congestion (or backlog) of all suppliers – something that we observe in practice. Finally, identify the coordination failure of this marketplace, if there is one, and how it manifests itself, and propose simple coordination mechanisms. Methodologically, we focus on large scale systems and develop a useful framework for studying

¹There are other applications that share the same salient feature of several firms competing in offering some type of substitutable service that is differentiated with respect to its price and delay. Perhaps one of the most pervasive comes from the US equities market, which comprises of many exchanges, such as the NYSE, NASDAQ, NYSE Arca, BATS Global Markets, etc. Exchanges typically function as electronic limit order books, operating under a "price-time" priority rule, and their high-frequency dynamics can be modeled as multi-class queueing systems. Exchanges offer a rebate to liquidity providers, i.e., traders that post limit orders that "make" markets when their orders get filled, and charge a fee to "takers" of liquidity that initiate trades using marketable orders that transact against posted limit orders. The magnitudes of these make-take fees vary across exchanges and are comparable to the spread plus a significant fraction of the overall trading costs. Exchanges often change their fees and rebates in an effort to attract liquidity. Market participants employ so called "smart order routers" that take into account real-time market data, including queue and trading rate information, and formulate an order routing problem to trade off between rebates and a notion of expected delay, fill probabilities, and/or expected adverse selection. Once again, prices are fixed but delays are state dependent.

the supplier game, which is otherwise intractable, by studying a much simpler approximating one dimensional diffusion process instead of a multi-dimensional, state dependent discrete Markov process.

In a bit more detail, we assume that potential buyers arrive according to a Poisson process and submit order requests, and that the suppliers (modeled as $M/GI/1$ queues) compete for these requests. Initially, suppliers decide the capacity (i.e., service rate) to offer, which is a static, long term decision. While operating, each supplier processes orders in a first-in-first-out manner, and submits a bid that comprises a fixed price and a target leadtime that depends on his own queue status. In fact, a key distinction of our work is that the delay quotations are dynamic rather than based on a steady-state assessment of the queue size. When suppliers submit bids, they face an economic tradeoff: a high price will lead to high revenues per order, but will reduce the total number of orders awarded, which will cause excessive idleness and implicit revenue loss; a low price will result in many awarded orders and large backlogs, that, in turn, will cause long delay quotations thus increasing the full cost of the respective bid. The arriving buyer then uses a scoring function to compute the net utility associated with her bid, and awards the order to the lowest-quote supplier in order to maximize her own surplus (provided that it is nonnegative).

11.1.2 Overview of Results

This appears to be one of the first papers to study competition in queues with substitutable products or services and state dependent congestion information. The discrete choice among substitutable products of potential consumers, the state dependent nature of the congestion signals, and the decentralized control among suppliers complicate the analysis of this system, rendering brute force analysis to be essentially intractable.²

The first contribution is to propose a tractable way for studying the decentralized market. As a preliminary step towards solving the capacity and pricing game, we analyze the queueing performance of this stochastic dynamic system assuming that the price vector is given. The solution to this problem allows suppliers to evaluate their revenues given their prices as well as the prices of the competitors, which is an essential subroutine in the equilibrium analysis of the supplier pricing game. So,

²The distinction regarding consumer choice model is important. With partially substitutable products, one could model consumer choice through a multi-product demand function, where the demand for one product depends on the price and delay of all products in a continuous manner. This is not the case with perfectly substitutable products, where demand may switch from one product to another in a discontinuous manner. For example, in a setting with two equally priced products, all consumers will select the one with the lower delay. This would increase the delay estimate of the faster option, causing all consumers to choose the alternative option. That is, small differences in price and delay, may lead to dramatic differences in demand for one of the suppliers.

given a specified price vector, our first set of results characterizes the behavior of the marketplace using an asymptotic analysis where the potential demand and the supplier processing capacities grow large simultaneously. This asymptotic analysis is motivated by the following observation: If this market were served by a unique supplier (modeled as an $M/GI/1$ queue as well), then it would be economically optimal for this supplier to set the price that induces the so-called “heavy-traffic” operating regime; i.e., rather than assuming that the system is operating in the heavy traffic regime, as is often done, this result provides a primitive economic foundation that this regime emerges naturally since it optimizes the system-wide revenues (e.g., see Besbes 2006; Maglaras and Zeevi 2003). Specifically, if Λ is the market size, then the above result states that the economically optimal price is of the form $p^* = \bar{p} + \pi/\sqrt{\Lambda}$, where \bar{p} is the price that induces full resource utilization in the absence of any congestion, and π is a constant.

With the above fact in mind, we formulate the performance analysis sub-problem in a novel way that becomes asymptotically tractable in settings with large capacities and large volume of transactions. Specifically, based on the above observation, the starting point of our analysis is to write the suppliers’ price bids as perturbations around the price \bar{p} of the form $p_i = \bar{p} + \pi_i/\sqrt{\Lambda}$, for a constant π_i , where Λ is viewed as a natural proxy for system size. Letting Λ grow large, we derive the corresponding fluid and diffusion approximations. The fluid model transient analysis is helpful in establishing an important state space collapse (SSC) result through a variation of an approach developed by Bramson (1998). The SSC property establishes that the suppliers’ dynamics are asymptotically coupled and can be described as a function of the aggregate (system-wide) workload process, and, moreover, SSC implies that a supplier is able to know his competitors’ bids by simply observing awaiting orders in his own buffer. We prove a weak convergence result of the workload process to a one-dimensional reflected Ornstein–Uhlenbeck (O-U) process, where interestingly the reflection point may be away from zero depending on the suppliers’ prices. The latter implies the nonintuitive property that the aggregate workload process can never drain even though some of the suppliers may be idling, and this happens if the suppliers differ in their pricing. The derivation of the diffusion model extends “standard” results to this setting with self-interested routing policies based on dynamic information, which is of independent interest.

Second, using the asymptotic characterization of system behavior at a fixed choice of prices and capacities, we characterize the revenue stream of each supplier using the steady state properties of the reflected O-U process. This is then used to study the resulting pricing game. We find that the pricing game does not admit a pure strategy equilibrium. We specify the structure of the supporting mixed strategy equilibria where suppliers randomize over their pricing decisions. We also prove that the second-order efficiency loss of the decentralized solution can be arbitrarily large. It is worth noting that our approximate analysis of the supplier game is internally consistent in the sense that the lower order price perturbations that essentially capture the supplier pricing game do not become unbounded, but rather stay finite. In essence, all suppliers choose to operate in the asymptotic regime we identified and used in our analysis. The framework of studying the appropriate asymptotic

formulation of the aggregated market in the context with self interested buyers and state dependent congestion information, and using the derived diffusion to study the supplier game is novel. Such problems had not been studied in the literature before, in part due to their inherent complexity, and their proposed roadmap advanced seems to be of broader interest.

Third, the discrepancy between the centralized solution and the decentralized equilibrium calls for the development of a mechanism to coordinate the marketplace, so that all suppliers would self-select to price according to the centralized solution. Our proposal relies on a transfer pricing scheme that compensates suppliers during idle periods. The existence of the intermediary in the motivating examples described above provides the natural support to implement it.

11.1.3 Literature Review

Our work touches on three related bodies of literature: (1) Economics of queues, (2) Competitive models in queueing contexts, and (3) Approximation schemes to analyze complex queueing models.

The literature that studies pricing in the context of single-server queues dates back to Naor (1969). The demand model that we consider here is inspired by Mendelson (1985): There is a single class of potential customers that arrive according to a Poisson arrival process, each having a private valuation that is an independent draw from a general distribution, and a delay sensitivity parameter that is common across all customers. Mendelson and Whang (1990) extends that model to multiple customer types, and Afèche (2013) extends it to a revenue maximization setting. In the context of queueing models with pricing and service competition, starting from the early papers by Levhari and Luski (1978) and Luski (1976), customers are commonly assumed to select their service provider on the basis of a “full cost” that consists of a fixed price plus a waiting cost. In both Levhari and Luski (1978) and Luski (1976), competition is modeled in a duopoly setting where firms operate as $M/M/1$ queueing systems. Relaxations of the early papers include Loch (1991), which studies a variant of Luski (1976) in which the providers are modeled as symmetric $M/GI/1$ systems. Lederer and Li (1997) generalizes Loch (1991) for arbitrary number of service providers. Allon and Federgruen (2007) treats the price and waiting time cost as separate firm attributes that can be traded off differently by each arriving customer. These papers focused on customers making decisions based on steady-state performance measures.

Our use of asymptotic approximations and heavy traffic analysis to study the supplier game is motivated by the results of Besbes (2006) and Maglaras and Zeevi (2003), who showed that in large scale systems the heavy traffic regime is the one induced by the revenue maximizing price. Our work implicitly assumes the validity of the heavy traffic regime in deriving its asymptotic approximation (as opposed to proving it as in the two papers above). The equilibrium pricing behavior of the competing suppliers supports the rationale of this assumption in the sense that no

supplier wishes to price in a way that would deviate from that operating regime. The derivation of our limit model makes heavy use of the work by Mandelbaum and Pats (1995) on queues with state dependent parameters, and of the framework developed by Bramson (1998) for proving state space collapse results. We also use technical results from Ata and Kumar (2005) and Williams (1998) in our analysis. However, the combination of all our model features does not fit neatly the technical requirements of the aforementioned papers, as shown in the proofs contained in the online appendix.

A queueing paper that studies a model that is similar to ours in a heavy traffic asymptotic regime is Stolyar (2005). The key differences are the following: (a) Stolyar (2005) assumes strictly convex delay cost functions as opposed to linear, (b) it does not consider pricing (or some term that could account for its effect in the routing rule), and (c) it does not allow for admission control decisions that can turn away users when the system is congested. The latter is captured by the behavior of self-interested users that differ in their valuations, and as a result will choose not join the market if the full cost exceeds their value. Taken together, these three elements necessitate a new analysis that leads to insights that differ than what was observed in Stolyar (2005). Perhaps the most notable difference is the fact that as a result of the pricing game, the workload process will not reflect at the origin, but instead it will reflect at some strictly positive quantity.

Other papers accounting for congestion pricing include DiPalantino et al. (2011), which studies two types of contractual agreements in oligopolistic service industries: service level guarantees (SLG) and best effort (BE), where firms provide the best possible service given their infrastructure. Allon and Gurvich (2010) appeals to an asymptotic analysis to study a competitive game of a queueing model, and propose a general recipe for relating the asymptotic outcome to that of the original system. They show that the pricing decisions and service level guarantees result in respectively first-order and second-order effects on the suppliers' payoffs. At a high level the approach of Allon and Gurvich (2010) is similar to ours, but the presence of real-time delay information and of perfectly substitutable lead to different models and results. However, and different from our mixed-strategy equilibrium analysis, they exploit the notion of ϵ -Nash equilibrium to explain the players' behavior. The follow-up paper Allon et al. (2012) studies a large-scale marketplace with a moderating firm and numerous service providers, but different from us the authors use a static measure of the waiting time standard (usually the expected value or some percentile of the steady state distribution). Finally, Maglaras et al. (2016) studied a mathematical model similar to ours, albeit, in a non-stationary environment, which was motivated by limit order routing in fragmented limit order book (financial equity) markets; apart from deriving a similar state space collapse result, they found empirical support for their findings in an analysis of a large scale US equities market data set.

The remainder is organized as follows. In Sect. 11.2, we describe the model details. Next, Sect. 11.3 derives the asymptotic characterization of the marketplace behavior, and Sect. 11.4 characterizes the equilibrium behavior of the supplier pricing game. Finally, Sect. 11.5 includes our concluding remarks. All the proofs and more details of the analysis can be found in the technical report Chen et al. (2008).

11.2 Model

11.2.1 Description of the Market

We consider an aggregated marketplace where a divisible homogeneous product (e.g., computer programming work) is exchanged. The market functions as follows:

Order arrivals Potential buyers arrive to the marketplace according to a Poisson arrival process with intensity Λ , and submit requests-for-quotes (RFQs). Each RFQ corresponds to the procurement of one unit of the product. Each buyer has a private valuation v for her order that is an independent draw from a general, continuously differentiable distribution $F(\cdot)$. Buyers are delay-sensitive and incur a cost c per unit of delay. Thus, buyers are homogeneous with respect to delay preferences, and heterogeneous with respect to valuations (though symmetric across the common c.d.f. $F(\cdot)$). A buyer that arrives at time t initiates a RFQ process to procure one unit of the product.

Suppliers The market is served by a set of suppliers $\mathcal{N} = \{1, \dots, n\}$. Each supplier i is modeled as an $M/GI/1$ queue with an infinite capacity buffer managed in a First-In-First-Out fashion. Service times at supplier i follow a general distribution with mean $1/\mu'_i$ and standard deviation σ_i . Let $\hat{\mu} := \sum_{i \in \mathcal{N}} \mu'_i / \Lambda$ be the (normalized) aggregated service rate of the market. We assume that the capacity vector $\mu' \equiv \{\mu'_i\}$ is common knowledge. This fact can also be sustained by information provided by intermediary entities like the ones discussed in Sect. 11.1.

Market mechanism Suppliers compete for this request by submitting bids that comprise a price p_i and a target leadtime $d_i(t)$. We assume that the price component of the bid is state-independent, i.e., supplier i always submits the same price bid p_i for all orders. The leadtime component of the bid submitted by each supplier i is state-dependent and equals the expected time it would take to complete that order; cf. Eq. 11.6 later on. We are assuming here that the supplier always submits a truthful estimate of the expected delay $d_i(t)$. In fact, in the presence of a market intermediary, the misreport of expected delays is discouraged through the display of past experiences of buyers with a given supplier, e.g. through publicly available ratings and reviews. This revealed information acts as a threatening device to favor the honest disclosure of suppliers' availabilities.

On their end, buyers are price- and delay-sensitive, and for each supplier i they associate a "full cost" given by $p_i + cd_i(t)$, where the delay sensitivity parameter c is assumed to be common for all buyers. Upon reception of the bids, the buyer awards her order to the lowest cost supplier, provided that her net utility is positive, i.e., $v \geq \min_{i \in \mathcal{N}} \{p_i + cd_i(t)\}$; otherwise, the buyer leaves without submitting any order. Whenever a tie occurs, the order is awarded by randomizing uniformly among the cheapest suppliers.³

³We could also allow other tie-breaking rules, and it can be verified that our results are not prone to the specific choice of tie-breaking rules.

Given vectors $p = (p_1, \dots, p_n)$, and $d(t) = (d_1(t), \dots, d_n(t))$, the instantaneous rate at which orders enter this aggregated market is given by

$$\lambda(p, d(t)) = \Lambda \bar{F} \left(\min_i \{p_i + c d_i(t)\} \right). \tag{11.1}$$

Focusing on the right-hand-side of the above expression, we note that the buyers' valuation distribution $F(\cdot)$ determines the nature of the aggregate demand rate function.⁴ Let $x = \min_i p_i$ and, with slight abuse of notation, write $\lambda(x)$ in place of $\lambda(p, 0)$, where 0 is the vector of zeros. We further define $\epsilon(x) = -(d\lambda(x)/dx) \cdot (x/\lambda(x))$. The expression $\epsilon(x)$ can be regarded as the price elasticity of the demand rate, as it measures the proportional change of demand rate in response to the price change. We will make the following intuitive economic assumption:

Assumption 1 *Given $x = \min_i p_i$, $\lambda(p, 0)$ is elastic in the sense that $\epsilon(x) > 1$ for all price vectors p in the set $\{p : 0 \leq \lambda(p, 0) \leq \sum_{i=1}^n \mu'_i\}$.*

The above assumption implies that in the absence of delays, a decrease in the minimum price would result in an increase in the market-wide aggregated revenue rate $p \cdot \lambda(p, 0)$.⁵ Of course, this would increase the utilization levels of the suppliers and lead to increased congestion and delays, thereby moderating the aggregate arrival rate $\lambda(p, d(t))$.

Let $A(t)$ be the cumulative number of orders awarded to all the competing suppliers up to time t ,

$$A(t) = N \left(\Lambda \int_0^t \bar{F} \left(\min_{i \in \mathcal{N}} \{p_i + c d_i(s)\} \right) ds \right), \tag{11.2}$$

where $N(t)$ is a unit rate Poisson process and the equality holds only in distribution. To represent the cumulative number of orders for each individual supplier, we first define

$$\mathcal{J}(t) \equiv \{i \in \mathcal{N} : p_i + c d_i(t) \leq p_j + c d_j(t), \forall j \in \mathcal{N}\} \tag{11.3}$$

as the set of cheapest suppliers at time t . Further, define $\Xi_{\mathcal{J}(t)}$ as the random variable that assigns the orders uniformly amongst the cheapest suppliers. That is, $\Xi_{\mathcal{J}(t)} = i$ with probability $1/|\mathcal{J}(t)|$ if $i \in \mathcal{J}(t)$, where $|\mathcal{J}(t)| > 0$ is

⁴For example, if $v \sim U[0, \Lambda/\alpha]$, then the demand function is linear, of the form $\lambda(x) = \Lambda - \alpha x$, where $x = \min_i \{p_i + c d_i(t)\}$; if $v \sim \exp(\alpha)$, then the demand is exponential, with $\lambda(x) = \Lambda e^{-\alpha x}$.

⁵This implication follows directly from the economics literature. When $\epsilon(x) > 1$, the proportional increase of demand rate is larger than the proportional decrease of price. As the revenue is the product of price and demand rate, the aggregated revenue rate ends up being higher. Suppose further that there are no congestion effects and that there exists a central planner who could select a common price p and an aggregate capacity $\hat{\mu}$. Under a linear capacity cost $h\hat{\mu}$ and an arrival rate Λ , the solution to the problem $\max_{p, \hat{\mu}} \{p\lambda(p, 0) - h\Lambda\hat{\mu} : 0 \leq \lambda(p, 0) \leq \Lambda\hat{\mu}\}$ results in a capacity decision that satisfies the above assumption.

the cardinality of $\mathcal{J}(t)$, and $\mathcal{E}_{\mathcal{J}(t)} = i$ with zero probability otherwise. For the ease of the exposition, we will assume that $\mathcal{J}(t)$ and $\mathcal{E}_{\mathcal{J}(t)}$ are defined as continuous processes for all $t \geq 0$, even if no actual arrival occurs at that time. This allows us to write the cumulative number of orders awarded to supplier i , denoted by $A_i(t)$, as

$$A_i(t) = \int_0^t \mathbf{1}\{\mathcal{E}_{\mathcal{J}(s)} = i\} dA(s), \tag{11.4}$$

where $\mathbf{1}\{\cdot\}$ is the indicator function. Note also that $A(t) = \sum_{i \in \mathcal{N}} A_i(t)$.

Supplier dynamics Let $Q_i(t)$ denote supplier i 's number of jobs in the system (i.e., in queue or in service) at time t , and $T_i(t)$ denote the cumulative time that supplier i has devoted into producing orders up to time t , with $T_i(0) = 0$. Let $Y_i(t)$ denote the idleness incurred by supplier i up to time t . Note that $T_i(t) + Y_i(t) = t$ for each supplier i ; moreover, $Y_i(t)$ can only increase at a time t when the queue $Q_i(t)$ is empty. Let $S_i(t)$ be the number of supplier i 's service completions when working continuously during t time units, and $D_i(t) = S_i(T_i(t))$ be the cumulative number of departures up to time t . The production dynamics at supplier i are summarized in the expression:

$$Q_i(t) = Q_i(0) + A_i(t) - D_i(t). \tag{11.5}$$

Given this notation, then

$$d_i(t) = \frac{Q_i(t) + 1}{\mu'_i}, \tag{11.6}$$

is the expected sojourn time of the new incoming order, given that it gets awarded to supplier i and the current queue length is $Q_i(t)$. Under our modeling assumptions, supplier i will therefore bid $(p_i, d_i(t))$, where $d_i(t)$ is given by Eq. 11.6. Each supplier knows his own system queue length $Q_i(t)$, but is not informed about his competitors' queue lengths.

11.2.2 Problems to Address

We study three problems related to the market model described above:

1. *Performance analysis for a given p and μ'* : Given a fixed price vector p and a vector of processing capacities μ' , the first task is to characterize the system performance, i.e., to characterize the behavior of the queue length processes $Q_i(t)$ at each supplier, and calculate the resulting revenue streams for each supplier. A supplier's long-run average revenue is

$$\Omega_i(p_i, p_{-i}) \equiv p_i \cdot \lim_{t \rightarrow \infty} \frac{S_i(T_i(t))}{t}, \tag{11.7}$$

where $p_{-i} \equiv (p_1, \dots, p_{i-1}, p_{i+1}, \dots, p_n)$ denotes other suppliers' price decisions. Our goal, therefore, is to analyze the performance of relevant system dynamics that leads to a tractable representation of these long-run average revenues.

2. *Characterization of market equilibrium:* The above problem serves as an input to study the competitive equilibrium that characterizes the supplier capacity and pricing games, both of which are one-shot games where each supplier selects his service rate and static price, successively. We further show that the capacity selections constitute the first-order effects on the suppliers' payoffs and the pricing decisions are of second order; thus, we can conveniently decouple the equilibrium analysis into two separate stages. For the first-stage, capacity game, since the capacities (service rates) are assumed to be publicly observable, we adopt the Nash equilibrium as our solution concept. For the second-stage, pricing game, the suppliers may be uninformed about the queue lengths of the competitors, which may potentially lead to information incompleteness. However, as we will show, the suppliers' competitive behavior is insensitive to this knowledge; consequently, we adopt again the standard Nash equilibrium (under complete information) as our solution concept. Given the revenue specified in Eq. 11.7, a Nash equilibrium $\{p_i^*\}$ requires that $p_i^* = \arg \max_{p_i} \Omega_i(p_i, p_{-i}^*), \forall i \in \mathcal{N}$.
3. *Market efficiency and market coordination:* Our objective here is to characterize the efficiency loss:

$$\max_p \left\{ \sum_{i \in \mathcal{N}} \Omega_i(p_i, p_{-i}) \right\} - \sum_{i \in \mathcal{N}} \Omega_i(p_i^*, p_{-i}^*), \quad (11.8)$$

i.e., the difference between the revenue of a system where a central planner would control the pricing decision of each supplier (i.e., the *first best solution*), and the sum of the revenues collected in the competitive framework. If the market equilibrium is inefficient, we would like to specify a simple market mechanism that coordinates the market and achieves the first best solution identified above. Such a mechanism could specify, for example, the rules according to which orders are allocated and payments are distributed among the suppliers.

11.3 Asymptotic Analysis of Marketplace Dynamics

This section focuses on the first problem described in Sect. 11.2. Despite the relatively simple structure of the suppliers' systems and the customer/supplier interaction, it is still fairly hard to study their dynamics due to the state-dependent delay quotations and the dynamic allocation of orders. Our approach is to develop an approximate model for the market dynamics that is rigorously validated in settings where the demand volume and processing capacities of the various suppliers are large. In this regime, the market and supplier dynamics simplify significantly, and

are essentially captured through a tractable one-dimensional diffusion process. This limiting model provides insights about the structural properties of this market, and provides a vehicle within which we are able to analyze the supplier game and the emerging market equilibrium. This is pursued in the next section.

11.3.1 Background: Revenue Maximization for an $M/M/1$ Monopolistic Supplier

As a motivation for our subsequent analysis, this subsection will summarize some known results regarding the behavior of a monopolistic supplier modeled as an $M/M/1$ queue that offers a product to a market of price and delay sensitive customers. The supplier posts a static price and dynamically announces the prevailing (state-dependent) expected sojourn time for orders arriving at time t , which is given by $d(t) = (Q(t) + 1)/\mu$. The assumptions on the customer purchase behavior are those described in the previous section. Given p and $d(t)$, the instantaneous demand rate into the system at time t is given by $\lambda(t) = \Lambda \bar{F}(p + cd(t))$. The supplier wants to select p to maximize his long-run expected revenue rate.

It is easy to characterize the structure of the revenue maximizing solution in settings where the potential market size Λ and the processing capacity μ grow large. Specifically, we will consider a sequence of problem instances indexed by r , where $\Lambda^r = r$ and $\mu^r = r\mu$; that is, r denotes the size of the market. The characteristics of the potential customers, namely their price sensitivity c and valuation distribution $F(\cdot)$, remain unchanged along this sequence. Let $\hat{p} = \arg \max p \bar{F}(p)$, and \bar{p} be the price such that relation $\Lambda^r \bar{F}(\bar{p}) = \mu^r$ holds. That is, neglecting congestion effects, \hat{p} is the price that maximizes the revenue rate and \bar{p} is the price that induces full resource utilization, and both of these quantities are independent of r . Assumption 1 implies that $\hat{p} < \bar{p}$ (or equivalently, $\Lambda^r \bar{F}(\hat{p}) > \mu^r$) thus accentuating the tension between revenue maximization and the resulting congestion effects. Besbes (2006) showed that the revenue maximizing price, denoted by $p^{*,r}$, is of the form

$$p^{*,r} = \bar{p} + \pi^*/\sqrt{r} + o(1/\sqrt{r}), \quad (11.9)$$

where π^* is a constant independent of r . Moreover, the resulting queue lengths $Q^r(t)$ are of order \sqrt{r} , or in a bit more detail, the normalized queue length process $\bar{Q}^r(t) = Q^r(t)/\sqrt{r}$ has a well defined stochastic process limit as $r \rightarrow \infty$. Since the processing time is itself of order $1/r$, the resulting delays are of order $1/\sqrt{r}$. Following Maglaras and Zeevi (2003), the delays are moderate in absolute terms (of order $1/\sqrt{r}$) but significant when compared to the actual service time (of order $1/r$). If the supplier can select the price and capacity μ , the latter assuming a linear capacity cost, then the optimal capacity choice is indeed such that $\hat{p} < \bar{p}$ (where \bar{p} is determined by μ), i.e., making the above regime the “interesting” one to consider. Finally, we note that the above results also hold for the case of generally distributed service times (Besbes 2006).

11.3.2 Setup for Asymptotic Analysis

Given a set of suppliers characterized by their prices and capacities p_i, μ'_i , we propose the following approximation:

1. Define the normalized parameters $\mu_i = \mu'_i/\Lambda$ for every supplier i .
2. Define \bar{p} to be the price such that $\Lambda\bar{F}(\bar{p}) = \sum_{i \in \mathcal{N}} \mu'_i$. Define $\pi_i = \sqrt{\Lambda}(p_i - \bar{p})$ so that the prices p_i can be represented as $p_i = \bar{p} + \pi_i/\sqrt{\Lambda}$.
3. Embed the system under consideration in the sequence of systems indexed by r and defined through the sequence of parameters:

$$\Lambda^r = r, \quad \mu_i^r = r\mu_i, \quad \forall i \in \mathcal{N}, \quad c^r = c, \quad v \sim F(\cdot), \quad (11.10)$$

and prices given by $p_i^r = \bar{p} + \pi_i/\sqrt{r}$ for all i .

Given the preceding discussion, one would expect that the market may operate in a manner that induces almost full resource utilization, and where the underlying set of prices takes the form assumed in item 3 above. This would be true if the market were managed by a central planner that could coordinate the supplier pricing and capacity decisions. The approach we pursue is to embed the system we wish to study in the sequence of systems indexed by r and described in Eq. 11.10, and subsequently approximate the performance of the original system with that of a limit system that is obtained as $r \rightarrow \infty$, which is more tractable. Note that for $r = \Lambda$ in Eq. 11.10, where Λ denotes the market size of the potential order flow as described in the previous section, we recover the exact system we wish to study. If Λ is sufficiently large, then the proposed approximation is expected to be fairly accurate.

The remainder of this section derives an asymptotic characterization of the performance of a market that operates under a set of parameters (p, μ') that are embedded in the sequence Eq. 11.10.

11.3.3 Transient Dynamics via a Fluid Model Analysis

The derivation of the asymptotic limit model (specifically, Proposition 2) will show that the following set of equations

$$\bar{Q}_i(t) = \bar{Q}_i(0) + \bar{A}_i(t) - \bar{D}_i(t), \quad \forall i \in \mathcal{N}, \quad (11.11)$$

$$\bar{A}_i(t) = \int_0^t \frac{1}{|\mathcal{J}(s)|} \mathbf{1}\{i \in \mathcal{J}(s)\} \Lambda\bar{F}(\bar{p}) ds, \quad \forall i \in \mathcal{N}, \quad (11.12)$$

$$\bar{D}_i(t) = \mu_i \bar{T}_i(t), \quad \forall i \in \mathcal{N}, \quad (11.13)$$

$$\int_0^t \bar{Q}_i(s) d\bar{Y}_i(s) = 0, \quad \forall i \in \mathcal{N}, \quad (11.14)$$

$$\bar{T}_i(t) + \bar{Y}_i(t) = t, \quad \forall i \in \mathcal{N}, \quad (11.15)$$

$$\bar{W}(t) = \sum_{i \in \mathcal{N}} \frac{\bar{Q}_i(t)}{\mu_i}. \quad (11.16)$$

captures the market's transient dynamics over short periods of length $1/\sqrt{r}$. This subsection studies the transient evolution of Eqs. 11.11, 11.12, 11.13, 11.14, 11.15, and 11.16 starting from arbitrary initial conditions.

In Eqs. 11.11, 11.12, 11.13, 11.14, 11.15, and 11.16, the processes \bar{Q} , \bar{A} , \bar{D} , \bar{T} , \bar{Y} are the fluid analogues of Q , A , D , T , Y defined in Sect. 11.2, and \bar{W} is the fluid analog of the system workload W . Equation 11.11 keeps track of the queue sizes. Equation 11.12 indicates how the arrivals are routed to these servers: An arrival walks away if her valuation is sufficiently low; otherwise, she joins server i based on the routing rule specified in Sect. 11.2. From Eq. 11.12, the orders get awarded to the various suppliers at a rate $\Lambda \bar{F}(\bar{p}) = \sum_{i \in \mathcal{N}} \mu'_i$, i.e., $\bar{F}(\bar{p}) = \sum_{i \in \mathcal{N}} \mu_i$ (as indicated by the aggregate counting process $N(\bar{F}(\bar{p})t)$). Equation 11.14 demonstrates the non-idling property: $\bar{Y}_i(t)$ cannot increase unless $\bar{Q}_i(t) = 0$. Equation 11.15 is a time-balance constraint. Finally, Eq. 11.16 establishes the connection between the total workload and the queue lengths.

The next proposition establishes that starting from any arbitrary initial condition, the transient evolution of the market (as captured through Eqs. 11.11, 11.12, 11.13, 11.14, 11.15), and 11.16 converges to a state configuration where all suppliers are equally costly in terms of the full cost of the bids given by (price + $c \times$ delay). This is, of course, a consequence of the market mechanism that awards orders to the cheapest supplier(s), until their queue lengths build up so that their full costs become equal. Simultaneously, expensive suppliers do not get any new orders and therefore drain their backlogs until their costs become equal to that of the cheapest suppliers. From then onwards, orders are distributed in a way that balances the load across suppliers. This result is robust with respect to the tie-breaking rule that one may apply when multiple suppliers share the same full cost.

Proposition 1 *Let sextuple $(\bar{Q}, \bar{A}, \bar{D}, \bar{T}, \bar{Y}, \bar{W})$ be the solution to Eqs. 11.11, 11.12, 11.13, 11.14, 11.15, and 11.16 with $\max\{|\bar{Q}(0)|, |\bar{W}(0)|\} \leq M_0$ for some constant M_0 . Then for all $\delta > 0$, there exists a continuous function $s(\delta, M_0) < \infty$ such that*

$$\max_{i, j \in \mathcal{N}} \left| \left(\pi_i + c \frac{\bar{Q}_i(s)}{\mu_i} \right) - \left(\pi_j + c \frac{\bar{Q}_j(s)}{\mu_j} \right) \right| < \delta, \quad \forall s > s(\delta, M_0). \quad (11.17)$$

11.3.4 State-Space Collapse and the Aggregate Marketplace Behavior

The next result shows that the transients studied above appear instantaneously in the natural time scale of the system, and as such the marketplace dynamics evolve as if all suppliers are equally costly at all times.

We use the superscript r to denote the performance parameters in the r -th system, e.g., $A_i^r(t)$, $S_i^r(t)$, $T_i^r(t)$, and $Q_i^r(t)$. The (expected) workload (i.e., the time needed to drain all current pending orders across all suppliers) is defined as $W^r(t) = \sum_{i \in \mathcal{N}} (Q_i^r(t)/\mu_i^r)$.

Motivated by the discussion in Sect. 11.3.1, we will optimistically assume (and later on validate) that the supplier queue lengths are of order \sqrt{r} , and accordingly define the re-scaled queue length processes for all suppliers according to

$$\tilde{Q}_i^r(t) = \frac{Q_i^r(t)}{\sqrt{r}}. \quad (11.18)$$

The corresponding re-scaled expected workload process is given by $\tilde{W}^r(t) = \sqrt{r}W^r(t) = \sum_{i \in \mathcal{N}} (\tilde{Q}_i^r(t)/\mu_i)$.

Define $\tilde{Z}^r(t) = \tilde{\pi} + \bar{c}\tilde{W}^r(t)$, where $\tilde{\pi} = \sum_{i \in \mathcal{N}} \pi_i/n$ and $\bar{c} = c/n$. $\tilde{Z}^r(t)$ can be regarded as a proxy for the average of the second-order terms of suppliers' bids since $\tilde{Z}^r(t) = (1/n) \sum_{i \in \mathcal{N}} (\pi_i + c(\tilde{Q}_i^r(t)/\mu_i))$. Note that the first-order term, \bar{p} , is common for all suppliers, and can be omitted while comparing suppliers' bids.

Proposition 2 (State Space Collapse) *Suppose $\pi_i + c(\tilde{Q}_i^r(0)/\mu_i) = \tilde{\pi} + \bar{c}\tilde{W}^r(0)$ in probability, $\forall i \in \mathcal{N}$. Then, for all $\tau > 0$, for all $\epsilon > 0$, as $r \rightarrow \infty$,*

$$\mathbb{P} \left\{ \sup_{0 \leq t \leq \tau} \max_{i, j \in \mathcal{N}} \left| \left(\pi_i + c \frac{\tilde{Q}_i^r(t)}{\mu_i} \right) - \left(\pi_j + c \frac{\tilde{Q}_j^r(t)}{\mu_j} \right) \right| > \epsilon \right\} \rightarrow 0,$$

$$\mathbb{P} \left\{ \sup_{0 \leq t \leq \tau} \max_{i \in \mathcal{N}} \left| \left(\pi_i + c \frac{\tilde{Q}_i^r(t)}{\mu_i} \right) - \tilde{Z}^r(t) \right| > \epsilon \right\} \rightarrow 0.$$

The proof applies the ‘‘hydrodynamic scaling’’ framework of Bramson (1998), which is introduced in the context of studying the heavy-traffic asymptotic behavior of multi-class queueing networks. Our model falls outside the class of problems studied in Bramson (1998), but as we show in the online appendix, his analysis can be extended to address our setting in a fairly straightforward manner.

11.3.5 Limit Model and Discussion

Proposition 2 shows that the supplier behavior can be inferred by analyzing an appropriately defined one-dimensional process $\tilde{Z}^r(t)$ that is related to the aggregated market workload. This also implies that although each supplier only observes his own backlog, he is capable of inferring the backlog (or at least the full cost) of all other competing suppliers.

The next theorem characterizes the limiting behavior of the one-dimensional process $\tilde{Z}^r(t)$, and as a result also those of $\tilde{W}^r(t)$ and $\tilde{Q}^r(t)$.

Theorem 1 (Weak Convergence) *Suppose $\pi_i + c(\tilde{Q}_i^r(0)/\mu_i) = \bar{\pi} + \bar{c}\tilde{W}^r(0)$ in probability, $\forall i \in \mathcal{N}$. Then $\tilde{Z}^r(t)$ weakly converges to a reflected Ornstein-Uhlenbeck process $\tilde{Z}(t)$ that satisfies*

$$\tilde{Z}(t) = \tilde{Z}(0) - \gamma c \int_0^t \tilde{Z}(s) ds + \tilde{U}(t) + \frac{c\sqrt{\sigma^2 + \hat{\mu}}}{\hat{\mu}} B(t), \tag{11.19}$$

where $B(t)$ is a standard Brownian motion, $\tilde{U}(0) = 0$, $\tilde{U}(t)$ is continuous and nondecreasing, and $\tilde{U}(t)$ increases only when $\tilde{Z}(t) = \hat{\pi}$. The parameters are $\gamma = f(\bar{p})/\bar{F}(\bar{p})$, and $\sigma \equiv \sqrt{\sum_{i \in \mathcal{N}} \sigma_i^2}$. In addition, $\tilde{W}^r(t) \Rightarrow (1/\bar{c})(\tilde{Z}(t) - \bar{\pi})$, and $\tilde{Q}_i^r(t) \Rightarrow (\mu_i/c)(\tilde{Z}(t) - \pi_i)$, $\forall i \in \mathcal{N}$.

The process $\tilde{U}(t)$ is the limiting process of $\tilde{U}^r(t) \equiv (c/\hat{\mu}) \sum_{i \in \mathcal{N}} \mu_i \tilde{Y}_i^r(t)$ (defined in the proof of Theorem 1), which can be regarded as the aggregate market idleness of the system.

This theorem characterizes the limiting marketplace behavior under a given price vector p . The market exhibits a form of “resource pooling” across suppliers. Given that $\tilde{Z}(t) \geq \hat{\pi}$, it follows that $\tilde{W}(t) \geq (1/\bar{c}) \max_{i \in \mathcal{N}} (\pi_i - \bar{\pi}) := \zeta$. This says that unless all the suppliers submit the same price bid, the aggregate workload in the marketplace will always be strictly positive and at a given time t , some suppliers will never incur any idleness. The intuition for this result is the following. When the queue of the most expensive supplier(s) gets depleted, and this supplier(s) starts to idle, the imbalance between the aggregate arrival rate and service rate force suppliers to build up their queue lengths instantaneously. Consequently, suppliers that price below $\hat{\pi}$ never deplete their queue lengths asymptotically. γ , that controls the speed of the reversion of the aggregate workload process, is extracted from the customer valuation distribution. γ measures the sensitivity of the demand function to changes to the full price “ $\pi + cd(t)$ ”, and it is proportional to the demand elasticity at \bar{p} .

To summarize, in the limit model, the suppliers’ queue length processes follow from $\tilde{Q}_i(t) = (\mu_i/c)(\tilde{Z}(t) - \pi_i)$, $\forall i \in \mathcal{N}$, where $\tilde{Z}(t)$ is defined through Eq. 11.19. The next section will use this result as an input to study the suppliers’ pricing game.

11.3.6 A Numerical Example

To demonstrate the system dynamics, we consider a system with two $M/M/1$ servers, delay sensitivity parameter $c = 0.5$, and arrival rate of buyers $\Lambda = 1$. The valuation v of each customer is assumed to follow an exponential distribution with mean 0.1. The aggregate and individual service rates are respectively $\hat{\mu} = e^{-1.3}$, $\mu_1 = 0.8\hat{\mu}$, and $\mu_2 = 0.2\hat{\mu}$. Moreover, suppose the price parameters are $\pi_1 = -1$, $\pi_2 = -2$, and therefore $\bar{\pi} = (\pi_1 + \pi_2)/2 = -1.5$.

We run simulations using Arena (a discrete-event simulation software) and at places supplement it with Matlab to compute the relevant parameters.

Fig. 11.1 An instance of workload process when $r = 30$

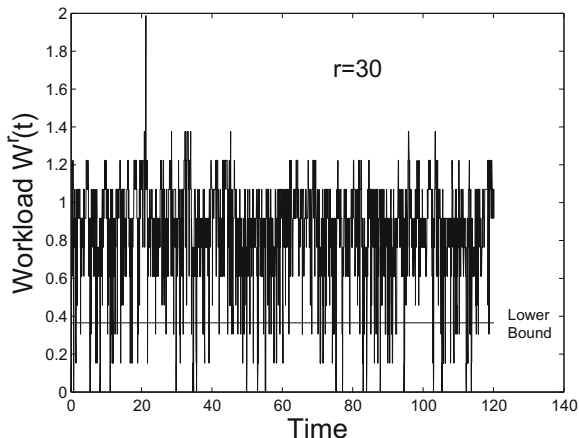
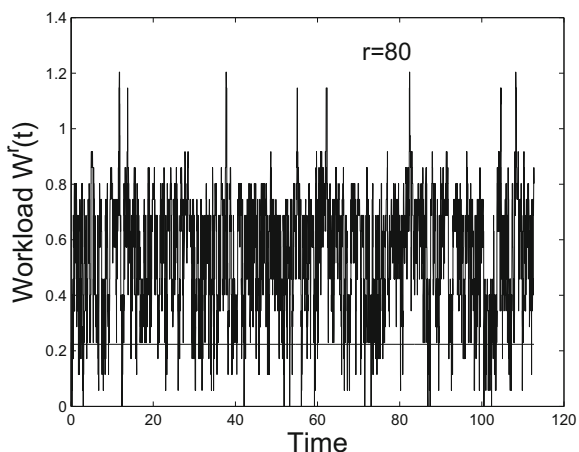


Fig. 11.2 An instance of workload process when $r = 80$



In Figs. 11.1 and 11.2, we illustrate the workload trajectories for $r = 30$ and $r = 80$, respectively, for one replication. Given these parameters, the respective boundaries are 0.365 and 0.224 for $r = 30$ and $r = 80$. Note that our mathematical statement is established on the steady-state workload process. We can either perform a very long run (say with 600,000 arrivals in expectation) and break each output record from the (single) run into a few large batches. Alternatively, we can run many replications and identify appropriate warm-up and run-length times (for example 3,000 replications, each of which generates roughly 200 arrivals). We find that both approaches lead to very similar statistical outcomes, and hence we report only the former. When $r = 80$ and the simulation time is 600,000/80 units, we find that 92.93% of time the system workload is above the boundary 0.224. When $r = 30$ (and the corresponding simulation time is 600,000/30 units), however, this proportion goes down to 71.22%. This suggests that our mathematical result of asymptotically negligible time below boundary is more applicable when the scaling factor is large.

11.4 Competitive Behavior and Market Efficiency

In this section, we discuss the equilibrium behavior of suppliers in the capacity and pricing games described in Sect. 11.2. We use the performance characterization of Sect. 11.3 and cast the suppliers' prices as "small" deviations around the market clearing price \bar{p} . This distinction results in a first-stage capacity game that affects the first-order revenues, and a second-stage, pricing game, that adjusts prices around the first-order price. Then, we briefly discuss the centralized solution that maximizes the aggregate payoffs. Next, we characterize the non-cooperative behavior of suppliers under the competitive environment. Finally, we propose a coordination scheme that achieves the aggregate payoff under the centralized solution, and describe how this coordination scheme can be implemented in the original system.

11.4.1 Suppliers' First-Order Payoffs and the Capacity Game

Let $R_i^r(t)$ denote supplier i 's cumulative revenue. Since the pricing is static, supplier i earns $R_i^r(t) = (\bar{p} + \pi_i/\sqrt{r})S_i^r(t - Y_i^r(t))$, where $S_i^r(\cdot)$ is the counting service completion process, and $Y_i^r(t)$ is the cumulative idleness process for supplier i . The next lemma shows that the "first-order" revenues of the suppliers only depend on the first-order price term \bar{p} and the service rates $\{\mu_i\}$'s.

Lemma 1 $R_i^r(t)/r \rightarrow \bar{p}\mu_i t$, as $r \rightarrow \infty$, $\forall i \in \mathcal{N}$.

Lemma 1 demonstrates that the capacity choice (μ_i) has a first-order effect on the suppliers' revenues. If we attach an appropriate capacity cost $c_i(\mu_i)$ to the suppliers, the capacity game can be explicitly posted. In the centralized system, a central planner decides the service rates (capacities) to maximize the net revenue:

$$\max_{\{\mu_i\}} \left\{ \sum_i \mu_i \bar{p} (\sum_i \mu_i) - \sum_i c_i(\mu_i) \right\}. \tag{11.20}$$

This centralized solution can be obtained via a two-stage problem in which we first find the optimal allocation for each individual to minimize the aggregate cost:

$$C(\hat{\mu}) = \min_{\{\mu_i\}} \left\{ \sum_i c_i(\mu_i), \text{ s.t. } \sum_i \mu_i = \hat{\mu}, \mu_i \geq 0, \forall i \right\}, \tag{11.21}$$

and then optimize over the aggregate service rate via $\max_{\hat{\mu}} \{\hat{\mu} \bar{p}(\hat{\mu}) - C(\hat{\mu}) \mid \hat{\mu} \geq 0\}$.

In a Nash equilibrium $\{\mu_i^*\}$'s, supplier i chooses a capacity such that

$$\mu_i^* = \arg \max_{\mu_i} \left\{ \mu_i \bar{p} (\mu_i + \sum_{j \neq i} \mu_j^*) - c_i(\mu_i) \right\}, \tag{11.22}$$

where $\bar{p} = (\bar{F})^{-1}(\sum_i \mu_i/\lambda)$ is the price that induces the full resource utilization. Furthermore, if $\hat{\mu} \bar{p}(\hat{\mu})$ is concave in $\hat{\mu}$ and $c_i(\mu_i)$ is convex in μ_i , $\forall i$,⁶ there exists a pure-strategy Nash equilibrium in the capacity game. Based on this existence result, we can characterize the pure-strategy equilibrium from the best responses of suppliers against others' strategies. Specifically, a Nash equilibrium $\{\mu_i^*\}_{i=1}^n$ satisfies

$$\mu_i^* \bar{p}'\left(\sum_i \mu_i^*\right) + \bar{p}\left(\sum_i \mu_i^*\right) - c'_i(\mu_i^*) = 0, \quad \forall i. \tag{11.23}$$

It can be verified that in the decentralized (Nash) equilibrium, each supplier intends to build a capacity μ_i^* higher than the centralized solution. This over-investment result follows from the ignorance of the negative externality a supplier brings to the entire system, as a supplier may benefit from over-investment since this allows him to capture a higher market share. This is reminiscent of the demand-stealing effect in the classical Cournot competition.

11.4.2 Suppliers' Second-Order Payoffs and the Pricing Game

To study the suppliers' pricing game we will focus on the second order correction around $R_i^r(t)$ defined as $r_i^r(t) \equiv (1/\sqrt{r})(R_i^r(t) - r \bar{p} \mu_i t)$, $\forall i \in \mathcal{N}$. The limiting processes of these corrected terms are characterized in the following lemma.

Lemma 2 $r_i^r(t) \Rightarrow r_i(t)$, as $r \rightarrow \infty, \forall i \in \mathcal{N}$, where $r_i(t) := \mu_i \pi_i t + \bar{p} \sigma_i B_{s,i}(t) - \mu_i \bar{p} \tilde{Y}_i(t)$ and $\tilde{Y}_i(t)$ is the limiting process of $\tilde{Y}_i^r(t)$ as $r \rightarrow \infty$.

Instead of using the revenue functions $\{\Omega_i(p_i, p_{-i})\}$'s defined in (11.7), we will study the suppliers' pricing game based on their (second-order) revenues given by

$$\Psi_i(\pi_i, \pi_{-i}) \equiv \lim_{t \rightarrow \infty} \frac{r_i(t)}{t} = \mu_i(\pi_i - \bar{p} \mathbb{E}[\tilde{Y}_i(\infty)]), \tag{11.24}$$

where $\pi_{-i} \equiv (\pi_1, \dots, \pi_{i-1}, \pi_{i+1}, \dots, \pi_n)$, and (with some abuse of notation) $\mathbb{E}[\tilde{Y}_i(\infty)] := \lim_{t \rightarrow \infty} (\tilde{Y}_i(t)/t)$. Define h_i as the (steady-state) proportion of the market idleness incurred by supplier i , i.e.,

$$\mathbb{E}[\tilde{Y}_i(\infty)] = h_i \mathbb{E}[\tilde{U}(\infty)], \tag{11.25}$$

⁶These assumptions are commonly adopted in revenue management, although in some context the arrival rate is used instead of the service rate, which makes no difference in heavy traffic regime, see e.g. Gallego and van Ryzin (1994).

where we again denote by $\mathbb{E}[\tilde{U}(\infty)] := \lim_{t \rightarrow \infty} (\tilde{U}(t)/t)$ the long-run average of the aggregate idleness $\tilde{U}(t)$ specified in Theorem 1.

Dividing Eq. 11.19 by t and letting $t \rightarrow \infty$, we obtain that

$$\mathbb{E}[\tilde{U}(\infty)] = \lim_{t \rightarrow \infty} \frac{\tilde{U}(t)}{t} = \gamma c \mathbb{E}[\tilde{Z}(\infty)] = \gamma c \beta \frac{\phi(\hat{\pi}/\beta)}{1 - \Phi(\hat{\pi}/\beta)}, \tag{11.26}$$

where

$$\beta = \sqrt{\frac{c(\hat{\mu} + \sigma^2)}{2\gamma\hat{\mu}^2}}, \tag{11.27}$$

and the closed-form expression follows from the fact that the reflected Ornstein–Uhlenbeck process $\tilde{Z}(t)$ has the stationary distribution as a truncated Normal random variable (Browne and Whitt 2003, Proposition 1).⁷ We do not derive the closed-form expressions of $\{h_i\}$'s since they are not needed for our equilibrium analysis.

Let $J = \{j \mid \pi_j = \hat{\pi}\}$, where $\hat{\pi} \equiv \max_{i \in \mathcal{N}} \pi_i$, denote the set of the most expensive suppliers (allowing for ties). From Theorem 1 we have that $\tilde{Z}(t) \geq \hat{\pi}$ for all $t \geq 0$ and that $\tilde{Q}_j^r(t) \Rightarrow (\mu_j/c)(\tilde{Z}(t) - \pi_j) > 0$, for all $t \geq 0$, $\forall j \notin J$. It follows that $\tilde{Y}_i(t) = 0$ for all $t \geq 0$, and therefore $h_j = 0$, $\forall j \notin J$. Given Eqs. 11.24 and 11.25, we obtain the suppliers' second-order long-run average revenue functions as follows:

$$\Psi_i(\pi_i, \pi_{-i}) = \begin{cases} \mu_i \pi_i - \mu_i \bar{p} h_i \gamma c \beta \frac{\phi(\hat{\pi}/\beta)}{1 - \Phi(\hat{\pi}/\beta)}, & \text{if } i \in J, \\ \mu_i \pi_i, & \text{otherwise.} \end{cases} \tag{11.28}$$

11.4.3 Centralized System Performance

In the centralized version of the system, a central planner makes the price decisions $\pi \equiv (\pi_1, \dots, \pi_n)$ in order to maximize the total aggregated revenue:

$$\max_{\pi} \left\{ \sum_{i \in \mathcal{N}} \mu_i \pi_i - \bar{p} \gamma \hat{\mu} \beta \frac{\phi(\hat{\pi}/\beta)}{1 - \Phi(\hat{\pi}/\beta)} : \pi_i \leq \hat{\pi} \right\}, \tag{11.29}$$

where we have applied $\sum_{i \in J} \mu_i h_i = \hat{\mu}/c$ to combine all the penalties imposed on the most expensive suppliers.

⁷We can verify that using their notation, the $\tilde{Z}(t)$ process corresponds to the following parameters: $a = \gamma c$, $m = 0$, and the process has only a left reflecting barrier $\hat{\pi}$.

For convenience, we define $\mathcal{L}(\hat{\pi}) \equiv \bar{p}\gamma\hat{\mu}\beta\phi(\hat{\pi}/\beta)/[1 - \Phi(\hat{\pi}/\beta)]$ as the revenue loss that the system suffers if $\hat{\pi}$ is the highest price offered. The optimal pricing decisions are summarized in the following lemma:

Lemma 3 *In a centralized system, all prices π_i 's are equal. The optimal static price is $\pi^C := \arg \max_{\pi} [\hat{\mu}\pi - \mathcal{L}(\pi)]$.*

11.4.4 Competitive Equilibrium

In a decentralized (competitive) system, each supplier is maximizing his own payoff, $\Psi_i(\pi_i, \pi_{-i})$, in a non-cooperative way: $\max_{\pi_i} \Psi_i(\pi_i, \pi_{-i})$. Recalling the definition of $\mathcal{L}(\hat{\pi})$, we can rewrite the supplier's payoff in Eq. 11.28 as

$$\Psi_i(\pi_i, \pi_{-i}) = \mu_i\pi_i - \mu_i h_i \frac{c}{\hat{\mu}} \mathcal{L}(\hat{\pi}) \mathbf{1}\{i \in J\} \quad (11.30)$$

In the following we characterize the equilibrium behavior. We will split our discussion in two cases, depending on whether suppliers are endowed with homogeneous or heterogeneous service rates.

11.4.4.1 Homogeneous Service Rate Case

We first consider the case where the service rates are the same across suppliers, i.e., $\mu_i = \mu_j \equiv \mu, \forall i, j \in \mathcal{N}$, and focus on symmetric equilibria. Define $\pi^* = \arg \max_{\pi} [\mu\pi - \mathcal{L}(\pi)]$ and $\Psi^* \equiv \mu\pi^* - \mathcal{L}(\pi^*)$. Note that we are charging all the idling penalty to a single supplier. In this way, a price π^* guarantees a lower bound for the payoff $\Psi_i(\pi_i, \pi_{-i})$. Hence, Ψ^* is the payoff that any supplier can guarantee for himself, i.e., his *minmax* level. We further let $\underline{\pi} := \Psi^*/\mu = \pi^* - \mathcal{L}(\pi^*)/\mu < \pi^*$ and observe that choosing price $\pi < \underline{\pi}$ is a dominated strategy. Thus, $\underline{\pi}$ can be regarded as a lower bound of suppliers' rational pricing strategies.

Although a standard approach is to look for a pure-strategy Nash equilibrium, in the next proposition we show that none exists. Instead, we shall focus on the mixed-strategy competitive equilibrium. Let $G(\pi)$ denote the mixing cumulative probability distribution of a supplier's pricing strategy π . The next proposition characterizes the structure of these mixing probabilities.

Proposition 3 *With homogeneous rates, there exists a unique symmetric equilibrium in which all suppliers randomize continuously over $[\underline{\pi}, \pi^*]$, and every supplier gets Ψ^* . The randomizing distribution is $G(\pi) = [\mu(\pi - \underline{\pi})/\mathcal{L}(\pi)]^{1/(n-1)}, \forall \pi \in [\underline{\pi}, \pi^*]$.*

The reason for not having any pure-strategy equilibrium is intuitively due to the discontinuity of suppliers' revenue functions. This creates an incentive for the cheap suppliers to increase their prices all the way to $\hat{\pi}$; however, they would also avoid to

reach $\hat{\pi}$ when themselves become the most expensive and incur a discontinuous penalty. Note that the range over which the price is randomized is completely determined by the individual's problem. In all generic cases, no tie of the highest static price may occur. In other words, the market idleness process is contributed by only one supplier. Moreover, any tie of two prices takes place with probability zero, which is in contrast to the centralized system where prices are always equal. Therefore, our homogeneous service model suggests that *price dispersion can be regarded as a sign of incoordination*. Also note that in equilibrium, the expected payoff of a supplier is identical to the case where he carries the entire market idleness, and hence he receives on average the minmax level. Competitive behavior drives away the possibility of extracting additional revenues.

Having characterized the competitive equilibrium, we now turn to the market efficiency issue. Define $\Pi^C \equiv \max_{\hat{\pi}} \{\hat{\mu}\hat{\pi} - \mathcal{L}(\hat{\pi})\}$ as the aggregate (second-order) revenue under the centralized solution and Π^* as the aggregate revenue among suppliers in the unique competitive equilibrium. The next proposition shows that the efficiency loss can be arbitrarily large when the number of suppliers explodes.

Proposition 4 *Suppose that the service rates are homogeneous. For any given aggregate service rate $\hat{\mu}$, for any given constant M , there exists a sufficiently large number N_M such that $|\Pi^C - \Pi^*| > M, \forall n > N_M$.*

Proposition 4 shows that as the number of suppliers grows, the competitive behavior among the suppliers may result in an unbounded efficiency loss. This demonstrates a significant inefficiency due to the market mechanism and it therefore calls for the need of a coordination scheme, as we investigate in Sect. 11.4.5. Note that this statement is asymptotic in the sense of the number of suppliers, which is different from the case in Sect. 11.3, and it is particularly relevant in the context of the large-scale systems discussed in Sect. 11.1. By restricting ourselves to the case of fixed aggregate service rate, we can then illustrate that the efficiency loss that results from the market idleness term also plays a pivotal role.

Note also that the first-order aggregate revenues of the centralized solution and the competitive equilibrium coincide; nevertheless, this is by construction of the asymptotic regime specified in Sect. 11.3.

11.4.4.2 Heterogeneous Service Rate Case

Now we consider the scenario where suppliers are endowed with different service rates. We again first define a global maximizers $\pi_1^*, \pi_2^*, \dots, \pi_n^*$, if supplier i ($1 \leq i \leq n$) is the one who proposes the highest price solely; i.e., we define $\pi_i^* = \arg \max_{\pi_i} [\mu_i \pi_i - \mathcal{L}(\pi_i)]$ and $\Psi_i^* \equiv \mu_i \pi_i^* - \mathcal{L}(\pi_i^*)$ as the global maximum revenue that supplier i can achieve as $J = \{i\}$. Next we let $\underline{\pi}_i = \Psi_i^* / \mu_i$ and recall that choosing price $\pi < \underline{\pi}_i$ is a dominated strategy for supplier i . The following proposition characterizes the relevant properties of an equilibrium needed for our purpose. $G_i(\cdot)$ denotes the mixing distribution that supplier i adopts in equilibrium.

Proposition 5 *Suppose suppliers are endowed with heterogeneous service rates. Then in a competitive equilibrium,*

- All $G_i(\cdot)$'s have the same left endpoint (denoted by \underline{s}) of their supports. Moreover, $\underline{s} \geq \max_{i \in \mathcal{N}} \underline{\pi}_i$.
- Suppliers' expected payoffs are proportional to their service rates $\{\mu_i\}$'s.
- If $n = 2$ and $\mu_1 > \mu_2$, then there exists a unique equilibrium in which supplier i 's revenue is $\mu_i \underline{\pi}_1$, $i = 1, 2$. The equilibrium mixing probabilities are respectively

$$\begin{aligned}
 G_2(\pi) &= \frac{\mu_1(\pi - \underline{\pi}_1)}{\mathcal{L}(\pi)}, \quad \forall \pi \in [\underline{\pi}_1, \pi_1^*], \\
 G_1(\pi) &= \frac{\mu_2(\pi - \underline{\pi}_1)}{\mathcal{L}(\pi)}, \quad \forall \pi \in [\underline{\pi}_1, \pi_1^*],
 \end{aligned}
 \tag{11.31}$$

and $G_1(\pi_1^*) = 1$. $G_1(\pi) = G_2(\pi) = 0, \forall \pi \leq \underline{\pi}_1$.

The first result on left endpoints is not surprising. This comes directly from an analogous argument for Proposition 3. The second result captures the *ex ante* difference between suppliers' payoff function: higher service rate brings higher equilibrium payoff. When we restrict to the duopoly setting, we know perfectly the range over which suppliers randomize their prices, and we can obtain closed-form expressions for their expected payoffs. They randomize the prices over the same range, and the supplier with a higher service rate tends to set a lower price: his mixing distribution stochastically dominates the other's in the usual, first-order sense. This implies that when a supplier has a capacity advantage, he can afford to price lower to capture more customers.

11.4.4.3 Numerical Results

In this section, our goal is to compare the performance between the centralized solution and the competitive equilibrium. We consider a system with n $M/M/1$ servers, delay sensitivity parameter $c = 0.5$, and arrival rate of buyers $\Lambda = 1$. The valuation v of each customer is assumed to follow an exponential distribution with mean 0.1, and \bar{p} is set such that the effective arrival rate $\mathbb{P}(v \geq \bar{p})$ matches the total service rate $\hat{\mu}$. As an example, if we let $\hat{\mu} = e^{-1.3}$, then \bar{p} can be obtained as follows: $\Lambda e^{-10\bar{p}} = \hat{\mu} \Leftrightarrow \bar{p} \approx 0.13$. The other relevant parameter is $\gamma = f(\bar{p})/(1 - F(\bar{p})) = 10$. Note that as we scale according to $\Lambda = r$ and $\hat{\mu} = \hat{\mu}^r$, \bar{p} stays unchanged.

The next two figures compare the centralized solution and the competitive equilibrium. Take $n = 2$ and assume $\hat{\mu} = \mu_1 + \mu_2 = e^{-1.3}$. Without loss of generality, we assume that supplier 1 has a higher capacity and let $a \equiv \mu_1/\hat{\mu} \in (0.5, 1)$ denote the heterogeneity of service rates between these two suppliers. Figure 11.3 demonstrates the mixing distributions of supplier 1 under a competitive equilibrium with different values of a . Figure 11.4 presents the upper and lower

Fig. 11.3 The mixing distribution of prices versus the heterogeneity of service rates

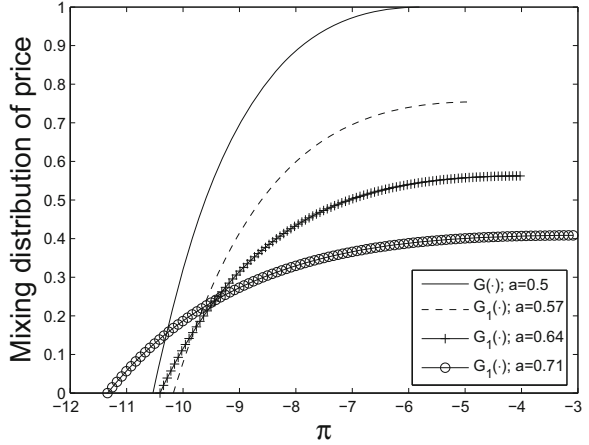
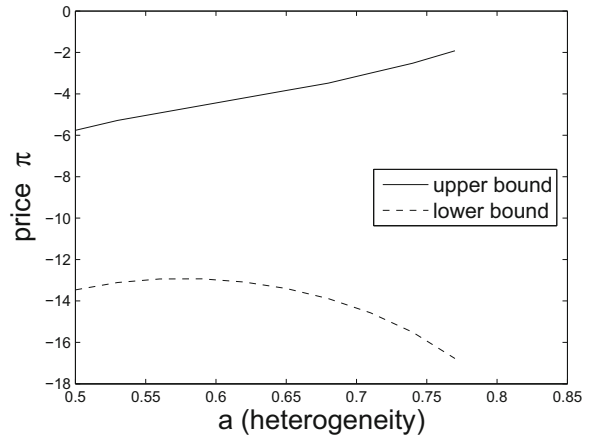


Fig. 11.4 The bounds of prices versus the heterogeneity of service rates



bounds of the price for the mixing distributions. Note that the mixing distribution may have a point mass at the upper bound π_1^* , in which case the mixing distribution jumps to 1 at π_1^* (e.g., $a = 0.57, 0.64, 0.71$ in Fig. 11.3). Although the mixing distributions of supplier 2 have no point mass, the comparison of the mixing distributions across different degrees of heterogeneity is qualitatively similar and therefore is omitted. Combining Figs. 11.3 and 11.4, there is no unambiguous prediction for the suppliers' pricing decisions when the capacity heterogeneity increases. The increase of the heterogeneity, a , has two effects. First, it mitigates the competition between the suppliers due to the difference of capacities. This might induce higher prices. Second, the increase of a also increases the variance of the service time (since $\sigma = [(1/(a\hat{\mu}))^2 + (1/((1-a)\hat{\mu}))^2]^{1/2}$ is increasing in $a \in (0.5, 1)$). This increases the magnitude of the second-order price through the parameter β . Since the second-order price is negative, it implies that the suppliers would set a lower price when the variance is higher. Because of these two conflicting

forces, no clear ranking of the mixing distribution can be obtained (as seen in Fig. 11.3), and the bounds are not monotonic (in the same direction) as the degree of heterogeneity increases (see Fig. 11.4). Note also that in the centralized solution, only one price is set:

$$\pi^C \equiv \arg \max_{\pi} \left\{ \hat{\mu}\pi - \bar{p}\gamma\hat{\mu}\beta \frac{\phi(\pi/\beta)}{1 - \Phi(\pi/\beta)} \right\} \in [4.5, 6.0]$$

when $a \in [0.5, 0.71]$. Since π^C is strictly positive, the prices in the competitive equilibrium are significantly lower than the price under the centralized control.

Finally, we investigate how the number of suppliers affects the efficiency gap between the centralized solution and the competitive equilibrium. To this end, we focus on the case with homogeneous suppliers. This allows us to fully characterize the equilibrium pricing strategies and the suppliers' expected (second-order) revenues. We first assume $\hat{\mu} = e^{-1.3}$ and increase n , the number of suppliers. The individual service rate is $\mu_i = \hat{\mu}/n, \forall i \in \mathcal{N}$. As demonstrated in Fig. 11.5, the range of prices becomes more negative when more suppliers participate in the market, due to a more severe competition among suppliers. In Fig. 11.6, we draw the expected aggregate (second-order) revenue of the market, $\sum_{i \in \mathcal{N}} \Psi_i(\pi_i, \pi_{-i})$, and vary the number of suppliers. We find that the expected aggregate revenue decreases when there are more suppliers due to the increasing price competition (as presented in Fig. 11.6). Thus, the mis-coordination problem becomes more serious when more suppliers participate in the market.

Fig. 11.5 The mixing distribution of prices versus the number of suppliers

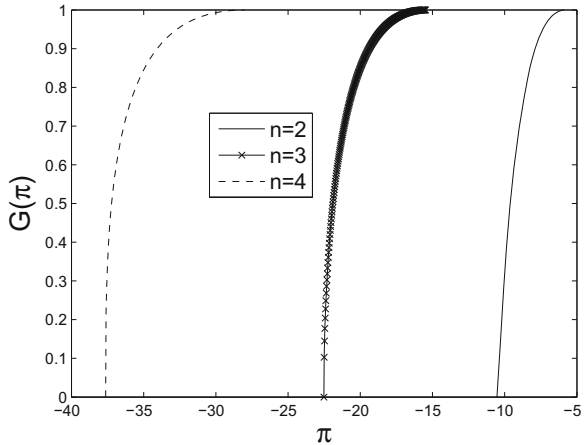
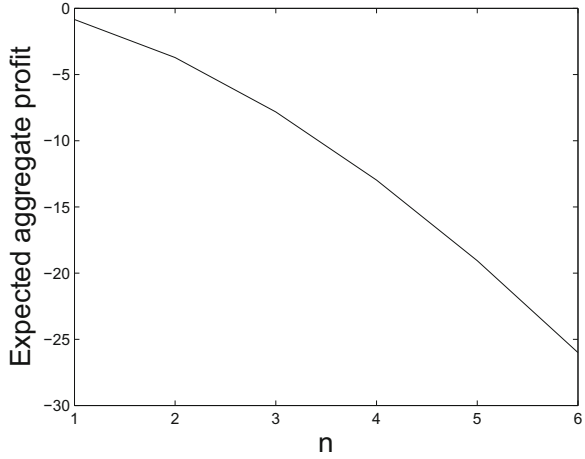


Fig. 11.6 The expected aggregate revenue versus the number of suppliers



11.4.4.4 A Remark on the Suppliers’ Participation

In characterizing the equilibrium behavior of the suppliers’ pricing game, we have neglected the suppliers’ participation decisions. Note that the pricing decisions $\{\pi_j\}$ ’s only affect the suppliers’ second-order revenues, which are simply small perturbations around the first-order revenues $\bar{p}\mu_i$. (as seen in Lemma 1). Thus, a supplier is willing to participate if and only if his first-order revenue $\bar{p}\mu_i$ is positive, which depends on the capacity (service rate) decisions rather than the pricing decisions.

To study the capacity game, we can assume that each supplier incurs a cost of capacity, $c_i(\mu_i)$. Since the pricing decisions do not affect the first-order term, the suppliers choose their capacities to maximize

$$\max_{\mu_i \geq 0} \{\bar{p}\mu_i - c_i(\mu_i)\}, \tag{11.32}$$

where the value of \bar{p} is endogenously determined through $\bar{F}(\bar{p}) = \sum_{j \in \mathcal{N}} \mu_j$, i.e., $\bar{p} = \bar{F}^{-1}(\sum_{j \in \mathcal{N}} \mu_j)$. The function $\bar{F}^{-1}(\sum_{j \in \mathcal{N}} \mu_j)$ can be interpreted as the inverse demand function, since it represents the customers’ effective arrival rate given the aggregate capacity $\sum_{j \in \mathcal{N}} \mu_j$. Notably, the above capacity game does not involve any stochastic term.

Moreover, according to Watts (1996, Corollary 1), this capacity game has a unique equilibrium if the following conditions are satisfied: (1) $\mu \bar{F}^{-1}(\mu)$ is concave in μ ; (2) $c_i(\mu_i)$ is weakly convex in μ_i ; and (3) there exists a sufficiently large μ^* such that $\mu \bar{F}^{-1}(\mu) - c_i(\mu)$ is decreasing in μ when $\mu > \mu^*$. The first condition is related to the price elasticity of the demand, the second condition implies a diseconomy of scale for the capacity investment, and the third condition simply ensures that the aggregate market payoff never explodes. These conditions are

widely adopted in many surplus sharing games, which contains the celebrated Cournot competition as a special case, to ensure that the competitive equilibrium is well-behaved (see Watts 1996, and the references therein).

Finally, a supplier is willing to participate in the market whenever in equilibrium $\max_{\mu_i \geq 0} \{\bar{p}\mu_i - c_i(\mu_i)\} \geq 0$. If we consider a special case in which the marginal cost of capacity is constant, i.e., $c_i(\mu_i) = c_i\mu_i, \forall \mu_i$, it is verifiable that only the suppliers that are more cost efficient will participate, i.e., the ones for whom $c_i < \bar{p}$.

11.4.5 Coordination Scheme

The above competitive equilibrium analysis reveals that each supplier receives an expected payoff lower than what he would obtain under the centralized solution. This implies that the centralized solution Pareto dominates all decentralized equilibria. Thus, implementing a coordination scheme results in no conflict of interests, even though the suppliers may be *ex ante* heterogeneous with respect to service rates. In addition, as the market size grows, the competitive behavior among suppliers may result in an unbounded efficiency loss. This demonstrates a significant inefficiency due to the market mechanism and motivates the search for a coordination scheme.⁸

11.4.5.1 Sufficient Condition for Coordination

We will first study the suppliers' behavior if they were "forced" to share the penalty, or revenue loss, that arises due to the market idleness. Under this scheme, supplier i 's payoff is

$$\Psi_i^{PS}(\pi_i, \pi_{-i}) \equiv \mu_i \pi_i - \frac{\mu_i}{\hat{\mu}} \mathcal{L}(\max_{j \in \mathcal{N}} \pi_j), \quad (11.33)$$

where the superscript *PS* refers to *penalty sharing* according to the service rates. The first term $\mu_i \pi_i$ is the gross revenue supplier i earns by serving customers, and the second term $(\mu_i/\hat{\mu})\mathcal{L}(\max_{j \in \mathcal{N}} \pi_j)$ corresponds to his penalty share that is proportional to his service rate μ_i . Note that this scheme is budget-balanced, i.e., no financing from outside parties is required. Let $\{\pi_i^{PS}\}'s$ denote the equilibrium prices under this sharing scheme. Under this sharing scheme, the centralized solution can be achieved.

⁸It is worth mentioning that the mixed-strategy equilibrium is studied mainly to demonstrate the discrepancy between the centralized system and the decentralized market equilibrium. It is conceivable that the implementation or identification of such a mixed-strategy equilibrium requires fairly sophisticated communication and consensus among the suppliers. Nevertheless, the Pareto dominance result justifies why such a coordination scheme is required and desired irrespective of the implementation issue of the mixed-strategy equilibrium.

Proposition 6 *Under the penalty sharing schemes that satisfy Eq. 11.33, $\{\pi_i^{PS} = \pi^C, \forall i \in \mathcal{N}\}$ is the unique equilibrium.*

Under the PS scheme, a supplier’s objective is in fact an affine function of the aggregate revenue Eq. 11.29. Hence, this coordination mechanism eliminates the wrong incentives of suppliers, regardless of the number of suppliers and their service rates. Since in both the competitive and the coordinated equilibria the suppliers’ expected payoffs are proportional to their service rates, all suppliers have a natural incentive to join.

11.4.5.2 “Compensation-While-Idling” Mechanism That Achieves Coordination

The natural question is whether we can implement a penalty-sharing mechanism based on observable quantities. We now show that this is achievable through an appropriate set of transfer prices between suppliers when one or more suppliers are idling.

Let η_{ij} be the transfer price per unit of time from supplier i to supplier j when supplier j is idle in the limit model, with $\eta_{ii} = 0$. The second-order revenue process for a supplier i under this compensation scheme becomes

$$\tilde{r}_i^{PS}(t) = \tilde{r}_i(t) + \tilde{\delta}_i(t), \tag{11.34}$$

where

$$\tilde{r}_i(t) \equiv \mu_i \pi_i t + \bar{p} \sigma_i B_{s,i}(t) - \mu_i \bar{p} \tilde{Y}_i(t), \tag{11.35}$$

$$\tilde{\delta}_i(t) \equiv \sum_{j \in \mathcal{N}, j \neq i} \eta_{ji} \tilde{Y}_i(t) - \sum_{j \in \mathcal{N}, j \neq i} \eta_{ij} \tilde{Y}_j(t). \tag{11.36}$$

According to Eq. 11.35, the three terms correspond to his revenue from serving customers. In Eq. 11.36, $\tilde{\delta}_i(t)$ corresponds to the net transfers for supplier i : $\sum_{j \in \mathcal{N}, j \neq i} \eta_{ji} \tilde{Y}_i(t)$ is the compensation he receives from other suppliers during the idle period, and $\sum_{j \in \mathcal{N}, j \neq i} \eta_{ij} \tilde{Y}_j(t)$ is the cash outflow to other suppliers while compensating their idleness.

Given Eq. 11.34, supplier i ’s long-run average revenue can be expressed as

$$\begin{aligned} \tilde{\Psi}_i(\pi_i, \pi_{-i}) &\equiv \lim_{t \rightarrow \infty} \frac{1}{t} [\tilde{r}_i(t) + \tilde{\delta}_i(t)] \\ &= \mu_i \pi_i - \mu_i \bar{p} \mathbb{E}[\tilde{Y}_i(\infty)] + \sum_{j \in \mathcal{N}, j \neq i} \eta_{ji} \mathbb{E}[\tilde{Y}_i(\infty)] \\ &\quad - \sum_{j \in \mathcal{N}, j \neq i} \eta_{ij} \mathbb{E}[\tilde{Y}_j(\infty)]. \end{aligned} \tag{11.37}$$

The next proposition specifies a set of transfer prices that implement the PS scheme.

Proposition 7 *The transfer prices*

$$\eta_{ij} = \frac{\mu_i \mu_j}{\hat{\mu}} \bar{p}, \quad \forall i \neq j, i, j \in \mathcal{N}, \quad \text{and} \quad \eta_{ii} = 0, \quad \forall i \in \mathcal{N}, \quad (11.38)$$

implement the PS rule, i.e., $\tilde{\Psi}_i(\pi_i, \pi_{-i}) = \Psi_i^{PS}(\pi_i, \pi_{-i})$.

The transfer prices proposed in Proposition 7 essentially eliminate the imbalance between the current share of the market idleness incurred by an individual supplier and his required share $((\mu_i/\hat{\mu})\mathcal{L}(\max_{j \in \mathcal{N}} \pi_j))$. Given these transfer prices, every supplier's objective is aligned with the centralized system (i.e., the objective is $\Psi_i^{PS}(\pi_i, \pi_{-i})$ in Eq. 11.33), and thus all suppliers are induced to set prices equal to π^C .

Proposition 7 shows that we are able to achieve coordination since given any chosen prices, we can align the suppliers' objectives with the planner's objective. To implement this compensation scheme in the original system, we can simply request each supplier make transfers according to Eq. 11.38. This mechanism can be implemented and monitored by the intermediary in the market.

Note that the coordination scheme is independent of the static prices $\{\pi_i\}$'s; it only requires the information of the service rates $\{\mu_i : 1 \leq i \leq n\}$, which are publicly available in our model. In fact, to facilitate the coordination scheme, the market intermediary needs to have access to the current queue lengths, and should be able to perfectly observe the idleness of suppliers.

11.4.6 Simulation Results

To close the loop, we shall return to the original system described in Sect. 11.3.2 and see how the competitive equilibrium and coordination scheme fare. To this end, we again run simulations using the Arena model. The parameters are the same as those in Sect. 11.4.4.3: $c = 0.5$, $\Lambda = 1$, $\hat{\mu} = e^{-1.3}$, and the valuation v is exponentially distributed with mean 0.1.

Mixing distributions of prices Compared with Sect. 11.3.6, a new challenge arises: we cannot arbitrarily assign prices, because now the suppliers determine their competitive prices as equilibrium outcomes. We shall use the results from Sect. 11.4.4.3 as inputs to our Arena model for both homogeneous and heterogeneous cases of suppliers. We note that the equilibrium pricing strategy is described by a continuous distribution without simple expressions (see Propositions 3 and 5). Thus, the usual inverse-transform method fails to apply (because the inverse of distribution function is not known). Furthermore, the acceptance-rejection method is also not suitable for this problem, because it requires an explicit expression of density function that is not available. Our treatment follows from a similar idea to Fig. 11.3. We first discretize them and record the cumulative distribution at discrete points. We choose the mesh sufficiently small and make linear interpolation to replicate approximately the original continuous distribution.

Second-order revenues In terms of suppliers' profits, we focus exclusively on the pricing game in which the second-order correction around $R_i^r(t)$ defined as $r_i^r(t) \equiv (1/\sqrt{r})(R_i^r(t) - r\bar{p}\mu_i t)$, $\forall i \in \mathcal{N}$. Note that given $\hat{\mu} = e^{-1.3}$, $\bar{p} \approx 0.13$. When there are two suppliers ($n = 2$), we let $a \equiv (\mu_1/\hat{\mu}) \in (0.5, 1)$ denote the heterogeneity of service rates between these two suppliers. We examine two scenarios: in the homogeneous case $a = 0.5$, these two suppliers are endowed with the same service rate (capacity). In the heterogeneous case, we choose $a = 0.71$.

Regarding the tie-breaking rule, here we examine two rules: the smallest index first rule by which supplier 1 gets the priority, and the random priority rule by which customers choose between tied suppliers with equal probabilities. For the random priority rule, we add a "two-way by chance" module to rout the customers randomly with 50–50 chances when there is a tie. For all the following simulations, we conduct 1,000 replications, each of which takes the warm up of 120 arrivals and regular simulation of 600 arrivals in expectation. The scaling factor r is fixed at 1,000. The confidence level is set at 5% when we make the statistical statements of hypothesis testing. We use two-sample-t two-tailed tests when we compare across different scenarios and paired-t tests when comparing between the two suppliers within each scenario.

Revenue comparison: Symmetric case First, we consider two symmetric suppliers ($a = 0.5$) and compare the suppliers' second-order revenues in the competitive equilibrium and under the coordination scheme. We find that the average difference of suppliers 1s and 2s revenues in these two scenarios are statistically significant. For supplier 1, the estimated revenue difference is 0.758 whereas the 95% confidence interval is 0 ± 0.00419 . Similarly, supplier 2s estimated revenue difference 0.76 and 0.76 falls outside ± 0.00421 . Therefore, the coordination scheme indeed leads to higher expected revenues for both suppliers.

Revenue comparison: Asymmetric case Second, we consider two asymmetric suppliers ($a = 0.71$). In this case, the coordination scheme again yields higher expected revenues for both suppliers that are statistically significant. The estimated revenue improvements for suppliers 1 and 2 are 0.992 and 0.351 respectively, and they fall outside the 95% confidence intervals ± 0.0026 and ± 0.00284 . We can also compare the two suppliers' revenues. Naturally, their (second-order) revenues are different due to heterogeneous service rates. Using paired-t tests, we observe that the revenue differences are statistically significant in both the competitive and coordinated scenarios (-0.435 and 0.206 on average), and their corresponding confidence intervals are ± 0.00295 and ± 0.00488 .

Tie-breaking rule Third, we can also examine the impact of tie-breaking rule. For this matter, we use the symmetric supplier case as illustration. We first start with the competitive equilibrium and compare the two tie-breaking rules: smallest index first and random rules. For the competitive equilibrium we fail to reject the null hypothesis, i.e., the suppliers' revenues are statistically indistinguishable under the two rules. In contrast, under the coordination scheme the tie-breaking rule matters substantially. The estimated revenue difference is 0.0857 and it falls outside the 95% confidence interval ± 0.00323 .

The above discrepancy can be explained intuitively. In the competitive equilibrium, both suppliers randomize their prices. Thus, the chance of seeing an actual tie is infinitesimal (zero probability in theory). Therefore, the tie-breaking rule rarely comes in action. However, under the coordination scheme, both suppliers are induced to set prices at π^C . Because their service rates are identical, often times ties actually happen and the priority rule goes in favor of supplier 1. In this case, random tie-breaking ensures the fair routing between suppliers and it leads to statistically significant consequences. Further to the above observation, we run additional comparisons between the two suppliers. Under the smallest index first rule, supplier 1 earns on average 0.176 more than supplier 2, which is outside the 95% confidence interval ± 0.00308 . Under the random priority rule, this difference is negligible (-0.000792 on average).

To summarize, our simulations suggest that (1) the coordination scheme is effective in both homogeneous and heterogeneous scenarios, and this benefit applies to all suppliers; (2) tie-breaking rules are inconsequential when suppliers adopt randomized pricing, but they do matter when instead deterministic prices are chosen.

11.5 Conclusions

We study an oligopolistic model in which suppliers compete for buyers that are both price and delay sensitive. We apply both fluid and diffusion approximations to simplify the multi-dimensional characteristics of the decoupled suppliers into a single-dimensional aggregated problem. Specifically, we establish the “state space collapse” result in this system: the multi-dimensional queue length processes at the suppliers can be captured by a single-dimensional workload process of the aggregate supply in the market, which can be expressed explicitly as a reflected Ornstein–Uhlenbeck process with analytical expressions. Based on this aggregated workload process, we derive the suppliers’ long-run average revenues and show that the suppliers’ competition results in a price randomization over bounded ranges, whereas under the centralized control suppliers should set identical and deterministic prices.

To eliminate the inefficiency due to the competition, we propose a novel compensation-while-idling mechanism that coordinates the system: each supplier gets monetary transfers from other suppliers during his idle periods. This mechanism alters suppliers’ objectives and implements the centralized solution at their own will. The implementation only requires a set of static transfer prices that are independent of the suppliers’ prices and the queueing dynamics such as the current queue lengths or the cumulative idleness. Its simplicity is an appealing feature to be considered for practical implementations in intermediary platforms such as online exchanges.

References

- Afèche P (2013) Incentive-compatible revenue management in queueing systems: optimal strategic delay. *Manuf Serv Oper Manag* 15(3):423–443
- Allon G, Federgruen A (2007) Competition in service industries. *Oper Res* 55(1):37–55
- Allon G, Gurvich I (2010) Pricing and dimensioning competing large-scale service providers. *Manuf Serv Oper Manag* 12(3):449–469
- Allon G, Bassamboo A, Cil EB (2012) Large-scale service marketplaces: the role of the moderating firm. *Manag Sci* 58(10):1854–1872
- Ata B, Kumar S (2005) Heavy traffic analysis of open processing networks with complete resource pooling: asymptotic optimality of discrete review policies. *Ann Appl Probab* 1:331–391
- Besbes O (2006) Revenue maximization for a queue that announces real-time delay information. Working paper, Graduate School of Business, Columbia University
- Bramson M (1998) State space collapse with applications to heavy-traffic limits for multiclass queueing networks. *Queueing Syst* 30:89–148
- Browne S, Whitt W (2003) Piecewise-linear diffusion processes. In: *Advances in queueing*. CRC Press, pp 463–480
- Chen YJ, Maglaras C, Vulcano G (2008) Design of an aggregated marketplace under congestion effects: asymptotic analysis and equilibrium characterization. Technical Report, Columbia University
- DiPalantino D, Johari R, Weintraub GY (2011) Competition and contracting in service industries. *Oper Res Lett* 39(5):390–396
- Gallego G, van Ryzin G (1994) Optimal dynamic pricing of inventories with stochastic demand over finite horizons. *Manag Sci* 40:999–1020
- Lederer P, Li L (1997) Pricing, production, scheduling and delivery -time competition. *Oper Res* 45:407–420
- Levhari D, Luski I (1978) Duopoly pricing and waiting lines. *Eur Econ Rev* 11:17–35
- Loch C (1991) Pricing in markets sensitive to delay. Ph.D. dissertation, Stanford University, Stanford
- Luski I (1976) On partial equilibrium in a queueing system with two servers. *Rev Econ Stud* 43:519–525
- Maglaras C, Zeevi A (2003) Pricing and capacity sizing for systems with shared resources: approximate solutions and scaling relations. *Manag Sci* 49:1018–1038
- Maglaras C, Moallemi C, Zheng H (2016) Queueing dynamics and state space collapse in fragmented limit order book markets. Working paper, Columbia University
- Mandelbaum A, Pats G (1995) State-dependent queues: approximations and applications. In: Kelly F, Williams R (eds) *Stochastic networks*. Proceedings of the IMA, vol 71. North-Holland, pp 239–282
- Mendelson H (1985) Pricing computer services: queueing effects. *Commun ACM* 28:312–321
- Mendelson H, Whang S (1990) Optimal incentive-compatible priority pricing for the $M/M/1$ queue. *Oper Res* 38:870–883
- Naor P (1969) On the regulation of queue size by levying tolls. *Econometrica* 37:15–24
- Stolyar AL (2005) Optimal routing in output-queued flexible server systems. *Probab Eng Inf Sci* 19:141–189
- Watts A (1996) On the uniqueness of equilibrium in Cournot oligopoly and other games. *Games Econ Behav* 13(2):269–285
- Williams RJ (1998) An invariance principle for semimartingale reflecting Brownian motions in an orthant. *Queueing Syst* 30:5–25

Chapter 12

Operations in the On-Demand Economy: Staffing Services with Self-Scheduling Capacity



Itai Gurvich, Martin Lariviere, and Antonio Moreno

Abstract Motivated by recent innovations in service delivery such as ride-sharing services and work-from-home call centers, we study capacity management when workers self-schedule. Our service provider chooses capacity to maximize its profit (revenue from served customers minus capacity costs) over a horizon. Because demand varies over the horizon, the provider benefits from flexibility to adjust its capacity from period to period. However, the firm controls its capacity only indirectly through compensation. The agents have the flexibility to choose when they will or will not work and they optimize their schedules based on the compensation offered and their individual availability. To guarantee adequate capacity, the firm must offer sufficiently high compensation. An augmented newsvendor formula captures the tradeoffs for the firm and the agents. If the firm could keep the flexibility but summon as many agents as it wants (i.e., have direct control) for the same wages it would not only generate higher profit, as is expected, but would also provide better service levels to its customers. If the agents require a “minimum wage” to remain in the agent pool they will have to relinquish some of their flexibility. To pay a minimum wage the firm *must* restrict the number of agents that can work in some time intervals. The costs to the firm are countered by the self-scheduling firm’s flexibility to match supply to varying demand. If the pool of agents is sufficiently large relative to peak demand, the firm earns more than it would if it had control of agents’ schedules but had to maintain a fixed staffing level over the horizon.

I. Gurvich (✉)
Cornell Tech, New York, NY, USA
e-mail: gurvich@cornell.edu

M. Lariviere
Kellogg School of Management, Evanston, IL, USA
e-mail: m-lariviere@kellogg.northwestern.edu

A. Moreno
Harvard Business School, Boston, MA, USA
e-mail: amoreno@hbs.edu

12.1 Introduction

Staffing in service environments is a challenging problem. Firms must control costs while assuring adequate capacity to serve demand. In tackling this problem, managers have always maintained an important trump card: the ability to order workers to work at specific times. The construction of the schedule might involve worker preference, union rules, or government regulations but, at the end of the day, each worker has been told when she is expected to begin and end her shift. Furthermore, these directives have been backed by implicit (and often explicit) consequences for not adhering to an assigned schedule.

In some novel service settings, however, firms are surrendering this power. Instead of ordering workers to punch in and out at appointed times, firms are allowing agents to create their own schedules, choosing whether and when to work based on personal preferences. We are not speaking here of professional knowledge workers who are given flexible schedules as long as projects are completed on time. Rather, we are focusing on industries such as ride-sharing services (e.g., Uber and Lyft), work-from-home call centers (e.g., Arise Virtual Solutions and LiveOps), or delivery services (e.g., Instacart and GrubHub) which must have capacity available to service demand as it arises.

These service providers have put themselves in a tenuous position. On the one hand, they need to provide their customers with good service. Ride-sharing services, for example, compete against conventional taxis and public transportation in part by emphasizing their availability. In the words of Uber's chief executive, "Uber is ALWAYS a reliable ride." (Kalanick 2012). Delivering on these commitments requires capacity; without adequate staffing, these service providers will fail to honor their obligations.

On the other hand, these service providers promise their agents¹ flexibility and cannot simply dictate when they should work. Food delivery service GrubHub, for example, promises that its "delivery partners" can pick when they want to work and that they will earn "competitive pay" when they do.² Work-from-home call center LiveOps makes a similar pitch:

As a LiveOps independent agent, you can benefit from a highly flexible and rewarding opportunity. ... As an independent contractor providing services to LiveOps' clients, you are your own boss!³

Flexibility and control of one's schedule are important to agents. In a study of Uber drivers (conducted for Uber), 85% of respondents cited the ability "to have more

¹Describing the people serving customers for these firms requires some finesse. Generally, those answering calls or driving customers are not employees. Rather, they are independent contractors whose continued relationship with the service provider is dependent on achieving a minimal level of performance (e.g., an Uber driver rating) over time. We will generally refer to those serving customers as agents.

²<http://driver.grubhub.com/>. Accessed May 24, 2016.

³<http://www.liveops.com/company/careers-jobs> Accessed May 24, 2016.

flexibility in my schedule” as a motivation to drive for the company (Hall and Krueger 2016). Additionally, Lyft and Uber have pointed to the fact that drivers set their own schedules in contesting lawsuits on whether drivers should be deemed employees or independent contractors (Levine and McBride 2015). Consequently, these firms cannot simply renege on allowing agents to self schedule. They must instead use incentives schemes to induce the right number of agents to be available at the right time.

A service provider must also assure that its agents have adequate earnings over time. Some of these firms aggressively recruit and compete with each other for agents.⁴ There are blogs that inform about work conditions in competing services.⁵ Firms are consequently rightly concerned when websites ask whether a particular company is a “work-from-home scam”⁶ or when former agents complain in public forums that a firm is “the worst company ever” offering “below average” pay.⁷

The provider’s problem can thus be understood as managing agent participation on two different time scales. On a longer-term basis (measured in weeks or months), the firm must maintain an adequate pool of eligible agents. In order to keep agents in the pool, the firm must ensure that agents earn enough to make collaborating with the firm an attractive opportunity. On a short-term basis (measured in hours or less), it must attract enough – but not too many – agents for each time interval over some horizon to maximize its profit while achieving a desired service level.

The goal of this paper is to examine how a firm that allows its agents to self schedule solves this problem. We consider a firm that staffs a service system facing time-varying arrivals over a horizon. The firm recruits a pool of agents who in each period choose whether or not to work. A given agent’s willingness to work varies with each period as she draws an *availability threshold* at the start of each period. Thus, given the terms the firm offers, an agent may want to work this morning but then be unavailable this afternoon.

The firm has three control levers at its disposal. First, it can set the *pool size* – that is, how many agents it recruits and qualifies to serve customers. Since training agents takes time, the pool size is set at the start of the horizon and cannot be adjusted based on the demand in a given period. The second lever is the *compensation* offered to agents who work in a period. This can vary from time period to time period. For most of our analysis, we assume the firm offers a fixed compensation for each time interval (e.g., \$15 per hour). However, we demonstrate that the firm can achieve identical results through alternative schemes, such as a piece-rate compensation, that depend on the number of customers an agent serves. Finally, we allow the firm to impose a *cap* on the number of agents that are active in

⁴<http://www.forbes.com/sites/ellenhuet/2014/05/30/how-uber-and-lyft-are-trying-to-kill-each-other/>

⁵<http://therideshareguy.com/category/lyft-vs-uber/>

⁶See workathomemoms.about.com/od/callcenterdataentry/a/arise.htm accessed on May 24, 2016.

⁷See www.glassdoor.com/Reviews/Employee-Review-LiveOps-RVW2743190.htm accessed on May 24, 2016.

a period. That is, the firm can tell an agent she cannot work in a given time interval even though she is willing to do so. Delivery service DoorDash, for example, limits the number of agents working in a given area at a given time (Campbell 2015).

We use a newsvendor setting in which demand and capacity can vary from period to period and find that the optimal decision is an elegant variant of the classical critical-fractile formula. Suppose that the firm offers agents a wage of η and receives revenue p from successfully serving a customer. Under conventional staffing (i.e., assuming that the firm can order any number of agents to work in a given period), the firm would employ enough agents so that the probability of turning customers away is η/p . The corresponding probability under self scheduling, however, is $\eta/p + F(\eta)/(pf(\eta))$, where F is a distribution governing agent availability and f its density.

This explicitly captures the cost of relinquishing direct control of capacity. If the firm could choose how many agents to summon for the same wage of η , rather than leaving this decision to the agents, it would set a higher staffing level and have higher profits. Further, customers would benefit from a higher service level. The drop in customer service from self scheduling is exacerbated when demand varies over the horizon. In a horizon with both high and low demand periods (in the sense of stochastically larger or smaller demand distributions), the provider offers agents higher pay in high demand periods and makes more capacity available, but the service level customers see falls in high demand periods.

We also demonstrate that the firm must use all three control levers – particularly capping the number of active agents – when it has to satisfy a nontrivial constraint on agent earnings. Absent an earnings constraint (i.e., when the firm only needs to consider gaining adequate agent participation in each period), the firm has an incentive to make its pool of agents as large as possible. It is then able to offer relatively low wages in both high and low demand periods and still induce a large number of agents to work. Once there is a constraint on agent earnings, however, the firm cannot slash wages. This drives up costs both because it pays more and because that higher pay induces too many agents to work. In particular, low-volume periods will be overstaffed. Capping the number of active agents addresses this problem. This, however, implies that *agents must sacrifice some scheduling flexibility in order to guarantee a minimum compensation level*.

The necessity of a cap does not go away if one replaces a per-period wage with a piece rate. However, its role changes. Under a per-period wage, a cap keeps the firm from paying for agents it does not want at the prevailing wage. Under a piece rate, the cap keeps excessive competition between agents from diluting agents' earnings.

Our work is related to the literature on principal-agent models (see Salanie 1997; Laffont and Martimort 2009, for reviews). The classical models focus on hiring an agent to exert effort for the benefit of the principal when the agent's actual effort cannot be observed. The principal is concerned with both directing the agent's action as well as gaining the agent's participation. We do not explicitly model effort. In effect, we assume monitoring is sufficient to assure that agents provide the appropriate level of effort. Our attention is squarely on assuring agent participation.

There has been some work in the operations literature looking at two-sided markets that match tasks with service providers (e.g., see Allon et al. 2012; Moreno and Terwiesch 2015). In these papers, individual clients look to buy a specific service (e.g., coding a smart phone app) that can be carried out by one individual. The question here is how different rules or information structures affect market performance. In our case, the service provider commits to serving customers with homogeneous requests that any available agent can handle. The question is then how the firm assures it has sufficient capacity to meet demand. Closer to our motivating applications, Riquelme et al. (2015) model a ride-sharing platform as a queue with customers requesting service and drivers executing rides. The paper studies how dynamic pricing benefits the platform.

To our knowledge there are only a handful of papers that explicitly deal with self-scheduling agents. Cachon et al. (2017) consider a one period problem in which agents join a platform's pool before knowing the market demand. Their analysis focuses on how the platform optimally adjusts agent payments and retail prices to efficiently allocate capacity after uncertainty is resolved. In our multi-period model temporal variation is predictable (i.e., everyone knows when "rush hour" is). Compensation is then dynamic but set in advance. In Taylor (2018), customers are delay sensitive so demand increases as waiting time decreases. In a queuing setting, it is shown that prices do not necessarily increase with congestion. Uncertainty is a key driver of these results since compensation and prices are set before uncertainty is resolved. In our setting, demand is not delay sensitive and all parties are informed about the level of demand.

A key distinction between our paper and the aforementioned two papers is that we consider not one but a sequence of periods with different but known statistical properties. The size of the agent pool is constant during the entire time horizon which creates a coupling between the different periods. The firm can, however, vary the compensation levels and the caps it puts on realized capacity. Ibrahim (2018) also considers a multi-period queueing model in which the pool size is optimized once in the beginning of the horizon and the decision maker has no further control of the number of active servers in each period. Our decision maker, in contrast, is endowed with the ability to affect the number of active servers by varying the compensation.

12.2 Model

We consider a service provider selling to customers over a horizon composed of T time intervals. In period t (for $1 \leq t \leq T$), the firm's revenue is determined by the number of available agents and market conditions. Let A_t denote the number of agents available in period t and $\mathbf{A} = (A_1, \dots, A_T)$. We assume that each agent can serve one customer per period making the firm's staffing level equivalent to its capacity. We assume that market conditions in period t are captured by a probability distribution G_t . That is, the actual demand in period t , D_t , is drawn from G_t . One

expects, for example, G_t to exhibit day of the week or time of day seasonality (e.g., call volume peaks in the late morning). Let $\mathbf{G} = (G_1, \dots, G_T)$.

Let $R(A_t, G_t)$ denote the firm's revenue in a period t with A_t available agents and market conditions G_t . Then,

$$R(A_t, G_t) = pS_t(A_t) = p \left(\int_0^{A_t} xg_t(x) dx + A_t \bar{G}_t(A_t) \right), \quad (12.1)$$

where g is the density of G , which we assume to be strictly positive, and $\bar{G} = 1 - G$. $S_t(A_t)$ is then expected unit sales in period t given staffing level A_t . The retail price p is fixed over the horizon. In Sect. 12.4, we allow the firm to choose p . The firm's revenue over the horizon is $R_T(\mathbf{A}, \mathbf{G}) = \sum_{t=1}^T R(A_t, G_t)$.

We assume that the firm pays agents η_t for being available in period t . For now, we assume that compensation is implemented through a per-interval compensation (e.g., paying \$15 per hour). We discuss alternative compensation schemes in Sect. 12.4. The firm's profit at period t is then given by

$$\Pi(A_t, G_t) = R(A_t, G_t) - \eta_t A_t,$$

and its profit over the horizon by $\Pi_T(\mathbf{A}, \mathbf{G}) = \sum_{t=1}^T \Pi(A_t, G_t)$.

Given a vector $\boldsymbol{\eta} = (\eta_1, \dots, \eta_T)$, the firm would like to use staffing levels \mathbf{A}^* that maximize Π_T and schedule A_t^* agents to be available in period t . However, under self scheduling it cannot directly order A_t agents to work. Instead it must offer sufficient compensation to induce that many agents to choose to work. We suppose that the firm has a pool of N qualified agents. Interpret N as the number of agents that are affiliated with (or belong to) the network of a firm, who have been trained to serve customers. Thus, N is the maximum number of agents that *could* potentially work in a given period. However, it is not the case that all pool members *will* work; some may find the firm's offered compensation in that period to be insufficient.

We model variation in agents' availability to work by assuming that each agent has an *availability threshold* for each period. An agent may thus be available for work this morning because they have drawn a low threshold but be unavailable this afternoon or tomorrow morning because they have drawn a significantly higher threshold. More formally, each agent draws an availability threshold τ from a distribution F at the start of each period. Agents are assumed to be statistically identical and independent of each other. The distribution does not vary over time, and a given agent's draw for period t is independent of her draw for any other period. We consider having F depend on the time period in Sect. 12.4. We assume that F is continuous with a strictly positive density f on a support $(0, \Phi)$. Let $\bar{F}(\tau) = 1 - F(\tau)$. We assume that F is log-concave, a condition that holds for many common distributions (see Bergstrom and Bagnoli 2005).

Agents are risk neutral and seek to maximize their earnings subject to only working in periods in which they expect to earn more than the availability threshold they have drawn for that period. Thus, an agent with realized availability threshold τ_t

in period t makes herself available to work if the firm offers compensation η_t greater than τ_t . The total number of agents *interested* in working in period t is then $NF(\eta_t)$. Note that we are implicitly appealing to the law of large numbers by assuming that the pool of qualified agents is sufficiently large that working with average number of available agents is a reasonable approximation of the actual number of available agents.

The firm's problem is then to maximize $\Pi_T(\mathbf{A}, \mathbf{G})$ by manipulating its available control levers. We consider three. The first is the pool size N . Since training agents takes time, this decision must be made up front. The pool size is thus constant over the horizon. The second variable is the agent compensation which is allowed to vary from period to period. Finally, the firm may *impose a cap* K_t on the number agents allowed to work in period t . If, under the offered compensation, the number of interested agents, $NF(\eta_t)$, exceeds the number the firm wants, it can choose to limit access only to the number it needs. With a cap K_t , the staffing level is $A_t = NF(\eta_t) \wedge K_t$. Allowing the possibility of an access cap requires some assumption regarding how the firm chooses from among interested agents. We will assume *random rationing*: the agents who work in an interval are chosen randomly from amongst those that are interested. Other rationing mechanisms are possible; for example, Netessine and Yakubovich (2012) discuss several settings in which better workers are given priority.

To this basic problem we can add a constraint related to agent welfare. As discussed in the introduction, firms have an interest in assuring that they are seen as providing good opportunities for workers. We model this by imposing a constraint on the agents' compensation. We consider a *per-period earnings* constraint that requires $\eta_t \geq \beta$ for all $t = 1, \dots, T$. This β is the "minimum wage." In a setting with a long repetitive horizon (as in virtual call centers where consecutive weeks are similar), this is equivalent to requiring that agents get sufficient earnings on any "type" of interval in which they work. If $\beta = 0$ the firm faces no earnings constraints.

Given these considerations, we can write the general form of the firm's optimization problem as

$$\begin{aligned} \max_{\eta, N, K} \quad & \Pi_T(\mathbf{A}, \mathbf{G}) \\ \text{s.t.} \quad & A_t = NF(\eta_t) \wedge K_t, \quad t = 1, \dots, T, \\ & \eta_t \geq \beta, \quad t = 1, \dots, T. \end{aligned} \tag{12.2}$$

We will consider variants of this problem where some decision variables are fixed (rather than adjustable). If any variable is fixed, this will be explicitly stated. For future reference, with one period this optimization problem is spelled out as

$$\begin{aligned} \max_{\eta, N, K} \quad & p \left(\int_0^A xg(x) dx + A\bar{G}(A) \right) - \eta A \\ \text{s.t.} \quad & A = NF(\eta) \wedge K, \\ & \eta \geq \beta. \end{aligned} \tag{12.3}$$

12.3 Analysis

In this section we establish the following three key results:

1. **Theorem 1:** Controlling capacity only indirectly imposes costs on the firm and its customers. If the firm could pay the same wages but have direct control of capacity it would have (obviously) a higher profit but it would also staff with more agents and, in turn, provide a higher service level.
2. **Theorem 2:** Under self-scheduling agents are better off in high demand periods (their compensation is higher) but customers are worse off (they experience a lower service level).
3. **Theorem 3:** In the presence of earnings constraints, the firm must use all of the tools in its toolbox. To maximize its profit, it is necessary for the firm to cap access in low demand periods but not necessarily in high demand period. Thus, an earnings guarantee comes at a cost to agents – their flexibility to self-schedule will be compromised.

The challenges brought about by self-scheduling should be weighed against the value of the flexibility self-scheduling brings. We explore this point in Sect. 12.3.4.

12.3.1 The Cost of Self Scheduling

To begin, we assume that customer arrivals are identically distributed in each period, i.e., $G_t \equiv G$, and hence drop the dependence on the time period from the notation. We first optimize the compensation level assuming that the pool size is fixed, that the firm does not consider earnings constraints (i.e., $\beta = 0$), and that no access cap is used.

The firm here solves Eq. 12.3 with $\beta = 0$, i.e., it maximizes $R(A, G) - \eta A$ where

$$R(A, G) = pS(A) = p \left(\int_0^A xg(x) dx + A\bar{G}(A) \right).$$

The following lemma is immediately derived from the first-order condition.

Lemma 1 *The unique optimal compensation level η^* satisfies*

$$G(NF(\eta^*)) = 1 - \frac{\eta^* + F(\eta^*)/f(\eta^*)}{p}. \quad (12.4)$$

The uniqueness of η^* follows from the log-concavity of F , which implies that the reversed hazard rate $f(\eta^*)/F(\eta^*)$ is monotonically decreasing. In the mechanism design literature, $\eta^* + F(\eta^*)/f(\eta^*)$ is known as the virtual cost. That is, the decision maker acts as if her marginal cost is $\eta^* + F(\eta^*)/f(\eta^*)$ even though she pays agents only η^* .

To capture the cost of indirect control, one can ask how many agents would the firm choose if, *for the same wages*, it could summon as many as agents as it wanted to maximize the profit:

$$\Pi^*(\eta^*) := \max_A \{pS(A) - \eta^*A\}. \quad (12.5)$$

With *direct control* its staffing level would be $A(\eta^*)$ given by

$$G(A(\eta^*)) = 1 - \frac{\eta^*}{p}. \quad (12.6)$$

The following theorem supports that, with the same wages, direct control of capacity not only increases profits (as intuitively expected) but also improves service level for the customer (All proofs appear in the appendix.)

Theorem 1 *For any given $N \geq 0$, direct-control of capacity leads to higher staffing level and profit relative to self scheduling, i.e., $NF(\eta^*) \leq A(\eta^*)$ and $\Pi(NF(\eta^*), G) \leq \Pi(A(\eta^*), G)$.*

We thus have a clear statement that direct control may be valuable to the firm and its customer. Of course, to obtain direct control of capacity the firm would likely have to pay the agents a premium for losing their flexibility. But for a sufficiently small premium, the service provider will opt for direct control. Suppose that the firm must offer $\eta^* + \Delta$ to compensate agents for giving up their flexibility. As long as $\Delta < F(\eta^*)/f(\eta^*)$, directly controlling capacity will result in a higher staffing level – and higher revenue – than self scheduling at η^* . The firm will prefer direct control as long as the increase in revenue is sufficient to offset higher staffing costs. The move to direct control would also benefit customers as they have a greater chance of receiving service.

The extent of the difference in outcomes between the self-scheduling setting and the benchmark depends on problem parameters. We next consider how the pool size and the distribution of the agent's availability threshold affect the firm's actions. To this end, we write $\eta_{F,N}^*$ to capture explicitly the dependence of the optimal agent compensation on the availability threshold distribution.

Lemma 2 *The firm's profit increases as either N increases or the availability threshold distribution decreases in the reversed hazard rate order, i.e.,*

$$\Pi(NF_1(\eta_{F_1,N}^*), G) \geq \Pi(NF_2(\eta_{F_2,N}^*), G) \text{ given } F_1 \text{ and } F_2 \text{ if } \frac{f_1(\tau)}{F_1(\tau)} \leq \frac{f_2(\tau)}{F_2(\tau)}.$$

The compensation rate for agents decreases as either N increases or the availability threshold distribution decreases in the reversed hazard rate order. Finally, the service level increases as N increases.

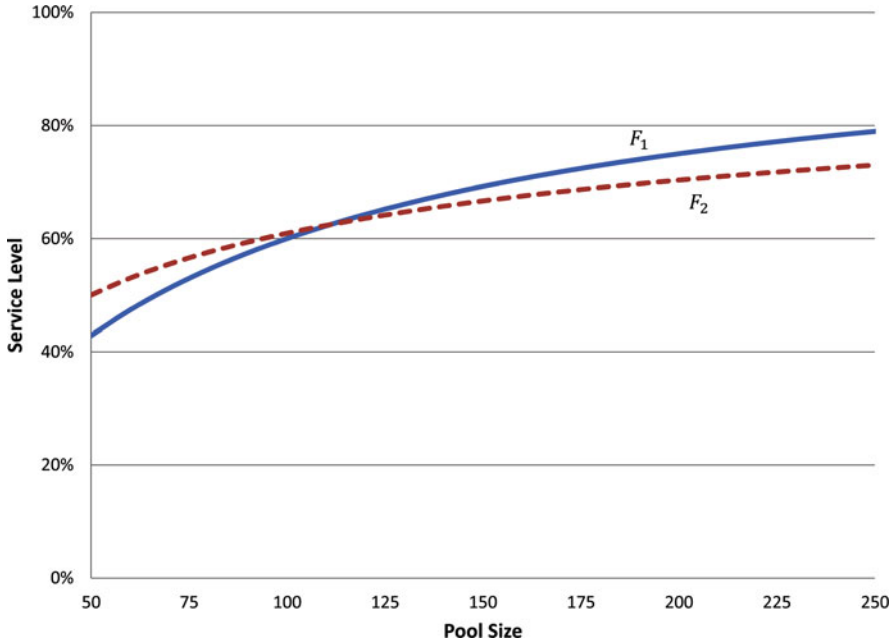


Fig. 12.1 The impact of pool size and availability threshold distribution on service level. F_2 is larger than F_1 in the reverse hazard rate ordering

A large pool promises more agents with low availability thresholds, which allows for a lower payment to agents. Additionally, since $F(\eta)/f(\eta)$ is increasing in η , the service level increases as η decreases and the negative effect of self-scheduling on service level decreases. Consequently, while self scheduling is less profitable, a self-scheduling firm sacrifices little when it has a large pool of agents.

Holding the pool size fixed, a smaller distribution of availability thresholds in the reverse hazard order means that a greater fraction of the pool is available at any compensation rate, which increases the firm’s profit. No simple monotonicity claims can be made about the dependence of the service level on the threshold distribution. Suppose demand is uniform over $[0, 100]$ while the availability threshold is one of two power function distributions, $F_1(x) = x$ for $0 < x < 1$ or $F_2(x) = x^2$ for $0 < x < 1$. Note that F_2 is larger than F_1 in the reverse hazard rate ordering. Figure 12.1 shows that the lower compensation rate under a smaller distribution does not necessarily translate to a higher service level. For smaller pool sizes, customers see a higher service level when agents have the larger availability threshold distribution F_2 . This relationship is reversed when the pool size is large. Intuitively, two countervailing forces are at play. On the one hand, a smaller threshold distribution means that more agents will be willing to work for any value of η , which favors a higher service level. On the other, a smaller threshold distribution means that the firm faces a higher virtual cost for any value of η , which argues for a lower service level. Larger pool sizes amplify the former effect causing it to be the dominating force as the number of agents grows.

12.3.2 Earnings Constraint and Agent Flexibility

Taking Lemma 2 to its logical conclusion, and assuming no costs of maintaining its pool (see Sect. 12.4.3), the firm would set N as large as possible, maximizing its earnings but squeezing the agents.

This is no longer feasible if there is a nontrivial constraint, $\beta > 0$, on agent earnings. It is obvious that if η^* , for a given N , is greater than β (where η^* is determined by Eq. 12.4), the earnings constraint is not binding and the firm can use η^* for compensation. If $\eta^* < \beta$, the firm would set the compensation per interval at β and would have $NF(\beta)$ interested agents. As stated in Lemma 2, the optimal compensation η^* is a decreasing with the pool size: for sufficiently large N , the earnings constraint will bind.

Lemma 3 *Any optimal solution to Eq. 12.3 with $\beta > 0$ has $\eta^* = \beta$ and $N^* \geq \bar{N} := \bar{G}^{-1}(\beta/p)/F(\beta)$. For all values of $N > \bar{N}$ such that $\eta_{N^*}^* < \beta$, the firm can strictly increase its profit by setting a cap and the optimal cap is set at $K^* = A(\beta)$.*

If $N = \bar{N}$ the optimal compensation will be precisely $\eta^* = \beta$. If the pool is larger compensation is fixed at β and the staffing is capped at $A(\beta)$. This cap ensures that the firm does not have to pay more agents than it needs at the prevailing wage: it is able to attract all the agents it needs at compensation β as if endowed with direct control. In other words, direct control would not affect the number of active agents, increase the firm's profit or improve the customers' service level. That not all interested agents are able to join (because of the cap at $A(\beta)$) implies that the agents give up some of their flexibility in return for a guaranteed compensation.

12.3.3 Time-Varying Demand

As we just saw, with stationary demand, if the firm is allowed to choose its pool size, it can set $N^* = \bar{G}^{-1}(\beta/p)/F(\beta)$ so that capping the number of agents that work is not necessary. The firm can choose N^* such that the number of interested agents is exactly the number it wants. However, the cap regains its relevance in an environment with time-varying demand.

Suppose that there are two types of intervals (low and high) with respective demand distributions G_l and G_h , such that G_h is stochastically greater than G_l in the sense of first order stochastic dominance. Let us assume that there are T_l intervals of low demand and T_h of high demand. The firm seeks to maximize

$$\Pi_T(\mathbf{A}, \mathbf{G}) := p(T_l S_l(A_l) + T_h S_h(A_h)) - (\eta_l T_l A_l + \eta_h T_h A_h),$$

where S_i ($i = l, h$) is as in Eq. 12.5 with G replaced by G_i . Let $\eta_{i,N}^*$ be the solution to Eq. 12.4 when the demand distribution is G_i , $i = l, h$, the pool size is N , and $K = \infty$. Let $A_i(\beta) = \bar{G}_i^{-1}(\beta/p)$ be the solution to Eq. 12.6 with $\eta = \beta$ and the demand distribution G_i , $i = l, h$.

Theorem 2 Consider Eq. 12.2 with fixed N and $\beta = 0$. Then, the optimal compensation is lower in low demand periods, i.e., $\eta_{l,N}^* \leq \eta_{h,N}^*$, and the staffing level is, consequently, lower. The service level is, however, higher in low demand periods.

Theorem 3 Every optimal solution to Eq. 12.2 has $\eta_{l,N}^* = \eta_{h,N}^* = \beta$ and $N^* \geq \bar{N} := A_h(\beta)/F(\beta) = \bar{G}_h^{-1}(\beta/p)/F(\beta)$. The assigned capacity satisfies $N^*F(\beta) \wedge K_l^* = A_l(\beta)$ and $N^*F(\beta) \wedge K_h^* = A_h(\beta)$. In particular, in any optimal solution, the firm uses a cap $K_l^* = A_l(\beta)$ in the low demand period.

Theorems 2 and 3 offer two important insights. First, the ability to cap the number of active agents is crucial to controlling the firm's costs when it must guarantee a minimum earning level under time-varying demand. Without it, low demand periods would be overstaffed. Second, regardless of earning constraints, it is never optimal for the firm to offer a higher service level in high demand periods. This continues to be true even when the retail price is a decision variable; see Sect. 12.4.2.

Before moving on, a brief discussion of alternative service models is due. We could consider other revenue models than the newsvendor. All that is needed for the insights to persist is that the expected unit sales be increasing and concave in the staffing level. One could, for example, suppose that sales are given by an Erlang loss model in which inducing more agents to work results in fewer lost sales.

That said, the newsvendor is applicable in a wide variety of settings. Facing significant uncertainty in call volume, a newsvendor model provides a good approximation for call-center optimization; see e.g. Bassamboo et al. (2010). To relate this specifically to our setup, consider a call center with a single group of servers serving a single type of customers with finite patience. If the number of agents is A and the call volume comes from a distribution G , then $S(A)$ provides a good approximation for the number of calls served. The average number of calls that abandon is the expected volume minus those served. If a contract with a client compensates the call center a p for each call served, the call center is optimizing capacity so as to maximize $pS(A) - \eta A$ just as in our study above. In this way, the self-scheduling newsvendor captures, at least in first order, the challenges faced by a call-center provider such as Arise Virtual Solutions or LiveOps.

12.3.4 The Benefit of Flexible Capacity

Despite of the challenges we have highlighted thus far with self scheduling – driven by the indirect control a self-scheduling firm has of its capacity – a significant advantage of self scheduling in a service setting is that it provides greater flexibility to respond to demand variation relative to traditional models. In call centers, this advantage corresponds to the fact that agents can work for very short periods (e.g., in half hour windows), eliminating the constraint of scheduling workers for shifts lasting several hours. The benefit of flexibility could overwhelm the challenges of self scheduling. We explore this issue through a numerical experiment that both supports this intuitive claim and highlights an important subtlety.

As an approximation for Poisson arrivals, the demand distribution is normal (truncated at 0) with a mean λ_i and standard deviation $\sqrt{\lambda_i}$ for $i \in \{l, h\}$. The agents' threshold distribution is normal with mean 7 and standard deviation of 3. The retail price is 10. We impose no earnings constraint (i.e., $\beta = 0$) and fix the pool size to $N = 400$.

We consider an eight-period horizon. We manipulate the extent of time variation by changing the number T_h of high demand periods. Given the mix of high and low demand periods, we solve for the optimal wages, η_l and η_h for low and high demand periods respectively, using (12.4) and compute the service provider's profit under self-scheduling.

We then compute the firm's profit if it relinquished its flexibility in return for direct control. This is the *inflexible + direct control* benchmark. So that we can compare against the self-scheduling firm, the benchmark still pays the wages η_l and η_h determined before. Direct control means that, given these wages, it can choose the number of servers it uses. However, this number is chosen once and remains the same on all periods.⁸

In the upper panel of Fig. 12.2, we consider the case where $\lambda_l = 100$ and $\lambda_h = 200$. As intuitively expected, the value of flexibility overwhelms the relative lack of capacity control. Self scheduling is close to the inflexible + direct control benchmark when there is little temporal variability in demand but performs substantially better when there is a mix of high and low demand periods.

In the lower panel, we increase the high demand volume to $\lambda_h = 400$. Here, when the number of high demand periods exceed five, the cost of incentivizing agents to work overwhelms the value of flexibility. The pool size ($N = 400$) is tight relative to the high demand rate of $\lambda_h = 400$. The self-scheduling firm is consequently significantly understaffed, leaving a substantial amount of money on the table. The benchmark will summon many more agents. Thus, when the pool size is small relative to peak demand, flexibility may be less important than control. Said another way, it is important for a self-scheduling firm to have a pool that is large relative to peak demand.

12.4 Variants of the Base Model

We consider four extensions of our base model: (i) alternative compensation schemes in which agent earnings are tied to the volume of customers served; (ii) a price dependent newsvendor setting in which the firms sets both agent compensation and the retail prices; (iii) A setting where there is a cost associated with maintaining a larger pool; and (iv) agents that have availability threshold distributions that vary over the horizon.

⁸Fixing the staffing level for the entire horizon is admittedly extreme. However, it implies that the result do not depend on the exact sequence of high and low demand periods.

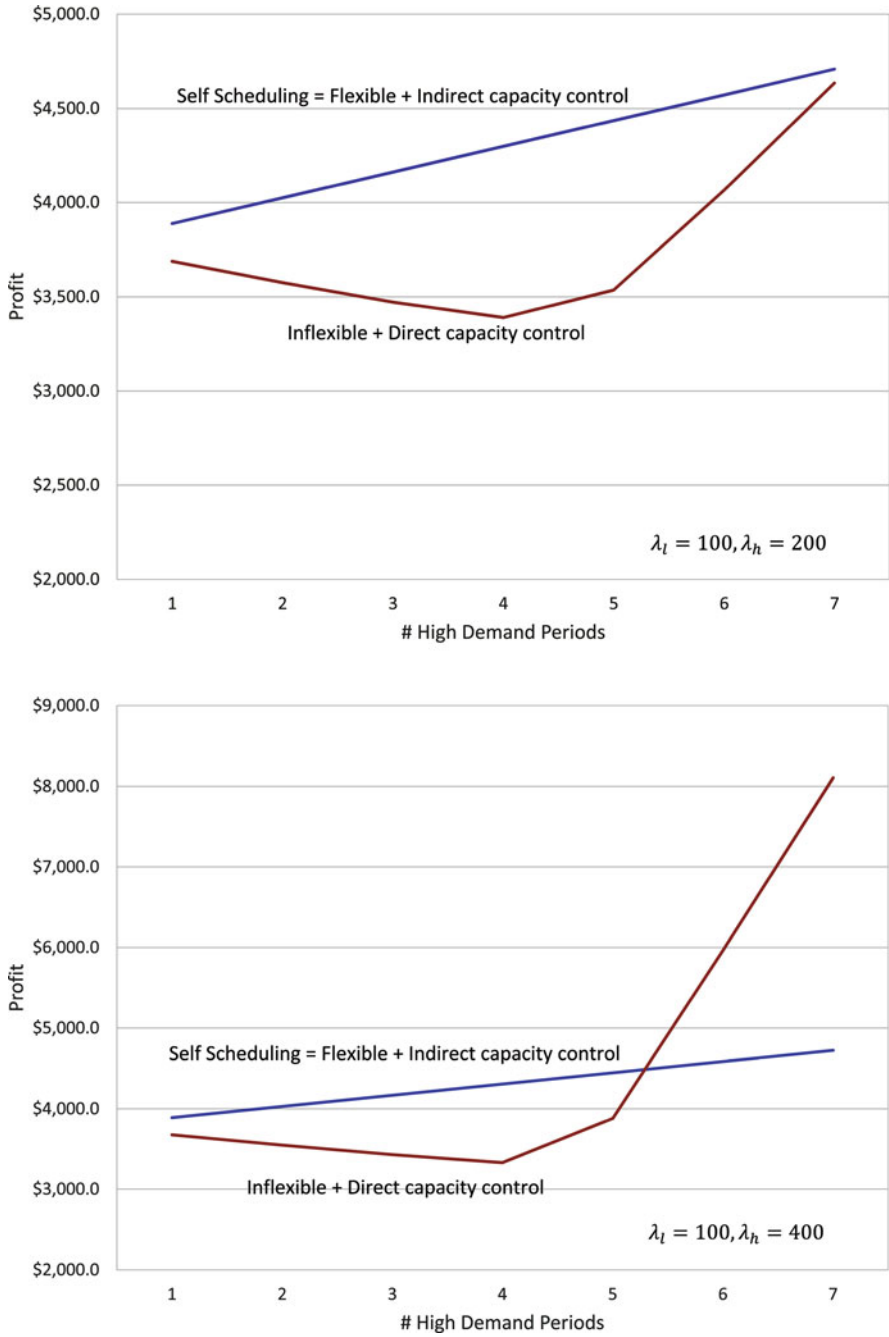


Fig. 12.2 How profit varies as the mix of high and low demand periods changes: The value of flexibility may be compromised if the pool size is not big enough relative to high demand

12.4.1 Volume-Dependent Compensation Schemes

Thus far we have assumed that the firm pays agents a fixed per-period amount η^* : an agent that signs-up to work Monday 10:00–10:30 gets η^* regardless of the number of customers served. However, many self-scheduling firms tie agent pay to utilization. Ride-sharing service such as Lyft and Uber compensate drivers by splitting fares with them. Similarly, call centers like LiveOps and Arise Virtual Solutions use piece-rate compensation or some combination of piece rate and a guaranteed per-interval minimum. The firm might reasonably prefer some sort of volume-dependent compensation. For example, a piece rate may address moral hazard issues (that we have left unmodeled) and induce agents to exert more effort. Additionally, a piece rate is easier on a firm's finances since it only pays agents when it has been paid by the client; such a consideration may be important for a nascent firm with limited resources.

Here we show that within our model with risk-neutral agents many reasonable compensation mechanisms are equivalent. Namely, that there exists a translation from one scheme to the other that generates the same outcomes in terms of staffing, service level and expected firm profit.

In particular, we will focus on *piece-rate* compensation in which an agent earns ϕ_t per completed transaction in period t . To determine how much an agent earns under such a scheme, we need to know how many transactions she completes. Let X_t^j denote the number of customers agent j serves in period t . The distribution of X_t^j depends on several factors including the number of active agents, the demand distribution G_t , and how the firm allocates work among agents. If some agents are given higher priority as demand is allocated among working agents, they will earn more money than those with lower priority. Here we assume that jobs are distributed uniformly among the active agents. (For a setting in which routing is non-uniform see Stouras et al. 2016). Thus, each agent will serve in expectation $S_t(A_t)/A_t$ customers where $S_t(\cdot)$ is as in Eq. 12.5 with G replaced by G_t . An agent's expected earnings in period t are then $\mu_t = \phi_t S(A_t)/A_t$, and the number of interested agents is $NF(\phi_t S(A_t)/A_t)$.

Recall that the number of agents interested in working in period t when the firm pays a fixed amount η_t is $NF(\eta_t)$. In comparing these values, note that under a fixed rate scheme, an agent can determine whether or not to work by considering only her own availability threshold. Under a piece rate scheme, however, an agent must consider both her threshold and what other agents are doing. We must in this case consider an equilibrium among the agents in which the number of interested agents is equal to the number that join, i.e., that $NF(\phi_t S_t(A_t)/A_t) = A_t$.

Lemma 4 Fix N , G_t , and ϕ_t . There then exists an equilibrium A_t^e , characterized by the unique solution to the equation

$$NF\left(\frac{\phi_t S_t(A_t^e)}{A_t^e}\right) = A_t^e. \quad (12.7)$$

It is a priori conceivable that this equilibrium structure introduces constraints into the firm's optimization problem or, in other words, that an optimal solution (N^*, η^*, K^*) to the firm's optimization problem is not implementable via a piece rate. The following simple argument is a proof to the contrary: the firm can move from per-interval compensation to piece-rate compensation without compromising its profits or customer service level.

Suppose that the firm is using a feasible solution (N, η, K) with $K_i \geq NF(\eta_i)$, $i = l, h$ (so that access is not really limited). The firm should offer the piece rate $\phi' = (\phi'_l, \phi'_h)$ such that

$$NF\left(\frac{\phi'_i S_i(A_i)}{A_i}\right) = A_i,$$

where $A_i = NF(\eta_i)$. Since (N, η, K) is a feasible solution to the firm's problem, it must be that $N \geq NF(\eta_i) = A_i$ so that the existence of ϕ' follows from the continuity of F . With this choice of ϕ' , the number of agents that sign-up in equilibrium is (using Lemma 4) the unique solution to $NF(\phi'_i S_i(A_i)/A_i) = A_i$ which must equal A_i by construction.

If (N^*, η^*, K^*) is an optimal solution to the firm's optimization problem Eq. 12.3, then the firm can set the piece rate at $\phi_i^* = \eta_i^* A_i / S_i(A_i)$ where $A_i = N^* F(\eta_i^*) \wedge K_i^*$. With this translation, the optimal solution (N^*, η^*, K^*) to the firm's problem with interval compensation is equivalent to the solution (N^*, ϕ^*, K^*) with piece rate compensation: (i) the number of agents interested for each interval is the same (and using the same cap, so is the number of people actually signing-up), (ii) the staffing level is the same and, hence, (iii) the expected firm profit and customer service level are the same.

There are, however, subtle differences between the piece rate and fixed compensation. First, the firm's staffing costs are deterministic under fixed compensation but these costs are variable under a piece rate. Thus a properly chosen piece rate delivers the same *expected* profit as the optimal fixed interval compensation but the *realized* profit for a given demand outcome differs.

Second, piece rate compensation lessens the impact of increasing the pool size. Under fixed rate compensation, doubling the pool size while holding the compensation rate constant will double the number of interested agents. The response to an increase in the pool size is less elastic when a piece rate is used. If the pool size is doubled while the piece rate is unchanged, the number of interested agents increases but does not double. Competition between agents could dissuade an agent, who was indifferent between working and not working under the original pool size, from participating as the pool grows.

Lastly, the cap on the number of active agents plays a different role under a piece rate than under a fixed per-period compensation. Under the latter, the firm's labor costs are fixed with regard to the demand realization. The cap serves to control this fixed cost and prevents the firm from paying for labor on which it would expect an inadequate return. From this perspective, a cap would seem to be unnecessary when the firm moves to a piece rate system. Labor is no longer a fixed cost and unutilized

capacity would appear to be costless to the firm. That unutilized capacity, however, does impose a burden on the firm. A large number of available agents reduces everyone's expected utilization and expected earnings for a given ϕ . Competition between agents undermines the firm's ability to compensate agents adequately with a relatively low piece rate. Without a cap, the firm would have to raise the piece rate, driving up its cost of serving customers.

Piece rate and fixed compensation are two extremes. Fixed compensation means agents do not face any volume risk while under a piece rate they carry all the risk. A two-part tariff (i.e., $v_t + \phi_t X_t^j$) or a piece rate with a minimum guarantee (i.e., $\max\{\kappa_t, \phi_t X_t^j\}$) offer intermediate mechanisms. A call center we have worked with pays agents a piece rate with a guaranteed minimum payment level. Uber has also been reported to guarantee an hourly rate at some time periods (Kirsner 2014). Given our analysis for piece rate it is not a surprise (and, indeed, can be easily shown) that these mechanisms are also equivalent, within our model, to fixed compensation.

12.4.2 Price-Dependent Newsvendor

Our assumption that the retail price is fixed regardless of whether demand is high or low is appropriate in some settings (e.g., a work-from-home call centers). Yet, other services (notably ride-sharing firms) raise their prices when demand increases. This calls into question one of our earlier results that customer experience a lower service in high-demand periods. In Theorem 2 we showed that if the retail price is fixed, an increase in staffing costs results in the firm picking a lower service level. If now the retail price also increases, it is not clear that it is still optimal to let customer service levels fall.

To examine these issues, we suppose that demand in a low-volume interval given a retail price p is a random variable ξ_p with distribution $G_l(x|p)$ and that ξ_p becomes smaller in the sense of first order stochastic dominance as p increases. That is, $G_l(x|p) \leq G_l(x|\hat{p})$ for all x for all $p \leq \hat{p}$. Next we assume that demand in a high-volume interval for a given price is $\theta\xi_p$ for some $\theta > 1$. The corresponding demand distribution is then $G_h(x|p) = G_l((x/\theta)|p)$. We assume that there is sufficient structure on $G_l(x|p)$ that the firm has unique pricing and staffing decisions for both low and high volume periods (see Petruzzi and Dada 1999).

Letting $S_i(A, p)$ be expected unit sales in a type $i \in \{l, h\}$ period given staffing level A and retail price p , one can show that

$$S_h(A, p) = \theta S_l(A/\theta, p).$$

For the moment, suppose the firm directly controls staffing, hiring as many agents as it wants at wage η in either type of period. Let $\hat{A}_i(\eta)$ and $\hat{p}_i(\eta)$ be the optimal staffing level and price, respectively, for a type $i = l, h$ interval. It is straightforward to show that $\hat{p}_h(\eta) = \hat{p}_l(\eta)$ and $\hat{A}_h(\eta) = \theta \hat{A}_l(\eta)$. Thus a firm which manages its

staff in a conventional fashion does not employ period-dependent pricing; it sticks with the same retail price and adjusts its staffing level to achieve the same service level in both high and low volume periods.

This will no longer hold if the firm allows agents to self schedule. Increasing staff above $\hat{A}_l(\eta)$ requires dipping further into the pool of agents which, in turn, drives up staffing cost. Higher costs then lead to a higher retail price. That is, prices surge higher in this framework not because of the market structure but because of higher costs.

The fact that the firm boosts the retail price does not yet tell us how service-level behaves and whether (or not) making the retail price endogenous leads to a departure from Theorem 2. To examine how the service level varies with demand under self scheduling, we work with a specific demand distribution. Suppose that $G_l(x|p) = x/D(p)$ where $D(p)$ is a non-negative strictly decreasing function. Since demand is uniformly distributed, the expected demand (given a price) p is $D(p)/2$. Let $\varepsilon(p) = -pD'(p)/D(p)$ be the elasticity of expected demand. We assume that $\varepsilon(p)$ is increasing, i.e., that demand becomes less elastic as the price falls.

If the firm directly controls staffing, its problem for high volume periods is

$$\max_{p,A} \{pS_h(A, p) - \eta A\}.$$

With the uniformly distributed demand, we have

$$S_h(A, p) = A - \frac{A^2}{2\theta D(p)},$$

which results in the following first order conditions:

$$\frac{A}{\theta D(p)} = 1 - \frac{\eta}{p}, \quad (12.8)$$

$$\varepsilon(p) = \frac{2}{A/(\theta D(p))} - 1. \quad (12.9)$$

Equation 12.8 is the classical critical fractile solution that ties the capacity A to a targeted service level. Equation 12.9 relates the service level to the elasticity of demand: the higher the service level, the lower the elasticity. To go the other way, a higher elasticity corresponds to a lower service level.

Substituting Eq. 12.8 into Eq. 12.9 yields an implicit expression for the optimal price \hat{p}

$$\hat{p} = \eta \frac{1 + \varepsilon(\hat{p})}{\varepsilon(\hat{p}) - 1}. \quad (12.10)$$

Equation 12.10 does not depend on θ . This is a specific instance of our more general argument above that the retail price is independent of the scale of demand when the

firm directly controls staffing. Further, the right-hand side of Eq. 12.10 is decreasing in p if $\varepsilon(p)$ is strictly increasing. Consequently, \hat{p} is increasing in η so that higher agent wages move the firm to a higher level of elasticity on the demand curve. Going back to Eq. 12.9, a higher elasticity corresponds to a lower service level. Thus when faced with higher staffing costs, the firm charges more but offers worst service.⁹

Turning to the self-scheduling setting, the firm’s problem for given a pool-size of N is

$$\max_{p,\eta} \{pS_h(NF(\eta), p) - \eta NF(\eta)\},$$

which yields the following first order conditions (the analogues to Eqs. 12.8 and 12.9)

$$\frac{NF(\eta)}{\theta D(p)} = 1 - \frac{\eta + F(\eta)/f(\eta)}{p},$$

$$\varepsilon(p) = \frac{2}{NF(\eta)/(\theta D(p))} - 1.$$

We again have (recall Eq. 12.4) that the self-scheduling firm works with an inflated marginal cost of capacity. That higher cost leads to higher price than one would have when capacity is directly controlled. If $\varepsilon(p)$ is strictly increasing, this leads to a lower service level. Further, the chosen compensation rate is now increasing in θ . Thus a firm that lets its agents self-schedule will charge more but, as in Theorem 2, offer worse service in high-demand periods.

12.4.3 When Maintaining a Larger Pool Costs More

In our basic model, the only mechanism keeping the firm from recruiting an infinite number of agents is a possible concern for agent welfare, as modeled through the earnings constraint. Suppose, alternatively, that the firm incurs a cost of $M(N)$ for maintaining a pool of size N agents. Before considering the cost of maintaining the pool, the profit in period i given N is then

$$\Pi_i(N) = p \left(\int_0^{NF(\eta_{i,N}^*)} x g_i(x) dx + NF(\eta_{i,N}^*) \bar{G}_i(NF(\eta_{i,N}^*)) \right) - \eta_{i,N}^* NF(\eta_{i,N}^*).$$

⁹This conclusion depends on $\varepsilon(p)$ being strictly increasing. If $D(p) = p^{-\bar{\varepsilon}}$, the elasticity of demand is constant at $\bar{\varepsilon}$, making \hat{p} proportional to η and the optimal service level independent of η .

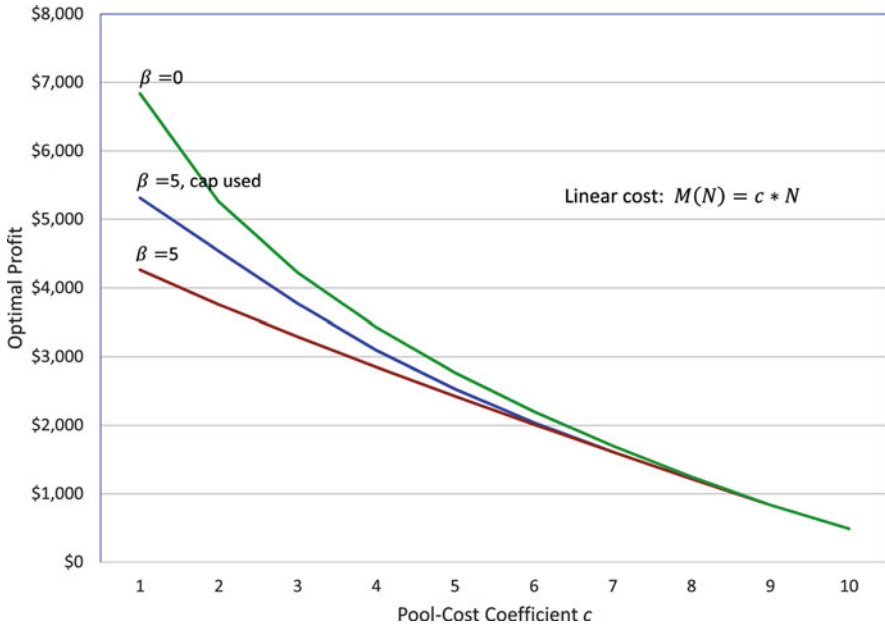


Fig. 12.3 Optimizing jointly pool size and wages with linear holding cost $M(N) = c * N$: Profits as a function of the cost coefficient c

Notice that $\eta_{i,N}^*$ depends on N . It is straightforward to show that profit is concave in N and that $\Pi_h'(N) \geq \Pi_l'(N)$. Additionally, $\Pi_l'(N) = 0$ if the earnings constraint would be binding in a type i period given pool size N .

If $M(N)$ is weakly convex the objective function (summing the profit over all periods and subtracting the pool cost) is concave implying that the maximization problem has an interior (and finite) solution N^* . Because the marginal return on agents is higher in high-demand periods, N^* is increasing in the number of high-demand periods.

A cost to maintaining the pool does not necessarily obviate the need to impose an access cap. If the cost of maintaining the pool does not increase too quickly, the service provider may still need to restrict access in low demand periods. Specifically, if $N^* > \bar{G}_l^{-1}(\beta) / F(\beta)$, the service provider would benefit from capping the number of agents working in low demand periods.

This is demonstrated in Fig. 12.3. Here we use a linear pool-cost function $M(N) = c * N$. For each value of the coefficient c , we solve the optimization problem to find the optimal staffing N^* and the optimal wages $\eta^*(N^*)$. We use a retail price $p = 10$, arrival rates $\lambda_l = 100$ and $\lambda_h = 200$. We have a total of eight periods in the horizon, with $T_h = 5$ and $T_l = 3$. Agent thresholds follow a normal distribution with mean 7 and standard deviation of 3.

We vary the coefficient c from 1 to 10 and consider three cases. In the first, there is no earnings constraint (i.e., $\beta = 0$). In the second, there is an earnings

constraint with $\beta = 5$ but the firm never limits access. The third case replicates the second but now allows for the use of an access cap. Across all three examples, we see that the higher values of c lower profitability as one would expect. Note that this is not driven just by increased pool costs; higher pool costs drive up per-period compensation which reduces staffing levels and thus revenues.

We also see that the earnings constraint decreases the provider's profit when the pool costs are not too high. This is particularly true when the firm does not limit access. (That is, an access cap still creates value.) However, when the marginal cost of increasing the pool is high, the differences between the three cases go away. At high values of c , the pool is so small that the optimal compensation is greater than β in low-demand periods even in the first scenario.

Finally, it is worth thinking about what drives the cost of maintaining the pool. Clearly, there is some overhead to tracking a large number of agents. But some costs could be endogenous and depend on how the firm manages its agents. For example, because the introduction of caps inevitably limit agent flexibility, imposing caps may decrease the willingness of agents to work for the firm in the first place.¹⁰ Within this model this could be captured by assuming that allowing caps increase the cost of recruiting and maintaining the pool, i.e., c may increase as we move from the second to the third case. In such a setting, we are over-stating the gain in moving between these settings.

12.4.4 *Period-Dependent Threshold Distributions*

In our base model, the distribution of agent threshold values is independent of the period. This is obviously unrealistic. Many people opt to work for a self-scheduling firm in part because existing obligations (e.g., having young children) make working a conventional schedule difficult. However, many of these obligations have a known schedule (e.g., the preschool gets out at the same time every weekday); this should be reflected in the distribution of threshold values.

Here, we suppose that the distribution of threshold values depends on the time interval of the horizon. Let Ω_T be the set of all time intervals. Let Ω_d and Ω_u be subsets of Ω_T such that $\Omega_d \cup \Omega_u = \Omega_T$ and $\Omega_d \cap \Omega_u = \emptyset$. An agent draws her threshold value for period t from F_d [F_u] if $t \in \Omega_d$ [$t \in \Omega_u$]. Further, we assume that F_u is stochastically larger than F_d , i.e., that $F_u(t) \leq F_d(t)$ for all t . Ω_d is

¹⁰For example, note the following complaint about UK on-demand delivery service Deliveroo: "I was working in a bar when a friend of mine started working for Deliveroo. I was sick of working until 2am and I really like cycling so I decided to join too. I wanted as many hours as possible, but cycle couriers mainly do a three-hour shift at lunchtime and three hours in the evening, because those are the busiest times for Deliveroo. . . . Its not flexible either. We used to have a system where you could swap shifts with people but they said it was too chaotic. Now you do the same shifts every week." (https://www.theguardian.com/money/2016/jun/15/he-truth-about-working-for-deliveroo-uber-and-the-on-demand-economy?CMP=share_btn_tw)

then the set of *desirable* time intervals in the sense that each agent has a higher probability of drawing a low threshold in these intervals than in the *undesirable* intervals of Ω_u . Stated another way, a given compensation rate will induce more agents to work in a desirable interval than in undesirable interval.

When there is no earnings constraint, agent preferences over periods only modestly complicate the firm's problem. Let Ω_h and Ω_l be the collection of time periods in which demand is high or low respectively. When the threshold distribution does not depend on the time period, the compensation offered in period t depends only on whether t falls in Ω_h or Ω_l . If the distribution varies with the time interval, compensation in period t depends on whether t lies in $\Omega_d \cap \Omega_h$ or $\Omega_d \cap \Omega_l$ and so on. The firm must then calculate four different compensation rates using the appropriate version of (12.4).

Things get more interesting when the firm must satisfy a non-trivial per-period earning. Now the firm offers β in every period and the question becomes how large a pool to recruit. Regardless of whether a high demand period falls in Ω_d and Ω_u , the firm would want $A_h(\beta) = \bar{G}_h^{-1}(\beta/p)$ working. The pool size necessary to achieve this staffing level is higher if the period in question is undesirable. Thus, if $\Omega_u \cap \Omega_h \neq \emptyset$, the firm sets $N = A_h(\beta)/F_u(\beta) > A_h(\beta)/F_d(\beta)$. However, assuming $\Omega_d \cap \Omega_h \neq \emptyset$, then desirable, high-volume periods will be overstaffed. Consequently, we conclude that the firm may cap access to high-volume periods if it must satisfy a non-trivial earnings constraint and the distribution of agent thresholds varies with the time period.

It is conceptually straightforward (albeit notationally cumbersome) to extend these results to having two types (say, A and B) of agents each of which has its own set of desirable and undesirable time periods. Assume that the firm offers the same compensation to both types of agents. This would be appropriate if the agents type (say, stay-at-home parents and college students) affects their availability but not their productivity. If there is no earnings constraint, the firm would now need to have eight different compensation rates that depend on the demand and whether or not the period is desirable for Type A agents, Type B agents or both. With a non-trivial earnings constraint, the action turns on how many agents of each type to recruit. The number of agents will be determined by a subset of the possible kinds of high-demand periods (e.g., periods that both types undesirable and those that Type A agents find undesirable but Type B desires but not those periods that both types desire). Caps will not be necessary in the periods that determine the staffing levels but will be in other high-demand periods.

12.5 Concluding Remarks

We studied a model in which a service provider allows its agents to choose when to work, a scheme being adopted in multiple service markets and stands in contrast to traditional models where the firm has direct control of capacity. The firm faces time varying demand over a horizon and must recruit and train a pool of agents.

The firm sets the pool size at the start of the horizon and the pool size cannot be adjusted dynamically. In each period, agents in the pool choose whether or not to work. An agent's willingness to work varies from period to period and the firm can offer a different compensation every period to attract enough workers to provide an adequate service level.

Our objective is to understand the costs of relinquishing direct control of capacity. These should be weighted against the value of flexibility that such self-scheduling models bring.

We frame the problem in terms of a newsvendor model in which demand and capacity can vary from period to period and show that the optimal decision is a variant of the classical critical-fractile formula. We show that the firm picks a lower service level than it would in a standard newsvendor setting, lowering its profits and making it harder for customers to get served. The drop in customer service is affected by the time-varying pattern of demand. In a horizon with both high and low demand periods, the provider offers agents higher pay in high demand periods and makes more capacity available, but the service level customers see is lower in high demand periods.

These issues are mitigated if the firm can recruit a large pool of agents. As the pool size grows, the firm can pay agents less and less and the gap between self scheduling and the benchmark newsvendor problem decreases. If the firm does not need to ensure a minimum earnings constraint to the agents, the firm has an incentive to make its pool of agents as large as possible. It is then able to offer relatively low wages in both high and low demand periods and still induce enough agents to work. Of course, agents are worse off in this case. If there is a constraint on agent earnings, however, the firm is limited in its ability to reduce wages. This is costly to the firm because it pays has to pay higher wages and because that higher pay may induce too many agents to work in some periods. Capping the number of active agents that can work in some periods addresses this problem. The firm chooses its pool large enough that it can get, on a period-by-period basis, exactly the number of agents it needs at the guaranteed wage. However, to control its cost, the firm caps the number of agents working in some periods. The implication is that agents must sacrifice some scheduling flexibility in order to guarantee a minimum compensation level.

For most of our analysis we have assumed the firm offers a fixed compensation for each time interval the agent is active. However, we have demonstrated that the firm can achieve identical results through alternative schemes, such as a piece-rate compensation, that depend on the number of customers an agent serves. We have also presented extensions that consider price-dependent demand and period-dependent threshold distributions.

Appendix

Proof of Theorem 1

The right hand side of Eq. 12.4 is smaller than that of Eq. 12.6 so that, since G is increasing in its argument, we must have that $NF(\eta^*) \leq A(\eta^*)$. For the profit comparison, notice that $NF(\eta^*)$, would generate in the benchmark problem the profit $\Pi(NF(\eta^*), G)$. Since $A(\eta^*)$ is, by definition, the optimal solution for the fixed wages η^* , it must be the case that $\Pi(A(\eta^*), G) \geq \Pi(NF(\eta^*), G)$. \square

Proof of Lemma 2

We first show the monotonicity results for the compensation followed by the service level and, finally, the profits. It is useful to re-write Eq. 12.4 as

$$\bar{G}(NF(\eta^*)) = \frac{\eta^* + F(\eta^*)/f(\eta^*)}{p}. \tag{12.11}$$

Compensation That compensation strictly increases with p and strictly decreases with N is evident from Eq. 12.11. Consider for instance p . Suppose to reach a contradiction that, as p increases, the compensation η^* actually decreases. Then, the right-hand-side of Eq. 12.11 decreases by the monotonicity of F/f so that the left-hand side $\bar{G}(NF(\eta))$ must also decrease. This would entail (since F is increasing and \bar{G} is decreasing) that η^* increases with p which is a contradiction.

To prove that the compensation increases with the agent availability distribution F notice that, assuming F_2 dominates F_1 in the reverse hazard rate ordering,

$$\eta_2^* + \frac{F_2(\eta_2^*)}{f_2(\eta_2^*)} \leq \eta_2^* + \frac{F_1(\eta_2^*)}{f_1(\eta_2^*)}. \tag{12.12}$$

Further, if F_1 is smaller than F_2 in the reverse hazard rate order, then it is also smaller in the regular stochastic ordering sense, $\bar{F}_1(x) \leq \bar{F}_2(x)$ (or $F_1(x) \geq F_2(x)$), so that, since \bar{G} is decreasing,

$$\bar{G}(NF_1(\eta_2^*)) \leq \bar{G}(NF_2(\eta_2^*)) \tag{12.13}$$

By Eq. 12.11,

$$\bar{G}(NF_1(\eta_1^*)) = \frac{\eta_1^* + F_1(\eta_1^*)/f_1(\eta_1^*)}{p} \text{ and } \bar{G}(NF_2(\eta_2^*)) = \frac{\eta_2^* + F_2(\eta_2^*)/f_2(\eta_2^*)}{p},$$

so that combining Eqs. 12.12 and 12.13 we have

$$\bar{G}(NF_1(\eta_2^*)) \leq \bar{G}(NF_2(\eta_2^*)) = \eta_2^* + \frac{F_2(\eta_2^*)}{f_2(\eta_2^*)} \leq \eta_2^* + \frac{F_1(\eta_2^*)}{f_1(\eta_2^*)}. \quad (12.14)$$

Assume now, to reach a contradiction, that $\eta_2^* < \eta_1^*$. Then, since the right hand side of Eq. 12.14 increases strictly in the compensation and the left hand side strictly decreases, we would get that

$$\bar{G}(NF_1(\eta_1^*)) < \bar{G}(NF_1(\eta_2^*)) \leq \eta_2^* + \frac{F_1(\eta_2^*)}{f_1(\eta_2^*)} < \eta_1^* + \frac{F_1(\eta_1^*)}{f_1(\eta_1^*)},$$

which contradicts Eq. 12.11. We conclude that $\eta_1^* \leq \eta_2^*$.

Service level The fact that the service level increases with p is evident from the optimal fractile formula Eq. 12.11 and from the fact, already proved, that the optimal compensation increases with p . Similar is the observation that the optimal service level increases with N .

Profits To show that profits increase with the pool size N , take $N_2 > N_1$. Let $\eta_{N_1}^*$ be the optimal compensation level at N_1 . Since $N_2 > N_1$, $N_2F(\eta_{N_1}^*) > N_1F(\eta_{N_1}^*)$. In particular, we can find $\bar{\eta} < \eta_{N_1}^*$ such that $N_2F(\bar{\eta}) = N_1F(\eta_{N_1}^*)$. With this $\bar{\eta}$, then, the firm gets the same staffing level under $(N_2, \bar{\eta})$ as under $(N_1, \eta_{N_1}^*)$ and, in turn, the same revenue. The staffing costs are smaller under $(N_2, \bar{\eta})$ since $\bar{\eta}N_2F(\bar{\eta}) = \bar{\eta}N_1F(\eta_{N_1}^*) < \eta_{N_1}^*N_1F(\eta_{N_1}^*)$. Thus, the pair $(N_2, \bar{\eta})$ generates a higher profit for the firm than the pair $(N_1, \eta_{N_1}^*)$. In particular, $(N_2, \eta_{N_2}^*)$ generates higher profits than $(N_1, \eta_{N_1}^*)$.

An identical argument is used to study the effect of an increase (in the sense of reverse hazard ordering) in the availability distribution starting with the observation that, since reverse hazard rate ordering implies stochastic ordering, $NF_1(\eta_{F_2}^*) \geq NF_2(\eta_{F_2}^*)$ where N is fixed and $\eta_{F_2}^*$ is the optimal compensation under F_2 . If $NF_1(\eta_{F_2}^*) = NF_1(\eta_{F_2}^*)$, then under F_1 , $\eta_{F_2}^*$ generates the same profit as the optimal solution for F_1 and, in particular, the optimal profit under F_1 is higher. If the inequality is strict, i.e., $NF_1(\eta_{F_2}^*) > NF_2(\eta_{F_2}^*)$ we can proceed, as before, by finding $\bar{\eta}$ that generates the same staffing and revenue but lower staffing cost. \square

Proof of Lemma 3

Let η_N^* be the optimal compensation in Eq. 12.4 when the pool size is N . Suppose that N is such that $\eta_N^* > \beta$. By Lemma 2, the firm's profits are strictly increasing in N and the compensation is decreasing in N . Thus, the firm will optimally increase N (and decrease η_N^*) until it hits β and we conclude that any optimal solution must have $\eta_N^* = \beta$. The firm's optimal N , is then given by maximizing (over N), the profits

$$\Pi(NF(\beta), G) = pS(NF(\beta)) - \beta NF(\beta).$$

This is a standard newsvendor problem so that the optimal level of N is given by the (unique) solution to $\bar{G}(NF(\beta)) = \beta/p$, or, equivalently, $N^* = \bar{G}^{-1}(\beta/p)/F(\beta)$. If the pool size is set at $N^* = \bar{N} := \bar{G}^{-1}(\beta/p)/F(\beta)$ no caps are needed.

The firm also has optimal solutions with $N^* > \bar{N}$ but in that case it must use a cap. To prove this take $N \neq N^*$ with $\eta_N^* < \beta$. The firm, to meet, the earnings constraint must increase the compensation to β in which case $NF(\beta)$ agents sign up and the firm's profit is given by

$$\Pi(NF(\beta), G) = p \int_0^{NF(\beta)} xg(x) dx + pNF(\beta) \bar{G}(NF(\beta)) - \beta NF(\beta).$$

Recall that

$$\Pi(A(\beta), G) = p \left(\int_0^{A(\beta)} xg(x) dx + A(\beta) \bar{G}(A(\beta)) \right) - \beta A(\beta) = p \int_0^{A(\beta)} xg(x) dx,$$

where we use the fact that, by definition, $\bar{G}(A(\beta)) = \beta/p$. There are two cases to consider depending on whether $NF(\beta) > A(\beta)$ or $NF(\beta) < A(\beta)$. The case that $NF(\beta) = A(\beta)$ is ruled out by the assumption that $N \neq N^*$. Suppose that $NF(\beta) > A(\beta)$ (the other case is argued identically).

$$\Pi(NF(\beta), G) - \Pi(A(\beta), G) = p \int_{A(\beta)}^{NF(\beta)} xg(x) dx + pNF(\beta) \bar{G}(NF(\beta)) - \beta NF(\beta).$$

Notice that

$$p \int_{A(\beta)}^{NF(\beta)} xg(x) dx \leq pNF(\beta) (\bar{G}(A(\beta)) - \bar{G}(NF(\beta))),$$

Thus,

$$\begin{aligned} \Pi(NF(\beta), G) - \Pi(A(\beta), G) \\ \leq pNF(\beta) (\bar{G}(A(\beta)) - \bar{G}(NF(\beta))) + pNF(\beta) \bar{G}(NF(\beta)) - \beta NF(\beta) = 0, \end{aligned}$$

where we used the fact that $\bar{G}(A(\beta)) = \beta/p$. In fact, since $NF(\beta) > A(\beta)$,

$$pNF(\beta) (\bar{G}(A(\beta)) - \bar{G}(NF(\beta))) > p \int_{A(\beta)}^{NF(\beta)} xg(x) dx,$$

we can conclude that

$$\Pi(NF(\beta), G) - \Pi(A(\beta), G) < 0,$$

so that the firm is better off with the cap. By the definition of $A(\beta)$ it is immediate that $A(\beta)$ is the optimal cap. \square

Proof of Theorem 2

Here we fix N and omit it from the subscript. Recall that η_l^* and η_h^* are characterized through the equations

$$\bar{G}_h(NF(\eta_h^*)) = \frac{\eta_h^* + F(\eta_h^*)/f(\eta_h^*)}{p} \quad \text{and} \quad \bar{G}_l(NF(\eta_l^*)) = \frac{\eta_l^* + F(\eta_l^*)/f(\eta_l^*)}{p}.$$

Suppose, to reach a contradiction, that $\eta_h^* < \eta_l^*$. Then, using the log-concavity of F (which implies, in particular, that F/f is increasing), we have that

$$\bar{G}_h(NF(\eta_h^*)) = \frac{\eta_h^* + F(\eta_h^*)/f(\eta_h^*)}{p} < \frac{\eta_l^* + F(\eta_l^*)/f(\eta_l^*)}{p} = \bar{G}_l(NF(\eta_l^*)). \tag{12.15}$$

Since F and G have strictly positive densities $F(\eta_h^*) < F(\eta_l^*)$ so that (since \bar{G} is strictly decreasing) $\bar{G}_l(NF(\eta_h^*)) > \bar{G}_l(NF(\eta_l^*))$. Using the assumed stochastic ordering we then have that

$$\bar{G}_h(NF(\eta_h^*)) \geq \bar{G}_l(NF(\eta_h^*)) > \bar{G}_l(NF(\eta_l^*)),$$

which is a contradiction to Eq. 12.15. It must be then that $\eta_h^* \geq \eta_l^*$. Consequently, the staffing levels satisfy $NF(\eta_h^*) \geq NF(\eta_l^*)$. Finally, since F/f is increasing, $\eta_h^* + F(\eta_h^*)/f(\eta_h^*) \geq \eta_l^* + F(\eta_l^*)/f(\eta_l^*)$ and

$$G_h(NF(\eta_h^*)) = 1 - \frac{\eta_h^* + F(\eta_h^*)/f(\eta_h^*)}{p} < 1 - \frac{\eta_l^* + F(\eta_l^*)/f(\eta_l^*)}{p} = G_h(NF(\eta_h^*)),$$

so that the service level is higher in low demand periods. \square

Proof of Theorem 3

Consider a pool size $N < A_h(\beta)/F(\beta) = \bar{G}_h^{-1}(\beta/p)/F(\beta)$. We will show that such a level cannot be optimal. There are two cases to consider depending on how $\eta_{h,N}^*$ in Eq. 12.4 relates to β .

Case I: $\eta_{h,N}^* < \beta$ In this case, in the absence of the earnings constraint the firm would optimally choose a compensation level below β . For a given level N , we can treat both types of periods (high and low) independently and, by Lemma 3, the firm sets its compensation levels at β and utilizes a cap $K_l = A_l(\beta)$ in the low demand periods where $A_l(\beta)$ is the solution to Eq. 12.6 with demand distribution G_l . In this case the cap in high demand periods is unnecessary because $NF(\beta) < A_h(\beta)$.

The active capacity is then $NF(\beta) \wedge A_h(\beta)$ and $NF(\beta) \wedge A_l(\beta)$ in the high and low demand periods. Thus, the firm's profits for values of $N < A_h(\beta)/F(\beta)$ with $\eta_{h,N}^* < \beta$ is given by

$$\bar{\Pi}(N) := T_h \Pi(NF(\beta) \wedge A_h(\beta), G_h) + T_l \Pi(NF(\beta) \wedge A_l(\beta), G_l).$$

Notice that $\bar{\Pi}(N)$ is increasing in N . Since $\eta_{h,N}^*$ is decreasing in N , it continues to hold that $\eta_{h,N}^* < \beta$ as we increase N . Thus, the firm's profit follows $\bar{\Pi}(N)$ and is increasing in N and, in particular, any optimal solution must have $N^* \geq A_h(\beta)/F(\beta)$. If the firm chooses $N = \bar{G}_h^{-1}(\beta/p)/F(\beta)$ no cap is needed at the high demand period because $NF(\beta) = A_h(\beta)$. A cap is needed in the low demand period unless $\eta_{l,N}^* = \eta_{h,N}^* = \beta$ (notice that by Theorem 2 it is always the case that $\eta_{l,N}^* \leq \eta_{h,N}^*$).

Case II: $\eta_{h,N}^* > \beta$ Since the earnings constraint is not binding the firm will use $\eta_{h,N}^*$ as the optimal compensation in the high demand period. By Lemma 2, the firm could increase its profit in high demand period by increasing N . If also, $\eta_{l,N}^* > \beta$, the same applies to low demand periods so that strictly increasing N is optimal. If $\eta_{l,N}^* < \beta$ the firm uses a cap in the low demand period and, as before, the firm's profit are increasing in N .

Thus, as long as N is such that $\eta_{h,N}^* > \beta$, the firm can increase its profits by increasing N . Let \tilde{N} be the smallest pool size such that $\eta_{h,N}^* = \beta$. (Recall that we treat N as continuous variable so that such a \tilde{N} exists). It must be the case that $\tilde{N} \geq A_h(\beta)/F(\beta) = \bar{G}_h^{-1}(\beta/p)/F(\beta)$. Otherwise, $\bar{G}_h(\tilde{N}F(\beta)) < \beta/p$ but, at the same time, being a solution to Eq. 12.4, $\eta_{h,N}^* = \beta$ satisfies $\bar{G}_h(\tilde{N}F(\beta)) = (\beta + F(\beta)/f(\beta))/p \geq \beta/p$ which is a contradiction. We conclude that $\tilde{N} \geq A_h(\beta)/F(\beta)$.

Finally, by Lemma 2 if \tilde{N} is strictly greater than $A_h(\beta)/F(\beta)$, the firm will set a cap to $K_l = A_l(\beta)$ and $K_h = A_h(\beta)$. As above, the firm can then decrease N until it hits $A_h(\beta)/F(\beta)$ without decreasing its profits. At this point no cap is needed in the high demand period because $N = A_h(\beta)/F(\beta)$ but it is needed in the low demand periods if $\eta_{l,N}^* < \eta_{h,N}^* = \beta$. □

Proof of Lemma 4

Consider the function $h(x) := NF(\phi S(x)/x) - x$. It is easily verified that $S(x)/x \rightarrow 1$ as $x \rightarrow 0$, so that $h(x) = NF(\phi S(x)/x) - x \rightarrow NF(\phi)$ as $x \rightarrow 0$. Since F is bounded by 1 we have, as $x \rightarrow \infty$, that $h(x) \rightarrow -\infty$. Combined, we just established that both $h(x) \rightarrow -\infty$ as $x \rightarrow \infty$ and $h(x) \rightarrow NF(\phi)$ as $x \rightarrow 0$. Since F and G have densities the function $h(x)$ is continuous on $(0, \infty)$ so that there must exist x_0 such that $h(x_0) = 0$. The fact that this point is unique then follows from the fact that h is monotone decreasing. Indeed, since $S'(x) = \bar{G}(x)$,

$$h'(x) = Nf\left(\phi \frac{S(x)}{x}\right) \phi \frac{S'(x) - 1}{x^2} - 1 = -Nf\left(\phi \frac{S(x)}{x}\right) \phi \frac{1 - \bar{G}(x)}{x^2} - 1 < 0.$$

□

References

- Allon G, Bassamboo A, Çil E (2012) Large-scale service marketplaces: the role of the moderating firm. *Manag Sci* 58(10):1854–1872
- Bassamboo A, Randhawa R, Zeevi A (2010) Capacity sizing under parameter uncertainty: safety staffing principles revisited. *Manag Sci* 56(10):1668–1686
- Bergstrom T, Bagnoli M (2005) Log-concave probability and its applications. *Econ Theory* 26:445–469
- Cachon G, Daniels K, Lobel R (2017) The role of surge pricing on a service platform with self-scheduling capacity. *Manuf Serv Oper Manag* 19(3):368–384
- Campbell H (2015) Delivering for doordash: what was my doordash orientation like? The ride share guy <http://therideshareguy.com/delivering-for-doordash-what-was-my-doordash-orientation-like>. Accessed 22 Feb 2018
- Hall J, Krueger A (2016) An analysis of the labor market for Uber’s driver-partners in the United States. National Bureau of Economic Research. Working paper No. 22843
- Ibrahim R (2018) Managing queueing systems where capacity is random and customers are impatient. *Prod Oper Manag* 27(2):234–250
- Kalanick T (2012) Surge pricing follow up. <http://blog.uber.com/2012/01/03/surge-pricing-followup/>. Accessed 12 Jan 2012
- Kirsner S (2014) What happened when Boloco founder John Pepper became an Uber driver. http://www.boston.com/business/technology/innoeco/2014/02/what_happened_when_a_boston_en.html. Accessed 22 Feb 2018
- Laffont J, Martimort D (2009) The theory of incentives: the principal-agent model. Princeton and Oxford University Press, Princeton/Oxford
- Levine D, McBride S (2015) Uber, Lyft face crucial courtroom test over driver benefits. Reuters <https://www.reuters.com/article/us-uber-lyft-workers/uber-lyft-face-crucial-courtroom-test-over-driver-benefits-idUSKBN0L11BN20150128>. Accessed 22 Feb 2018
- Moreno A, Terwiesch C (2015) Doing business with strangers: reputation in online service marketplaces. *Inf Syst Res* 25(4):865–886
- Netessine S, Yakubovich V (2012) The Darwinian workplace. *Harv Bus Rev* 90(5):25
- Petruzzi N, Dada M (1999) Pricing and the newsvendor problem: a review with extensions. *Oper Res* 47(2):183–194

- Riquelme C, Banerjee S, Johari R (2015) Pricing in ride-share platforms: a queueing-theoretic approach. SSRN Working paper, https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2568258. Accessed 22 Feb 2018
- Salanie B (1997) The economics of contracts: a primer. The MIT Press, Cambridge
- Stouras K, Girotra K, Netessine S (2016) First ranked first to serve: strategic agents in a service contest. INSEAD. Working paper
- Taylor T (2018) On-demand service platforms. Manufacturing and service operations management. Published online: 23 July 2018. <https://doi.org/10.1287/msom.2017.0678>

Chapter 13

On Queues with a Random Capacity: Some Theory, and an Application



Rouba Ibrahim

Abstract One standard assumption in workforce management is that the firm can dictate to workers when to show up to work. However, that assumption is challenged in modern business environments, such as those arising in the sharing economy, where workers enjoy various degrees of flexibility, including the right to decide when to work. For example, a ride-sharing service cannot impose on its drivers to be on the road at specific times; similarly, a virtual call-center manager cannot direct her agents to be available for select shifts. When self-scheduling is allowed, the number of workers available in any time period is uncertain. In this chapter, we are concerned with the effective management of service systems where capacity, i.e., the number of available agents, is random. We rely on a queueing-theoretic framework, because customers are time-sensitive and delays are ubiquitous in the services industry, and focus on the performance analysis and control of a queueing system with a random number of servers. In particular, we begin by surveying some theoretical results on the control of queueing systems with uncertainty in parameters (here, the number of servers). Then, we illustrate how to apply those theoretical results to study the problems of staffing and controlling queueing systems with self-scheduling servers and impatient, time-sensitive, customers.

13.1 Introduction

Nowadays, there seems to be “an Uber for everything,” e.g., for food delivery, parking, haircuts, domestic cleaning, etc. Such on-demand services are typically provided by online service platforms, which match consumers with workers who are willing to perform said services for a certain fee. The rise of those on-demand platforms has transformed the ways in which services are sought and delivered. For example, rather than managing a team of employees, an on-demand service provider

R. Ibrahim (✉)
University College London, London, UK
e-mail: rouba.ibrahim@ucl.ac.uk

must manage a virtual pool of independent contractors who have a legal right to various degrees of flexibility, e.g., in deciding which tasks to perform, or in setting their own schedules.

The effective management of innovative businesses in the sharing economy presents new challenges for practitioners and academics alike. To illustrate, let us consider the call-center industry. In recent years, virtual call centers have become increasingly prevalent (Vocalcom 2014). In virtual call centers, such as LiveOps (liveops.com) or Arise (arise.com), self-scheduling agents are usually independent contractors who have no requirement on the least number of working hours to be fulfilled, and are free to choose their own working periods in 30-min intervals. The standard staffing question (how many agents should the firm have?) is especially relevant in virtual call centers. This is because managing those systems involves two different time scales, and the agent pool size cannot be easily adjusted at very short time notice: (i) weeks ahead of time, typically 4–10 weeks, the system manager selects the total staffing level in the system to allow sufficient time for agent training and qualification; and, (ii) days or, in many cases, just hours ahead of time, agents select their own schedules. Since the agent population is both remote and large, up to hundreds of agents, system managers cannot simply solicit their agents' scheduling preferences ahead of time. For example, hiring decisions in virtual call centers often do not even involve a face-to-face interview. Moreover, the promised scheduling flexibility constitutes the main appeal of these jobs, and cannot be simply restricted by the firm. Therefore, hiring the right number of self-scheduling agents which scales appropriately to fit customer needs is a fundamental challenge.

Similar operational challenges arise in other service contexts as well. Amazon Flex (flex.amazon.com) relies on independent contractors to deliver Amazon Prime Now packages, which have a short delivery deadline, usually 1–2 h. Those delivery workers enjoy the flexibility of setting their preferred delivery times. Ride-sharing services, such as Uber (Uber.com) or Lyft (lyft.com), also allow their drivers to self-schedule. They use “surge pricing” to ensure the participation of a sufficient number of drivers in different time periods. While each of those settings poses unique operational challenges, agents may be viewed as being strategic in each. That is, they are decision makers who choose whether or not to be available for work in a given shift based on their individual preferences or availabilities.

We are interested in studying the effective operational management of such service systems. To do so, we adopt a queueing-theoretic framework. Relying on queueing models is natural, in our setting, because customers are time-sensitive and delays are ubiquitous in the services industry. To wit, there is a broad literature in queueing theory which studies the problems of staffing and controlling large-scale service systems; e.g., for surveys of applications in call-center management, see Gans et al. (2003) and Akşin et al. (2007). Much of that body of research formulates recommendations based on queueing models with several realistic features, such as time-varying parameters and non-standard network structures. However, one prevalent assumption in those models is that the number of servers is deterministic. As such, the realized staffing level in any given time period is assumed to be equal to the planned staffing level for that period. In contrast, with self-scheduling agents,

a firm cannot simply impose on its workers to show up to work at a given time. In other words, the number of agents in any given shift is uncertain, i.e., it must be modelled as a random variable instead.

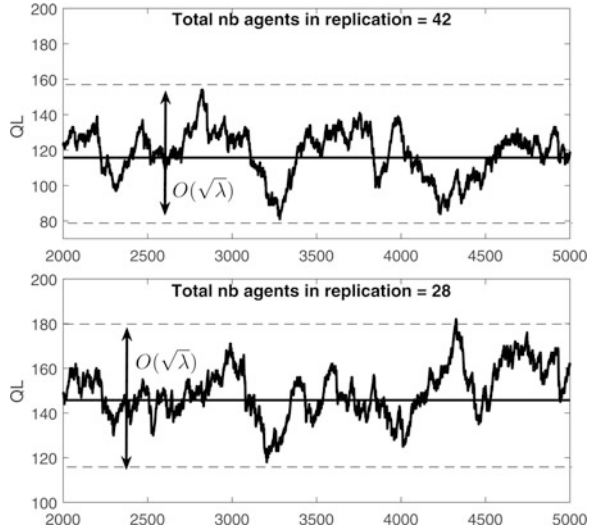
This chapter studies optimal staffing and control decisions in queueing systems with a random number of servers. Thus, we can position our work, broadly, as being part of the literature on controlling queueing systems with model-parameter uncertainty. This literature can be classified into two main categories: The first category aims at reducing parameter uncertainty through better forecasting (Shen and Huang 2008; Aldor-Noiman et al. 2009; Ibrahim and L'Ecuyer 2013, etc.). The second category, which is more closely related to our approach, investigates effective decision-making in the context of queueing systems with uncertain parameters (Harrison and Zeevi 2005; Whitt 2006b; Bassamboo and Zeevi 2009; Gans et al. 2015, etc.). With that in mind, our aim in this chapter is two-fold. First, we provide some required theoretical background. Specifically, in Sect. 13.2, we survey recent papers which propose approximations to queueing systems with uncertain parameters; those approximations are grounded in many-server heavy-traffic limits. Second, we illustrate how those theoretical results may be applied; we devote all remaining sections to that aim. Specifically, we describe some results from Ibrahim (2017a) who studies the operational management of queueing systems with self-scheduling agents, using both short-term and long-term controls.

13.2 Theoretical Background: Queues with Uncertain Parameters

At a high level, the analysis of queueing systems with uncertain parameters is complicated, for the most part, because it involves two “layers” of variability: (i) *stochastic variability*, for any realized value of the underlying uncertain parameter, since e.g., interarrival, service, and patience times are random; and (ii) *parameter uncertainty*, since the parameter itself, e.g., the number of servers in our setting, is random. Because of that analytical complexity, and because we are primarily interested in studying the operations of large service systems, it is useful to rely on many-server heavy-traffic limiting regimes, which typically simplify the analysis and yield valuable insight. In particular, performance measures of interest, e.g., the expected queue length in our setting, are approximated by limits of appropriate sequences, where the arrival rate is allowed to grow without bound. To rigorously justify the appropriateness of such approximations, we have to quantify their corresponding errors, asymptotically in large systems. Specifically, we have to determine how the orders of magnitudes of those errors grow as the system size (or the arrival rate) increases. Here, we focus on two problem formulations, corresponding to two regimes, which have been proposed in the literature.

Stochastic-fluid approximation The first formulation assumes that uncertainty effects dominate stochastic fluctuations. That is, stochastic fluctuations may be ignored, in large systems, and one can focus solely on uncertainty effects. The

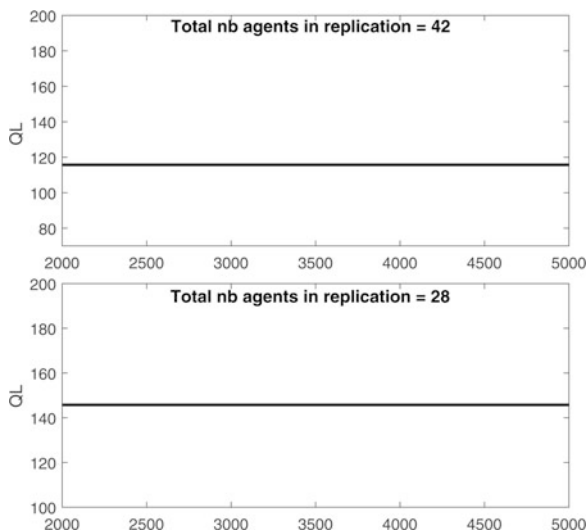
Fig. 13.1 Queue lengths seen by arriving customers in two independent simulation replications of an $M/M/N + M$ model where N follows a truncated $Nor(72, 30)$ distribution and $\lambda = 100$



intuition is that detailed fluctuations of stochastic processes (item i above) are typically realized on a short time scale, whereas parameter uncertainty (item ii above) is typically realized on a longer time scale. So, if uncertainty effects are “large,” then stochastic fluctuations become less critical in describing performance in the system. For example, with self-scheduling servers, variability in the numbers of servers is realized from shift to shift, i.e., over the hourly time scale. In contrast, variability due to the randomness in arrivals, service times, and times to abandon, is realized on a shorter time scale, e.g., over the course of a few seconds or minutes. In this case, if the uncertainty in the number of servers is larger than the order of stochastic fluctuations, then one can derive e.g., cost-minimizing staffing levels, by solving a *stochastic-fluid* optimization problem which is a fluid-type problem where only parameter uncertainty is accounted for. This type of approximation was first proposed in Harrison and Zeevi (2005) for the control of a multi-class queueing system, and considered in e.g., Bassamboo et al. (2010) for random arrival rates, and Dong and Ibrahim (2017) for a random number of servers.

For further emphasis, that difference in time scales can be visualised in e.g., computer simulations of the system. In particular, to simulate the system, one first draws a random variate from the distribution of the number of servers, and then simulates a queueing model with that realization for the number of servers. In another (independent) simulation run, one draws another realization from that same distribution, simulates the system for that new realization, and so on. The variability in the number of servers across multiple simulation runs is due to the long-run parameter uncertainty, whereas the variability within a given run, e.g., fluctuations of the queue length around its average in that run, is due to short-term stochastic fluctuations. We illustrate stochastic-fluid approximations in Figs. 13.1 and 13.2. In Fig. 13.1, we plot simulation sample paths of the queue lengths seen

Fig. 13.2 Average queue lengths in two independent simulation replications of an $M/M/N + M$ model where N follows a truncated $\text{Nor}(72, 30)$ distribution and $\lambda = 100$



by arriving customers in a queueing system where the number of servers is random (here, assumed to follow a truncated normal distribution), and where the system is, on average, overloaded. The solid horizontal line in each subfigure corresponds to the average queue length in that simulation run. As expected, conditional on the realized number of servers, the magnitude of stochastic fluctuations around the average is on the order of the square-root of the arrival rate (Garnett et al. 2002), while the average queue-length itself is on the order of the arrival rate. In Fig. 13.2, we “ignore” stochastic fluctuations, and approximate the queue length seen upon arrival by the average queue length in that run. We note that this average changes depending on the run, i.e., depending on the specific realization for the number of servers in that particular run. Relying on a stochastic-fluid approximation amounts to focusing on that cross-run variation in the means, while ignoring the more refined stochastic fluctuations within a given simulation run.

Fluid approximation The second formulation assumes that both uncertainty effects and stochastic fluctuations are negligible, and can be ignored. In this regime, we derive, e.g., the optimal staffing policy, by solving a deterministic *fluid* optimization problem instead. To illustrate, this corresponds to additionally ignoring the variations in the average queue-lengths across the different runs in Fig. 13.2. Whitt (2006a) conjectured the existence of a deterministic fluid limit for general overloaded queueing systems. That fluid limit was later established in Kang et al. (2010) and Zhang (2013). While crude fluid approximations are generally less accurate than their stochastic-fluid counterparts, they remain very useful because they usually have a remarkably simple form. Moreover, they are extremely accurate in many cases, so that there may be no tangible advantage from considering more refined approximations.

Dong and Ibrahim (2017) compare the asymptotic accuracies, i.e., the orders of magnitude of errors, corresponding to both stochastic fluid and fluid approximations in a system with a random number of servers. To summarize their main result, which parallels the result in Bassamboo et al. (2010) for random arrival rates: When the variance of the number of servers is asymptotically large, in particular larger than the square-root order of stochastic fluctuations, the system may be considered to be in an *uncertainty-dominated regime* where stochastic-fluid approximations are remarkably accurate. Moreover, the more variable the number of servers, the more accurate are those stochastic-fluid approximations. In contrast, if that variance is asymptotically small, in particular at most equal to the square-root order of stochastic fluctuations, then the system may be considered to be in a *variability-dominated regime*, where there is no tangible benefit from using stochastic-fluid approximations over fluid approximations.

13.2.1 Self-Scheduling Servers: A Binomial Distribution

To model self-scheduling agent behavior, it is natural to assume that there is a pool of agents, of size n , and that each agent from that pool makes an independent decision to join a shift j with a given probability, r_j . In this case, the random number of servers in shift j , which we denote by N_j , has a binomial distribution, $\text{Bin}(n, r_j)$, where n is the number of trials and r_j is the success probability. Because $\text{Var}[N_j] = \sqrt{nr_j(1-r_j)}$, the variance is on the square-root order, i.e., it is of the same order as stochastic fluctuations in the system. Thus, there should be no advantage in using the stochastic-fluid model, over the fluid model, when the system is large. Intuitively, this is because the binomial distribution “concentrates” around its mean, nr_j , when n is large. We can formally prove this intuition (here, we focus on a single shift since the results easily extend to multiple shifts). To do so, we restrict attention to exponentially-distributed service times and a Poisson arrival process. In particular, service times are independent and identically distributed (i.i.d.) random variables with an exponential distribution and mean $1/\mu$. We assume, without loss of generality, that $\mu = 1$; this amounts to measuring time in units of mean service times. We assume that each customer will abandon if he is unable to start service before a random amount of time, which we refer to as his patience time. Abandonment makes the system stable, irrespective of the realized numbers of servers. There is unlimited waiting space, and we use the first-come-first-served (FCFS) service discipline.

We consider a sequence of queueing models indexed by the arrival rate λ , and study system performance as λ increases without bound. The number of servers in the λ th system is $N_\lambda \sim \text{Bin}(n_\lambda, r)$. We assume that the traffic intensity $\rho \equiv \lambda/\mathbb{E}[N_\lambda] = \lambda/rn_\lambda$ remains fixed as λ increases. Let Q_{N_λ} denote the steady-state queue length and α_{N_λ} the net customer abandonment rate in the $M/M/N_\lambda + GI$ queue (abandonment makes the system stable). We refer to the cases with $\rho > 1$, $\rho < 1$, and $\rho = 1$ as the overloaded, underloaded, and quality-and-efficiency driven

(QED) regimes, respectively. Since N_λ is random, an $M/M/N_\lambda + GI$ system with e.g., $\rho > 1$ may or may not be overloaded, i.e., having $\lambda > N_\lambda$. Let \bar{q}_ρ and $\bar{\alpha}_\rho$ be the fluid approximations for the queue length and net abandonment rates with a traffic intensity ρ . The following theorem establishes the asymptotic accuracy of fluid approximations with a binomially-distributed number of servers.

Theorem 1 Consider an $M/M/N_\lambda + GI$ queueing model with $N_\lambda \sim \text{Bin}(n_\lambda, r)$,

1. If $\rho > 1$ (overloaded regime), then there exists a finite constant $K > 0$ such that

$$\limsup_{\lambda \rightarrow \infty} |\mathbb{E}[Q_{N_\lambda}] - rn_\lambda \bar{q}_\rho| \leq K \quad \text{and} \quad \lim_{\lambda \rightarrow \infty} |\mathbb{E}[\alpha_{N_\lambda}] - rn_\lambda \bar{\alpha}_\rho| \rightarrow 0.$$

2. If $\rho = 1$ (critically-loaded regime), then there exist finite constants $K'_1, K'_2 > 0$ such that

$$\limsup_{\lambda \rightarrow \infty} \mathbb{E}[Q_{N_\lambda}] \leq K'_1 \sqrt{\lambda} \quad \text{and} \quad \limsup_{\lambda \rightarrow \infty} \mathbb{E}[\alpha_{N_\lambda}] \leq K'_2 \sqrt{\lambda}.$$

3. If $\rho < 1$ (underloaded regime), then

$$\lim_{\lambda \rightarrow \infty} \mathbb{E}[Q_{N_\lambda}] \rightarrow 0 \quad \text{and} \quad \lim_{\lambda \rightarrow \infty} \mathbb{E}[\alpha_{N_\lambda}] \rightarrow 0.$$

Theorem 1 shows that, in the overloaded system, the fluid approximation for the expected queue length is asymptotically accurate up to $\mathcal{O}(1)$,¹ and the fluid approximation for the net abandonment rate is asymptotically accurate up to $o(1)$, i.e., the corresponding error is asymptotically bounded in the former case, and it decreases with the arrival rate in the latter case. In other words, fluid approximations are “extremely accurate” in the overloaded regime. In the critically-loaded system, those fluid-approximation errors are $\mathcal{O}(\sqrt{\lambda})$, i.e., they grow in the square-root of the size of the system. In the underloaded regime, fluid approximations are $o(1)$ -accurate since errors for both performance measures decrease with the arrival rate. In other words, relying on fluid approximations is justifiable when the number of servers follows a binomial distribution, which is a reasonable model for self-scheduling server behavior.

13.2.2 What Do the Asymptotic Results Mean?

To interpret the asymptotic results of Theorem 1, it is important to clearly distinguish between *stochastic fluctuations* for the queue-length process (which can be observed, e.g., in a given simulation run), and the *accuracy of many-server fluid*

¹Let f and g be two functions defined on some subset of \mathbb{R} . Then, as $n \rightarrow \infty$, (1) $f(n) = \mathcal{O}(g(n))$ if there exists $M > 0$ and $C > 0$ such that $|f(n)| \leq M|g(n)|$ for $n \geq C$; (2) $f(n) = o(g(n))$ if for all $\epsilon > 0$, there exists N such that $|f(n)| \leq \epsilon|g(n)|$ for all $n \geq N$.

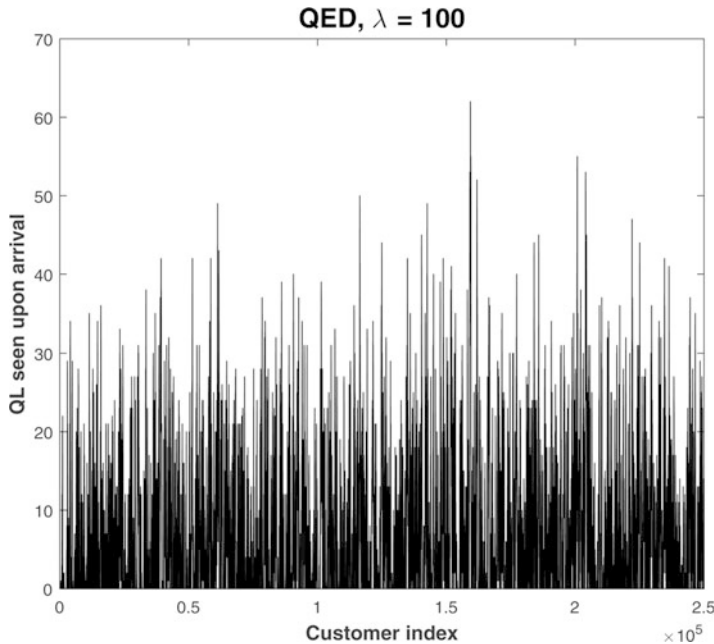


Fig. 13.3 Queue length seen upon arrival in a single simulation run with $\rho = 1$ and $n = 100$. The fluid limit is equal to 0

approximations for the expected queue-length (obtained by averaging over multiple simulation runs, and quantified by letting λ increase). Figures 13.3, 13.4, 13.5, 13.6, 13.7, and 13.8 are based on simulations of an $M/M/n + M$ queueing model with service rate $\mu = 1$ and abandonment rate $\theta = 0.5$. We consider a deterministic number of servers in these simulations because the same intuitions continue to hold when the number of servers is binomially distributed instead.

In Figs. 13.3 and 13.4, we present queue-length sample paths based each on a single simulation run in a system where $n = \lambda$ i.e., $\rho = 1$. For such parameter values, critical-loading approximations (Garnett et al. 2002) are known to describe the system well. For Fig. 13.3, we let $\lambda = 100$, and for Fig. 13.4, we let $\lambda = 1000$. In each figure, the fluid limit is identically equal to 0. It is clear from Figs. 13.3 and 13.4 that the magnitude of stochastic fluctuations in the system is on the order of $\sqrt{\lambda}$, as expected. The same continues to hold in an overloaded system, as illustrated in Figs. 13.5 and 13.6, where we let $\rho = 1.4$ instead.

The asymptotic accuracy results of Theorem 1 describe how the expected queue length differs from its fluid limit as λ increases. Figure 13.7 considers different λ values, ranging from $\lambda = 100$ to $\lambda = 1000$, in critically-loaded systems where $n = \lambda$. As a function of λ , we plot estimates of the expected queue length which are based on averages over 10 independent simulation runs of length 10 million arrivals each. In Fig. 13.8, we do the same but let $\rho = 1.4$ instead, i.e., we consider an overloaded system. Contrasting Figs. 13.7 and 13.8 illustrates our asymptotic

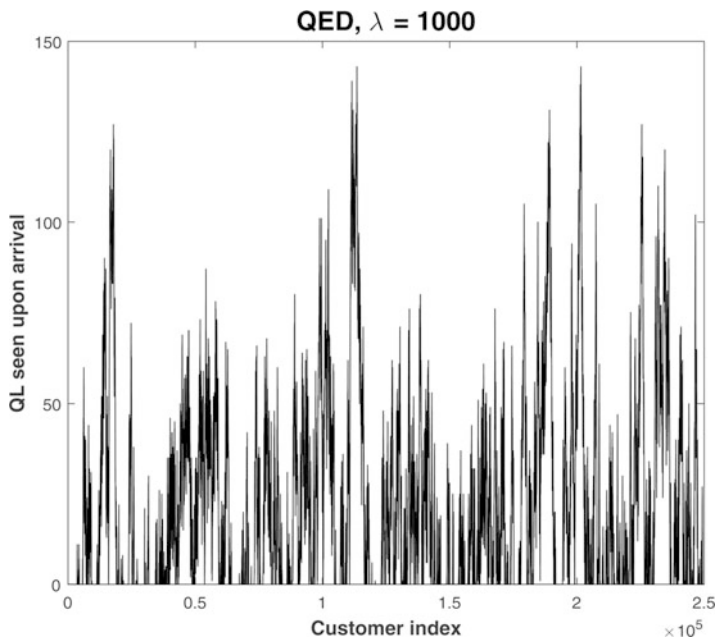


Fig. 13.4 Queue length seen upon arrival in a single simulation run with $\rho = 1$ and $n = 1000$. The fluid limit is equal to 0

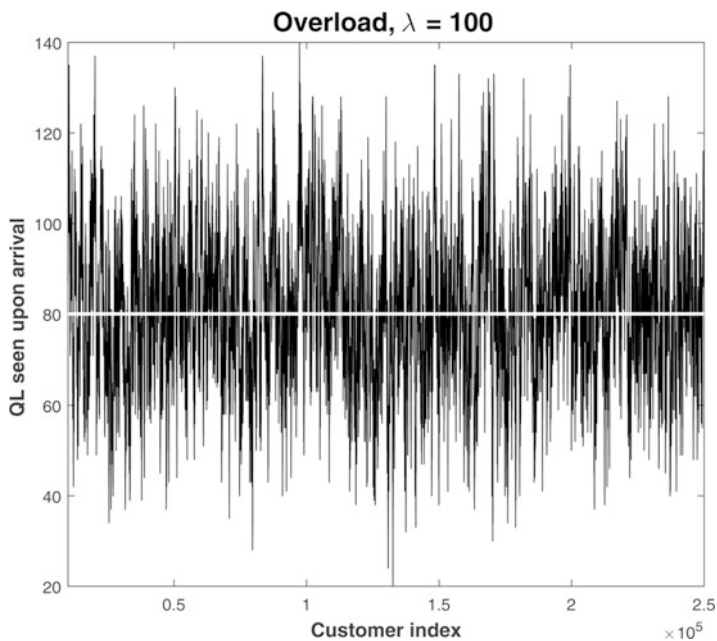


Fig. 13.5 Queue length seen upon arrival in a single simulation run with $\rho = 1.4$ and $n = 100$. The fluid limit is equal to 80

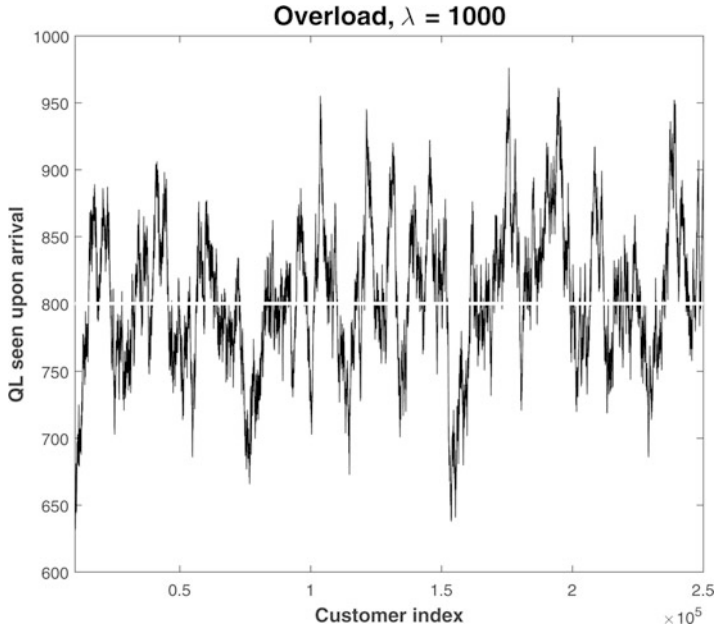


Fig. 13.6 Queue length seen upon arrival in a single simulation run with $\rho = 1.4$ and $n = 1000$. The fluid limit is equal to 800

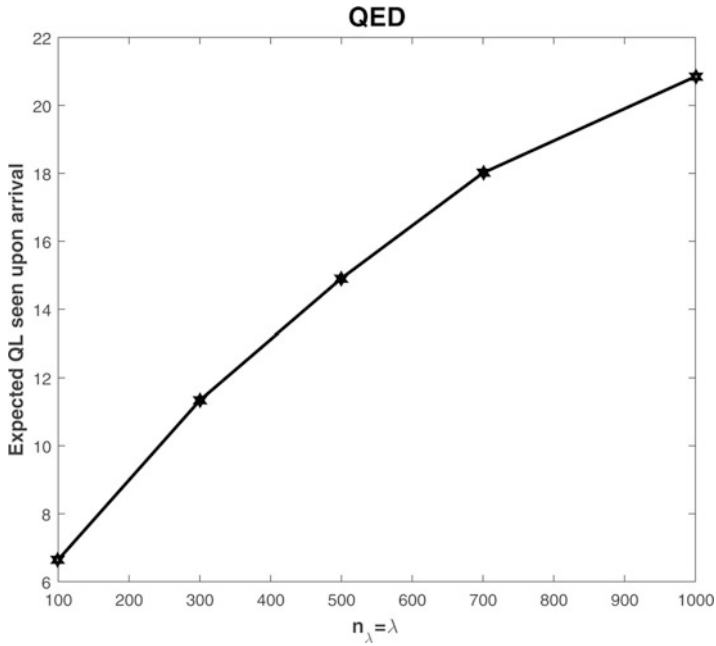


Fig. 13.7 Averages of queue lengths for varying λ where $n_\lambda = \lambda$

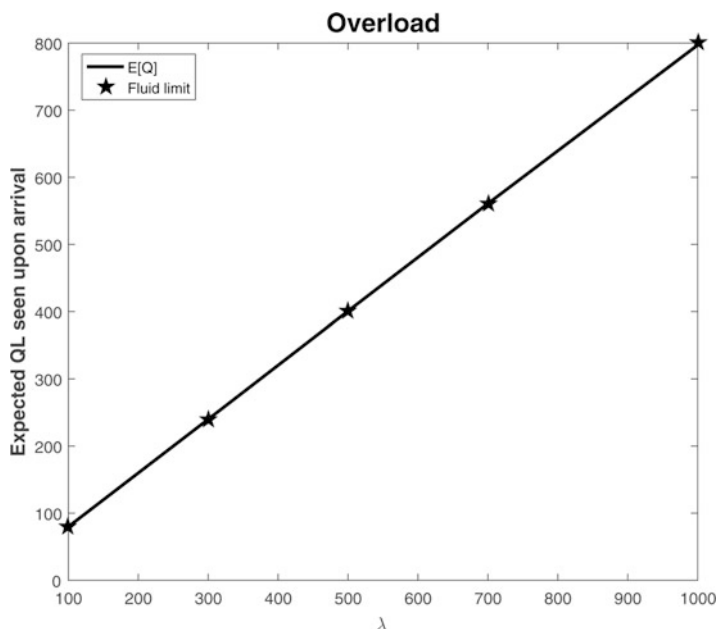


Fig. 13.8 Averages of queue lengths for varying λ where $\lambda/n_\lambda = 1.4$

accuracy results in Theorem 1. In the critically-loaded case, i.e., Fig. 13.7, stochastic fluctuations are consistent with those suggested by the *central limit theorem*, i.e., they are on the order of $\sqrt{\lambda}$. In the overloaded case, i.e., Fig. 13.8, stochastic fluctuations are better explained by *large deviations theory*: Fluid approximations are practically indistinguishable from the estimates for average queue-lengths; see Bassamboo et al. (2010) for related additional discussion.

This section presented the theoretical background needed to study performance in queueing systems with randomness in capacity. In the remainder of this chapter, we apply that theoretical framework to formulate managerial recommendations on the operational management of queueing systems with self-scheduling servers.

13.3 Self-Scheduling Agents: A Long-Term Staffing Decision

13.3.1 The Model

There are k shifts, and agents show up at random for these shifts. In particular, an agent shows up for work in shift j with a given probability, r_j , independently of other agents. We denote the total pool size by n . The number of agents in shift j is a random variable with a binomial distribution, $N_j \sim \text{Bin}(n, r_j)$, where n is the

number of trials and r_j is the success probability. That is, nr_j is the expected number of agents that show up to work in shift j . Customers arrive to the system according to general stationary processes with rates λ_j . We assume that there is no service overlap between the different shifts, i.e., customers who arrive during a shift must be served by agents who are assigned to that shift. This assumption is reasonable when the system is large enough: In this case, processing any customers who remain in queue at the end of a given shift would not take too long, because there are many servers working in parallel. In each shift, we consider a $G/G/N_j + GI$ model. Patience times are i.i.d. across customers, and have a cdf F , complementary cdf (ccdf) \bar{F} , density function f , hazard-rate function h_a , and mean $1/\theta$ for some $\theta > 0$. Service times are assumed to be i.i.d with a general distribution. The arrival, service, and abandonment processes are all mutually independent, also independent of the number of servers. We continue to use the FCFS service discipline.

The system manager must select an appropriate total staffing level, n , to effectively balance staffing and operational costs. With a binomially-distributed number of servers, relying on a fluid approximation is justified in large systems (Sect. 13.2.1). In other words, both stochastic variability and parameter uncertainty may be assumed to be of second order, relative to average performance measures in the system. By ignoring stochasticity, the key challenge in managing a random capacity reduces to the salient *heterogeneity across shifts*, in both demand rates and agent availabilities. To illustrate this point, let us compare the settings with a single and two shifts. With a single shift, assume that each agent has a probability 0.25 of showing up for that shift, and that 100 agents are needed. Then, with a pool of $n = 400$ self-scheduling agents, 100 agents will show up on average. Thus, by staffing a large enough agent pool, the manager could induce the desired number of agents to show up, on average. Now, with two shifts, e.g., morning and afternoon, assume that each agent has probabilities 0.25 to show up in the morning shift, and 0.5 to show up in the afternoon shift. Also, assume that 100 agents are needed for the morning shift, and 150 for the afternoon shift. Then, staffing a pool that is large enough to meet demand in one of the two shifts, on average, will lead to either overstaffing or understaffing the other shift. With multiple shifts and heterogeneous arrival rates and show-up probabilities, it is not clear, a priori, how the manager should staff her system: Should she aim to match demand in some shifts, but not others? What should this decision depend on? We investigate such questions in what follows.

13.3.2 Fluid Formulation

As in Bassamboo and Randhawa (2010), we consider two quality-of-service costs, indexed by the shift j : (i) A delay cost, h_j , per customer for each unit of time that this customer spends waiting to be served, and (ii) an abandonment penalty cost, p_j , incurred per customer who abandons before being served. Each agent is paid c_j per unit time in shift j , if she is available for work in that shift. The system manager

must decide on the staffing level, n . Let \bar{q}_{ρ_n} and $\bar{\alpha}_{\rho_n}$ be the fluid approximations for the queue-length and net abandonment rates with a traffic intensity $\rho_n \equiv \lambda/n\mu$.

Staffing decisions in systems with self-scheduling agents cannot, usually, be made “on the fly”. Because of this, a manager must make her staffing decision in advance: She may not know, with certainty, the availability of each agent in her workforce and she cannot enforce attendance but, given current technological advances, she should be able to obtain *historical estimates* of joining probabilities of agents; these are the r_j in our model. For example, based on analyzing human resources data in her firm, she may know that stay-at-home parents usually prefer to work in the morning, while children are at school, but may not know whether a specific work-from-home parent will show up to work on a given morning. She may know that higher compensations are typically offered during certain times of day (ride-sharing) or for certain client companies (virtual call centers), but may not know the exact “surge prices” that are going to be in effect, if any. To capture such challenges, we assume that the compensation, c_j , and the show-up probability, r_j , are fixed (we will relax those assumptions later). In other words, we begin by solving the problem:

$$\text{minimize}_{n \in \mathbb{N}} C(n) \equiv \sum_{1 \leq j \leq k} (nr_j c_j + p_j \cdot \bar{\alpha}_{\rho_n}/r_j + h_j \cdot \bar{q}_{\rho_n}/r_j), \quad (13.1)$$

where \mathbb{N} denotes the set of natural integers.

13.3.3 Optimal Staffing Policy

13.3.3.1 No Self-Scheduling

To quantify the impact of self-scheduling, we need to choose a useful benchmark. Without self-scheduling, the system manager can select the optimal staffing levels, n_j^* , independently for each shift j . To specify n_j^* , we need to make additional assumptions. The density of the fluid that has been waiting for exactly u time units, in shift j , is equal to $\lambda_j \bar{F}(u)$. Therefore, the corresponding (unscaled) queue length is given by $q_j = \int_0^{w_j} \lambda_j \bar{F}(u) du$, where w_j denotes the waiting time given service. The net abandonment rate (unscaled) in shift j is equal to $\lambda_j \bar{F}(w_j)$. In the absence of self-scheduling, we must have that $n_j^* = \lambda_j \bar{F}(w_j^*) \leq \lambda_j$ where w_j^* is the optimal waiting time in shift j ; this is because it is suboptimal to staff more than λ_j agents in shift j , i.e., underload that shift. In shift j , w_j^* is determined by solving:

$$\min_{w_j \geq 0} \lambda_j \left((c_j - p_j) \bar{F}(w_j) + h_j \int_0^{w_j} \bar{F}(u) du \right). \quad (13.2)$$

Hereafter, we make the following assumption, which states that staffing costs are sufficiently inexpensive.

Assumption 1 For all j , $c_j < \min\{h_j/h_a(0) + p_j, h_j/\theta + p_j\}$.

It is useful to offer a brief comment on the validity of Assumption 1. To do so, let us assume that c_j corresponds to the minimum wage which is close to \$7 per hour in the United states. Let us also take the time unit to be one hour. Assume that $\theta = 4$, i.e., a customer is willing to wait on average for 15 min before abandoning. For a numerical value of h_j , we make use of existing empirical evidence from the call-center literature, e.g., Akşin et al. (2013). Based on their results (see Table 4 in that paper), customers attribute a waiting cost of roughly 1 \$ per minute. This translates into 60 \$ per hour. For such values, the assumption that we make on the staffing cost is satisfied irrespective of the value of p_j . Under Assumption 1, it is easy to establish the following result for the solution to problem Eq. 13.2.

Proposition 1 Under Assumption 1, in a system with no self-scheduling servers, it is optimal to match the supply and demand rates in every shift, i.e., $n_j^* = \lambda_j$.

In other words, we choose as benchmark a setting where it is optimal to match demand and supply in each of the shifts. This is the case when staffing costs are not too high, as per Assumption 1. When demand and supply are matched in each shift, there is no delay, at fluid scale. Thus, the customer abandonment distribution does not play any role, since customers do not abandon. With self-scheduling, the manager is no longer able to set n_j^* independently for each shift and must decide, instead, on the total pool size n . Because of the ensuing imbalance between demand and supply, some shifts may be congested. Thus, because of self-scheduling, the customer abandonment distribution will now play an important role in congested shifts. Here, we study how to exploit that role to mitigate the cost of self-scheduling.

13.3.3.2 Self-Scheduling Capacity

To capture the heterogeneity across different shifts, we define the *augmented arrival rate* $\Gamma_j \equiv \lambda_j/r_j$, and let $\Gamma_0 \equiv 0$. This will allow us to characterize, in a simple manner, the optimal solution to the staffing problem in Eq. 13.1. Letting $n = \Gamma_j$ amounts to matching the supply and demand rates in shift j . This is because the number of agents who show up in shift j is then equal to $n \cdot r_j = \Gamma_j \cdot r_j = \lambda_j$. In a sense, the respective values of Γ_j , across shifts, quantify the degree of self-scheduling imbalance in the system. In particular, if $\Gamma_j \equiv \Gamma$ are identical across all shifts, then it is easy to see that staffing $n = \Gamma$ would eliminate the cost of self-scheduling, on average. However, if the Γ_j 's are “very different” across the different shifts, then managing self-scheduling agents becomes increasingly difficult, i.e., leading to a higher cost. In an overloaded shift j , we have that $\Gamma_j \bar{F}(w_j) = n$, i.e., $w_j = \bar{F}^{-1}(n/\Gamma_j)$. Since it is never optimal to strictly underload all shifts, the staffing problem in Eq. 13.1 can be defined piecewise:

$$\underset{0 \leq n \leq \Gamma_k}{\text{minimize}} C(n) \equiv \sum_{j=1}^k \mathbf{1}(\Gamma_{j-1} \leq n < \Gamma_j) u_j(n), \tag{13.3}$$

where $\mathbf{1}(n \in A)$ denotes the indicator function over the set A , and $u_j(n)$ is given by:

$$u_j(n) \equiv \sum_{i=1}^k c_i n r_i + \sum_{i=j}^k \left(p_i (\lambda_i - n r_i) + h_i \lambda_i \int_0^{\bar{F}^{-1}(n/\Gamma_i)} \bar{F}(u) du \right), \quad (13.4)$$

i.e., $u_j(n)$ is the *total* cost incurred if n is chosen in the interval $[\Gamma_{j-1}, \Gamma_j)$. It turns out that the solution to Eq. 13.3 depends on the monotonicity of the hazard rate of the abandonment distribution. Here is how.

Monotonically increasing hazard rate A monotonically increasing hazard rate corresponds to customer patience “wearing out” as the customers waits longer in queue. For abandonment distributions with a monotonically increasing hazard rate (including the exponential distribution with a constant hazard rate), we find that it is optimal to match the supply and demand rates in *one* of the k shifts when servers self schedule, with the remaining shifts being either over or under staffed; this lies in contrast to matching supply and demand in *all* shifts without self-scheduling, as per Proposition 1. In particular, the following proposition holds:

Proposition 2 *For abandonment distributions with a monotonically non-decreasing hazard rate, there is one shift i_0 where the supply and demand rates must be matched, i.e., $n^* = \Gamma_{i_0}$.*

The optimality of overstaffing certain shifts lends some support to the staffing policies adopted in virtual call centers such as LiveOps or Arise, where agents regularly complain about the fact that there are “too many other agents on board” and, consequently, “too few calls to answer”. However, the compensation structure in those settings is different: There, the manager typically uses volume-dependent pay, e.g., agents earn a piece-rate compensation in addition to some base salary. Under our fixed compensation structure, we find that overstaffing certain shifts can minimize costs, but that this is not true across all shifts.

Monotonically decreasing hazard rate We now consider abandonment distributions with a monotonically decreasing hazard rate, which is consistent with the way call-center customers abandon in practice. A monotonically decreasing hazard rate for abandonment corresponds to customers becoming increasingly patient as they wait long in queue, e.g., because they feel that “they have waited already for so long, so why not wait a little longer?”.

Proposition 3 *For abandonment distributions with a monotonically decreasing hazard rate, it is optimal to either under or over staff every shift (no matching), or to match the supply and demand rates in one of the shifts.*

Interestingly, Proposition 3 shows that it may be optimal for the manager to not match the supply and demand rates anywhere, i.e., to effectively under or over load every shift. In practical terms, Proposition 3 shows that it may be optimal for the manager to maintain an *imbalance* between the average supply and demand

rates in each of the shifts. In other words, the conventional wisdom for workforce management in call centers, which is to staff just enough agents to meet projected incoming demand, may no longer be the right approach with self-scheduling agents, since it may be optimal *not* to meet the established service level in any shift, but rather to exceed or fall below it.

To summarize, the optimal staffing policy in a system with self-scheduling agents is not straightforward, and strongly depends on both the show-up behavior of agents and the impatience distribution of customers. In particular, it may be optimal to match supply and demand in exactly one shift (monotonically increasing hazard rate), or no shift at all (monotonically decreasing hazard rate). This lies in contrast with the benchmark solution where the abandonment distribution played no role, and it is optimal to match supply and demand across all shifts. Of course, having both understaffed and overstaffed shifts in the optimal staffing policy means that the manager cannot eliminate the imbalance between supply and demand, which is due to self-scheduling, by adjusting the staffing level in her system. Thus, we need to investigate short-term controls in the system as well, in addition to the long-run staffing decision. We do so in what follows.

13.4 Short-Term Controls

It is natural to investigate how to control the **compensation**, c_j , for each shift j . In particular, we assumed in Eq. 13.1 that the agent show-up probability, r_j , was exogenously specified. In practice, r_j usually depends on the compensation offered in shift j . In order to capture how changes in c_j may impact agent show-up behavior, we assume, as in Gurvich et al. (2017), that agents are statistically identical and have an availability threshold (opportunity cost) T for showing up in shift j . Letting $G(\cdot)$ denote the cumulative distribution function (cdf) of T , an agent shows up in shift j with probability $r_j \equiv G(c_j)$. We also assume that $G(\cdot)$ is log-concave with positive density function $g(\cdot)$; this will be used later to ensure uniqueness of solutions in our optimization problems.

Nevertheless, there is also a need to consider alternative tools, besides compensation and staffing, to control the system. First, the manager may be restricted in how much and how often she can modify compensation. This is certainly the case in virtual call centers because of market transparency and fierce competition between providers. Also, in virtual call centers, compensations are often set in advance by client companies rather than by the virtual call-center platform itself. In this case, the responsibility of the platform is to staff and train agents, and act as an intermediary between client companies and their agents. Second, while pricing influences agents, it cannot always be used to influence the behaviour of customers, e.g., in service-oriented virtual call centers; thus, there is a need to consider customer-side controls. Third, there is considerable concern about the extent to which pricing should be used as a control in on-demand service platforms, because of extreme and frequent fluctuations.

In some settings, it may be possible to cap the participation of agents in certain shifts; e.g., this is the case in virtual call centers where the manager can easily choose which shifts to make available for self-scheduling agents to choose from. However, capping agent participation is restrictive, and agents usually complain when too few shifts are available. Moreover, capping may not be possible in certain settings, e.g., with ride-sharing services where drivers may already be on the road, so that it would be difficult to prohibit them from driving at different times. Thus, we do not consider such a control in this chapter. Instead, we investigate controls on the customer side, as follows.

In Sect. 13.3, we characterized the role played by the abandonment distribution in a system with randomness in capacity. Because of this, it is natural to investigate ways of controlling customer impatience to alleviate the system's cost. Here, we propose to do so via *delay announcements* in the system. We assume that the provision of delay announcements is costless for the manager, and that a single announcement is given to each delayed customer immediately upon arrival. The idea of using delay announcements as a control of customer impatience is not new. Indeed, it has been explored both empirically and analytically in several papers, albeit in contexts different from ours; for a survey of those papers, see Ibrahim (2017b). At a high level, while compensation is used to control agent joining behavior in our setting, the announcements are used to control customer behavior instead. However, it should be noted at the onset that the announcements cannot be used to restore balance in the system, i.e., entirely eliminate the cost of self-scheduling. Indeed, while delay information incites impatient customers (who would have abandoned anyway) to abandon earlier, thus leading to a reduction in waiting costs, it does not impact the overall abandonment rate in the system.

In the remainder of this chapter, we study how a manager should decide on staffing, compensation, and announcements both separately and jointly. It is unclear, a priori, what the interaction between our three controls will be, i.e., how one would affect the other. It is also unclear how a manager who has the options of making announcements and controlling compensations would staff her system: Would she use these three controls to consistently match supply and demand in all shifts? Or, would she continue to overstaff/understaff some shifts? If so, then when?

13.4.1 Delay Announcements: Performance Impact

In this work, we contend that the announcements made must be truthful and accurate, for otherwise customers will learn to mistrust them. Studying the performance impact of delay announcements, when customers respond to these announcements by updating their abandonment distributions, is challenging. Indeed, changes in customer impatience affect system dynamics and, in turn, the future announcements made. For example, if customers abandon faster because of high announcements, then future waiting times, and future announcements which depend on those waiting times, should be shorter. When waiting times decrease, customers are inclined to

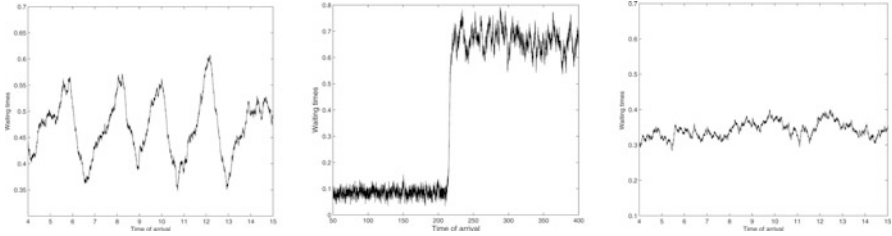


Fig. 13.9 No equilibrium, multiple equilibria, and a single equilibrium under different announcement-dependent abandonment distributions

be patient and wait for service, thus making delays longer again, and so on. In a nutshell, studying the impact of the announcements involves characterizing an equilibrium in the system. At a high level, an equilibrium must correspond to the long-run performance in the system, where the average announced delay coincides with the average experienced delay.

First, it is not clear whether such an equilibrium exists, or if it is unique; indeed, there may be multiple equilibria and the system may exhibit oscillations between those equilibria. We illustrate these possibilities in Fig. 13.9, where we assume different customer-response functions, in each subfigure, and plot simulation-based sample paths of waiting times for each such function. (The specifics of the model are not important, and are therefore omitted.) The three subfigures, from left to right, correspond to having no equilibrium, multiple equilibria, and a unique equilibrium, respectively. Second, even when a unique equilibrium exists, it is not clear how to specify that the announcement and the corresponding delay, which are both random variables, coincide in that equilibrium, e.g., this could be in expectation, in distribution, or asymptotically when scaled in an appropriate way. Third, it is not clear how stochastic fluctuations around the equilibrium affect the system's performance. Even under Markovian assumptions, explicit analysis of the underlying birth-and-death process is analytically complex. This is so because the transition rates of the birth-and-death chain would all be dependent on the announcements. Therefore, analysis is typically done in an asymptotic regime instead. Here, we do so in the context of a fluid model, as in Armony et al. (2009) who consider a single shift instead, and a context different from ours.

Because we consider a system with multiple shifts, and different shifts have different congestion levels and therefore different delay announcements, we obtain in each shift a different announcement-dependent abandonment distribution. Herein lies the complexity of considering multiple shifts: The announcements may lead to shorter delays in some shifts, but not in others, and the aggregate effect of those announcements is unclear. To derive insights, we focus hereafter on an exponential abandonment distribution with an announcement-dependent rate. In particular, letting w be the announcement made, customers abandon according to an

exponential abandonment distribution with rate $\theta(w)$. For a fixed agent pool size, n , we let $w_j^e(n)$ denote the equilibrium delay in shift j , which is dependent on n . We must have:

$$\lambda_j e^{-w_j^e(n)\theta(w_j^e(n))} = nG(c_j), \quad \text{i.e.,} \quad e^{-w_j^e(n)\theta(w_j^e(n))} = \frac{n}{\Gamma_j}, \quad (13.5)$$

by conservation of flow in shift j . The total cost in the system, with the announcements, is

$$C_a(n) \equiv \sum_{i=1}^k c_j n G(c_j) + \sum_{i=1}^k \left(p_j + \frac{h_j}{\theta(w_j^e(n))} \right) (\lambda_j - nG(c_j))^+. \quad (13.6)$$

Assuming that $\theta(\cdot)$ is continuous and strictly increasing, consistently with the empirical evidence in e.g., Mandelbaum and Zeltyn (2013) and Aksin et al. (2016), guarantees the existence and uniqueness of an equilibrium $w_j^e(n)$ in every shift j . In what follows, we also assume that $\theta(w)$ is a differentiable function of w and that $\lim_{w \rightarrow \infty} \theta(w) > 0$.

13.4.1.1 When Do the Announcements Reduce the Cost of Self-Scheduling?

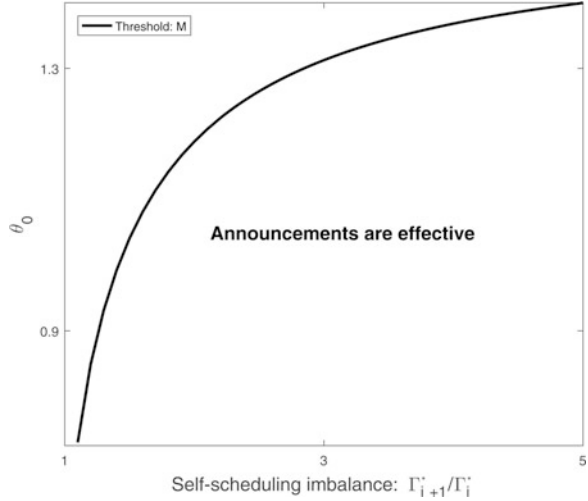
To gain a deep understanding, it is useful to begin by exploring the performance impact of each short-term control separately. Therefore, we first assume that the manager communicates delay announcements in all congested shifts, and that compensations and the staffing level are fixed. We then ask the question: Do the announcements help in reducing the cost of self-scheduling? Naturally, the announcements are effective if they incite customers to abandon faster than they would have otherwise. In our problem, we have different Γ_j values and, consequently, different announcement-dependent abandonment rates given by Eq. 13.5. We let θ_0 denote the abandonment rate without the announcements, which is constant across all shifts. By Proposition 2, because the times to abandon are assumed to be exponentially distributed, it is optimal to critically load one shift, call it i_c , i.e., $n^* = \Gamma_{i_c}$ without the announcements. We now derive a simple sufficient condition under which the announcements lead to an overall decrease in the system's cost.

Proposition 4 *With exponential abandonment with an announcement-dependent rate $\theta(w)$, if*

$$\theta_0 \cdot \theta^{-1}(\theta_0) < \ln \left(\frac{\Gamma_{i_c+1}}{\Gamma_{i_c}} \right), \quad (13.7)$$

then $C_a(n^) < C^*$ for $C_a(\cdot)$ in Eq. 13.6, where C^* is the optimal solution to Eq. 13.1 with $n^* = \Gamma_{i_c}$.*

Fig. 13.10 Threshold on θ_0 as given by Proposition 4



The condition in Eq. 13.7 means that customers do not abandon “too fast” in the absence of the announcements. This is because it can be shown, under our assumption on $\theta(\cdot)$, that the function on the left-hand-side of Eq. 13.7 is increasing in $\theta_0 \geq 0$. Thus, the condition may be equivalently interpreted as imposing a threshold, say M , on θ_0 .

In Fig. 13.10, we plot how the value of the threshold M varies with the degree of self-scheduling “imbalance,” as measured by $\Gamma_{i_c+1}/\Gamma_{i_c}$.² The area under the threshold curve corresponds to values of θ_0 for which the provision of delay announcements decreases the overall cost in the system. In other words, this is when the announcements are effective. It is interesting to note that this area increases as $\Gamma_{i_c+1}/\Gamma_{i_c}$ increases, i.e., the announcements are *increasingly* effective as self-scheduling causes a *greater* imbalance in the system.

13.4.1.2 A New Staffing Problem

Since the announcements lead to a decrease in waiting times, it is natural to investigate whether it is optimal for the manager to create additional congestion by understaffing her system. This is because this increased congestion would, subsequently, be reduced by the announcements. To explore this, we now assume that the manager can jointly optimize the staffing level in her system, along with the announcements. The manager’s staffing problem, assuming that she makes announcements in all shifts at a later stage, is given by:

²We assume that $k = 10$; $c = 1.1$; $h = 0.5$; $p = 1.0$; avg. $\lambda = 55$; $r = 0.4$; with announcements: $\theta(w) = 1.5 - e^{-2w}$.

$$\min_{n \in \mathbb{N}} \sum_{j=1}^k \left(c_j n G(c_j) + \left(p_j + \frac{h_j}{\theta(w_j^e(n))} \right) (\lambda_j - n G(c_j))^+ \right), \quad (13.8)$$

where we replace the constant abandonment rate θ_0 by different announcement-dependent rates, $\theta(w_j^e(n))$, depending on both the shift and the staffing level n . That is, in setting her optimal staffing level, the manager needs to consider the subsequent dependence of customer abandonment behavior on the selected pool size. Let n_a^* denote the optimal solution to Eq. 13.8, with the announcements, and n^* denote the optimal solution to Eq. 13.1, without the announcements.

Proposition 5 *With exponential abandonment with an announcement-dependent rate $\theta(w)$, if*

$$\theta_0 \cdot \theta^{-1}(\theta_0) < \min_{1 \leq i \leq k-1} \ln(\Gamma_{i+1}/\Gamma_i), \quad (13.9)$$

then $n_a^* < n^*$.

That is, under Eq. 13.9, it is optimal for the manager to hire a smaller agent pool than without the announcements. Condition Eq. 13.9 could also be interpreted as an upper bound on θ_0 , albeit a tighter one than in Proposition 4. Proposition 5 does not give an indication of the extent to which cost can be reduced by staffing a smaller pool. In Fig. 13.11, we plot the percent decrease in the cost of self-scheduling as a function of the self-scheduling imbalance in the system, as measured through Γ_k/Γ_1 : The higher Γ_k/Γ_1 , the larger the imbalance. Figure 13.11 shows that the announcements become *increasingly* effective as the self-scheduling imbalance in the system *increases*.

To summarize, we find that delay announcements are most effective when there is a significant imbalance which arises from having “very” heterogeneous augmented arrival rates, e.g., because the agents have “very” different show-up probabilities across the different periods. This is a desirable property, because we want the announcements to be able to control the system in that case. The announcements are also effective when customers are relatively patient in the benchmark system, because delay information incites them to abandon faster.

13.5 Joint Control of Compensation and Delay Announcements

In a ride-sharing platform, the manager may know that a concert will end shortly in a given region (say in the next 15–30 min), and anticipate a surge in demand for Uber cars. She would then use short-term controls, e.g., pricing, to incite more agents to go to that region. She would not, however, be able to control her overall pool of drivers, i.e., hire and train more Uber drivers, because such a decision must

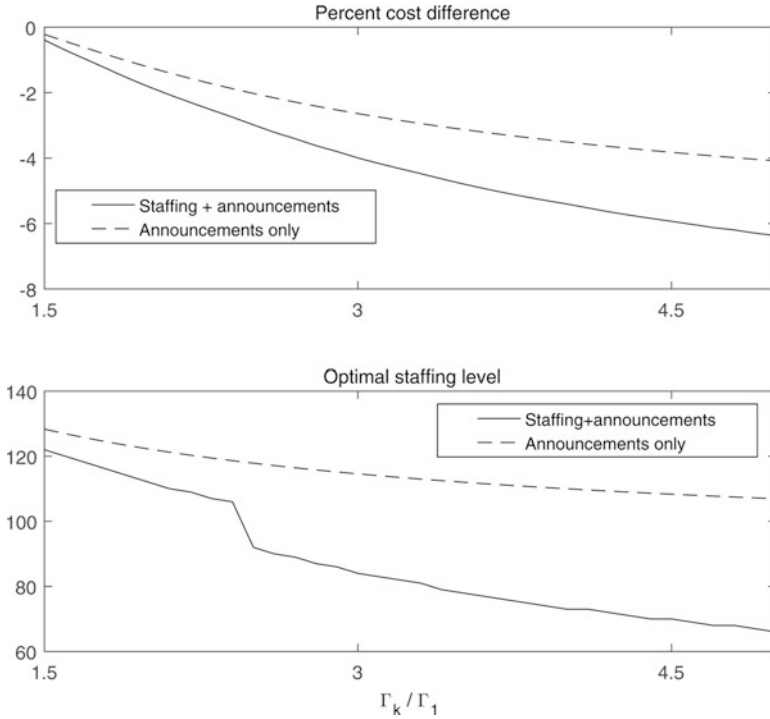


Fig. 13.11 Percent decrease in the system’s cost by solving Eq. 13.8

be made weeks in advance. To mimic this situation, we assume that the staffing level is equal to n , and investigate the optimal compensation to be offered in shift k . We let l denote the minimum wage allowed in any shift. The manager may decide on c_k separately for each shift k :

$$\min_{c_k \geq l} \left\{ c_k n G(c_k) + \left(p_k + \frac{h_k}{\theta} \right) (\lambda_k - n G(c_k))^+ \right\}. \tag{13.10}$$

For expositional ease, we let $L_k \equiv p_k + h_k/\theta$ capture customer-related costs, and denote $\psi_k^n \equiv G^{-1}(\lambda_k/n)$. If $c_k = \psi_k^n$, then $nG(c_k) = \lambda_k$. In other words, using compensation $c_k = \psi_k^n$ in shift k incites just enough agents to meet demand in that shift. It will also be convenient to define $a_k < L_k$ as follows:

$$G(a_k) \left(1 + (a_k - L_k) \frac{g(a_k)}{G(a_k)} \right) = 0. \tag{13.11}$$

The optimal compensation in problem Eq. 13.10 is given by the following lemma.

Lemma 1 *The optimal compensation in shift k , solution to Eq. 13.10, depends on n as follows:*

1. *If $n \geq \lambda_k/G(l)$, then $c_k^* = l$ and shift k is overstaffed;*
2. *If $\lambda_k/G(a_k) \leq n < \lambda_k/G(l)$, then $c_k^* = \psi_k^n$ and demand and supply are matched in shift k ;*
3. *If $n < \lambda_k/G(a_k) < \lambda_k/G(l)$, then $c_k^* = a_k$ and shift k is understaffed.*

It is not surprising that the optimal compensation in a given shift depends on the total pool size, n , and that the larger the pool, the smaller the compensation needed to incite agents to participate. In particular, we find that the manager uses the minimum wage in shift k when the agent pool size is very large (case (a)). In this case, the manager need not use high compensation to incite sufficient agent participation in the shift. For moderate values of the agent pool size (case (b)), the manager sets compensation to match demand and supply in the shift, i.e., $c_k^* = \psi_k^n$. Finally, when the pool size is very small (case (c)), inciting sufficient agent participation is too costly for the manager, so she sets a compensation that leads to an *understaffed* shift k .

We now turn to the more interesting case where the manager may jointly control both the provision of delay announcements and the compensation offered, in each shift. For tractability, we assume that the announcement-dependent abandonment rate is constant and equal to $\tilde{\theta} > \theta$, where θ is the rate without the announcements. We denote $\tilde{L}_k \equiv p_k + h_k/\tilde{\theta}$ and note that $\tilde{L}_k < L_k \equiv h_k + p_k/\theta$. Thus, it is optimal for the manager to make announcements in every overloaded shift, since doing so would reduce the cost of congestion in that shift. While it is clear that making announcements is beneficial to the manager in that case, it is unclear whether agents will be better or worse off because of the announcements. We continue to assume that the staffing level n is fixed, and we investigate the optimal compensation in a shift where the manager is allowed to make delay announcements. Since the announcements are only relevant when the system is congested, we focus on case (c) in Lemma 1, i.e., we assume that $n < \lambda_k/G(a_k)$ for a_k in Eq. 13.11. We let \tilde{c}_k^* denote the optimal compensation in shift k , assuming that the manager makes announcements in that shift; i.e., \tilde{c}_k^* minimizes $c_k n G(c_k) + \tilde{L}_k (\lambda_k - n G(c_k))^+$. In the following lemma, we show that $\tilde{c}_k^* = \tilde{a}_k < c_k^* = a_k$ where

$$G(\tilde{a}_k) \left(1 + (\tilde{a}_k - \tilde{L}_k) \frac{g(\tilde{a}_k)}{G(\tilde{a}_k)} \right) = 0. \quad (13.12)$$

Lemma 2 *If the manager has the option of making delay announcements, then agents receive lower compensation, i.e., they are worse off.*

Intuitively, because the manager is able to reduce congestion in the system by resorting to the announcements, she does not need to incite too many agents to participate. Thus, she offers lower compensation. In other words, instead of using high compensation to incite higher agent participation, she uses the announcements to disincentivize customer waiting instead, thereby relieving the congestion caused by self-scheduling.

13.6 Jointly Optimizing Long and Short-Term Controls

We are now ready to investigate the joint optimization problem, where the manager may use all three controls, staffing, compensation, and the announcements, at once. Here is the manager's problem when she can optimize all controls:

$$\underset{c_j \geq l, n \in \mathbb{N}}{\text{minimize}} \quad \Pi(n, \mathbf{c}) \equiv \sum_{1 \leq j \leq k} \left(c_j \cdot nG(c_j) + \tilde{L}_j(\lambda_j - nG(c_j))^+ \right), \quad (13.13)$$

where $\mathbf{c} \equiv (c_1, c_2, \dots, c_k)$ is the k -dimensional vector of compensations and, as before, $\tilde{L}_j \equiv p_j + h_j/\theta$ is the adjusted congestion cost which accounts for the effect of the announcements. To better position our results, we recall that when capping agents is allowed, the optimal compensation is set equal to the minimum wage in all shifts (Gurvich et al. 2017), irrespective of the value of that wage, and the staffing level high enough to match demand in the highest-demand shift (with the offered minimum wage). Supply is capped in the overstaffed shifts. In our context, because capping is deemed undesirable and not allowed, we find that the optimal compensation depends on the *value* of the minimum wage, in particular whether it is “low” or “high,” the manager may offer higher compensation than the minimum wage in some shifts, and may still either understaff or overstaff some shifts. In understaffed shifts, she uses the announcements. Problem Eq. 13.13 may be solved in two stages, first fixing the staffing level n (assuming that the announcements are made in every congested shift) and solving for the optimal compensation, as a function of n and, second, determining the optimal staffing level by exploiting that structure for the optimal compensation. However, the solution to Eq. 13.13 is algebraically complex with multiple shifts. Thus, we focus on two special cases: (i) the minimum wage is sufficiently low, and (ii) the minimum wage is sufficiently high.

13.6.1 Low Minimum Wage

We begin by considering the case where the minimum wage is “sufficiently low”. In particular, we define:

$$l_0 = G^{-1} \left(\frac{\min_i \{\lambda_i\}}{\max_i \{\lambda_i / G(\tilde{a}_i)\}} \right) \quad \text{where } \tilde{a}_k \text{ is given in Eq. 13.12.} \quad (13.14)$$

Then, the following lemma holds for $l < l_0$.

Lemma 3 *If the minimum wage is sufficiently low, then all shifts are either overstaffed or have matched supply and demand. Moreover, there exists at least one shift where demand and supply are matched. The manager does not resort to using delay announcements.*

Lemma 3 shows that the manager need not always resort to using the announcements. In particular, if the minimum wage is “low enough,” then she will hire enough agents and offer high compensations so that $n^*G(c_i^*) \geq \lambda_i$ in each shift i , i.e., no period is congested and there is no need to resort to the announcements. Moreover, she will offer a compensation that is strictly higher than the minimum wage in at least one of the shifts (with highest demand rates). Intuitively, because the minimum wage is small, the manager is less restricted in the compensation that she has to pay her agents. Therefore, she can afford to staff a larger pool and eliminate congestion in her system. This also explains why she is then able to pay her agents a compensation which is strictly larger than the minimum wage. Because no shift is congested, the manager does not resort to making delay announcements.

13.6.2 High Minimum Wage

We now explore the case where the minimum wage is “sufficiently high”. In particular, we assume that $\tilde{a}_i < l < \tilde{L}_i$ for all i . In the following lemma, we show that the manager would then make announcements.

Lemma 4 *If the minimum wage is sufficiently high, then the manager uses the minimum wage in all shifts. Moreover, there exists a shift where supply and demand are matched, with the remaining shifts either under or over staffed; announcements are made in every congested shift.*

Lemma 4 shows that the manager must set compensation equal to the minimum wage in every shift, if that minimum wage is sufficiently high. In this case, the manager must staff a smaller agent pool (because it would be too costly to employ many agents), and she will use the announcements to alleviate congestion in understaffed periods.

13.7 Conclusions

The recent and ongoing growth of the sharing economy has motivated several recent papers in the academic literature; indeed, this book is testament to that growing interest. In this chapter, we surveyed some theoretical results on the analysis of queueing systems with uncertain parameters, and described how such results may be applied for the effective management of queueing systems with self-scheduling agents. Because of the analytical complexity in such settings, queueing-theoretic approximations, which are grounded in many-server heavy-traffic limits, are useful in generating valuable insight.

Nevertheless, there remains numerous directions that are interesting to explore. For example, several modelling extensions (e.g., multiplicity of customer classes) remain to be explored. Our modelling approach was based on approximating system dynamics by using a fluid model. This is justifiable when the number of servers is binomially-distributed. In general, e.g., when there is considerable variability in agent show-up behavior or when the binomial distribution is not appropriate, there is a need to explore more refined approximations, jointly with dynamic compensation decisions and other controls. In studying the effect of delay announcements, we focused on a setting where the announcement-dependent abandonment rate is constant. In practice, customer response to the announcements tends to be non-regular, exhibiting jumps at the epochs of announcements. Developing tools to study such a response, in a setting where there is randomness in capacity, would be interesting to explore as well.

Technical Appendix

Proof of Theorem 1

The Overloaded Regime

$\mathcal{O}(1)$ -Accuracy for the Fluid Queue Length

We begin by establishing the asymptotic $\mathcal{O}(1)$ -accuracy for the expected queue length. Let $0 < \epsilon < r$ and define $k_1 \equiv r - \epsilon$ and $k_2 \equiv r + \epsilon$. Assume that ϵ is small enough so that $\rho r / (r + \epsilon) > 1$. Denote $\mathbb{E}[Q_{N_\lambda} | N_\lambda = s] \equiv \mathbb{E}[Q_s]$ where Q_s is the steady-state queue length in the corresponding $M/M/s + GI$ queue with the same arrival rate.

Conditioning and unconditioning on N_λ Conditioning on N_λ , we can write:

$$\begin{aligned} & |\mathbb{E}[Q_{N_\lambda}] - rn_\lambda \bar{q}_\rho| \\ &= \left| \sum_{s \geq 0} \mathbb{E}[Q_s] \mathbb{P}(N_\lambda = s) - rn_\lambda \bar{q}_\rho \right| \\ &= \left| \sum_{s \geq 0} (\mathbb{E}[Q_s] - s \bar{q}_\rho) \mathbb{P}(N_\lambda = s) \right| \quad \text{since } \mathbb{E}[N_\lambda] = rn_\lambda = \sum_{s \geq 0} s \mathbb{P}(N_\lambda = s), \\ &\leq \left| \sum_{s < k_1 n_\lambda \text{ or } s > k_2 n_\lambda} (\mathbb{E}[Q_s] - s \bar{q}_\rho) \mathbb{P}(N_\lambda = s) \right| \\ &\quad + \left| \sum_{k_1 n_\lambda \leq s \leq k_2 n_\lambda} (\mathbb{E}[Q_s] - s \bar{q}_\rho) \mathbb{P}(N_\lambda = s) \right|. \end{aligned}$$

We now turn to establishing asymptotic bounds for A_λ and B_λ , defined as follows:

$$A_\lambda \equiv \left| \sum_{s < k_1 n_\lambda \text{ or } s > k_2 n_\lambda} (\mathbb{E}[Q_s] - s\bar{q}_\rho) \mathbb{P}(N_\lambda = s) \right| \quad \text{and}$$

$$B_\lambda \equiv \left| \sum_{k_1 n_\lambda \leq s \leq k_2 n_\lambda} (\mathbb{E}[Q_s] - s\bar{q}_\rho) \mathbb{P}(N_\lambda = s) \right|.$$

Asymptotic bound for N_λ far from $n_\lambda r$ We begin by showing that A_λ is asymptotically negligible.

Lemma 5 $\lim_{\lambda \rightarrow \infty} A_\lambda = 0$.

Proof We can write,

$$A_\lambda = \left| \sum_{s > k_2 n_\lambda \text{ or } s < k_1 n_\lambda} \mathbb{E}[Q_s] \mathbb{P}(N_\lambda = s) - \sum_{s > k_2 n_\lambda \text{ or } s < k_1 n_\lambda} s\bar{q}_\rho \mathbb{P}(N_\lambda = s) \right|,$$

$$\leq \mathbb{E}[Q_0] \sum_{s > k_2 n_\lambda \text{ or } s < k_1 n_\lambda} \mathbb{P}(N_\lambda = s) + \sum_{s > k_2 n_\lambda \text{ or } s < k_1 n_\lambda} s\bar{q}_\rho \mathbb{P}(N_\lambda = s).$$

Also, define $A_\lambda^{(1)} \equiv \mathbb{E}[Q_0] \sum_{s > k_2 n_\lambda \text{ or } s < k_1 n_\lambda} \mathbb{P}(N_\lambda = s)$ and $A_\lambda^{(2)} \equiv \sum_{s > k_2 n_\lambda \text{ or } s < k_1 n_\lambda} s \cdot \bar{q}_\rho \mathbb{P}(N_\lambda = s)$. Note that Q_0 has the same distribution as the steady-state number in the system in an $M/GI/\infty$ model with Poisson arrivals at rate $\lambda = rn_\lambda \rho$ and i.i.d. generally distributed service times having the same distribution, F , as the abandonment times in our original model. Therefore, exploiting standard results for the infinite-server queue, Q_0 has a Poisson distribution with mean $\lambda/\theta = rn_\lambda \rho/\theta$, i.e., $\mathbb{E}[Q_0] = \mathcal{O}(\lambda)$. Applying Hoeffding's inequality to the binomial distribution: $\mathbb{P}(k_1 n_\lambda \leq N_\lambda \leq k_2 n_\lambda) \geq 1 - 2e^{-2\epsilon^2 n_\lambda}$; equivalently, $\mathbb{P}(k_1 n_\lambda > N_\lambda \text{ or } N_\lambda > k_2 n_\lambda) \leq 2e^{-2\epsilon^2 n_\lambda}$. Thus,

$$A_\lambda^{(1)} = \mathbb{E}[Q_0] \sum_{s > k_2 n_\lambda \text{ or } s < k_1 n_\lambda} \mathbb{P}(N_\lambda = s) = \mathbb{E}[Q_0] \cdot \mathbb{P}(k_1 n_\lambda > N_\lambda \text{ or } N_\lambda > k_2 n_\lambda) \rightarrow 0$$

as $\lambda \rightarrow \infty$.

We now turn to showing that $A_\lambda^{(2)}$ is asymptotically negligible as well. Note that:

$$A_\lambda^{(2)} = \bar{q}_\rho \sum_{s > k_2 n_\lambda \text{ or } s < k_1 n_\lambda} s \mathbb{P}(N_\lambda = s) = \bar{q}_\rho \mathbb{E}[N_\lambda \mathbb{1}\{N_\lambda > k_2 n_\lambda \text{ or } N_\lambda < k_1 n_\lambda\}],$$

where $\mathbb{1}\{\cdot\}$ denotes an indicator random variable. By the Cauchy-Schwarz inequality:

$$\begin{aligned} & \mathbb{E}[N_\lambda \mathbb{1}\{N_\lambda > k_2 n_\lambda \text{ or } N_\lambda < k_1 n_\lambda\}] \\ & \leq \sqrt{\mathbb{E}[N_\lambda^2] \mathbb{P}(N_\lambda > k_2 n_\lambda \text{ or } N_\lambda < k_1 n_\lambda)} \\ & = \sqrt{(n_\lambda r(1-r) + n_\lambda^2 r^2) \mathbb{P}(N_\lambda > k_2 n_\lambda \text{ or } N_\lambda < k_1 n_\lambda)} \rightarrow 0 \text{ as } \lambda \rightarrow \infty. \end{aligned}$$

Therefore, $A_\lambda^{(2)} \rightarrow 0$ as $\lambda \rightarrow \infty$. Combining the above, we obtain that $A_\lambda \rightarrow 0$ as well.

Asymptotic bound for N_λ close to $n_\lambda r$ We now characterize B_λ for large λ .

Lemma 6 *There exists a finite constant $C > 0$ such that $\limsup_{\lambda \rightarrow \infty} B_\lambda \leq C$.*

Proof We begin by writing B_λ as follows,

$$B_\lambda \leq \sum_{k_1 n_\lambda \leq s \leq k_2 n_\lambda} \left| \mathbb{E}[Q_s] - s \bar{q}_{\rho_s} \right| \mathbb{P}(N_\lambda = s) + \left| \sum_{k_1 n_\lambda \leq s \leq k_2 n_\lambda} s(\bar{q}_{\rho_s} - \bar{q}_\rho) \mathbb{P}(N_\lambda = s) \right|, \tag{13.15}$$

where $\rho_s \equiv n_\lambda r \rho / s$ and \bar{q}_{ρ_s} is the fluid limit for the queue length in the $M/M/s + GI$ queue with traffic intensity ρ_s (the arrival rate is $\lambda = r n_\lambda \rho$ and the number of servers is s). Let,

$$\begin{aligned} B_\lambda^{(1)} & \equiv \sum_{k_1 n_\lambda \leq s \leq k_2 n_\lambda} \left| \mathbb{E}[Q_s] - s \bar{q}_{\rho_s} \right| \mathbb{P}(N_\lambda = s) \quad \text{and} \\ B_\lambda^{(2)} & \equiv \left| \sum_{k_1 n_\lambda \leq s \leq k_2 n_\lambda} s(\bar{q}_{\rho_s} - \bar{q}_\rho) \mathbb{P}(N_\lambda = s) \right|. \end{aligned}$$

First, we consider $B_\lambda^{(1)}$ and show that it is asymptotically bounded. Fix n_λ and note that to each $k_1 n_\lambda \leq s \leq k_2 n_\lambda$ corresponds a traffic intensity ρ_s in the $M/M/s + GI$ system, where $\rho_s = n_\lambda r \rho / s$ and $1 < \rho r / (r + \epsilon) \leq \rho_s \leq \rho r / (r - \epsilon)$. By Theorem 5 of Bassamboo and Randhawa (2010), assuming that f is strictly positive and continuously differentiable,

$$\limsup_{\lambda \rightarrow \infty} \left| \mathbb{E}[Q_s] - s \bar{q}_{\rho_s} \right| \leq \sqrt{f(\bar{w}_{\rho_s})} \left(\frac{3|f'(\bar{w}_{\rho_s})|}{\rho_s f^2(\bar{w}_{\rho_s})} + \frac{1}{2} \right), \tag{13.16}$$

where \bar{w}_{ρ_s} is the fluid limit for the steady-state waiting time in the overloaded $M/M/s + GI$ queue with traffic intensity ρ_s . Note that for $\rho r / (r + \epsilon) \leq \rho_s \leq \rho r / (r - \epsilon)$, we have that $\bar{w}_{\rho r / (r + \epsilon)} \leq \bar{w}_{\rho_s} \leq \bar{w}_{\rho r / (r - \epsilon)}$. By the continuity of the bounding function in Eq. 13.16 and the boundedness theorem, we conclude that there exists a finite constant $C_1 > 0$ such that

$$\sup_{k_1 n_\lambda \leq s \leq k_2 n_\lambda} \sqrt{f(\bar{w}_{\rho_s})} \left(\frac{3|f'(\bar{w}_{\rho_s})|}{\rho' f^2(\bar{w}_{\rho_s})} + \frac{1}{2} \right) \leq C_1. \quad (13.17)$$

Since

$$\begin{aligned} B_\lambda^{(1)} &= \sum_{k_1 n_\lambda \leq s \leq k_2 n_\lambda} |\mathbb{E}[Q_s] - s\bar{q}_{\rho_s}| \mathbb{P}(N_\lambda = s) \\ &\leq \sup_{k_1 n_\lambda \leq s \leq k_2 n_\lambda} |\mathbb{E}[Q_s] - s\bar{q}_{\rho_s}| \sum_{k_1 n_\lambda \leq s \leq k_2 n_\lambda} \mathbb{P}(N_\lambda = s) \\ &\leq \sup_{k_1 n_\lambda \leq s \leq k_2 n_\lambda} |\mathbb{E}[Q_s] - s\bar{q}_{\rho_s}|, \end{aligned}$$

combining Eqs. 13.16 and 13.17 yields that $\limsup_{\lambda \rightarrow \infty} B_\lambda^{(1)} \leq C_1$ by taking limits on both sides. There remains to study the asymptotic behaviour of $B_\lambda^{(2)}$. Note that $\bar{q}_{\rho_s} = \rho_s \int_0^{(\bar{F})^{-1}(1/\rho_s)} \bar{F}(u) du$, e.g., by equations (3.6) and (3.7) in Whitt (2006a). Consider,

$$\begin{aligned} &\left| \sum_{s \geq 0} s \left(\rho_s \int_0^{(\bar{F})^{-1}(1/\rho_s)} \bar{F}(x) dx - \rho \int_0^{(\bar{F})^{-1}(1/\rho)} \bar{F}(u) du \right) \mathbb{P}(N_\lambda = s) \right| \\ &= \left| \sum_{s \geq 0} \left(n_\lambda r \rho \int_0^{(\bar{F})^{-1}(s/n_\lambda r \rho)} \bar{F}(u) du - s \rho \int_0^{(\bar{F})^{-1}(1/\rho)} \bar{F}(u) du \right) \mathbb{P}(N_\lambda = s) \right|, \\ &= \left| \mathbb{E} \left[\left(n_\lambda r \rho \int_0^{(\bar{F})^{-1}(N_\lambda/n_\lambda r \rho)} \bar{F}(u) du - N_\lambda \rho \int_0^{(\bar{F})^{-1}(1/\rho)} \bar{F}(u) du \right) \right] \right|, \\ &= \left| n_\lambda \rho r \mathbb{E} \left[\left(\int_{(\bar{F})^{-1}(1/\rho)}^{(\bar{F})^{-1}(N_\lambda/n_\lambda r \rho)} \bar{F}(u) du \right) \right] \right|. \end{aligned}$$

We now show that there must exist a finite constant $C_2 > 0$ such that

$$\left| n_\lambda \rho r \mathbb{E} \left[\left(\int_{(\bar{F})^{-1}(1/\rho)}^{(\bar{F})^{-1}(N_\lambda/n_\lambda r \rho)} \bar{F}(u) du \right) \right] \right| \leq C_2$$

for λ large enough. To this aim, define the function

$$g_\lambda(x) = n_\lambda \rho r \int_{(\bar{F})^{-1}(1/\rho)}^{(\bar{F})^{-1}(x/n_\lambda r \rho)} \bar{F}(u) du \text{ for } x \geq 0.$$

For a given λ , we use a Taylor-series expansion of $\mathbb{E}[g_\lambda(N_\lambda)]$ around $\mathbb{E}[N_\lambda] = n_\lambda r$ (we can do this since g_λ is sufficiently differentiable and the moments of N_λ are finite):

$$|\mathbb{E}[g_\lambda(N_\lambda)]| = \left| \mathbb{E} \left[g_\lambda(n_\lambda r) + g'_\lambda(n_\lambda r) (N_\lambda - n_\lambda r) + \frac{1}{2} g''_\lambda(n_\lambda r) (N_\lambda - n_\lambda r)^2 \right] \right| + \mathcal{O}(1/\lambda).$$

Indeed, by computing the centralized moments of N_λ and higher-order derivatives of g_λ , it can be shown that the remainder term in the Taylor series is $\mathcal{O}(1/\lambda)$. Also, $g_\lambda(n_\lambda r) = 0$ and

$$g'_\lambda(n_\lambda r) = -\frac{1/\rho}{f(\bar{F}^{-1}(1/\rho))} \quad \text{and} \quad g''_\lambda(n_\lambda r) = -\frac{1}{rn_\lambda \rho} \frac{h_1(\rho) + (1/\rho)h_2(\rho)/h_1(\rho)}{h_1^2(\rho)},$$

where $h_1(\rho) = f(\bar{F}^{-1}(1/\rho))$ and $h_2(\rho) = f'(\bar{F}^{-1}(1/\rho))$. Thus, there exists $C_2 > 0$ such that:

$$|\mathbb{E}[g_\lambda(N_\lambda)]| \approx \left| \frac{1}{2} g''_\lambda(n_\lambda r) n_\lambda r (1 - r) \right| \leq C_2 \quad \text{for } \lambda \text{ large enough.}$$

We now turn to the asymptotic behaviour of $B_\lambda^{(2)}$. Note that:

$$B_\lambda^{(2)} = |\mathbb{E}[g_\lambda(N_\lambda) \mathbb{1}\{N_\lambda \in [k_1 n_\lambda, k_2 n_\lambda]\}]|, \quad \text{and} \\ |\mathbb{E}[g_\lambda(N_\lambda)]| = |\mathbb{E}[g_\lambda(N_\lambda) \mathbb{1}\{N_\lambda \in [k_1 n_\lambda, k_2 n_\lambda]\}] + \mathbb{E}[g_\lambda(N_\lambda) \mathbb{1}\{N_\lambda \notin [k_1 n_\lambda, k_2 n_\lambda]\}]|.$$

Bounding the second term in the last equality,

$$\begin{aligned} & \mathbb{E}[g_\lambda(N_\lambda) \mathbb{1}\{N_\lambda \notin [k_1 n_\lambda, k_2 n_\lambda]\}] \\ & \leq |\mathbb{E}[g_\lambda(N_\lambda) \mathbb{1}\{N_\lambda \notin [k_1 n_\lambda, k_2 n_\lambda]\}]| \\ & \leq \sqrt{\mathbb{E}[g_\lambda^2(N_\lambda)] \mathbb{P}(N_\lambda \notin [k_1 n_\lambda, k_2 n_\lambda])} \quad (\text{Cauchy Schwarz inequality}) \\ & \rightarrow 0, \end{aligned}$$

since $\mathbb{P}(N_\lambda \notin [k_1 n_\lambda, k_2 n_\lambda])$ vanishes exponentially fast as $\lambda \rightarrow \infty$, and $\mathbb{E}[g_\lambda^2(N_\lambda)] = \mathcal{O}(\lambda^2)$ since

$$\int_{(\bar{F})^{-1}(1/\rho)}^{(\bar{F})^{-1}(N_\lambda/n_\lambda r \rho)} \bar{F}(u) \, du \leq 1/\theta.$$

Thus,

$$\limsup_{\lambda \rightarrow \infty} B_{\lambda}^{(2)} = \limsup_{\lambda \rightarrow \infty} |\mathbb{E}[g_{\lambda}(N_{\lambda}) \mathbb{1}\{N_{\lambda} \in [k_1 n_{\lambda}, k_2 n_{\lambda}]\}]| \leq C_2.$$

Combining the above, there exists $C > 0$ such that $\limsup_{\lambda \rightarrow \infty} B_{\lambda} \leq C$.

$\mathcal{O}(1)$ -accuracy Since both A_{λ} and B_{λ} are asymptotically bounded, there must exist $K > 0$ such that, as desired:

$$\limsup_{\lambda \rightarrow \infty} |\mathbb{E}[Q_{N_{\lambda}}] - r n_{\lambda} \bar{q}_{\rho}| \leq K.$$

$\mathcal{o}(1)$ -Accuracy for the Fluid Net Abandonment Rate

The proof for the net abandonment rate proceeds along similar lines, so we will be brief. Paralleling Eq. 13.16, and denoting $\mathbb{E}[\alpha_{N_{\lambda}} | N_{\lambda} = s] \equiv \mathbb{E}[\alpha_s]$, we can exploit Theorem 5 in Bassamboo and Randhawa (2010) to show that $\sum_{k_1 n_{\lambda} \leq s \leq k_2 n_{\lambda}} (\mathbb{E}[\alpha_s] - s \bar{\alpha}_{\rho_s}) \mathbb{P}(N_{\lambda} = s) \rightarrow 0$ as $\lambda \rightarrow \infty$. Moreover, by equation (3.3) in Whitt (2006a): $\bar{\alpha}_{\rho_s} = \rho_s - 1$; thus, $s(\bar{\alpha}_{\rho_s} - \bar{\alpha}_{\rho}) = \rho(n_{\lambda} r - s)$. We can then write:

$$\sum_{k_1 n_{\lambda} \leq s \leq k_2 n_{\lambda}} s(\bar{\alpha}_{\rho_s} - \bar{\alpha}_{\rho}) \mathbb{P}(N_{\lambda} = s) = \rho \mathbb{E}[(nr - N_{\lambda}) \mathbb{1}(k_1 n_{\lambda} \leq N_{\lambda} \leq k_2 n_{\lambda})],$$

and deduce that $\mathbb{E}[(nr - N_{\lambda}) \mathbb{1}(k_1 n_{\lambda} \leq N_{\lambda} \leq k_2 n_{\lambda})] \rightarrow 0$ since $\mathbb{E}[N_{\lambda}] = r n_{\lambda}$.

The Underloaded Regime

Let $0 < \epsilon < r$ be small enough so that $\rho r / (r - \epsilon) < 1$, and recall that $k_1 \equiv r - \epsilon$ and $k_2 \equiv r + \epsilon$. Then, conditioning on N_{λ} :

$$\begin{aligned} \mathbb{E}[Q_{N_{\lambda}}] &= \sum_{k_1 n_{\lambda} \leq s \leq k_2 n_{\lambda}} \mathbb{E}[Q_s] \mathbb{P}(N_{\lambda} = s) + \sum_{k_1 n_{\lambda} > s \text{ or } s > k_2 n_{\lambda}} \mathbb{E}[Q_s] \mathbb{P}(N_{\lambda} = s), \\ &\leq \sum_{k_1 n_{\lambda} \leq s \leq k_2 n_{\lambda}} \mathbb{E}[Q_s] \mathbb{P}(N_{\lambda} = s) + \mathbb{E}[Q_0] \sum_{k_1 n_{\lambda} > s \text{ or } s > k_2 n_{\lambda}} \mathbb{P}(N_{\lambda} = s). \end{aligned}$$

As in the proof of Theorem 1, we can show that: $\mathbb{E}[Q_0] \sum_{k_1 n_{\lambda} > s \text{ or } s > k_2 n_{\lambda}} \mathbb{P}(N_{\lambda} = s) \rightarrow 0$ as $\lambda \rightarrow \infty$. Also, $\sum_{k_1 n_{\lambda} \leq s \leq k_2 n_{\lambda}} \mathbb{E}[Q_s] \mathbb{P}(N_{\lambda} = s) \leq \mathbb{E}[Q(k_1 n_{\lambda})] \sum_{k_1 n_{\lambda} \leq s \leq k_2 n_{\lambda}} \mathbb{P}(N_{\lambda} = s)$. Since $\mathbb{E}[Q(k_1 n_{\lambda})]$ is the expected steady-state queue length in an underloaded queue, it converges to 0 as $\lambda \rightarrow \infty$, e.g., see Theorem 5.1 in Zeltyn and Mandelbaum (2005). The limit for the net abandonment follows similarly.

The Critically-Loaded Regime

We condition on N_λ :

$$\begin{aligned}
 \mathbb{E}[Q_{N_\lambda}] &= \sum_{k_1 n_\lambda \leq s < n_\lambda r} \mathbb{E}[Q_s] \mathbb{P}(N_\lambda = s) \\
 &+ \sum_{n_\lambda r < s \leq k_2 n_\lambda} \mathbb{E}[Q_s] \mathbb{P}(N_\lambda = s) + \mathbb{E}[Q_{n_\lambda r}] \mathbb{P}(N_\lambda = n_\lambda r), \\
 &\leq \sum_{k_1 n_\lambda \leq s < n_\lambda r} |\mathbb{E}[Q_s] - s \bar{q}_{\rho_s}| \mathbb{P}(N_\lambda = s) + \sum_{k_1 n_\lambda \leq s < n_\lambda r} s \bar{q}_{\rho_s} \mathbb{P}(N_\lambda = s) \\
 &+ \sum_{n_\lambda r < s \leq k_2 n_\lambda} \mathbb{E}[Q_s] \mathbb{P}(N_\lambda = s) + \mathbb{E}[Q(n_\lambda r)] \mathbb{P}(N_\lambda = n_\lambda r), \quad (13.18)
 \end{aligned}$$

where $\rho_s = r \rho n_\lambda / s$. Paralleling Eqs. 13.16 and 13.17, we can show that there exists a finite constant C'_1 such that for large λ : $\sum_{k_1 n_\lambda \leq s < n_\lambda r} |\mathbb{E}[Q_s] - s \bar{q}_{\rho_s}| \mathbb{P}(N_\lambda = s) \leq C'_1$ since $\rho_s > 1$ for all $k_1 n_\lambda \leq s < n_\lambda r$. Also,

$$\begin{aligned}
 &\sum_{k_1 n_\lambda \leq s < n_\lambda r} s \bar{q}_{\rho_s} \mathbb{P}(N_\lambda = s) \\
 &= \sum_{k_1 n_\lambda \leq s < n_\lambda r} n_\lambda r \left(\int_0^{(\bar{F})^{-1}(s/n_\lambda r)} \bar{F}(x) dx \right) \mathbb{P}(N_\lambda = s) \\
 &= \mathbb{E} \left[\left(n_\lambda r \int_0^{(\bar{F})^{-1}(N_\lambda/n_\lambda r)} \bar{F}(x) dx \right) \mathbb{1}(N_\lambda \in [k_1 n_\lambda, n_\lambda r)) \right]. \quad (13.19)
 \end{aligned}$$

Using arguments as in Theorem 1 (noting e.g., that $g_\lambda(n_\lambda r) = \int_0^{(\bar{F})^{-1}(1)} \bar{F}(x) dx = 0$), we can show that there exists a finite $C'_2 > 0$ such that

$$\limsup_{\lambda \rightarrow \infty} \mathbb{E} \left[\left(n_\lambda r \int_0^{(\bar{F})^{-1}(N_\lambda/n_\lambda r)} \bar{F}(x) dx \right) \mathbb{1}(N_\lambda \in [k_1 n_\lambda, n_\lambda r)) \right] \leq C'_2.$$

By Theorem 4.1 of Zeltyn and Mandelbaum (2005), there exists $K' > 0$ such that $\mathbb{E}[Q_{n_\lambda r}] \leq K' \sqrt{\lambda}$ for large enough λ . Given that $\sum_{n_\lambda r < s \leq k_2 n_\lambda} \mathbb{E}[Q_s] \mathbb{P}(N_\lambda = s) \rightarrow 0$ as $\lambda \rightarrow \infty$ (underloaded regime), we obtain that the entire expression in Eq. 13.18 is $\mathcal{O}(\sqrt{\lambda})$. The proof for the abandonment rate follows along similar lines, so we omit the relevant details.

Proofs of Propositions

Proposition 2 If the abandonment distribution is exponential, then for $\Gamma_{i-1} \leq n < \Gamma_i$, $u_i(n) = \sum_{j=1}^k c_j r_j n + \sum_{j=i}^k (p_j + h_j/\theta)(\lambda_j - nr_j)$. We assume:

$$\sum_{j=1}^k c_j r_j - \sum_{j=i_0}^k L_j r_j < 0 \quad \text{and} \quad \sum_{j=1}^k c_j r_j - \sum_{j=i_0+1}^k L_j r_j > 0. \quad (13.20)$$

Clearly, under condition Eq. 13.20, $C(n)$ is piecewise linear with piecewise negative slopes for $n \leq \Gamma_{i_0}$, and strictly positive slopes for $n > \Gamma_{i_0}$.

With a monotonically increasing hazard rate, we have

$$u_i(n) = \sum_{j=1}^k c_j r_j n + \sum_{j=i}^k \left(p_j (\lambda_j - nr_j) + h_j \lambda_j \int_0^{\bar{F}^{-1}(nr_j/\lambda_j)} \bar{F}(u) du \right).$$

Thus,

$$u_i'(n) = \sum_{j=1}^k c_j r_j - \sum_{j=i}^k r_j \left[p_j + \frac{h_j}{h_a(\bar{F}^{-1}(n/\Gamma_j))} \right],$$

which is strictly decreasing in n , i.e., $u_i''(n) < 0$. Thus, the objective is piecewise strictly concave. The minimum must be achieved at some Γ_i , at which we critically load shift i' .

Proposition 3 In $[\Gamma_{i-1}, \Gamma_i)$, $u_i'(n)$ is as in the proof of Proposition 2, so that $u_i''(n) > 0$ and the function is piecewise convex. It also follows that $u_i'(n_1) < u_{i+1}'(n_2)$ for $n_1 \in [\Gamma_i, \Gamma_{i+1})$ and $n_2 \in [\Gamma_{i+1}, \Gamma_{i+2})$. In other words, if $C'(x) > 0$, then $C'(y) > 0$ for $y \geq x$. Thus, the minimum n^* will be at the interior of an interval $(\Gamma_{i_0-1}, \Gamma_{i_0})$ if $u_{i_0}'(\Gamma_{i_0-1}) < 0$ and $u_{i_0}'(\Gamma_{i_0}-) > 0$. Here is a sufficient condition for this to be the case.

Sufficient condition There exists $i_0, \beta, \gamma > 0$ such that:

$$\begin{aligned} \frac{\Gamma_{i_0-1}}{\Gamma_{i_0}} < \beta; & \quad \sum_{i=1}^k c_i r_i - \sum_{i=i_0}^k r_i \left(p_i + \frac{h_i}{f_a(\bar{F}^{-1}(\beta))} \right) < 0; \\ \frac{\Gamma_{i_0}}{\Gamma_k} > \gamma; & \quad \sum_{i=1}^k c_i r_i - \sum_{i=i_0}^k r_i \left(p_i + \frac{h_i}{f_a(\bar{F}^{-1}(\gamma))} \right) > 0. \end{aligned}$$

To see why this implies an interior point solution, note that:

$$\begin{aligned}
 u'_{i_0}(\Gamma_{i_0-1}) &= \sum_{i=1}^k c_i r_i - \sum_{i=i_0}^k r_i \left(p_i + \frac{h_i}{f_a(\bar{F}^{-1}(\Gamma_{i_0-1}/\Gamma_i))} \right) \\
 &< \sum_{i=1}^k c_i r_i - \sum_{i=i_0}^k r_i \left(p_i + \frac{h_i}{f_a(\bar{F}^{-1}(\Gamma_{i_0-1}/\Gamma_{i_0}))} \right) \\
 &< \sum_{i=1}^k c_i r_i - \sum_{i=i_0}^k r_i \left(p_i + \frac{h_i}{f_a(\bar{F}^{-1}(\beta))} \right) < 0 \quad \text{by assumption.}
 \end{aligned}$$

Furthermore,

$$\begin{aligned}
 u'_{i_0}(\Gamma_{i_0}^-) &= \sum_{i=1}^k c_i r_i - \sum_{i=i_0}^k r_i \left(p_i + \frac{h_i}{f_a(\bar{F}^{-1}(\Gamma_{i_0}^-/\Gamma_i))} \right) \\
 &> \sum_{i=1}^k c_i r_i - \sum_{i=i_0}^k r_i \left(p_i + \frac{h_i}{f_a(\bar{F}^{-1}(\Gamma_{i_0}^-/\Gamma_k))} \right) \\
 &> \sum_{i=1}^k c_i r_i - \sum_{i=i_0}^k r_i \left(p_i + \frac{h_i}{f_a(\bar{F}^{-1}(\gamma))} \right) > 0.
 \end{aligned}$$

Combining both, we get that $u'_{i_0}(\Gamma_{i_0-1}) < 0$ and $u'_{i_0}(\Gamma_{i_0}^-) > 0$ which, combined with the fact that $C'(\cdot)$ increases across intervals, implies that the minimizer must lie strictly in the interval $(\Gamma_{i_0-1}, \Gamma_{i_0})$. In words, if the imbalance between the augmented arrival rates $\Gamma_{i_0}/\Gamma_{i_0+1}$ is small enough, it is optimal to “strike a balance” between the two shifts, i.e., underloading a shift, while overloading the other.

Proposition 4 It suffices to show that $\theta(w_{i+1}^e(\Gamma_{i_c})) > \theta_0$ for $i \geq i_c$. To see this, note that: $\lambda_{i+1} e^{-w_{i+1}^e \theta(w_{i+1}^e)} = \Gamma_{i_c} r_{i+1}$. This implies: $e^{-w_{i+1}^e \theta(w_{i+1}^e)} = \Gamma_{i_c}/\Gamma_{i+1}$, for $i \geq i_c$, i.e., $w_{i+1}^e \theta(w_{i+1}^e) = \ln(\Gamma_{i+1}/\Gamma_{i_c})$. Assume that $\theta_0 \cdot \theta^{-1}(\theta_0) < \ln(\Gamma_{i_c+1}/\Gamma_{i_c})$. Then, $\theta_0 \cdot \theta^{-1}(\theta_0) < w_{i+1}^e \theta(w_{i+1}^e)$ for $i \geq i_c$ since $\Gamma_{i_c+1} \leq \Gamma_{i+1}$ for $i \geq i_c$. Since $x\theta^{-1}(x)$ is increasing in x , we obtain that $w_{i+1}^e > \theta^{-1}(\theta_0)$, which implies that $\theta(w_{i+1}^e) > \theta_0$ for $i \geq i_c$, as desired. Then, $C_a(\Gamma_{i_c}) < C(\Gamma_{i_c}) = C^*$, and we get strict reduction in cost due to the announcements.

Proposition 5 It suffices to show that, for all n , $C'_a(n) > C'(n)$. If this holds, then $C'(n) > 0$ would imply $C'_a(n) > 0$, so that $C_a(\cdot)$ increases whenever $C(\cdot)$ increases, which leads to $n_a^* < n^*$. Fix $n \in [\Gamma_{i_0-1}, \Gamma_{i_0})$, for some i_0 . Then, $C_a(n) = \sum_{i=1}^k c_i n r_i + \sum_{i=i_0}^k (p_i + h_i/\theta(w_i^e))(\lambda_i - n r_i)$ where $e^{-w_i^e \theta(w_i^e)} = n/\Gamma_i$. That is, for $i \geq i_0$,

$$w_i^e \theta(w_i^e) = \ln\left(\frac{\Gamma_i}{n}\right) > \ln\left(\frac{\Gamma_i}{\Gamma_{i_0}}\right) > \min_{1 \leq i \leq k-1} \ln\left(\frac{\Gamma_{i+1}}{\Gamma_i}\right) > \theta_0 \cdot \theta^{-1}(\theta_0),$$

under condition Eq. 13.9. This implies that $\theta(w_i^e) > \theta_0$ for all $i \geq i_0$. Note that for $n \in [\Gamma_{i_0-1}, \Gamma_{i_0})$,

$$\begin{aligned} C'_a(n) &= \sum_{i=1}^k c_i r_i - \sum_{i=i_0}^k p_i r_i - \sum_{i=i_0}^k \frac{r_i h_i}{\theta(w_i^e)} \\ &\quad + \sum_{i=i_0}^k h_i (\lambda_i - n r_i) \frac{\theta'(w_i^e)}{n \theta^2(w_i^e) (\theta'(w_i^e) w_i^e + \theta(w_i^e))}. \end{aligned}$$

Thus, assuming condition Eq. 13.9 implies that $C'_a(n) > C'(n)$, for every $n \in [\Gamma_{i_0-1}, \Gamma_{i_0})$. Since we can let i_0 denote any period index, we obtain that $n_a^* < n^*$.

Lemma 1 We derive the optimal compensation for a fixed value of the pool size n . Since c_i^* can be decided upon separately for each shift, we focus on a single shift setting in what follows, i.e., we fix the shift i . The solution depends on the specific value of n .

1. $n \geq \lambda_i / G(l)$. $c_i^* = l$, i.e., offer minimum wage and overstaff shift i (under-loaded).
2. $n < \lambda_i / G(l)$. Note that we must have that $\lambda_i \geq nG(c_i)$ i.e., $c_i \leq G^{-1}(\lambda_i/n)$ because it will not be cost effective for the manager to incite more supply than the demand in the shift.

Subcase 1 We assume that $L_i \leq l$. In this case, the problem becomes:

$$\min_{L_i \leq l \leq c_i \leq G^{-1}(\lambda_i/n)} \{n c_i G(c_i) + L_i (\lambda_i - n G(c_i))\}$$

which is equivalent to

$$\text{minimize}_{L_i \leq l \leq c_i \leq G^{-1}(\lambda_i/n)} t_i(c_i) \equiv (c_i - L_i) G(c_i).$$

Since $c_i > L_i$, it is readily seen that the objective is increasing in c_i . Thus, we must have that $c_i^* = l$. That is, we offer minimum wage and understaff shift i (over-loaded).

Subcase 2 We now assume that $L_i > l$. In this case, $\lambda_i / G(L_i) < \lambda_i / G(l)$. We then consider the two intervals: (a) $n \leq \lambda_i / G(L_i) < \lambda_i / G(l)$ and (b) $\lambda_i / G(L_i) < n < \lambda_i / G(l)$.

(a) $n \leq \lambda_i / G(L_i) < \lambda_i / G(l)$. The problem is now:

$$\min_{l \leq c_i \leq \min\{G^{-1}(\lambda_i/n), L_i\}} \{n c_i G(c_i) + L_i (\lambda_i - n G(c_i))\}$$

which is equivalent to solving:

$$\underset{l \leq c_i \leq \min\{L_i, G^{-1}(\lambda_i/n)\}}{\text{minimize}} \quad t_i(c_i) \equiv (c_i - L_i)G(c_i).$$

Note that $t'_i(c_i) = G(c_i)(1 + (c_i - L_i)g(c_i)/G(c_i))$. In this case, we have $L_i \leq G^{-1}(\lambda_i/n)$. Since $t'(L_i) \geq 0$, and $t_i(\cdot)$ is convex under log-concavity of G , we obtain that:

- (i) If $t'_i(l) < 0$ i.e., $(1 + (l - L_i)g(l)/G(l)) < 0$, then there exists an optimal $c_i^* = a_i \in (l, L_i)$ where $t'(a_i) = 0$;
- (ii) If $t'_i(l) \geq 0$ i.e., $(1 + (l - L_i)g(l)/G(l)) \geq 0$, then we have $c_i^* = l$.

In both cases (i) and (ii), the system is overloaded, i.e., the manager incites a smaller supply than the demand in shift i .

- (b) Now, consider: $\lambda_i/G(L_i) < n < \lambda_i/G(l)$. Let $0 < a_i < L_i$ be such that $t'_i(a_i) = 0$ i.e.,

$$G(a_i) \left(1 + (a_i - L_i) \frac{g(a_i)}{G(a_i)} \right) = 0.$$

The optimization problem is

$$\min_{l \leq c_i \leq G^{-1}(\lambda_i/n) < L_i} t_i(c_i).$$

Note that if $a_i < l$, then $c_i^* = l$ (by the convexity of the objective); in other words, the manager offers the minimum wage and runs shift i overloaded. Now, assume that $a_i \geq l$. We then have the following two cases:

- (i) $t'(G^{-1}(\lambda_i/n)) \leq 0$ i.e., $G^{-1}(\lambda_i/n) \leq a_i$, i.e., $\lambda_i/G(L_i) < \lambda_i/G(a_i) \leq n < \lambda_i/G(l)$. In this case, $c_i^* = G^{-1}(\lambda_i/n)$ which means that the manager incites a supply equal to the demand, i.e., she critically loads her shift.
- (ii) $t'(G^{-1}(\lambda_i/n)) > 0$, i.e., $G^{-1}(\lambda_i/n) > a_i$, i.e., $\lambda_i/G(L_i) < n < \lambda_i/G(a_i) \leq \lambda_i/G(l)$. In this case, $c_i^* = a_i$ and the manager incites a supply that is smaller than the demand, i.e., she overloads her shift.

Lemma 2 We let \tilde{a}_k be the solution to Eq. 13.12. Then, $\tilde{t}'(x) \equiv G(x)(1 + (x - \tilde{L}_k)g(x)/G(x))$ is increasing for $x \leq \tilde{L}_k$ by the log-concavity of $G(\cdot)$. If $a_k > \tilde{L}_k$, then it must be that $a_k > \tilde{a}_k$ since $\tilde{a}_k < \tilde{L}_k$. Let us now assume that $a_k \leq \tilde{L}_k$. Since $\tilde{L}_k < L_k$, we must have that

$$\begin{aligned} G(a_k) \left(1 + (a_k - \tilde{L}_k) \frac{g(a_k)}{G(a_k)} \right) &> G(a_k) \left(1 + (a_k - L_k) \frac{g(a_k)}{G(a_k)} \right) \\ &= G(\tilde{a}_k) \left(1 + (\tilde{a}_k - \tilde{L}_k) \frac{g(\tilde{a}_k)}{G(\tilde{a}_k)} \right) 0. \end{aligned}$$

Because $\tilde{r}'(x)$ is increasing in x for $x \leq \tilde{L}_k$, and we have both $a_k, \tilde{a}_k \leq \tilde{L}_k$, we also obtain that $a_k > \tilde{a}_k$. If $n < \lambda_k/G(a_k)$, then we must also have that $n < \lambda_k/G(\tilde{a}_k)$, so that the optimal compensation as per Lemma 1 is to set $\tilde{c}_k^* = \tilde{a}_k < c_k^* = a_k$. We note that if n is as in cases (a) and (b) of Lemma 1, then the compensation offered to agents is unchanged since compensation is set so that there is no congestion in the shift. We also note that if $\tilde{a}_k < l < a_k$ then $\tilde{c}_k^* = l$ so that $\tilde{c}_k^* < c_k^*$ as well. In other words, agents are worse off in all cases.

Lemma 3 Note that if $l < l_0$ then $\max\{\lambda_i/G(\tilde{a}_i)\} < \min\{\lambda_i/G(l)\}$. For $\max\{\lambda_i/G(\tilde{a}_i)\} < n < \min\{\lambda_i/G(l)\}$, we must have that $\Pi'(n) < 0$. Thus, $n^* \geq \min\{\lambda_i/G(l)\} > \max\{\lambda_i/G(\tilde{a}_i)\}$, and we do not overload or use the announcements in any shift (since $\Pi'(n)$ is strictly increasing in n). It is readily seen that we cannot, for an optimal n^* , have all shifts strictly underloaded. Thus, there must exist i_0 as specified in the lemma.

Lemma 4 In this case, problem Eq. 13.13 simplifies to:

$$\begin{aligned} \underset{n \geq 0}{\text{minimize}} \Pi(n) \equiv & \sum_{\{i:n \geq \lambda_i/G(l)\}} nlG(l) \quad (\text{underloaded}) \\ & + \sum_{\{i:n < \lambda_i/G(l)\}} \ln G(l) + \tilde{L}_i(\lambda_i - nG(l)) \quad (\text{overload+} \\ & \text{announcements}) \end{aligned}$$

Note that $\Pi(n)$ is piecewise linear. Then, $\Pi'(n) = klG(l) - \sum_{\{i:n < \lambda_i/G(l)\}} \tilde{L}_i G(l)$. Clearly, as n increases, $\Pi'(n)$ increases too. Under our assumptions, there must exist a unique k_0 such that $\Pi'(n) < 0$ for $n < \lambda_{k_0}/G(l)$ and $\Pi'(n) > 0$ for $n > \lambda_{k_0}/G(l)$. The optimal solution is to set $n^* = \lambda_{k_0}/G(l)$.

References

- Akşın OZ, Armony M, Mehrotra V (2007) The modern call center: a multi-disciplinary perspective on operations management research. *Prod Oper Manag* 16(6):665–688
- Akşın Z, Ata B, Emadi SM, Su CL (2013) Structural estimation of callers' delay sensitivity in call centers. *Manag Sci* 59(12):2727–2746
- Akşın Z, Ata B, Emadi S, Su C (2016) Impact of delay announcements in call centers: an empirical approach. *Oper Res* 65(1):242–265
- Aldor-Noiman S, Feigin PD, Mandelbaum A (2009) Workload forecasting for a call center: methodology and a case study. *Ann Appl Stat* 3(4):1403–1447
- Armony M, Shimkin N, Whitt W (2009) The impact of delay announcements in many-server queues with abandonment. *Oper Res* 57(1):66–81
- Bassamboo A, Randhawa R (2010) On the accuracy of fluid models for capacity sizing in queueing systems with impatient customers. *Oper Res* 58(5):1398–1413
- Bassamboo A, Zeevi A (2009) On a data-driven method for staffing large call centers. *Oper Res* 57(3):714–726
- Bassamboo A, Randhawa RS, Zeevi A (2010) Capacity sizing under parameter uncertainty: safety staffing principles revisited. *Manag Sci* 56(10):1668–1686

- Dong J, Ibrahim R (2017) Flexible workers or full-time employees? On staffing service systems with a blended workforce. Working paper, Columbia University
- Gans N, Koole G, Mandelbaum A (2003) Telephone call centers: tutorial, review, and research prospects. *Manuf Serv Oper Manag* 5:79–141
- Gans N, Shen H, Zhou YP, Korolev N, McCord A, Ristock H (2015) Parametric forecasting and stochastic programming models for call-center workforce scheduling. *Manuf Serv Oper Manag* 17(4):571–588
- Garnett O, Mandelbaum A, Reiman M (2002) Designing a call center with impatient customers. *Manuf Serv Oper Manag* 4(3):208–227
- Gurvich I, Lariviere M, Moreno-Garcia T (2017) Operations in the on-demand economy: staffing services with self-scheduling capacity. Working paper, Northwestern University
- Harrison JM, Zeevi A (2005) A method for staffing large call centers based on stochastic fluid models. *Manuf Serv Oper Manag* 7(1):20–36
- Ibrahim R (2017a) Managing queueing systems where capacity is random and customers are impatient. *Prod Oper Manag* 27(2):234–250
- Ibrahim R (2017b) Sharing delay information in service systems: a literature survey. *Queue Syst.* <https://doi.org/10.1007/s11134-018-9577-y>
- Ibrahim R, L'Ecuyer P (2013) Forecasting call center arrivals: fixed-effects, mixed-effects, and bivariate models. *Manuf Serv Oper Manag* 15(1):72–85
- Kang W, Ramanan K et al (2010) Fluid limits of many-server queues with reneging. *Ann Appl Probab* 20(6):2204–2260
- Mandelbaum A, Zeltyn S (2013) Data-stories about (im)patient customers in tele-queues. *Queue Syst* 75(2–4):115–146
- Shen H, Huang JZ (2008) Interday forecasting and intraday updating of call center arrivals. *Manuf Serv Oper Manag* 10(3):391–410
- Vocalcom (2014) Virtual call center industry projected to more than quadruple – are you ready? <http://www.vocalcom.com/en/blog/customer-service/virtual-call-center-industry-projected-to-more-than-quadruple-are-you-ready/>
- Whitt W (2006a) Fluid models for multiserver queues with abandonments. *Oper Res* 54(1):37–54
- Whitt W (2006b) Staffing a call center with uncertain arrival rate and absenteeism. *Prod Oper Manag* 15:88–102
- Zeltyn S, Mandelbaum A (2005) Call centers with impatient customers: many-server asymptotics of the $M/M/n + G$ queue. *Queue Syst Theory Appl* 51(3–4):361–402
- Zhang J (2013) Fluid models of many-server queues with abandonment. *Queue Syst* 73(2):147–193

Part III
Crowdsourcing Management

Chapter 14

Online Group Buying and Crowdfunding: Two Cases of All-or-Nothing Mechanisms



Ming Hu, Mengze Shi, and Jiahua Wu

Abstract This chapter focuses on the two popular business models, namely, online group buying and crowdfunding. Both models use variations of all-or-nothing mechanisms, where transactions will take place only if the total number of committed purchases/pledges exceeds a specified threshold within a certain period. We seek to understand the impact of all-or-nothing mechanisms on consumer behavior, as well as the optimal design of such mechanisms, from the perspective of third-party platforms like Groupon and Kickstarter. First, using a dataset from the online group buying industry, we empirically identify two types of threshold-induced effects on consumer behavior. Next, we study optimal information disclosure and pricing strategies under all-or-nothing mechanisms. We show that it is always beneficial for the firm to disclose the cumulative number of sign-ups to reduce the uncertainty for later arrivals. Regarding pricing, we show that the introduction of a price menu for the same product can be a win-win for both the creator and buyers.

14.1 Introduction

The motivation of this study stems from two popular business models, namely, online group buying and crowdfunding. Online group buying is a scheme designed to help coordinate a group of interested buyers so that they can reach their common purchase goals. In a typical group-buying deal, no transaction will take place unless the total number of committed purchases exceeds a specified threshold within a certain period. Online group-buying websites first appeared in the late 1990s, as part of the wave of innovative online market-based mechanisms. Usually, the consumers had to make the purchase commitment through escrow payment systems. Most of

M. Hu (✉) · M. Shi

Rotman School of Management, University of Toronto, Toronto, ON, Canada
e-mail: ming.hu@rotman.utoronto.ca; mshi@rotman.utoronto.ca

J. Wu

Imperial College Business School, Imperial College London, London, UK
e-mail: j.wu@imperial.ac.uk

© Springer Nature Switzerland AG 2019

M. Hu (ed.), *Sharing Economy*, Springer Series in Supply Chain Management 6,
https://doi.org/10.1007/978-3-030-01863-4_14

319

the representative group-buying websites that became popular in the late 1990s, including Mercata, Mobshop, and Letsbuyit, either ceased operating or changed their business models a few years later (Kauffman and Wang 2002). Interestingly, despite the failure of these pioneering group-buying sites, a decade later another generation of social buying websites like Groupon and LivingSocial emerged. Led by the market leader, Groupon, these newcomers started their business with offering “a deal a day” tailored to each local market.

Crowdfunding allows creators to raise funds from potential buyers to start their ventures and in return, the creators offer products to the buyers.¹ The creators can be designers, musicians, software developers, or any kind of inventors. Unlike in the conventional retail setting, here a buyer not only commits to purchasing the product but also prepays to fund the project. A project will be successfully funded only if the total value of committed purchases exceeds a specified goal within a certain time. The crowdfunding industry has experienced tremendous growth in recent years. For instance, Kickstarter is one of the leading online non-equity crowdfunding sites that match people (the “crowd”) with projects. It has raised more than \$800 million and supported more than 50,000 projects in the four years since its inception in 2009. With the passing in the US of the Jumpstart Our Business Startups Act in September 2013, the crowdfunding industry acquired legitimacy and is expected to lead a new era of entrepreneurship.

What these two business models share in common is the underlying mechanism, namely, the all-or-nothing mechanism, where transactions will take place only if the total number of committed purchases/pledges exceeds a specified threshold within a certain period. There are many reasons why a project needs a certain number of buyers and a certain amount of funds to start. In the context of crowdfunding, for instance, there may be *economies of scale* due to high initial setup costs on the supply side. Most digital products fall into this category. On the demand side, the product may exhibit *positive externality* and it requires enough users for the product to be valuable enough. In this chapter, we seek to have a thorough understanding of the impact of all-or-nothing mechanisms on consumer behavior, as well as their optimal design, from the perspective of third-party platforms like Groupon and Kickstarter.

In the context of online group buying, we first empirically investigate the effect of thresholds on consumer’s sign-up behavior. Our investigation utilizes a dataset collected from Groupon.com, during a period when the company predominantly used “a deal a day” format for each local market. For each Groupon deal, deal characteristics, threshold level, and *real-time* updated number of sign-ups were posted on the website. These data provide us with an opportunity to infer the effects of thresholds from the sign-up patterns over time. Our study finds two types of threshold-induced behavior. The first type of threshold effects refers to a *substantial increase* in the number of sign-ups around the time when the threshold is reached.

¹On equity crowdfunding sites, funders are also investors for the creators’ financial endeavors. This chapter focuses on non-equity crowdfunding where buyers are also funders but not investors.

The second threshold effect refers to a *stronger* positive relation between the number of new sign-ups and the cumulative number of sign-ups *before* the thresholds are met than afterward. We discuss several mechanisms compatible with the findings that may serve as potential hypotheses for future research.

After establishing the impact of all-or-nothing mechanisms on consumer behavior, we will study the optimal design of all-or-nothing mechanisms from the perspective of third-party platforms like Groupon and Kickstarter. Under all-or-nothing mechanisms, consumers are linked together by the common goal of project success. Thus, it is crucial for the firm to facilitate the coordination among consumers to improve the chance of project success.

We first investigate the information disclosure strategy, specifically, whether or not the sponsor should ask participants to make decisions without knowing the choices of others. We take the perspective of group-buying platforms like Groupon and investigate the impact of alternative information disclosure strategies on deal success rates. Knowing how to improve the success rates is vital because group-buying firms typically earn revenues from successful listings only and not all group-buying deals succeed. Without careful analysis, the firm's decision does not appear straightforward because of the uncertainty about the number of consumer arrivals and their individual valuations. Looking forward, one can find it beneficial to post the number of sign-ups if a large cohort of consumers with high individual valuations turn out in the early stage, but it can be detrimental if the first cohort of consumers turns out to be small and have low individual valuations.

To investigate the influence of information disclosure strategies on deal success rates, we develop a two-period model where two cohorts of consumers arrive at the deal sequentially. The two-period model is a stylized capture of the fact that earlier arrivals are faced with more uncertainty in the deal's success rate than later arrivals. The firm being studied chooses between a "sequential mechanism" where the firm posts the number of sign-ups at the end of the first period, and a "simultaneous mechanism" where the firm does not post the first-period outcome. Somewhat surprisingly, our analysis shows that the deal's success rate is always higher under the sequential mechanism. To understand this result requires backward-inductive reasoning, starting with the second period and then moving back to the first period. A sequential mechanism increases the ex-ante expected sign-up rates of the second cohort of consumers by eliminating the uncertainty facing them. The increased expected sign-up rates of the second cohort enhance the confidence of the first cohort of consumers, thereby increasing the ex-ante expected sign-up rates of the first cohort. This result underscores the importance of modeling and investigating the dynamics of sign-up behavior under online group buying. The result also offers a potential explanation for why firms like Groupon and Kickstarter display the updated number of sign-ups along with the minimum number required to unlock the deals.

Last, we study various pricing strategies in facilitating coordination among consumers in the context of crowdfunding. In crowdfunding, a creator may design a menu of price options when buyers have heterogeneous valuations. A crucial consideration in the menu design is incentive compatibility; that is, each type of

consumers should be better off by choosing the option designed for them than by choosing any other options. In addition, a creator needs to consider the anticipated project success rate. In the design of the price menu for a crowdfunding project, how would the creator's consideration of project success affect the optimal pricing decisions?

To investigate this question, we propose and analyze a two-period model where a creator may charge different prices over time to sequentially arriving buyers. The decisions of the earlier buyers are posted to the later arrivals. The buyers decide whether or not to purchase and which price option to choose according to their valuations. It turns out that, given the same product but at different prices, the buyer with a high product valuation may choose the high-price option in crowdfunding to ensure the success of the project. This will be the case as long as a buyer perceives that other buyers may have low product valuations. Interestingly, this "over-pay" behavior can also improve total buyer surplus, and hence the introduction of discriminatory pricing strategies, such as a menu, can be win-win for both the creator and buyers. This result is in stark contrast to the traditional situation where consumers would generally prefer the low-price option given *no* difference in quality.

14.2 Consumer Behavior Under All-or-Nothing Mechanisms

We first empirically investigate the effect of all-or-nothing mechanisms on consumer's sign-up behavior utilizing a dataset collected from Groupon.com. The online group-buying industry has witnessed phenomenal growth since the début of Groupon in 2008. Group-buying firms are third-party intermediaries that facilitate the coordination among a large group of consumers. Such coordination permits consumers to enjoy the quantity discounts offered by the sellers collectively. Groupon, since its début in 2008, increased its total number of subscribers to over 200 million as of March 2013. Groupon extended its coverage to more than 500 markets in 48 countries, up from just 28 U.S. markets in 2009.

We hired a research assistant at a major university to build a data crawler on the Google App Engine platform. The data crawler extracted deal information, such as deal description, deal price, discount level, and threshold, whenever a new deal was posted. The program updated the cumulative number of sign-ups with the interval of every five minutes. We use this real-time dataset to keep track of consumers' responses to various group buying deals during the lifetime of each deal and to uncover the patterns of sign-up accumulation.

Our data includes a total of 4,208 deals from 86 cities or regions covered by Groupon between September 28th, 2010 and December 07th, 2010. The duration of the observation period was 71 days in total. For each deal, we recorded a set of deal attributes and monitored the inter-temporal sign-up process. Table 14.1 presents the summary statistics for all these 4,208 deals. The average deal price in the sample was \$30.68, with an average discount level of 56% off. Each of these deals contained

Table 14.1 Summary statistics of all deals

	Mean	Std. Dev.	Minimum	Maximum
Deal attributes				
Deal price (\$)	30.68	30.53	2	250
Discount level (%)	56.35	9.96	19	96
Threshold	55.40	68.50	3	800
Market population (thousands)	854.30	1,332.94	56	8,364
Outcome				
Total amount purchased	784.80	1,331.57	5	29,380

Note: *Deal Price* denotes the net price a consumer needed to pay if the deal tips. *Discount Level* denotes the markdown of deal price relative to the regular price. *Threshold* denotes the minimum number of committed purchases for the deal to succeed. *Market Population* is the population of the local market where the deal was posted. *Total Amount Purchased* denotes the number of consumers who purchased the product or service by the end of the sign-up process

a threshold of sign-up numbers for the deal to succeed. A group-buying deal would be off if the total number of committed consumers did not reach the threshold. The average threshold value specified by Groupon was around 55. The average number of coupons purchased for each deal was around 785. In our sample, all deals reached the thresholds before expiration.

During our data-collection period, deals were posted daily for 24 h from Mondays to Thursdays. However, the duration of deals posted on Fridays and weekends could vary from 24 to 72 h. In some relatively small markets, Groupon would post 72-h deals on Friday. We also saw a transition from 72-h deals to 48-h deals during our data-collection period.

14.2.1 Empirical Model

Our primary objective is to investigate the effects of thresholds on the rate of signing up to the deals. However, many confounding factors, such as varying online traffic to the websites at different times of the day and unobserved deal heterogeneity, may contribute to the sign-up pattern. We seek to establish the threshold effects through formal statistical analyses rigorously.

14.2.1.1 The Base Model

We start our analysis with a flexible model specification to distill the sign-up pattern around the thresholds. We include a series of time dummy variables, with each variable capturing the sign-up pattern during a 5-min time interval around the period when the threshold was reached. The dependent variable in our model, denoted

by $y_{i,t}$, is the number of new sign-ups during the t th time interval for deal i . To control the unobserved deal heterogeneity, we apply a deal-fixed effect model with the following specification:

$$y_{i,t} = \sum_{j=-T}^T \alpha_j I_{\{s_{i,t}=j\}} + \psi_t + \mu_i + \epsilon_{i,t}, \quad (14.1)$$

where t represents the time index before re-aligning the deals, and $s_{i,t}$ indicates the time index after re-aligning the deals at the period when the threshold was reached. Recall that with the re-aligned data, time 0 is the period when a deal reaches its threshold. Consequently, $s_{i,t}$ is equal to 0 if deal i reaches its threshold at time period t . Similarly, $s_{i,t} = j$ for all $j > 0$ (resp., $j < 0$) represents that time period t is the j th period after (resp., before) deal i reaches its threshold, and $I_{\{s_{i,t}=j\}}$ for all j is a dummy variable which is equal to 1 if $s_{i,t} = j$, and zero otherwise. The set of time dummy variables, $I_{\{s_{i,t}=j\}}$, $j = -T, \dots, T$, is used to capture the sign-up pattern around the time when the thresholds are reached, where T reflects the width of the time window. In addition, ψ_t measures the time-of-the-day fixed effect using a 5-min time dummy, and μ_i measures the deal fixed effects. The term $\epsilon_{i,t}$ is the error component.

To estimate the base model, we can apply standard approaches for estimating fixed-effects panel models. The fixed effects can be eliminated by either taking differences between adjacent observations from the same deal or subtracting the average over time from every variable, i.e., time-demeaning. Then, we can apply the generalized least squared (GLS) estimator to the transformed data.

14.2.1.2 The Extended Model with Lagged Variables

Though the base model allows us to capture the sign-up pattern around thresholds in a flexible way, it does not reflect the dependency of the new sign-ups on the cumulative number of sign-ups. Consequently, we extend the base model by introducing the lagged cumulative sign-ups, $Y_{i,t-1}$, into the model. The extended model can be formulated as follows:

$$y_{i,t} = \sum_{j=-T}^T \alpha_j I_{\{s_{i,t}=j\}} + \sum_{j=-T}^T \beta_j I_{\{s_{i,t}=j\}} Y_{i,t-1} + \psi_t + \mu_i + \epsilon_{i,t}. \quad (14.2)$$

Similar to the base model, we include the interactions between time dummies after re-alignment, $I_{\{s_{i,t}=j\}}$, and the lagged cumulative number of sign-ups, $Y_{i,t-1}$, to capture the relationship between the new sign-ups and the cumulative number of sign-ups in a flexible way.

However, unlike the base model which can be estimated consistently using GLS estimator, the estimation of a fixed effects model with lagged variables is more technically involved. The lagged regressor is likely to be correlated with the fixed

effects, which gives rise to “dynamic panel bias” (Nickell 1981). To solve this problem, we apply the generalized method of moments (GMM) approach proposed by Arellano and Bond (1991). First, we take a difference of Eq. (14.2) to eliminate the deal fixed effects:

$$\begin{aligned}
 y_{i,t} - y_{i,t-1} &= \sum_{j=-T}^T \alpha_j (I_{\{s_{i,t}=j\}} - I_{\{s_{i,t-1}=j\}}) \\
 &+ \sum_{j=-T}^T \beta_j (I_{\{s_{i,t}=j\}} Y_{i,t-1} - I_{\{s_{i,t-1}=j\}} Y_{i,t-2}) \\
 &+ (\psi_t - \psi_{t-1}) + (\epsilon_{i,t} - \epsilon_{i,t-1}). \tag{14.3}
 \end{aligned}$$

As $Y_{i,t-1}$ is correlated with the error term, specifically $\epsilon_{i,t-1}$, GLS yields inconsistent estimates after the first-difference transformation. However, if there is no serial correlation in the error term $\epsilon_{i,t}$, then the longer lags of the regressors, i.e., $Y_{i,k}$, $k = t - 2, \dots, 1$, which are correlated with $y_{i,t-1}$ (see Eq. (14.2)), and thus $Y_{i,t-1}$, but not with the error term $\epsilon_{i,t-1}$, can serve as instruments for the model after the first-difference transformation. In the case of our model, $Y_{i,t-2}$, and $I_{\{s_{i,t-1}=j\}} Y_{i,t-2}$, together with their longer lags can serve as GMM instruments for Eq. (14.3). The differences of the strictly exogenous variables, i.e., $I_{\{s_{i,t-1}=j\}}$ and ψ_t , can serve as standard instruments.

We capture the unobserved heterogeneity across deals with deal fixed effects. The observed deal variations as described by product/service categories, deal prices and discounts, and city characteristics are unlikely to capture all sources of deal heterogeneities. For example, restaurants in a city can have different locations, offer different cuisines, and enjoy different reputations. For Eq. (14.2) to identify threshold effects on group-buying deals, we implicitly assume that unobserved deal attributes are accounted for by a time-invariant component, i.e., μ_i . This fixed component controls for unobserved deal attributes, which may positively correlate with both lagged cumulative number of sign-ups, $Y_{i,t-1}$, and the number of new sign-ups, $y_{i,t}$, and thus solves an “errors in variables” type of endogeneity problem (Villas-Boas and Winer 1999). Given the panel data structure, we can use deal-specific fixed effects to control the variations across deals.

14.2.2 Results

We next present the empirical results for the base model and the extended model. Our analysis uses a one-and-a-half-hour time window before and after the threshold was reached for regression analysis, i.e., $T = 17$. The usage of a relatively small time window around the time when thresholds were reached eliminates other unrelated factors and allows us to focus on the effects of thresholds on consumers’ sign-up behavior.

Since the earliest observation serves as the reference level, we have a total of 34 5-min time dummy variables, i.e., $I_{\{s_{i,t}=j\}}$, $j = -16, \dots, 17$, to capture the sign-up pattern around thresholds. To focus on the main findings, we do not interact every single time dummy variable with the lagged cumulative number of sign-ups. Instead, we divide the three-hour time window into four non-overlapping time periods of equal length, create four new time dummies, with each representing a 45-min time interval, and interact the newly created time dummies with the lagged cumulative number of sign-ups. The interaction terms between these four time dummies and the lagged cumulative number of sign-ups are sufficient to capture the level shift in the ratio between new sign-ups and the cumulative number of sign-ups when the thresholds were passed, as well as the trend of the ratio both before and after reaching the thresholds.

Table 14.2 shows the regression result of our base model. The coefficient of the time dummy when the thresholds were reached ($\hat{\alpha}_0 = 3.582$, $p < 0.01$) is significantly greater than the coefficients of other time dummies, which indicates a clear spike in the number of sign-ups during the time interval when the thresholds were reached after we control for heterogeneous time traffic and unobserved deal heterogeneity. This establishes the first-type of threshold effects: a surge in the number of new sign-ups around the time when the threshold was reached.

Table 14.3 shows the regression results of our extended model. The result from the GLS estimator is presented as a benchmark. We show the results from the

Table 14.2 Regression results of the base model using three-hour data

Estimates	Estimates	Estimates	Estimates	Estimates	Estimates
α_{-16} -0.094** (0.037)	α_{-10} -0.437*** (0.063)	α_{-4} -0.076 (0.101)	α_1 0.131 (0.185)	α_7 0.355** (0.178)	α_{13} 0.953*** (0.212)
α_{-15} -0.152*** (0.041)	α_{-9} -0.447*** (0.067)	α_{-3} -0.266** (0.106)	α_2 -0.085 (0.152)	α_8 0.623*** (0.193)	α_{14} 1.140*** (0.232)
α_{-14} -0.257*** (0.043)	α_{-8} -0.385*** (0.074)	α_{-2} -0.198* (0.112)	α_3 0.104 (0.150)	α_9 0.642*** (0.187)	α_{15} 1.277*** (0.237)
α_{-13} -0.288*** (0.051)	α_{-7} -0.396*** (0.080)	α_{-1} 0.023 (0.130)	α_4 0.400** (0.161)	α_{10} 0.634*** (0.198)	α_{16} 1.403*** (0.252)
α_{-12} -0.351*** (0.052)	α_{-6} -0.394*** (0.087)	α_0 3.582*** (0.315)	α_5 0.223 (0.156)	α_{11} 0.881*** (0.216)	α_{17} 1.477*** (0.296)
α_{-11} -0.357*** (0.058)	α_{-5} -0.347*** (0.094)		α_6 0.211 (0.181)	α_{12} 0.942*** (0.212)	
Time-of-the-day fixed effects			Yes		
Deal fixed effects			Yes		
Number of observations			147,112		
Number of deals			4,208		
Adjusted R-squared			0.026		

Note: The dependent variable is the number of new sign-ups per 5-min time interval. Standard errors are clustered by deal and reported in parentheses

* $p < 0.1$; ** $p < 0.05$; *** $p < 0.01$

Table 14.3 Regression results of the extended model using three-hour data

	GLS	GMM (2 lags)	GMM (3 lags)	GMM (4 lags)	α_1	GLS	GMM (2 lags)	GMM (3 lags)	GMM (4 lags)
α_{-16}	-0.069* (0.037)	-0.115*** (0.044)	-0.123*** (0.044)	-0.126*** (0.044)		-0.172 (0.174)	-1.044* (0.536)	-1.029* (0.550)	-0.920 (0.609)
α_{-15}	-0.101** (0.040)	-0.196*** (0.063)	-0.213*** (0.063)	-0.220*** (0.064)	α_2	-0.473*** (0.188)	-1.403** (0.610)	-1.377** (0.624)	-1.265* (0.681)
α_{-14}	-0.183*** (0.042)	-0.329*** (0.083)	-0.354*** (0.083)	-0.365*** (0.084)	α_3	-0.366* (0.204)	-1.346** (0.619)	-1.309** (0.633)	-1.193* (0.698)
α_{-13}	-0.195*** (0.049)	-0.391*** (0.104)	-0.424*** (0.103)	-0.438*** (0.105)	α_4	-0.166 (0.212)	-1.202* (0.632)	-1.153* (0.647)	-1.034 (0.710)
α_{-12}	-0.245*** (0.053)	-0.489*** (0.123)	-0.528*** (0.121)	-0.546*** (0.123)	α_5	-0.459** (0.207)	-1.555** (0.633)	-1.494** (0.645)	-1.371* (0.710)
α_{-11}	-0.242*** (0.061)	-0.537*** (0.142)	-0.581*** (0.139)	-0.603*** (0.141)	α_6	-0.583*** (0.213)	-1.735*** (0.649)	-1.659** (0.651)	-1.533** (0.711)
α_{-10}	-0.319*** (0.069)	-0.664*** (0.158)	-0.713*** (0.155)	-0.739*** (0.158)	α_7	-0.555** (0.226)	-1.761*** (0.649)	-1.669** (0.651)	-1.539** (0.717)
α_{-9}	-0.328*** (0.078)	-0.725*** (0.173)	-0.778*** (0.169)	-0.808*** (0.172)	α_8	-0.414* (0.229)	-1.677** (0.652)	-1.570** (0.649)	-1.437** (0.710)
α_{-8}	-0.468*** (0.128)	-0.840** (0.335)	-0.832** (0.347)	-0.802** (0.368)	α_9	-0.540** (0.239)	-1.864*** (0.665)	-1.739*** (0.660)	-1.602** (0.721)
α_{-7}	-0.510*** (0.131)	-0.931*** (0.349)	-0.919** (0.361)	-0.890** (0.383)	α_{10}	-0.309 (0.257)	-1.004 (0.685)	-0.941 (0.747)	-0.669 (0.852)
α_{-6}	-0.544*** (0.127)	-1.017*** (0.367)	-1.003*** (0.376)	-0.974** (0.399)	α_{11}	-0.196 (0.266)	-0.920 (0.687)	-0.843 (0.750)	-0.561 (0.860)
α_{-5}	-0.544*** (0.125)	-1.066*** (0.375)	-1.046*** (0.385)	-1.018** (0.407)	α_{12}	-0.284 (0.257)	-1.032 (0.689)	-0.938 (0.740)	-0.647 (0.852)

(continued)

Table 14.3 (continued)

	GLS	GMM (2 lags)	GMM (3 lags)	GMM (4 lags)	α	GLS	GMM (2 lags)	GMM (3 lags)	GMM (4 lags)
α_{-4}	-0.327*** (0.121)	-0.898** (0.382)	-0.872** (0.390)	-0.845** (0.415)	α_{13}	-0.427* (0.238)	-1.198* (0.672)	-1.085 (0.725)	-0.785 (0.835)
α_{-3}	-0.594*** (0.120)	-1.220*** (0.387)	-1.184*** (0.392)	-1.158*** (0.417)	α_{14}	-0.400 (0.251)	-1.198* (0.679)	-1.066 (0.726)	-0.757 (0.840)
α_{-2}	-0.600*** (0.116)	-1.276*** (0.388)	-1.231*** (0.390)	-1.206*** (0.415)	α_{15}	-0.433* (0.240)	-1.255* (0.671)	-1.103 (0.707)	-0.786 (0.819)
α_{-1}	-0.463*** (0.113)	-1.195*** (0.384)	-1.139*** (0.383)	-1.116*** (0.405)	α_{16}	-0.486* (0.250)	-1.339** (0.677)	-1.168* (0.708)	-0.841 (0.818)
α_0	2.992*** (0.272)	2.204*** (0.506)	2.273*** (0.478)	2.294*** (0.466)	α_{17}	-0.601*** (0.229)	-1.484** (0.728)	-1.294* (0.745)	-0.958 (0.833)
Lag cumulative sign-ups									
Lag cumulative sign-ups *									
within 45 min before threshold dummy (β_1)									
Lag cumulative sign-ups *									
within 45 min after threshold dummy (β_2)									
Lag cumulative sign-ups *									
from 45 to 90 min after threshold dummy (β_3)									
Time-of-the-day fixed effects									
Deal fixed effects									
Number of observations									
Number of deals									
Adjusted R-squared									
						Yes	Yes	Yes	Yes
						147,112	142,904	142,904	142,904
						4,208	4,208	4,208	4,208
						0.389			

Note: The dependent variable is the number of new sign-ups per 5-min time interval. Standard errors are clustered by deal and reported in parentheses

* $p < 0.1$; ** $p < 0.05$; *** $p < 0.01$

Table 14.4 Statistical tests using estimates in Table 14.3

	GLS	GMM (2 lags)	GMM (3 lags)	GMM (4 lags)		GLS	GMM (2 lags)	GMM (3 lags)	GMM (4 lags)
$\alpha_0 - \alpha_{-5}$	3.535*** (0.330)	3.270*** (0.335)	3.319*** (0.335)	3.312*** (0.327)	$\alpha_0 - \alpha_1$	3.164*** (0.303)	3.248*** (0.348)	3.302*** (0.369)	3.215*** (0.388)
$\alpha_0 - \alpha_{-4}$	3.319*** (0.323)	3.102*** (0.325)	3.145*** (0.326)	3.139*** (0.320)	$\alpha_0 - \alpha_2$	3.465*** (0.330)	3.607*** (0.392)	3.651*** (0.417)	3.559*** (0.444)
$\alpha_0 - \alpha_{-3}$	3.586*** (0.312)	3.423*** (0.312)	3.457*** (0.314)	3.452*** (0.309)	$\alpha_0 - \alpha_3$	3.358*** (0.348)	3.550*** (0.418)	3.582*** (0.437)	3.487*** (0.466)
$\alpha_0 - \alpha_{-2}$	3.591*** (0.304)	3.480*** (0.303)	3.504*** (0.305)	3.500*** (0.301)	$\alpha_0 - \alpha_4$	3.158*** (0.337)	3.406*** (0.416)	3.427*** (0.437)	3.328*** (0.463)
$\alpha_0 - \alpha_{-1}$	3.455*** (0.296)	3.399*** (0.295)	3.412*** (0.296)	3.410*** (0.295)	$\alpha_0 - \alpha_5$	3.451*** (0.328)	3.759*** (0.412)	3.767*** (0.431)	3.665*** (0.461)
$\beta_1 - \beta_2$	0.010*** (0.003)	0.011* (0.006)	0.009 (0.006)	0.010* (0.007)	$\beta_2 - \beta_3$	0.004 (0.003)	0.011*** (0.003)	0.010*** (0.003)	0.011*** (0.003)

Note: Standard errors are clustered by deal and reported in parentheses. Significance levels are related to the null hypothesis H_0 : combination of coefficients = 0

* $p < 0.1$; ** $p < 0.05$; *** $p < 0.01$

GMM estimator using two, three and four lags of GMM instrumental variables. Specifically, GMM instruments include the lagged cumulative number of sign-ups, and the lagged interaction terms between the 45-min time dummies and the cumulative number of sign-ups. In theory, we can use all valid lagged regressors, i.e., those with lags of two and more. However, the number of instruments would be quadratic in the time dimension of the panel, and the GMM estimator may perform poorly with a large number of instruments (Roodman 2009). Too many instruments may overfit endogenous variables, bias coefficient estimates, and thus the results from the finite sample may be far from the asymptotic ideal. In our analysis, we apply the GMM estimator using two, three and four lags of instrumental variables. As shown in columns 2, 3 and 4 in Table 14.3, the results are robust with respect to the number of lags used.

After controlling for the lagged cumulative number of sign-ups, the surge in the number of sign-ups during the interval when the threshold was reached remains. We further verify this finding by conducting Wald tests on the differences between the estimated coefficients of the time dummy when the threshold was reached, and other time dummies within the half-an-hour time window either before or after the threshold was reached. As shown in the first part of Table 14.4, the differences between the coefficients are all statistically significant. To quantify this effect, we define a *spike index* as the measure for the additional number of sign-ups due to the thresholds. Specifically, the *spike index* is equal to $\sum_{j=-5}^5 (\hat{\alpha}_0 - \hat{\alpha}_j) / 10$. Using the results from the GMM estimator with 2 lags, the spike index across all deals is 3.424. That is, on average, around 3.424 more consumers would sign up to the deal during the 5-min time interval when the threshold was reached (statistically significant, $p < 0.01$). To assess the magnitude of this effect, note that during the

half-an-hour time window before and after the threshold was reached, on average there were approximately 3.6 consumers signing up ever 5 min. Thus, the existence of thresholds produces a substantial boost in sales.

We also observe a level shift in the ratio between the new sign-ups and the cumulative number of sign-ups. Specifically, our test results show that the ratio in the 45-min time window before the thresholds is higher than that in the 45-min time window after the thresholds by around 0.01, which is consistent among the results from the GLS estimator and GMM estimator with various lags (see test results of $\beta_1 - \beta_2$ in Table 14.4). The difference in the ratio before and after the thresholds is significant from the GLS estimator and GMM estimator using either 2 lags or 4 lags. However, the difference is insignificant from the GMM estimator using 3 lags. The underlying reason might be the weak instruments for 3 and 4 lags to be discussed below. The estimated coefficients from the GMM estimator also indicate that the ratio between new sign-ups and the cumulative number of sign-ups continues decreasing after reaching the thresholds (see test results on $\beta_2 - \beta_3$ in Table 14.4). After the thresholds, the ratio in the second 45-min period is significantly lower than that in the first 45-min time window.

It is worth noting that the validity of the GMM estimator will be violated if the error component $\epsilon_{i,t}$ is serially correlated over time. To address this concern, we apply post-estimation tools of the GMM estimator and examine the serial correlation structure of the new error component $\epsilon_{i,t} - \epsilon_{i,t-1}$. The second-order serial correlation is -0.662 ($p = 0.508$), suggesting that the error components in Eq. (14.2), i.e., $\epsilon_{i,t}$, are indeed uncorrelated over time.

Another potential concern with the GMM estimator is weak instruments. When the correlation between instrumental variables and the endogenous variable is low, the asymptotic distribution of the coefficients breaks down, and GMM estimates may not be consistent (Bound et al. 1995). In this case, the standard errors on GMM estimates are likely to be larger than those on GLS estimates. For our model, the concern of weak instruments may become important if the lagged cumulative number of sign-ups is not informative in predicting the new sign-ups.

To test the existence of weak instruments, we regress the endogenous variable after the first-difference transformation, i.e., $y_{i,t}$, on various lags of the cumulative number of sign-ups. The regression analysis is conducted using 2 lags to 4 lags, and the F-statistics are summarized in Table 14.5. We refer to the rule of thumb suggested by Staiger and Stock (1997) that the finite-sample bias of instrumental variables would not be a serious problem when the F-statistic is greater than 10.

Table 14.5 F-statistics for the instrumental variable regressions

	2 lags	3 lags	4 lags
F-statistic of $Y_{i,t-1}$	18.03	16.32	12.56
F-statistic of $Y_{i,t-2}$	11.37	0.40	0.36
F-statistic of $Y_{i,t-3}$		13.21	0.24
F-statistic of $Y_{i,t-4}$			7.23

When we use only two lags, the F-statistics on both lags is greater than the cut-off value of 10. However, the results with more than 2 lags show that the correlation between the number of new sign-ups and some lagged cumulative number of sign-ups is low. Since the incremental number of sign-ups within a short period is likely to be highly correlated, utilizing a more lagged cumulative number of sign-ups may not increase the power in predicting the number of new sign-ups, rendering some lagged variables as weak instruments.

14.2.3 Potential Mechanisms Behind Threshold Effects

Our analysis has documented and substantiated two types of threshold effects in online group buying. However, the aggregate nature of the data prevents us from identifying the exact mechanisms contributing to these effects. In this section, we discuss several mechanisms compatible with the findings that may serve as potential hypotheses for future research.

For the first type of threshold effects, i.e., the sudden surge of sign-ups around the time when thresholds were reached, we consider three possible mechanisms: value enhancement, postponed decision making, and higher consumer awareness. First, a consumer may derive positive psychological value from beating a target. When the cumulative number of sign-ups approaches a threshold, consumers may experience an urge to beat the target and sign up in a “frenzy” fashion. Such “frenzy”, which is similar with “bidding frenzy” phenomenon widely observed towards the end of auctions, may reflect a mental state “characterized by a high level of excitement, a strong sense of competition, and an intense desire to win” (see, e.g., Ku et al. 2005; Heyman et al. 2004; Häubl and Popkowski Leszczyc 2004). Second, some consumers may postpone sign-up decisions until the deals are about to be on. When the number of sign-ups is small, consumers face the uncertainty on deal success and the risk of not receiving the discount at the end. Thus, some consumers may choose to postpone their sign-up decisions if the cost to track the sign-up numbers is sufficiently low. Such postponement of action could lead to a surge in the number of sign-ups around the time when the thresholds were reached. Third, the number of sign-ups may surge around the thresholds because of increased consumer awareness to the deals generated by a firm’s communication strategy. When a deal comes close to its threshold, a group-buying firm may feature the deal on its front page, highlight the deal in its email to the subscribers, or coordinate with third-party deal aggregators to enhance the placement of the deal on their websites.

For the second type of threshold effects, i.e., the level shift of the ratio between the number of new sign-ups and the cumulative number of sign-ups before and after the thresholds were reached, we discuss four alternative mechanisms: word-of-mouth referral, observational learning, consumer heterogeneity, and demand satiation. First, people who have already signed up to deals play an active role in disseminating deal information. In our model, the positive relationship between the number of new sign-ups and the cumulative number of sign-ups may capture

the intensity of such referrals (e.g., Bass et al. 1994). Following this logic, our regression results would suggest a stronger intensity of referrals before the thresholds were reached. This result is consistent with the view of group buying business as a marketing tool to exploit social interactions between consumers. For instance, Jing and Xie (2011) show that some informed consumers can be motivated to persuade their social contacts to join the group-buying deals to ensure deal success.

Second, herding behavior or observational learning may explain the positive relationship between the number of new sign-ups and the cumulative number of sign-ups, as demonstrated in Zhang and Liu (2012). In the context of group buying, some consumers can be uncertain about the actual quality of the suppliers and hence the value of the deals. These consumers may infer the quality of suppliers from the number of consumers who have already signed up to the deal. With this rationale, the level shift of β_j indicates that, when making their own purchase decisions, individual consumers would be more likely to resort to the choices of others before the thresholds were reached and when the deals were uncertain.

Third, different types of consumers may arrive at the deal site before and after thresholds are reached. For instance, tech-savvy consumers may learn about the deals and sign up earlier. These consumers can be more capable of engaging in referrals through online social networks. The early arrivals may also perceive the referrals to be more valuable because their contacts are unlikely to be aware of the deal. As a result, word-of-mouth referrals should be stronger before the thresholds were reached. Finally, there might exist satiation in both the market demand and the reach of word-of-mouth referrals. If the satiation levels happen to be reached around the same time when thresholds were achieved, then we would observe more sign-ups per 5-min interval before the thresholds than afterward.

The mechanisms discussed above have different implications on the economic impact of threshold effects. Based on the GMM estimator with 2 lags, on average there were 3.4 additional sign-ups during the periods when thresholds were reached. Moreover, the relation between the new sign-ups and the cumulative number of sign-ups experienced a drop right after reaching the thresholds and continued the decreasing trend afterward. To accurately quantify the economic implications of threshold effects, it is necessary to account for specific underlying mechanisms. For example, for the first-type of threshold effects, one needs to know the extent of inter-temporal substitutions, i.e., whether some consumers might sign up sooner or later if thresholds did not exist. Similarly, for the second-type of threshold effects, one needs to confirm if word-of-mouth referral was the driving force. As data limitation prevents us from identifying such specific mechanisms, we leave to future research to investigate the economic implications of threshold effects. For both types of threshold effects, the magnitude of the effect is likely to be greater when the success rate is lower, or the stake is higher. Groupon's daily deals studied in this paper were generally expected to be on. The success rate can be much lower in other markets with all-or-nothing mechanisms. For instance, the success rate was 43% at Kickstarter (blog.kickstarter.com) and 12.3% on Prosper.com (Zhang and Liu 2012). Economic implications of threshold effects are thus expected to be higher in these markets.

14.3 Coordination Under All-or-Nothing Mechanisms

Having established the impact of all-or-nothing mechanisms on consumers' behavior, in this section, we study how the firm could improve the chances of project success by facilitating coordination among consumers. Knowing how to improve the success rates is important because firms typically earn revenues from successful listings only and not all listings succeed. For example, from the launch of Kickstarter.com in April 2009 to March 2011, a total of 20,371 projects were proposed. Among them, 7,496 projects attracted enough funds, leading to a success rate of 43% (blog.kickstarter.com). We study the problem from two different perspectives, namely, information disclosure and pricing strategies.

14.3.1 Information Disclosure

We first investigate the optimal information disclosure policy under all-or-nothing mechanisms. We consider a market where a firm uses the group-buying format to promote a product or service to consumers. To illustrate the critical insights of the firm's design of the mechanism and subsequent consumer responses, we resort to a stylized model with two consumers and two periods. A two-person model is an efficient way to capture the sequential nature of consumer arrivals and the interdependence between purchase decisions of early and late arrivals. In practice, firms may use variations of all-or-nothing mechanisms depending on the specific products or services being promoted. For instance, Kickstarter sets the total amount of demanded dollar commitment as the threshold, and the project sign-up horizon typically can last for several weeks. We base our description of the model on group-buying websites like Groupon, where the individual decision is binary.

14.3.1.1 Model Setup

At the beginning of the first period, the firm posts its group-buying deal, which is characterized by three elements: group-buying price p , a minimum number of buyers N required, and a time horizon of two periods for signing up. In the two-person model, we assume that one consumer arrives in the first period, and the other in the second period. Each consumer demands and may purchase up to one unit of the product or service. We limit the threshold N to be 2. That is, the deal succeeds if and only if both consumers sign up for the deal, and they receive the product or service at a price of p . Otherwise, the deal is off and no transaction takes place. The individual product valuation for consumer i is denoted by V_i , which is drawn from a cumulative distribution function $F_i(\cdot)$.

A consumer decides to sign on to the group-buying deal if and only if she expects a discounted utility from the deal at least as high as that of not signing on to the deal. To sign on to the deal is a commitment to purchase if the deal is on. In making the decision, a consumer takes into consideration her own valuation for the deal as well as her expectation about the success rate of the group-buying deal. In the two-person case, one consumer uses the other consumer's sign-up likelihood as the belief in making sign-up decisions. Conditional on signing on to the deal by herself, the deal's success is a direct consequence of the other consumer's sign-up decision. We assume that consumers form rational expectations and focus on the pure strategy equilibria of the game. Specifically, we denote by $H_i(q_{-i})$ the likelihood of an individual consumer i 's signing on to a group-buying deal, where q_{-i} is consumer $-i$'s sign-up likelihood that consumer i expects. This notation emphasizes that the likelihood of signing up depends on the other consumer's sign-up likelihood, but suppresses its dependence on the characteristics of the deal. We assume the following relation between $H_i(q_{-i})$ and q_{-i} .

Assumption 1 (Sign-Up Likelihood) For all i and $q_{-i} \in [0, 1]$, we assume

1. $H_i(q_{-i} = 0) = 0$;
2. $H_i(q_{-i})$ is non-decreasing in q_{-i} .

Assumption 1(i) states that if a consumer expects with certainty that the other consumer will not sign up, then she expects zero benefits from signing on to the deal and will not sign up. Assumption 1(ii) indicates that the higher the likelihood that a consumer expects the other consumer to sign up, the higher the probability that the consumer will sign up.

If both consumers can adequately communicate with each other before making sign-up decisions, the first-best solution of the deal's success rate can be obtained at $H_1(1) \cdot H_2(1)$. However, a full information structure is uncommon in reality where typical consumers do not reveal their own private information to strangers in online group-buying settings. Instead, the firm considers two alternative group-buying mechanisms: a sequential mechanism and a simultaneous mechanism. The distinction between these two mechanisms is created by the firm's decision on whether to reveal information to the consumer who arrives in the second period; specifically, the second consumer can see the choice of the consumer who arrives in the first period under a sequential mechanism, but not under a simultaneous mechanism. When a firm adopts the simultaneous mechanism, although the second consumer arrives and makes a decision later than the first consumer, she does not gain any information advantage by arriving late. Not disclosing the choice of the first consumer makes the process equivalent to one where consumers simultaneously make decisions under appropriate time discounting adjustment. Thus, we define the simultaneous and sequential mechanisms from an information perspective, not based on the sequence of arrivals. Our approach follows the tradition of Varian (1994). The distinction between simultaneous and sequential mechanisms in this paper is similar in nature to that between sealed and sequential auctions. We

assume away observational learning where consumers may draw quality inferences from observation of peer choices, though the sign-up pattern under the sequential mechanism is similar to information cascades (Zhang 2010).

The firm's goal is to achieve a higher success rate. Higher success rates lead to higher expected profits when the firm receives a fixed lump-sum gain for each successful deal or a variable profit for each committed purchase given the deal is successful. Given any group-buying mechanism, the game consumers play is in nature a coordination game. To achieve a higher success rate, the firm aims at attaining better coordination between consumers. We interpret the success rate in the following sense.

Definition 1 (Success Rate) The success rate measures the ex-ante likelihood that a group-buying deal is successful before the consumers arrive.

The higher the success rate, the higher the expected payoffs for the firm and the higher the expected individual and total surpluses for consumers.

14.3.1.2 Equilibrium Analysis Under Simultaneous Mechanism

When the firm adopts the simultaneous mechanism, neither consumer is informed of the sign-up decision of the other. As a result, each consumer bases her sign-up decision on her belief in the valuation of the other one. The game is a Bayesian game in the Harsanyi sense (Harsanyi 1968) where “types” are defined by valuations. Specifically, the realized type of consumer i is defined by its realized valuation v_i , which is private information known to the consumer herself but not to the other consumer. Similarly, we denote by V_i the corresponding random variable of consumer i 's types, whose distribution is public information. For any consumer i , $s_i(v_i) \in \{0, 1\}$ denotes the decision rule that consumer i takes to decide whether to sign up or not, given her realized type being v_i . Given consumer i 's sign-up decision $s_i(v_i)$ and the other consumer $-i$'s sign-up decision $s_{-i}(v_{-i})$, the surplus of consumer i is denoted as $\pi_i(s_i(v_i), s_{-i}(v_{-i}))$.

Definition 2 (Bayesian Equilibrium) The Bayesian equilibrium strategy $\{s_i^*(v_i)\}$ for the simultaneous game is defined by the best-response strategy played by each consumer i , $s_i^*(v_i) \in \arg \max_{s_i(v_i) \in \{0,1\}} E_{V_{-i}} \{\pi_i(s_i(v_i), s_{-i}^*(V_{-i}))\}$ for all v_i .

The existence of Bayesian equilibria follows standard arguments (Tabarrok 1998). When consumers make sign-up decisions under the simultaneous mechanism, the consumer who arrives in the second period faces the same uncertainty in the sign-up probability of the other consumer as the one arrives in the first period. The solution scheme for this type of game has the following structure: there exists a valuation range for consumer i such that the consumer signs up if and only if her valuation falls into such a range. At the beginning of the game, the ex-ante belief of the firm on the success rate under a Bayesian equilibrium strategy $\{s_i^*(v_i)\}$

can be characterized by $q^* = P\left(\sum_{i=1}^2 s_i^*(V_i) = 2\right)$, which is the success rate at equilibrium. Since the revenue of a group-buying site depends on the success of deals, from this point on we compare different mechanisms by the seller’s expected deal success rate.

Denote q_i the belief of consumer i ’s sign-up likelihood held by consumer $-i$. As a result, consumer 1 signs up with probability $H_1(q_2)$, and consumer 2 signs up with probability $H_2(q_1)$. Equilibrium is characterized by a pair of beliefs (q_1^*, q_2^*) that satisfies the following conditions

$$H_1(q_2^*) = q_1^*, \quad H_2(q_1^*) = q_2^*.$$

Equivalently, q_i^* is given by

$$q = H_i(H_{-i}(q)), \quad i = 1, 2.$$

The above equations characterize equilibria in the sense that any Bayesian equilibrium results in a pair of beliefs (q_1^*, q_2^*) which satisfies the equations, and any pair of beliefs (q_1^*, q_2^*) that satisfies the equations correspond to a Bayesian equilibrium where consumer i behaves as if consumer $-i$ signs up for the deal with probability q_{-i}^* . The existence of such pair of beliefs (q_1^*, q_2^*) , and equivalently, the existence of a Bayesian equilibrium in pure strategies, is a direct consequence of Tarski’s Fixed Point Theorem, since $H_i(H_{-i}(q))$ is non-decreasing in $q \in [0, 1]$, $i = 1, 2$, by Assumption 1.

Following the preceding discussion, the success rate q^* at equilibrium can be characterized by

$$q^* = P(X_1(q_2^*) + X_2(q_1^*) = 2) = H_1(q_2^*) \cdot H_2(q_1^*), \tag{14.4}$$

where $X_i(q)$ is a Bernoulli random variable with success probability $H_i(q)$, $i = 1, 2$.

Notice that q_i^* is given by $q = H_i(H_{-i}(q))$, $i = 1, 2$, so there may exist multiple equilibria. Similar to Jackson and Yariv (2007), we categorize equilibria into two types, stable equilibria, and tipping points, depending on their sensitivity to minor perturbation in belief. Next, we provide a formal definition.

Definition 3 (Stable Equilibrium and Tipping Point) A pair of beliefs (q_1^*, q_2^*) is a stable equilibrium (tipping point) if for all $i = 1, 2$, there exists $\epsilon' > 0$ such that for all $\epsilon \in (0, \epsilon')$, when consumer i expects that the other consumer makes her decision with the belief $q_i^* - \epsilon$, consumer i ’s sign-up likelihood will be higher (lower) than $q_i^* - \epsilon$; when consumer i expects that the other consumer makes her decision with the belief $q_i^* + \epsilon$, consumer i ’s sign-up likelihood will be lower (higher) than $q_i^* + \epsilon$.

14.3.1.3 Equilibrium Analysis Under Sequential Mechanism

Under the sequential mechanism, at the beginning of the second period, the firm posts the decision of the consumer who arrives in the first period. Since two consumers make decisions sequentially, the first consumer needs to predict the sign-up probability of the second consumer. The sequential game we analyze follows the concept of rational expectations (RE) equilibrium.

Definition 4 (RE Equilibrium) For any realization v_i of the valuation for consumer i who moves first, an RE equilibrium $q_{-i}^*(v_i)$ in the sequential game satisfies: (i) Consumer i plays an optimal strategy of whether to sign up $S_i^*(v_i) \in \{0, 1\}$, given belief q_{-i} about the sign-up likelihood of consumer $-i$; (ii) Given the decision $S_i^*(v_i)$ from consumer i and any realization v_{-i} of the valuation of consumer $-i$, consumer $-i$ plays a best-response strategy $S_{-i}^*(v_{-i}; S_i^*(v_i)) \in \{0, 1\}$; (iii) The belief is consistent with the sign-up likelihood of consumer $-i$: $q_{-i} = q_{-i}^*(v_i) = P(S_{-i}^*(V_{-i}; S_i^*(v_i)) = 1)$.

Suppose consumer i arrives in the first period and consumer $-i$ arrives in the second period. When the second consumer makes her decision, there is no more uncertainty about the future. Given the decision of the first consumer, the optimal strategy for the second consumer can be characterized by a valuation range: to sign up if and only if (a) the first consumer signs up and (b) her own valuation falls into the range for signing up with the likelihood of the first consumer's sign-up being 1. After solving the best response from the second consumer, we move backward to the first consumer. Suppose consumer i expects that consumer $-i$ will sign up with probability q_{-i} . Then consumer i signs up with probability $H_i(q_{-i})$ in the first period. At equilibrium, q_{-i}^* is consistent with consumer $-i$'s sign-up likelihood, i.e., $q_{-i}^* = H_{-i}(1)$.

From the seller's perspective, the success rate Q_i^* at equilibrium can be characterized by

$$Q_i^* = P(X_i(q_{-i}^*) + X_{-i}(1) = 2) = H_i(q_{-i}^*) \cdot H_{-i}(1), \quad (14.5)$$

where $X_i(q)$, $X_{-i}(1)$ are Bernoulli random variables with success probability $H_i(q)$ and $H_{-i}(1)$, respectively. It is noteworthy that the equilibrium under the sequential mechanism is guaranteed to be unique, and hence it is stable.

14.3.1.4 Mechanism Design: Simultaneous or Sequential?

Given the potential presence of multiple equilibria under the simultaneous mechanism, how can one compare the success rates under simultaneous and sequential mechanisms? We adopt the approach proposed in Jackson and Yariv (2007) because it allows us to compare the set of equilibria regardless of equilibrium multiplicity. First, we formalize our criteria for comparisons of each consumer's belief regarding the other consumer's sign-up likelihood by the following definition.

Definition 5 (Higher Beliefs) For each individual consumer, one scenario or mechanism generates a higher belief than another if, for any belief at a stable equilibrium of the latter, there exists a higher belief at a stable equilibrium of the former, and for any belief at a tipping point of the latter there is a lower belief at a tipping point of the former or no lower tipping points of the former at all.

The reason why the sequential mechanism yields higher beliefs for both consumers is two-fold. First, with no tipping point, it is more likely for consumers' expectations to move upwards to the beliefs at the stable equilibrium. Second, when the stable equilibrium is reached, the beliefs regarding the other consumer's sign-up likelihood under the sequential mechanism is higher than that under the simultaneous mechanism for both consumers. If both consumers have higher expectations at equilibrium, then each individual is more likely to sign up, and thus the deal is more likely to succeed. Consequently, we can compare the deal's success rate by comparing the belief held by each individual consumer.

Definition 6 (Higher Success Rates) One scenario or mechanism generates a higher success rate than another if the belief of each consumer in the former is higher than the belief of the consumer in the latter in the sense of Definition 5.

To formally compare the belief held by each individual under alternative mechanisms, consider $f(q)$ and $g(q)$ as two functions parameterized by $q \in [0, 1]$. Suppose $f(q) \geq g(q)$ for any $q \in [0, 1]$. That is $f(q)$ is always higher than $g(q)$ point-wisely for any belief $q \in [0, 1]$. Then, for any left-to-right crossing point of curve $g(q)$ with the 45° line, there always exists a higher crossing point of curve $f(q)$; for any right-to-left crossing point of curve $g(q)$, there exists no lower crossing point of curve $f(q)$. Consequently, if $f(q) \geq g(q)$ for all $q \in [0, 1]$, we can conclude that mechanism with equilibrium characterization $f(q) = q$ yields higher belief than mechanism with equilibrium characterization $g(q) = q$.

Recall that, when consumer i arrives first, and consumer $-i$ arrives later, the belief held by consumer $-i$, q_i^* , under the simultaneous mechanism is given by $q = H_i(H_{-i})(q)$, and the belief held by consumer $-i$ under the sequential mechanism is simply 1. As $1 \geq H_i(H_{-i})(q)$ for any $q \in [0, 1]$, the sequential mechanism yields higher belief than the simultaneous mechanism for consumer i . Similarly, for the first mover consumer i , as $H_{-i}(1) \geq H_{-i}(H_i(q))$ for any $q \in [0, 1]$, the sequential mechanism yields higher belief than the simultaneous mechanism for consumer i as well. Given the preceding discussions, we have the following proposition.

Proposition 1 (Mechanism Comparison for the Two-Person Game) *Everything else being equal, the sequential mechanism always yields higher success rates than the simultaneous mechanism.*

The driving force behind this result is that in the simultaneous mechanism, each consumer faces uncertainty about the other consumer when making decisions. However, in the sequential mechanism, the second consumer decides only after the uncertainty about the first consumer has been resolved. Moreover, in anticipation of

that the second consumer will make the sign-up decision without facing uncertainty in the success rate, the first consumer's confidence about the second consumer's sign-up likelihood is consequently boosted.

14.3.2 Pricing

Next, we study how various pricing options could be used to facilitate consumer coordination under all-or-nothing mechanisms. To investigate this question, we use a similar two-period model, but now prices are endogenized. We base our description of the model on crowdfunding websites like Kickstarter.

14.3.2.1 Model Setup

Consider a risk-neutral creator who adopts a sequential crowdfunding mechanism for selling products. The creator posts a proposed project, with specific price information, on a crowdfunding platform. The sign-up process expires after two periods. In each period t , one buyer arrives at the proposed project. We denote the buyer at time t as B_t , with $t = 1, 2$. For the proposed project to succeed, both buyers would have to sign up.

Buyers may have different product valuations. To model this heterogeneity, we assume their valuations are i.i.d. with the following two-point distribution:

$$V_t = \begin{cases} H & \text{with probability } \alpha, \\ L & \text{with probability } 1 - \alpha, \end{cases}$$

where $H > L > 0$.

Upon arrival at the project in period 1, buyer B_1 realizes a private product valuation, makes the purchase decision, and leaves the site. The creator observes and announces the purchase decision of B_1 . Then, buyer B_2 arrives at the project, realizes a private product valuation, observes the purchase decision of buyer B_1 , and makes her own purchase decision. Both buyers are fully rational and make purchase decisions to maximize their own expected utility.

At the beginning of the game, the creator makes the pricing decision and posts it on a crowdfunding website. In addition, the creator decides the funding target, denoted by T . The creator commits to an all-or-nothing mechanism such that the project succeeds only if the total amount pledged reaches or exceeds T . Otherwise, the project fails. When making decisions, the creator knows the distribution of product valuations of two buyers but does not know the exactly realized valuations. The creator's goal is to maximize the expected profit from the proposed project. We assume, without loss of generality, that there is no transaction cost associated with pledging or rewarding, and there is no time discounting over the sign-up horizon.

Next, we define alternative pricing strategies with the target endogenized to be consistent with the pricing strategy and analyze their profitability.

14.3.2.2 Alternative Pricing Policies

Uniform Pricing With a uniform pricing strategy, the creator posts a single price of p for her product. Since the project succeeds only if both buyers sign up for the project and pay p , the creator can effectively set $T = 2p$. Price p can take any positive value; however, given the two-point distribution of product valuations, the optimal price should be either $p = H$ or $p = L$. Thus, it suffices to consider the following two cases:

Margin Strategy (H) With this strategy, the creator sets the price at $p^H = H$ and the target at $T^H = 2H$: any target beyond $2H$ would doom the project to failure; any target in $(H + L, 2H]$ is equivalent to $2H$, which sells only to high-type buyers, consistent with the term of “margin strategy.” Under this strategy, a high-type buyer will sign up, but a low-type buyer will decline. This strategy has a success rate of $s^H = \alpha^2$ and the creator has an expected profit of $\pi^H = 2\alpha^2 H$.

Volume Strategy (L) With this strategy, the creator sets the price at $p^L = L$ and the target at $T^L = 2L$: any target below $2L$ is equivalent to $2L$, with only the low price being paid, consistent with the term of “volume strategy.” Under this strategy, both buyers will sign up, regardless of their types, and the project always succeeds; i.e., $s^L = 1$. The creator’s profit is $\pi^L = 2L$.

Compared to the margin strategy, the volume strategy gives the creator a higher chance of project success; however, given that the project succeeds, the volume strategy yields a lower margin.

Intertemporal Pricing (D) With this strategy, a creator sets different prices for different periods, denoted by p_t^D for period t , $t = 1, 2$. Following the same logic as before, given the two-point distribution of product valuations, the optimal price in each period must be either L or H . That leads to two candidate strategies, either $(p_1^D, p_2^D) = (H, L)$ or $(p_1^D, p_2^D) = (L, H)$.

The creator’s goal is $T^D = H + L$: any target in $(2L, H + L]$ is equivalent to $H + L$, which leads to different prices charged in different periods, consistent with the term of “intertemporal pricing strategy.” The success rate is $s^D = \alpha$ and the creator’s expected profit is $\pi^D = \alpha(H + L)$. Note that since the buyers arriving at different periods have the same distribution of product valuations, the creator is indifferent between these two intertemporal pricing strategies.

Menu Pricing (M) With this strategy, the creator posts a menu of two prices, a high price of p_h^M and a low price of p_l^M , where $p_l^M \leq L \leq p_h^M \leq H$. Unlike intertemporal pricing, the optimal prices on the menu may not be equal to the two valuation points H and L (see Lemma 1 below). The creator sets the target at the sum of the high and low prices, i.e., $T^M = p_h^M + p_l^M$: any target in $(2p_l^M, p_h^M + p_l^M]$ is equivalent to $p_h^M + p_l^M$, which requires at least one buyer to pay the high price.

The menu pricing strategy proposed here may not appear to be a valid menu because the products of the same quality have different prices. That is done intentionally to tease out a buyer's incentive to overpay in crowdfunding.² With the traditional selling mechanism, such a menu would not work because each buyer would always choose the lower price option. However, in crowdfunding, a buyer, who is also a funder, may select the higher price if such a choice could substantially increase the likelihood of the project success. This is because that one buyer's behavior affects another buyer's expected utility; that is, positive externality arises from the common goal of project success.

We solve the optimal menu strategy with the backward induction method. In period 2, buyer B_2 observes the contribution of the earlier buyer B_1 and makes her own decision accordingly. Specifically, if B_1 has signed up at p_h^M , B_2 always signs up at the low price $p_l^M \leq L$, regardless of product valuations. On the other hand, if B_1 has signed up at p_l^M , B_2 either pledges p_h^M or does not sign up at all: for B_2 to sign up at p_l^M is meaningless because the project will certainly fail. Hence, buyer B_2 should sign up at p_h^M if her product valuation is H and otherwise not sign up at all. Recall that the probability is α for the product valuation to be H .

Next, we move back to the first period and consider buyer B_1 . If her product valuation is L , she always signs up at p_l^M . Otherwise, she can choose between two options. By choosing the low-price option p_l^M , B_1 expects a larger surplus $H - p_l^M$ but a lower success rate at α . Alternatively, by choosing the high-price option p_h^M , the buyer expects a smaller surplus $H - p_h^M$ but a higher success rate at 1. A high-type B_1 would prefer the high-price option p_h^M over the low-price option p_l^M if and only if the following incentive-compatibility (IC) condition is satisfied:

$$\alpha(H - p_l^M) \leq H - p_h^M. \quad (\text{IC})$$

The creator decides the optimal menu of prices to maximize the expected profit, subject to the condition (IC). Analyzing the creator's problem leads to the following lemma.

Lemma 1 (Optimal Menu Strategy) *With the menu strategy, the creator's optimal prices are $p_h^M = (1 - \alpha)H + \alpha L$, $p_l^M = L$ and the optimal target is $T^M = (1 - \alpha)H + (1 + \alpha)L$. The corresponding success rate is $s^M = \alpha(2 - \alpha)$, and the expected profit is $\pi^M = \alpha(2 - \alpha)((1 - \alpha)H + (1 + \alpha)L)$.*

Lemma 1 indicates that with the menu pricing strategy, everything else being equal, as long as the high price p_h^M is low enough such that the (IC) condition is satisfied, a high-type buyer B_1 prefers to pay the high price, even though a lower price option is available. The amount of overpayment is $p_h^M - p_l^M = (1 - \alpha)(H - L)$,

²The price menu may not be far stretched from practice. Many projects on Kickstarter contain levels with minute quality differences. For each interpretation, one may consider a trivial quality for the price menu in this section.

which increases with product valuation gap $H - L$ and decreases with α . Thus, when product valuation is more heterogeneous between different types of buyers or when buyer B_1 is more pessimistic about the product valuation of buyer B_2 , a high-type buyer B_1 has a greater incentive to overpay. With a larger gap $H - L$, a high-type buyer B_1 derives more utility from the project and is thus more willing to take a sacrifice to ensure the project's success. Similarly, with a smaller α , the risk of project failure is high, and thus the incentive to overpay is higher. As α decreases from 1 to a small value $\epsilon > 0$, the amount of overpayment increases from 0 to $(1 - \epsilon)(H - L)$, close to $H - L$.

When a high-type B_1 chooses price option p_h^M , this buyer enjoys a surplus of $\alpha(H - L)$. Since the project is sure to succeed, the buyer B_2 will choose the low price of $p_l^M = L$. Otherwise, had the buyer B_1 chosen p_l^M , a high-type buyer B_2 would have to pay p_h^M and incur a *reduction* of surplus $(1 - \alpha)(H - L)$. Thus, when the high-type buyer B_1 pays the higher price, the overpayment has a *positive externality* effect on the second buyer's surplus. Moreover, the size of the positive externality effect increases with the high price option in the menu. Thus, the crowdfunding mechanism not only artificially creates an externality effect among the decisions of buyers, but also determines the size of the externality effect endogenously.

14.3.2.3 Optimal Pricing Strategy

The creator determines the optimal pricing strategy by comparing the expected profits from each of the alternative pricing strategies. We summarize our analysis in the proposition below.

Proposition 2 (Optimal Pricing Strategy) *The creator's optimal strategy is*

1. *volume strategy, if*

$$\frac{H}{L} \leq \frac{2 - \alpha^2}{\alpha(2 - \alpha)};$$

2. *menu strategy, if*

$$\frac{2 - \alpha^2}{\alpha(2 - \alpha)} \leq \frac{H}{L} \text{ and } \alpha \leq \frac{3 - \sqrt{5}}{2}, \text{ or}$$

$$\frac{2 - \alpha^2}{\alpha(2 - \alpha)} \leq \frac{H}{L} \leq \frac{1 + \alpha - \alpha^2}{3\alpha - \alpha^2 - 1};$$

3. *intertemporal strategy, if*

$$\frac{1}{2} \leq \alpha \text{ and } \frac{1 + \alpha - \alpha^2}{3\alpha - \alpha^2 - 1} \leq \frac{H}{L} \leq \frac{1}{2\alpha - 1}, \text{ or}$$

$$\frac{3 - \sqrt{5}}{2} \leq \alpha \leq \frac{1}{2} \text{ and } \frac{1 + \alpha - \alpha^2}{3\alpha - \alpha^2 - 1} \leq \frac{H}{L};$$

4. *margin strategy*, if

$$\frac{1}{2} \leq \alpha \text{ and } \frac{1}{2\alpha - 1} \leq \frac{H}{L}.$$

Proposition 2 indicates that each of the four pricing strategies can be optimal within certain parameter subspaces. Uniform pricing strategies are optimal in two extreme cases. Specifically, given the valuation of L , the volume strategy is optimal when buyers are unlikely to have high product valuation (i.e., small α) or when high- and low-type buyers have a narrow valuation gap (i.e., H/L is close to 1). Intuitively, it can be seen that if both buyers are very likely to have a low product valuation, then the creator should pursue the volume strategy. The expected gain from targeting at only high-type buyers is not worth the risk that the project may fail. At the other extreme, if buyers are very likely to have a high valuation (i.e., large α) and the valuation gap between high- and low-type buyers is large (i.e., large H/L), then the creator should go for the margin strategy. In other words, if both buyers are very likely to have a high product valuation which is much higher than the low valuation, the creator should pursue the margin strategy. The gain from targeting the high-type buyers is significant, and the risk is bearable.

Two types of discriminatory pricing strategy, namely, the intertemporal and menu pricing strategies, are more profitable than the uniform pricing strategies when the fraction of high-type buyers (α) is not very large and the valuation ratio (H/L) is large enough, in other words, when buyers believe that others may have low product valuations and the product valuations are heterogeneous sufficiently. The difference between these two discriminatory pricing strategies is the timing: while the same menu of two options exists in both periods with a menu pricing strategy, a unique option is available in each period with an intertemporal pricing strategy. It is important to note that, first, with our model, an optimal discriminatory pricing strategy degenerates into a uniform pricing strategy in the traditional selling setting. The intertemporal, or menu, pricing strategy becomes optimal with a crowdfunding mechanism because a common target links two buyers. As a result, a high-type buyer may choose the high-price option to compensate for the small contribution from a low-type buyer. Second, the buyers can achieve coordination without any explicit interactions. In our model, self-interested buyers do not incorporate other buyers' utilities into their objectives (as in Chen and Li 2013), nor do some buyers communicate with others to increase their valuations (as in Jing and Xie 2011).

Table 14.6 shows the target amount and project success rate for each type of pricing strategy. Specifically, the target amount increases in the order of volume strategy, menu strategy, intertemporal strategy and margin strategy, and the project

Table 14.6 Success rate and target

	Success rate	Target
Volume strategy	1	$2L$
Menu strategy	$2\alpha - \alpha^2$	$(1 - \alpha)H + (1 + \alpha)L$
Intertemporal strategy	α	$H + L$
Margin strategy	α^2	$2H$

success rate decreases in the same order. When the fraction of high-type buyers (α) increases, the differences in the success rate among different strategies diminish, but the differences in the target level remain large. A high target leads to a low project success rate, and vice versa.³

To sum up, in a model without product quality differentiation, menu pricing strategy can be optimal in a crowdfunding mechanism. When buyers are sufficiently heterogeneous in product valuation, offering a menu of prices can help achieve a better balance between volume (or success rate) and margin. A novel insight is that a self-interested high-type buyer might be willing to pay extra to ensure the project's success. Thus, by turning buyers into funders, the crowdfunding mechanism enhances coordination among different buyers. Moreover, the menu pricing strategy, by choosing a proper level of the high-price option below H , moderates the high-type buyer's sacrifice and optimizes the coordination incentive. Crowdfunding creates compatibility between the purchases of two buyers who share the common goal of the project success. As a result, each buyer's behavior has an external effect on other buyers' utilities. The extent of the externality effect is regulated by the price difference in the menu.

14.4 Conclusion

Motivated by online group buying and crowdfunding, this chapter studies the common underlying mechanism behind the two business models, namely, all-or-nothing mechanisms. First, using a dataset from the online group buying industry, we empirically identify two types of threshold effects on consumer behavior induced by the mechanism. The first type of threshold effects refers to a *substantial increase* in the number of sign-ups around the time when the threshold was reached, and the second type of threshold effects refers to a stronger positive relation between the number of new sign-ups and the cumulative number of sign-ups before reaching the thresholds than afterward. For a more detailed discussion, please see Wu et al. (2014), who also discuss the heterogeneity of thresholds effects across different product categories and geographic locations.

Next, we discuss the optimal design of all-or-nothing mechanisms from the perspective of third-party platforms like Groupon and Kickstarter. We first study information disclosure strategies in the context of online group buying. In particular, we examine the success rate of a group-buying deal under two alternative mechanisms: a sequential mechanism and a simultaneous mechanism. Our analysis shows that all other things being the same, a sequential mechanism dominates a simultaneous mechanism. Interestingly, posting the cumulative number of sign-ups

³Our model can be easily adapted to situations where a not-for-profit creator wants to maximize the success rate, subject to raising enough funds to cover setup costs in advance. The fixed setup costs can become the exogenous target. Table 14.6 can be used as a guide for the optimal pricing strategy given the exogenous target. For example, if the exogenous target T falls in the range $(2L, (1 - \alpha)H + (1 + \alpha)L]$, the menu strategy maximizes the success rate.

from the first period can reduce uncertainty and thus increase the expected sign-ups among the second cohort of consumers. The increased second-period sign-ups can, in turn, improve the confidence among the first cohort of consumers and lead to a higher expected number of sign-ups in the first period, thus further increasing the deal's success rate. This backward-inductive perspective, starting from the second period and going back to the first period, is crucial to understanding the intuition behind our result. The driving force behind our result is that consumers essentially play a coordination game, and the sequential mechanism, by revealing information from one cohort to the other, allows for better coordination. We then study, in the context of crowdfunding, how all-or-nothing mechanisms may affect a creator's pricing decisions on the basis of a two-period game where cohorts of buyers arrive at a crowdfunding project and make sign-up decisions sequentially. Our results show that, even when product options are the same, high-type buyers may still choose the high-price option. This result is unique to the crowdfunding mechanism, where tacit coordination among buyers is necessary to ensure the project success. For more detailed analysis of the preceding two strategies, interested readers are referred to Hu et al. (2013) and Hu et al. (2015).

While this chapter offers useful insights into the impact and design of online group buying and crowdfunding, future research is required to further our understanding of the issue. For instance, traditionally, the effort to improve the success rates of projects concentrates on optimizing the upfront design of project characteristics. However, because of the inherent uncertainty and all-or-nothing mechanism of online group buying and crowdfunding projects, contingently providing incentives or adjusting project features over the course of a campaign can be as crucial as the upfront design. As a first step towards the understanding of dynamic policies, Du et al. (2017) study the optimal design of contingent stimulus policies for crowdfunding campaigns.

References

- Arellano M, Bond S (1991) Some tests of specification for panel data: Monte Carlo evidence and an application to employment equations. *Rev Econ Stud* 58(2):277–297
- Bass F, Krishnan T, Jain D (1994) Why the bass model fits without decision variables. *Market Sci* 13(3):203–223
- Bound J, Jaeger D, Baker R (1995) Problems with instrumental variables estimation when the correlation between the instruments and the endogeneous explanatory variable is weak. *J Am Stat Assoc* 90(430):443–450
- Chen Y, Li X (2013) Group buying commitment and fundraisers' competitive advantages. *J Econ Manag Strateg* 22(1):164–183
- Du L, Ming H, Wu J (2017) Contingent stimulus in crowdfunding, Working paper. <http://ssrn.com/abstract=2925962>
- Harsanyi J (1968) Games with incomplete information played by “Bayesian” players parts II. Bayesian equilibrium points. *Manag Sci* 14(5):320–334
- Häubel G, Popkowski Leszczyc P (2004) Bidding frenzy: intensity of competitive interaction among bidders and product valuation in auctions. *Adv Consum Res* 31:90–93

- Heyman J, Orhun Y, Ariely D (2004) Auction fever: the effect of opponents and quasi-endowment on product valuations. *J Interact Market* 18(4):7–21
- Hu M, Shi M, Wu J (2013) Simultaneous vs. sequential group-buying mechanisms. *Manag Sci* 59(12):2805–2822
- Hu M, Li X, Shi M (2015) Product and pricing decisions in crowdfunding. *Market Sci* 34(3): 331–345
- Jackson M, Yariv L (2007) Diffusion of behavior and equilibrium properties in network games. *Am Econ Rev* 97(2):92–98
- Jing X, Xie J (2011) Group buying: a new mechanism for selling through social interactions. *Manag Sci* 57(8):1354–1372
- Kauffman R, Wang B (2002) Bid together, buy together: on the efficacy of group-buying business models in Internet-based selling. In: Lowry PB, Cherrington JO, Watson RR (eds) *Handbook of electronic commerce in business and society*, CRC Press, Boca Raton
- Ku G, Malhotra D, Murnighan J (2005) Towards a competitive arousal model of decision-making: a study of auction fever in live and internet auctions. *Organ Behav Hum Decis Process* 96(2): 89–103
- Nickell S (1981) Biases in dynamic models with fixed effects. *Econometrica* 49(6):1417–1426
- Roodman D (2009) A note on the theme of too many instruments. *Oxf Bull Econ Stat* 71(1): 135–158
- Staiger D, Stock J (1997) Instrumental variables regressions with weak instruments. *Econometrica* 65(3):557–586
- Tabarrok A (1998) The private provision of public goods via dominant assurance contracts. *Public Choice* 96(3–4):345–362
- Varian H (1994) Sequential contributions to public goods. *J Public Econ* 53(2):165–186
- Villas-Boas J, Winer R (1999) Endogeneity in brand choice models. *Manag Sci* 45(10):1324–1338
- Wu J, Shi M, Hu M (2014) Threshold effects in online group buying. *Manag Sci* 61(9):2025–2040
- Zhang J (2010) The sound of silence: observational learning in the US kidney market. *Market Sci* 29(2):315–335
- Zhang J, Liu P (2012) Rational herding in microloan markets. *Manag Sci* 58(5):892–912

Chapter 15

Threshold Discounting: Operational Benefits, Potential Drawbacks, and Optimal Design



Simone Marinesi, Karan Girotra, and Serguei Netessine

Abstract We study the use of threshold discounting, the practice of offering service at a discounted price only if at least a given number of customers show interest in it, pioneered by Groupon. We model a capacity-constrained firm offering service to a random-sized population of strategic customers in two representative time periods, a desirable hot period and a less desirable slow period. A comparison with the traditional approach typically employed in such circumstances (slow period discounting or closing) reveals that threshold discounting boosts the firm's operational performance on account of two advantages. First, the contingent discount incentivizes slow period consumption when the market for the service is large and reduces supply of the service when the market is small, allowing the firm to respond to the service's unobserved market potential. Second, activation of the threshold discount signals the market size to strategic customers, supplying them with information on service availability, and inducing them into self-selecting the consumption period to one that improves the firm's capacity utilization and profit. Unlike typical settings with strategic customers, their strategic behavior in our setting increases the firm's profits. When threshold discounts are offered through an intermediary, arrangements often used in practice distort the incentives of the intermediary, and typically result in a higher discount and a lower activation threshold relative to what would be optimal for the service firm. We consider alternate deal designs, and we find that the best designs compromise the service provider's flexibility in order to provide customers with clear offer terms.

S. Marinesi (✉) · S. Netessine
The Wharton School, University of Pennsylvania, Philadelphia, PA, USA
e-mail: marinesi@wharton.upenn.edu; netessin@wharton.upenn.edu

K. Girotra
Cornell Tech, New York, NY, USA
e-mail: karan@girotra.com

15.1 Introduction

Firms often operate in environments in which they must serve highly variable demand with capacity that is fixed in the short term. Demand seasonality makes it even more difficult for service firms to maintain high capacity utilization, as spare capacity in low-demand periods typically cannot be used to serve customers in high-demand periods. Among the industries that struggle with this problem are the movie theater industry (\$10.2 Billion of revenues in 2016), the restaurant industry (\$783 Billion in 2016),¹ and a broad variety of retail services, such as spas, bowling clubs, circuses, museums, etc.

Over the last decade, advances in online technology have allowed firms to reach an unprecedented level of engagement with their customers. In particular, around 2010 there has been a sudden surge in online discounted deals. The most famous online deal website is certainly Groupon: founded in late 2008 and the first of an innovative breed of firms, it grew by 2,241% in its second year of operation – faster than celebrated firms like Amazon or Ebay – and went public in 2011, raising \$700 Million to become the largest IPO by a US Internet company after Google.² Groupon's growth was fueled by the use of an innovative discount structure, in which customers could purchase retail services with a substantial discount, but the discount was valid only if at least a certain number of customers showed interest in the offering. From here on, we refer to deals where discounted service is contingent on a threshold number of customers purchasing it as *threshold discounting* offers.

The popularity of threshold discounting increased dramatically since Groupon started using these offers, together with the industry that grew around them, becoming an almost essential feature for the hundreds of websites that spawned all over the world trying to imitate Groupon's business model.³ These offers have been highly praised by the business press, and have more recently received the attention of the academic community as well. Their celebrated advantages can be summarized in the ability to leverage “network effects” and economies of scale.

From the existing studies that exist on threshold discounting offers emerges a picture that is incomplete and controversial. Incomplete, because the advantages of such offers have so far been studied only from a marketing perspective – ignoring, for example, capacity constraints – and little is known about their impact on a firm's operations. Controversial, because it is difficult to reconcile the celebrated advantages of threshold discounting offers with their progressive discontinuation by many online discounters, including Groupon. Moreover, no analysis has been attempted so far to determine how threshold discounting offers should be implemented in details.

¹Sources: <https://goo.gl/emx1jR>, <https://goo.gl/7DvBqD>. Expected revenues for the United States.

²See “Groupon's IPO biggest by U.S. Web company since Google”, Reuters, November 4, 2011, <http://goo.gl/h8VFj>, and “Groupon IPO: Growth Rate Is 2,241%”, The Wall Street Journal, June 2, 2011, <http://goo.gl/UFFwK>

³See “LivingSocial aims to be different from Groupon”, Sep 23, 2011, <https://bloom.bg/2wcGtnt>, and “Psyched to Buy, in Groups”, New York Times, Feb 9, 2011, <http://nyti.ms/2wXVcmK>

Hence, three important questions remain open. What is the advantage of threshold discounting from an operational perspective? How can we reconcile the advantage that such offers provide to merchants with their progressive discontinuation? What is the best way to design threshold discounting offers in order to best leverage their operational advantage?

The objective of this paper is to provide answers to these three questions. We consider a typical situation in which threshold discounting is used, i.e., a capacity-constrained firm offers his services to a random-sized population of strategically-acting customers who prefer to be served on a desirable “hot” time period over a less desirable “slow” time period (e.g., a movie theater on Saturday versus Monday evening), with the degree of preference for the hot period over the slow period varying across the population. Demand is thus variable but substitutable between two vertically differentiated services, the hot-period service and the slow-period service.

To answer the first question on the operational advantages of threshold discounting offers, we compare them with a more traditional approach typically used by firms in this context, that is, to either close on the slow period, or open on the slow period at a discounted price. The comparison reveals that threshold discounting outperforms the traditional approach on account of two effects. First, by setting an activation threshold, threshold discounting endows the firm with a built-in, demand-responsive mechanism (which we refer to as *responsive duality*) that matches different market states with appropriate pricing/closing decisions, resulting in higher capacity utilization and better fixed costs management. Second, threshold discounting induces a *strategic scarcity effect* that increases customers’ responsiveness to slow-period discounts by exploiting their strategic behavior: inducing them into self-selecting their consumption period to one that better serves the firm’s interests of managing capacity and margins. Notably, the superior performance of threshold discounting over the traditional approach holds *even in the absence of economies of scale or network consumption effects*, the only benefits of threshold discounting that have been studied so far.

Our comparison also highlights an (unexpected) impact of strategic customers on the effectiveness of threshold discounting. The literature on strategic customers has largely found that they reduce a firm’s profit because they time their purchases in order to get lower prices, thus reducing margins for the firm. In contrast with the literature, we find that in our context having more strategic customers is *beneficial* for the firm, that is, the higher the fraction of strategic customers in the population, the higher the firm’s profit.

To answer the question on why threshold discounting offers may have been discontinued, we consider the case in which threshold discounting is offered through an intermediary with high negotiating power (such as Groupon) as it often happens in practice. We find that in these cases the intermediary has strong incentives to prefer a lower activation threshold (often equal to zero) a higher discount relative to what is in the interest of the service provider. This result suggests that disappearance of the activation threshold and the very high discounts that are usually featured

with threshold discounting offers may be due to incentive misalignment between the powerful intermediaries that dominate the industry and the service providers that need their services.

To answer our last question on how to best design threshold discounting offers, we study various dimensions of their design. We first consider if the firm should commit upfront to the number of customers required for the deal to be active, or instead decide on the deal activation only after observing the subscription level. Interestingly, we find that postponing the deal-activation decision to incorporate the market information contained in the subscription level is harmful to the firm. We then investigate the best time for the firm to reveal whether the threshold has been reached – specifically, if this should be before both periods begin or not – and find that early disclosure is the superior design. We also compare time-restricted threshold discounts with the often-used unrestricted discounts, and find that time-restricted discounts are superior. Finally, we show that under certain conditions the preferred design with *committed threshold, early disclosure, and time-restricted discount* can be further improved by offering targeted discounts that reduce lost margins while retaining all the operational advantages of these deals.

Our work makes several contributions. This is the first study to examine the operational advantages of threshold discounting, a popular phenomenon that has spawned a multi-billion-dollar online deals industry. In contrast with the strategic customer behavior literature in operations, we show that in our setting, strategic customer behavior is beneficial. Further, we provide clear prescriptions about the design and use of threshold discounts – specifically, that a service provider is always better off offering a *committed threshold, early disclosure, time-restricted* threshold discount as opposed to the variety of other designs touted by deal intermediaries and designers, and we show that focused threshold discounts can improve profits even further. Our analysis cautions potential users of these discounts to the incentive conflicts inherent in the current modes of offering these discounts through intermediaries, and provides an explanation for why threshold discounts have been discontinued despite their potential benefits for service providers. Overall, our work argues that threshold discounting is an overlooked strategy for managing capacity that, if used correctly, can significantly improve a firm's operational performance and profit.

15.2 Literature Review

Our work is related to three different streams of literature: group-buying and quantity discounts, strategic consumers, and demand manipulation via pricing.

Group-Buying and Quantity Discounts In the early 2000s, several group buying websites like Mercata.com, LetsBuyIt.com, and Mobshop.com were founded with the objective of aggregating the buying power of customers to obtain quantity discounts. Anand and Aron (2003) model these group buying practices with a firm

offering a price-quantity menu to customers, and find that group buying is better than a simple fixed price only when either demand uncertainty satisfies certain conditions, or economies of scale are coupled with production postponement. In a similar spirit, Chen et al. (2007) study group buying auctions, where a firm commits to a price-quantity function and customers arrive stochastically and bid their reservation price; they find that group buying is better than simple fixed pricing only in the presence of economies of scale, or when the firm is risk-seeking. Unlike our paper, the above works do not consider vertically differentiated services or threshold discounting schemes.

Closer to our work are recent papers that specifically investigate threshold discounting offers. Jing and Xie (2011) show that in a threshold discounting offer, informed players reach out to their uninformed friends in an attempt to reach the threshold and obtain the discount, to the firm's benefit. Chen and Zhang (2014) analytically show that threshold discounts are the optimal mechanism to price discriminate a population of customers, under some conditions on the size uncertainty of different customer segments. Li and Wu (2018) and Wu et al. (2014) empirically study the evolution of customers' subscriptions over time, with the former studying herding and word-of-mouth effects, and the latter examining effects driven by the activation threshold. Hu et al. (2013) models the impact of two modes of pledging – sequential or simultaneous – on customers' pledging decisions and firm profit. The success of Groupon has also spurred studies on broader issues other than threshold discounts – see Edelman et al. (2016) and references therein. While similarly inspired, none of these papers take an *operational* perspective on threshold discounting offers, consider inter-temporal demand substitution effects, provide explanations for why major players have discontinued threshold discounting offers, or studies how to best design threshold discounting offers.

Strategic Customers A few decades ago, Coase (1972) conjectured that a monopolist selling a durable good would eventually lower its price down to marginal cost when facing infinitely patient consumers. Recent years have seen renewed interest in the operational implications of customer strategic behavior. Most of the work has focused on strategic purchasing delay on the part of customers when a firm sells a finite inventory of a durable good and may change the price over time. Su (2007) considers customers with different valuations for the product and degrees of patience, and develops insights on how the interplay of these characteristics affects the firm's pricing policy and profit. Liu and van Ryzin (2008) study how the capacity choice of a firm can be used to induce a rationing risk on risk-averse strategic customers and limit their strategic purchasing delay. Cachon and Swinney (2009) consider a setting in which the firm cannot commit in advance to prices, and they study the value of quick response strategies to mitigate the negative effect of strategic purchase delays on the part of customers when there are different classes of customers, while Cachon and Swinney (2011) explore the interplay of quick response and enhanced design in fast fashion systems. Aviv and Pazgal (2008) consider both pre-announced and contingent pricing strategies, and they provide recommendations for when these different approaches should be

used if both are viable. Su and Zhang (2008) study the value of both quantity and price commitment, and show how a decentralized supply chain can exploit the inefficiencies of decentralization as a commitment device to indirectly implement price and quantity commitment strategies, even when commitments are not credible. Strategic customers have also been studied in other situations, see Netessine and Tang (2009) for more references.

Like many of the above work, our customers time their purchases taking into account the strategic behavior of other players. Unlike the above papers, however, we explore the consequences of such strategic behavior in a novel setting, a firm that employs threshold discounting while offering service in two vertically differentiated service periods. The impact of strategic customer behavior in our context are unexpected and in contrast with the main findings from this large literature.

Demand Manipulation via Pricing This body of literature deals with situations in which a capacity-constrained (or inventory-constrained) firm can use the pricing decision to reduce the supply-demand mismatch.

All the literature on revenue management, for instance, focuses on this topic (see Talluri and Van Ryzin 2006, for a survey) including all papers on peak load pricing (see Crew et al. 1995, for a survey on this topic). In his paper on price dispersion, Dana (1999) shows the operational benefit of shifting demand across time periods by rationing the number of seats offered at a lower price, even when firms cannot predict the peak time. In other settings, Lus and Muriel (2009) find that pricing is more effective than technology choices at balancing supply and demand when a firm sells highly substitutable products, and Boyacı and Özer (2010) show how advanced selling and pricing can be jointly used to reduce the demand-supply mismatch.

Our paper departs from the existing literature in that we study a way to reduce the supply-demand mismatch through a novel pricing approach: namely, we study the use of correctly designed threshold discounting offers in the presence of strategically-acting customers.

15.3 The Model

15.3.1 Preliminaries

Consider a capacity-constrained service provider that offers his services in two time periods to a random-sized population of strategically acting customers. Customers prefer to be served in a desirable “hot” time period over a less desirable “slow” time period (e.g., a movie theater on Saturday versus Monday night) and have varying degrees of preference. After briefly describing the model, we first examine the traditional approach typically employed by firms in similar circumstances – a choice between closing down or discounting on the slow period – and then we compare the result with threshold discounting offers as popularized by online deal sites such as Groupon, LetsGroop, BigDeal, etc.

Service Economics We model provision of the service in two representative time periods: in a hot period preferred by customers, at a price r_h , and in a less preferred slow period.⁴ The service provider has capacity to serve at most k customers during each service period, but has the flexibility to shut down or choose any price in the slow period. When offering the service in a given period, the service provider incurs fixed costs c_F (for employees, utilities, etc. . .) plus an additional expense of c for every customer served. The costs are not too prohibitive to preclude profits, $c_F < k(r_h - c)$. When demand exceeds capacity, the provider rations capacity randomly amongst customers.

Customers The service is made available to a market comprised of infinitesimal customers of aggregated size \tilde{x} , where \tilde{x} is an unobserved random variable with support \mathbb{R}_+ , cumulative distribution function G , and survival function $\bar{G} = 1 - G$. Customers value the service in the hot period at v_h , $v_h > r_h$, higher than their value in the slow period, \tilde{v}_s , which varies across customers – i.e. customers differ in the degree to which they prefer the hot period over the slow period. Each customer's slow period valuation, \tilde{v}_s , is privately known only by the customer herself; it is drawn from a continuous distribution, with cumulative density H , survival function \bar{H} , and support $[\underline{v}, v_h)$. Customers desire to consume the service in at most one period, and they can choose their time of consumption strategically – i.e. each customer takes into account the choices of other customers, thus forming expectations of the service availability in different periods. Customers use these refined beliefs, in addition to the provider's announced shutdown and pricing decisions, and the private information on the slow period valuation to make their consumption timing decisions.

The setup described above corresponds to a wide variety of consumer services such as movie theaters, spas, opera houses, etc. Each of these services share the key characteristics of our setup – desirable and less desirable service periods, single consumption, and per-period capacity that is fixed in the short run.

15.3.2 *The Traditional Approach: Seasonal Closure or Regular Discounting*

Traditionally, service providers either shut down in slow periods or remain open but try to attract customers by offering a discounted price. For example, in many cities of mainland Europe where fixed costs of operation are substantial, restaurants and museums typically close on Monday.⁵ On the other hand, in London, service providers are often open on Mondays, but offer discounts and promotions to attract

⁴We place no restrictions on which period comes first.

⁵See <http://goo.gl/M52do>

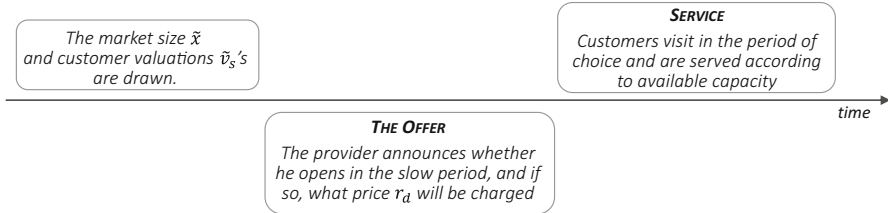


Fig. 15.1 Timeline for the traditional approach

customers.⁶ Formally, the service provider decides whether to offer the service in the slow period and, in case service is provided, what price to charge. The service provider makes these decisions knowing the probability distribution of the market size, G , but *not its realization* x . The sequence of events is provided in Fig. 15.1.

In the traditional approach, the service provider compares the profits from closing on the slow period (seasonal closure) and the profits from opening and offering a discount (regular discounting) and chooses to act so to maximize expected profit. Closing the business in the slow period implies that the service provider gives up some of his capacity – capacity available during the slow period – in order to save on the fixed cost c_F . In this case, customers visit during the hot period, and the service provider serves them up to capacity. The expected profit with this approach is

$$\Pi_c = (r_h - c) \int_0^{+\infty} \min(k, x) dG(x) - c_F, \tag{15.1}$$

where the subscript c stands for closure.

Alternatively, the service provider may offer the service in both hot and slow periods, albeit at a lower price $r_d \leq r_h$ in the slow period – where the subscript d stands for regular discounting. Under this strategy, a customer’s consumption timing best-response is driven by a trade-off between the higher utility she derives from the hot period on the one hand, and the better prices in the slow period on the other, both adjusted by her rational expectation regarding service availability in each time period.

Formally, a customer visits during the slow period iff her slow period valuation for the service is higher than a threshold valuation $\hat{v}_d(r_d)$, which is the valuation that makes a customer indifferent between the two service periods, and is given by

$$(v_h - r_h) \int_0^{+\infty} \min\left(1, \frac{k}{H(\hat{v}_d)x}\right) dG_c(x) = (\hat{v}_d - r_d) \int_0^{+\infty} \min\left(1, \frac{k}{\bar{H}(\hat{v}_d)x}\right) dG_c(x), \tag{15.2}$$

⁶For example, [Maxwell’s](#), [The Lexi Cinema](#), and [Cavendish Conference Venues](#) run “Monday madness” promotions, reducing their prices on Mondays, when they expect fewer customers.

where the LHS (RHS) represents the expected surplus of the customer from visiting during the hot (slow) period, obtained as the product of the service surplus times the expected availability of the service, and where $G_c(x) = \int_0^x u dG(u) / \int_0^{+\infty} dG(u)$ is the cdf of the market size from the perspective of an individual customer, i.e. conditional on her existence in the market (see Deneckere and Peck 1995, for the derivation of customer posterior beliefs in these cases). Note that customers' best-response visit strategy $\hat{v}_d(r_d)$ is unique since the LHS and RHS of Eq. 15.2 are respectively increasing and decreasing in \hat{v}_d for every price r_d . The expected profit for the provider is given by

$$\begin{aligned} \Pi_d^* = \max_{r_d} & \left[(r_h - c) \int_0^{+\infty} \min(k, H(\hat{v}_d(r_d))x) dG(x) \right. \\ & \left. + (r_d - c) \int_0^{+\infty} \min(k, \tilde{H}(\hat{v}_d(r_d))x) dG(x) - 2c_F \right]. \end{aligned} \quad (15.3)$$

The potential advantage of regular discounting is best explained by rewriting the profit as the product of the expected unit margin M_d (expected profit over expected sales) times expected capacity utilization U_d (expected sales over available capacity) times capacity k , minus fixed costs $2c_F$:

$$\Pi_d(r_d) = \max_{r_d} [M_d(r_d) \cdot U_d(r_d)] k - 2c_F. \quad (15.4)$$

A price reduction has two consequences for the firm: it always reduces the expected margin – both because customers pay less during the slow period and, because of the lower price, more customers visit in the slow period – and it shifts some demand from the hot to the slow period. While the reduction in margin is always harmful, the shift in demand can be beneficial; specifically, if the discount is not excessive, the demand shift improves capacity utilization by balancing demand across the two service periods, with the highest capacity utilization being at $r_d = \hat{v}_d^{-1}(H^{-1}(1/2))$, when demand is the same in both periods. Formally, $M_d(r_d)$ and $U_d(r_d)$ are respectively increasing and unimodal in r_d . The optimal price is the one that optimally trades off higher margins and higher capacity utilization.

The firm's expected profit with the traditional approach is therefore $\Pi_a^* = \max(\Pi_c, \Pi_d^*)$, where Π_c and Π_d^* are the profits if the firm closes on the slow day or not, defined in Eqs. 15.1 and 15.3, respectively. While the firm must decide ex-ante whether to close in the slow period (seasonal closure) or open and offer a discount (regular discounting), it is instructive to compare ex-post profits as a function of market size realization, to examine when it would have been better to open, and when it would have been better to close. Put differently, the next Lemma provides the strategy that would be followed by an omniscient firm, a firm that could observe the market size from the start.

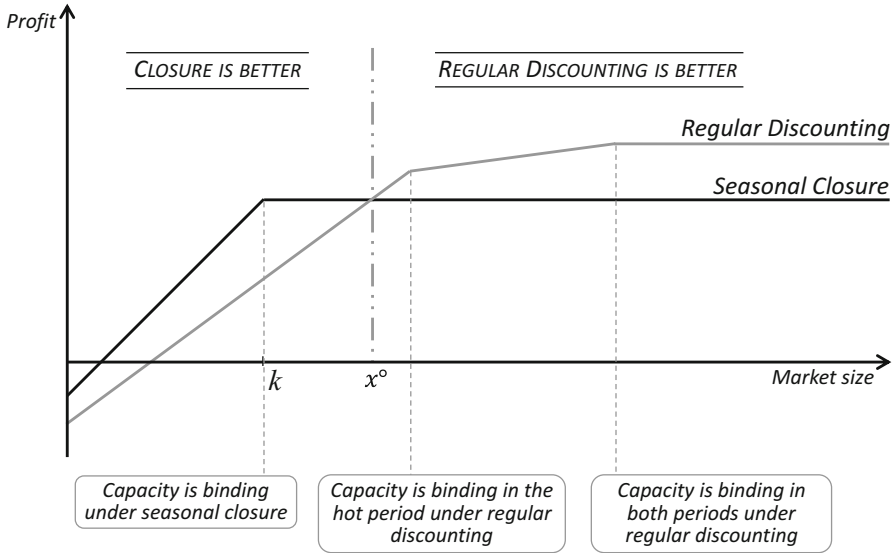


Fig. 15.2 Profits from regular discounting and closure

Lemma 1 *The realized profits under regular discounting are higher than under seasonal closure iff the market size realization is higher than a critical level, x° .*

All proof can be found in the online compendium.⁷ Figure 15.2 shows the realized, ex-post profit of regular discounting and closure, together with their difference, as a function of realized market size. When market size is low, closing on the slow day is preferred, as it both saves on fixed costs and keeps margins high; the advantage of closing is the highest when market size is equal to capacity k . Any higher market size results in lost sales with closure, but corresponds to higher sales if the firm is open in the slow period. Eventually, this makes regular discounting preferred to closure when market size is higher than x° .

Taken together, the above discussion highlights the key weakness of the traditional approach. The service provider, not knowing the actual size of the market, is forced to make an ex-ante trade-off: choosing the preferred strategy for low market sizes, closure, and bearing the risk of losing sales if the market size is high – due to limited capacity – or choosing the preferred strategy when the market size is high, regular discounting, and bearing the risk of not repaying the augmented fixed costs if the market size is low – due to thinner margins. We next examine a threshold discounting scheme, which alleviates this trade-off.

⁷<http://ssrn.com/abstract=3031173>

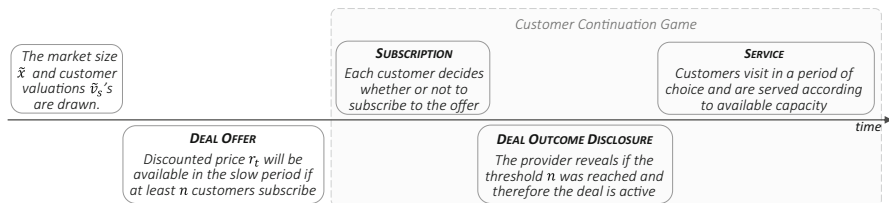


Fig. 15.3 Timeline of threshold discounting

15.3.3 Threshold Discounting

Threshold discounting allows customers to visit a firm and avail themselves of the service in the slow period at a discounted price, contingent on enough other customers showing interest in doing the same. In this section, we analyze the benefits of offering such a deal to strategically acting customers.

15.3.3.1 Sequence of Events

Figure 15.3 illustrates the sequence of events for a threshold discounting scheme. At the beginning, the market size is realized, but is not observed by either the service provider or the customers. Then, the service provider announces a deal: the service will be offered at a discounted price $r_t < r_h$ to all customers who subscribe to the offer, but only if at least n_t of them end up signing up for it. If less than n_t customers sign up for the deal, the service provider will close during the slow period.⁸ Each customer then decides whether to subscribe to the deal or not.⁹ At this point, the firm announces whether the number of subscriptions is at least equal to the activation threshold (and the deal is active) or not. Customers then choose a period to visit, and consume the service.

In order to study this game, it is convenient to break it into two parts: the initial deal offer, and the following continuation game (Fudenberg and Tirole 1991, page 331), in which customers subscribe, the firm reveals the deal outcome, and then customers visit the firm in a period of their choice. This continuation game, which follows the initial deal offer, is strategically played only among customers: in fact, in the deal outcome disclosure stage the firm simply reveals whether the pre-announced threshold was reached (in Sect. 15.3.7.1 we consider an alternate design

⁸A weaker form of threshold discounting is one where only the pricing decision is determined by subscribers, and the firm remains open in both periods. In this case the results are very similar, as briefly discussed in Sect. 15.4.

⁹We assume that subscriptions are not binding for customers; if subscriptions are binding, i.e., customers are pre-charged the slow period price r_d upon subscription, all our results are identical and in fact the analysis is simpler.

where the firm can freely make his activation decision after observing the number of subscribers). In what follows, we first proceed with the analysis of the customer continuation game for a given deal offer (r_t, n_t) , and then we include the deal offer decision made by the firm to find the equilibrium of the full game.

15.3.3.2 Customer Continuation Game

We examine the best-response strategies of an individual customer starting from the last stage of the continuation game, when she must decide her visit strategy, that is, in which period to visit the firm. When the deal is not active, the service is not available during the slow period and therefore she visits in the hot period. When instead the deal is active, she visits in the period in which she expects to obtain the highest surplus. Specifically, the visit strategy v_i of customer i is a function of her service valuation for the slow period $v_{s,i}$, the price that she will be charged during the slow time period, i.e. r_t if she has previously subscribed to the deal and r_h otherwise, and the expected service availability in each time period. To compute expected service availability, the customer takes into account the vector of visit strategies of all other customers, v_{-i} , and updates her belief on the realized market size *conditional on the information that the deal is active*. That is, her posterior distribution of the market size accounts for the fact that at least n_t customers subscribed employing the subscription strategy σ_{-i} , and is computed using Bayes' rule.

Next is the subscription stage, in which we assume that a customer subscribes iff this increases her expected future payoff (or alternatively that the frictional cost to subscribing is small). Specifically, the best-response subscription strategy of customer i is a function of her valuation for the service during the slow period $v_{s,i}$, of the announced deal price r_t and of the threshold n_t , as well as the vectors of subscription and visit strategies of all other players, σ_{-i} , v_{-i} . The subscription stage belongs to the class of Coordination games first defined by Schelling (1960): in this type of game there are typically multiple equilibrium outcomes, where if enough customers coordinate on a certain decision, a single customer has no incentives to deviate from what the majority does. The multiple equilibria that arise can be meaningfully grouped into two types. In type I equilibria, customers subscribe if their valuation for the slow period is sufficiently high, while in type II equilibria, customers never subscribe to the deal: as a consequence, the deal is never active, and therefore not subscribing is optimal. We discard type II equilibria from our analysis and restrict our attention hereafter to type I equilibria because a type I equilibrium Pareto Dominates a type II equilibrium.¹⁰ Type I equilibria are characterized in the next lemma

¹⁰Customers are better off coordinating on a type I equilibrium than on a type II equilibrium. The reason is that by subscribing to the deal, they have a chance to get a discount and visit in their preferred slow period, while at the same time increasing availability for those who did not subscribe.

Lemma 2 (Equilibrium strategies in Customer Continuation Game) *In a type I equilibrium:*

1. *A customer subscribes to the deal iff her valuation for the service in the slow period is higher than a certain threshold.*
2. *A customer visits during the slow period iff the deal is active and her valuation is higher than a threshold; she visits during the hot period otherwise.*
3. *The subscription and the visit thresholds coincide.*

In equilibrium, customer subscription and visit strategies are of a threshold type, and the thresholds for the two strategies coincide, since customers who would visit in the slow period are the same as those who subscribe to the deal. Thus, customer behavior can be fully summarized by just one threshold, \hat{v}_t , such that a customer with a slow-period valuation lower than \hat{v}_t does not subscribe to the offer and visits in the hot period, whereas a customer with a slow-period valuation higher than \hat{v}_t subscribes to the offer, and then visits in the slow period if the deal is active and in the hot period when the deal is not active. The threshold valuation \hat{v}_t is the one that, conditional on the deal being active, makes a customer indifferent between visiting in the slow and in the hot period – since when the deal is not active both subscribers and non-subscribers visit on the hot period and earn the same surplus. The following equation compares the threshold customer's surplus in each period, when the deal is active, for any deal offer (r_t, n_t) :

$$\begin{aligned} (\hat{v}_t - r_t) \int_{n_t \bar{H}(\hat{v}_t)^{-1}}^{+\infty} \min\left(1, \frac{k}{\bar{H}(\hat{v}_t)x}\right) dG_c(x) \\ = (v_h - r_h) \int_{n_t \bar{H}(\hat{v}_t)^{-1}}^{+\infty} \min\left(1, \frac{k}{Ht(\hat{v}_t)x}\right) dG_c(x). \end{aligned} \quad (15.5)$$

The LHS represents customer surplus when she visits in the slow period, and the RHS when she visits in the hot period. Unfortunately, for a general deal offer (r_t, n_t) there can be multiple solutions to Eq. 15.5, and consequently multiple type I equilibria. An increase in the threshold, \hat{v}_t , implies an increase in the number of visitors in the hot period and a corresponding decrease in the number of visitors in the slow period. While the hot period surplus decreases when \hat{v}_t increases because a higher fraction of customers visiting in the hot period reduces availability, the slow period surplus generally does not increase in \hat{v}_t . A higher threshold \hat{v}_t implies fewer customers visiting in the slow period, which should increase availability (the visit effect), but it also means that a smaller fraction of customers subscribe to the deal, which implies that the deal is active only when demand is higher, which in turn implies lower availability (the subscription effect). The overall effect is therefore ambiguous.

However, we can show that there exists a unique solution to Eq. 15.5 when the announced price of the deal, r_t , is higher than a certain level \bar{r} . To understand the

drivers of this effect, it is instructive to rewrite Eq. 15.5 in terms of a comparison between relative availability and relative surplus in the two periods:

$$\frac{\hat{v}_t - r_t}{v_h - r_h} = \frac{\int_{n_t \bar{H}(\hat{v}_t)^{-1}}^{+\infty} \min(1, k(H(\hat{v}_t)x)^{-1}) dG_c(x)}{\int_{n_t \bar{H}(\hat{v}_t)^{-1}}^{+\infty} \min(1, k(\bar{H}(\hat{v}_t)x)^{-1}) dG_c(x)}. \tag{15.6}$$

The LHS of the rewritten equation is the ratio of the service surplus in the slow period to that in the hot period, whereas the RHS is the ratio of service availability in the hot period to that in the slow period. The ratio of the service surplus (LHS) is always increasing in the threshold, \hat{v}_t . When the deal price r_t is higher than $\bar{r} = H^{-1}(1/2) - v_h + r_h$, a higher fraction of customers visit in the hot period, i.e., $H(\hat{v}_t) \geq 1/2$, which ensures that the ratio of service availability always decreases in the customer threshold, \hat{v}_t . To see why, note that, as before, a higher threshold implies a smaller fraction of visitors in the slow period and a higher fraction in the hot period, thus decreasing the service availability ratio (the visit effect). Also as before, a higher \hat{v}_t implies fewer subscribers, which means that the deal is active only when demand is higher (the subscription effect): however, since $r_t \geq \bar{r} \Leftrightarrow H(\hat{v}_t) \geq 1/2$, this implies that the impact of higher demand in the hot period is more severe than in the slow period. Hence, a price $r_t \geq \bar{r}$ ensures that there exists a unique equilibrium for the customer continuation game. We will show that this is always the case for the full game.

15.3.3.3 Optimal Announcement and Equilibrium Outcome

The service provider chooses the slow period price r_t and the activation threshold n_t that maximize expected profit, taking into account customer best-response strategy $\hat{v}_t(r_t, n_t)$ characterized in Eq. 15.5. The expected profit of the firm is then

$$\begin{aligned} \Pi_t^* = \max_{r_t, n_t} & \left[(r_h - c) \int_0^{n_t/\bar{\alpha}_t} (\min(k, x) - c_F) dG(x) \right. \\ & \left. + \int_{n_t/\bar{\alpha}_t}^{+\infty} (\min(k, \alpha_t x)(r_h - c) + \min(k, \bar{\alpha}_t x)(r_t - c) - 2c_F) dG(x) \right], \end{aligned}$$

where $\alpha_t = H(\hat{v}_t(r_t, n_t))$ and $\bar{\alpha}_t = \bar{H}(\hat{v}_t(r_t, n_t))$ are the fractions of customers that visit during the hot and slow periods, respectively, when the firm announces the deal (r_t, n_t) , and where the dependence of α_t and $\bar{\alpha}_t$ from (r_t, n_t) is omitted to improve readability.

Lemma 3 *The firm can restrict to deal offers with a discounted price higher than \bar{r} without any reduction in his expected profit.*

This result states that the firm needs to consider only announcements with a discounted price higher than \bar{r} , because it is never optimal to discount so much that more than half of the customers visit in the slow period when the deal is active.

This Lemma shows that even though there might be multiple type I equilibria in the customer continuation game that ensues after the deal is announced, there is a unique equilibrium for the full game, because the firm is always better off announcing a deal for which there exists a unique customer best response $\hat{v}_t(r_t, n_t)$. We next compare the profits under the unique equilibrium outcome of the threshold discounting game with those from the traditional approach, that is, with the highest profit between closure and regular discounting.

15.3.4 Comparing Threshold Discounting with the Traditional Approach

Theorem 1 *Threshold discounting leads to higher expected profit than the traditional approach, i.e. $\Pi_t^* > \Pi_a^*$.*

The superior performance of threshold discounting arises from its most characteristic feature, i.e., the activation threshold, which gives rise to two independent sources of advantage: a responsive duality effect and a strategic scarcity effect.¹¹

15.3.4.1 Responsive Duality

Lemma 1 showed that closing in the slow period ends up earning a higher profit than opening and discounting if and only if market size is below a threshold. Unfortunately, a firm considering the traditional approach needs to decide whether to employ regular discounting or seasonal closure ex-ante, without knowing the market size, and the choice that maximizes the expected profit may turn out to be wrong in retrospect once market size is realized and customers visit the firm. With threshold discounting, the firm does not have to trade off the relative strengths of seasonal closure and regular discounting, because he can get the best of both worlds.

An appropriately designed threshold discounting offer allows the firm to ensure that the deal gets activated only when the market size turns out to be above a threshold of his choice (see Lemma 5 in the online compendium). In such a contingency, the firm is balancing demand by effectively imitating the demand-shifting effect of regular discounting. On the other hand, when the market size is below this threshold, the deal is not activated and the service is not offered in the slow period, so that the firm achieves fixed-cost optimization and full margins by effectively using the seasonal closure approach. From Lemma 1, we know that regular discounting is better than closure if and only if the market size is high enough. This means that the activation threshold endows threshold discounting with

¹¹The main results from this section and the following Sects. 15.3.5 and 15.3.6 also appear, or are mentioned, in Marinesi et al. (2017), although the exact formulation and some of the underlying assumptions may differ.

a *responsive duality*, i.e. a built-in, market-responsive dual mechanism that allows the firm to use the information supplied by customers to choose the best demand manipulation technique to employ (closing or regular discounting) an advantage unavailable with the traditional approach.

This responsive duality is not the only advantage of threshold discounting, the benefits go further. Even more interesting is a strategic scarcity effect created by threshold discounting, which allows a firm to better price discriminate strategic customers than does regular discounting, thus improving capacity utilization even further.

15.3.4.2 Strategic Scarcity Effect

Customers strategically think about price and availability and they react differently to a slow-period discount that is active contingent on high market size, as opposed to a discount that is always active. In particular, we find that a discount conditional on a high enough market size – as the one employed by threshold discounting – increases the fraction of demand diverted from the hot to the slow period compared to the same level of a non-contingent discount. We call this observation the *strategic scarcity effect*, which we formalize in the next theorem.

Theorem 2 *For any potentially optimal slow period price $r > \bar{r}$, and for any positive activation threshold $n_t > 0$, threshold discounting diverts more demand from the hot to the slow period compared to regular discounting. Formally, $\forall r > \bar{r}$ and $\forall n_t > 0$, we have that $H(\hat{v}_t(r, n_t)) < H(\hat{v}_d(r))$.*

Remember from Eq. 15.4 that the advantage of discounting the slow period service lies in shifting demand from the hot to the slow period in order to achieve a more equitable allocation of demand across periods – but it comes at the cost of reducing margins. Strategic scarcity is beneficial because it accomplishes the same result as would a additional price reduction in the slow period, that is, diverting more customers to the slow period, but it comes as a free lunch, i.e., the provider enjoys the additional demand shift without paying through higher discounts or lost margins. Put differently, strategic scarcity is beneficial because it magnifies the returns from any discounting level by increasing strategic customers' elasticity to price reductions compared to regular discounting.

The key cause of this effect lies in the difference in service availability between the hot and slow periods under the two discount schemes. Under threshold discounting, the fact that the deal is active signals to the customers that the market size is high enough. This implies that availability will be lower in both time periods, but *more so in the hot period*, making the slow period more desirable to customers.¹² Thus,

¹²A simple example can clarify this property: suppose that capacity is 10, that 60% of customers visit in the hot period, and that the market size is either 10, 20 or 30 with equal odds; then the expected availability of the hot period relative to the slow period is $(1 + 10/12 + 10/18)/(1 + 1 +$

a customer who is indifferent between the two periods under regular discounting is instead willing to visit during the slow period under threshold discounting when the deal is active, because the active deal signals that the market size is higher than average, hence the odds of being served shift further in favor of the slow period. Overall, the higher effectiveness of threshold discounting due to strategic scarcity effectively implies that the service provider can achieve the same level of capacity utilization that a regular discounting strategy would, while keeping higher margins and thus earning a higher profit.

15.3.4.3 A Novel Operational Advantage

To summarize, the advantages of threshold discounting stem from (1) its responsive nature, imitating the fixed cost savings of seasonal closure when market size is low and the demand-balancing effect of regular discounting when market size is high; and (2) increasing customer responsiveness to slow-period discounts by signaling the market size – via the deal activation – which enables the customer to use this information in estimating service availability and self-selecting the consumption period, thus increasing capacity utilization for any discount level offered. Put simply, threshold discounting combines closure with an improved version of regular discounting, and takes the best of each. Interestingly, both the aforementioned advantages rely on information transmission, but while responsive duality exploits the information that customers send to the firm (by choosing whether to subscribe to the deal), strategic scarcity responds to information that the firm sends back to customers (by announcing whether the threshold was reached or not).

As pointed out in the introduction, anecdotal popular press discussions of the benefit of threshold discounts have focused on their network effects and a consequent demand increase. Note that our model deliberately leaves out network effects to focus on operational performance, and our effects stem solely from the better demand-supply matching enabled by threshold discounting. Further, all the results presented above continue to hold even when there are *no economies of scale* ($c_F = 0$). In fact, even when there are no fixed costs, threshold discounting is still better than regular discounting at servicing customers due to the strategic scarcity effect. This suggests that these innovative and profit-enhancing schemes do not need to be the exclusive prerogative of high-volume businesses, but can instead be employed by small businesses – such as those featured by Groupon and its competitors – with equally beneficial results.

$10/12) = 43/51 \simeq 0.84$ over all market states, $(10/12 + 10/18)/(1 + 10/12) = 25/33 \simeq 0.76$ over the two higher states, and $(10/18)/(10/12) = 2/3 \simeq 0.67$ for the highest state; that is, the expected service availability of the hot period relative to the slow period decreases as we consider only increasingly higher states – as strategic customers do when they learn that the activation threshold has been reached.

15.3.5 Impact of Strategic Customers on Threshold Discounting Performance

Our analysis has so far assumed that all customers are strategic, in the sense that they all account for other customers’ subscription and visit responses to the discounting scheme offered by the firm when they make their decisions. Arguably, not all customers are sophisticated enough to do this: Li et al. (2014), for example, estimate the percentage of strategic consumers in the airline industry to be between 5.2% and 19.2%. In this section, we extend our analysis to consider a mixed population in which a fraction γ of customers are strategic, and the remaining fraction $1 - \gamma$ are nonstrategic – they do not account for the decisions of other customers. This means that in making her decision, a nonstrategic customer naively ignores both the odds of the deal being active and the expected availability in each service period, since these depend respectively on the subscription and visit strategies of the other customers. A nonstrategic customer subscribes/visits in the slow period iff her service surplus is higher than in the hot period, i.e. iff $v_s - r_t > v_h - r_h$, where v_s is her slow-period valuation. The profit of threshold discounting when only a fraction γ of the population is strategic is given by (the profit expression for the traditional approach is easily updated following the same logic)

$$\begin{aligned} \Pi_t^{\gamma*} = \max_{r_t, n_t} & \left[(r_h - c) \int_0^{n_t \bar{\alpha}_{t,\gamma}(r_t, n_t)^{-1}} (\min(k, x) - c_F) dG(x) \right. \\ & + \int_{n_t \bar{\alpha}_{t,\gamma}(r_t, n_t)^{-1}}^{+\infty} [\min(k, \alpha_{t,\gamma}(r_t, n_t)x)(r_h - c) \\ & \left. + \min(k, \bar{\alpha}_{t,\gamma}(r_t, n_t)x)(r_t - c) - 2c_F] dG(x) \right], \end{aligned} \tag{15.7}$$

where $\alpha_{t,\gamma} = \gamma H(\hat{v}_t(r_t, n_t, \gamma)) + (1 - \gamma)H(v_h - r_h + r_t)$ is the fraction of the population that in equilibrium visits the firm during the slow period when the deal is active, given by the mix of strategic and nonstrategic customers, $\bar{\alpha}_{t,\gamma} = 1 - \alpha_{t,\gamma}$, and where $\hat{v}_t(r_t, n_t, \gamma)$ is defined as in Eq. 15.6, with $H(\hat{v}_t)$ and $\bar{H}(\hat{v}_t)$ being replaced by $\alpha_{t,\gamma}$ and $\bar{\alpha}_{t,\gamma}$ respectively.

We now reevaluate the superiority of threshold discounting when nonstrategic customers are also present in the population.

Theorem 3 *Threshold discounting outperforms the traditional approach for any composition of strategic and nonstrategic customers in the population, including when the population comprises entirely of nonstrategic customers.*

As explained above, the advantage of threshold discounting is driven both by its ability to mimic closure and regular discounting when most appropriate, as well as from the strategic scarcity effect it creates. While the strategic scarcity effect relies on customers’ ability to account for the decisions of other customers when making

their decisions, the responsive duality advantage exploits the information contained in the number of subscribers that does *not* require customers to be strategic, but rather to signal if they are planning to visit during the slow period. Thus, even when there are no strategic customers in the population, the operational advantages of threshold discounting persist.

Next, we study the impact that the proportion of strategic customers in the population has on the profits of a service provider employing threshold discounting. Most of the existing literature on strategic customers (Su and Zhang 2008; Liu and van Ryzin 2008; Cachon and Swinney 2009, 2011) has either proven that strategic customers are a threat to a firm's profit, or has taken it as granted and developed countermeasures to reduce their harmful effect.¹³ The typical setting often evoked is one in which an apparel retailer sells a finite inventory over a finite season, and may resort to price markdowns at the end of the season in order to dispose of leftover inventory. By anticipating price markdowns, strategic customers can decide to postpone their purchases until the end of the season, thus buying at a discount and reducing profits for the firm. Our setting shares several characteristics with this typical setting. In Cachon and Swinney (2009), for instance, strategic customers can decide to purchase in two different periods – during the season, when their valuation for the product is higher, or at the end of the season, when their valuation is lower – which maps exactly to the hot and slow periods in our framework. As in our paper, in Cachon and Swinney (2009) the firm offers a reduced price in the period that customers value the least. Finally, as in our paper, strategic customers take into account the actions of other customers and act to maximize their expected surplus. Despite these similarities, the effects of strategic customers in our setting are in stark contrast with those in the classic settings studied in the literature.

Theorem 4 *The profits under threshold discounting are higher with more strategic customers in the population. Formally, $\Pi_i^{\gamma_1^*} > \Pi_i^{\gamma_2^*}$ for every $\gamma_1, \gamma_2 \in [0, 1]$ such that $\gamma_1 > \gamma_2$.*

Strategic customers differ from nonstrategic ones in that, by accounting for the actions of the other players, they can better account for future prices and availability, and act accordingly. In the classic setting, this leads strategic customers to wait for otherwise unanticipated price markdowns, and this is always harmful for the firm. In our setting, strategic behavior has different implications. First, strategic customers account for the visit decision of the other customers, which allows them to form expectations on the service availability of each period, accounting for the odds

¹³There are two exceptions. One is the empirical work by Li et al. (2014), which argues that if, on the one hand, strategic customers reduce margins, on the other hand they increase demand, either by forcing the firm to reduce prices, which raises demand in itself, or by making consumers postpone purchases, thus having a second purchasing opportunity. As a result, the effect on profit may go either way. The second exception is the working paper Chun and Ovchinnikov (2017), who show that airlines can use loyalty programs to exploit customers' strategic behavior, inducing them to fly more than they need, thereby increasing demand and profit. Our result arises due to very different underlying dynamics.

of getting a unit of service before they visit, which is in the interest of the firm. Second, they also account for the subscription decision of other customers, which allows them to refine their expectation on service availability upon knowing that the deal is active (strategic scarcity effect) which also goes in the interest of the firm, as already discussed. In our context, there is no difference in how strategic and nonstrategic customers account for price reductions, since the firm clearly announces them upfront before the subscription stage – and with good reason, as discussed below. Hence, the sophisticated decision process of strategic customers always has a beneficial impact for the firm.

It should be noted that the firm's commitment to a price reduction has nothing in common with the use of price commitment strategies as a countermeasure to strategic customers, as in Su and Zhang (2008). In their setting, the firm commits to high enough prices at the end of the season to deter strategic customers from purchasing at the end of the season, i.e., in the "slow" period. In our setting, the firm announces price reductions to achieve the opposite effect, i.e., to induce customers to purchase in the slow period. The difference arises because they consider a firm selling inventory of a durable good, while we consider a service firm selling capacity. For a firm selling a physical product, a customer who decides to purchase in the low season rather than in the high season is always harmful, because it reduces margins: hence, the firm commits to high prices in the low season to prevent such behavior from occurring. For a service firm selling capacity, a customer who decides to purchase in the slow period rather than in the hot period may instead be beneficial, because it increases sales whenever capacity in the hot period is sold out, but there is still spare capacity in the slow period: hence, the firm commits to (appropriate) price reductions in the slow period to incentivize such behavior. Basically, the perishable nature of capacity transforms strategic customer's intertemporal purchasing decisions from a threat to margins into an opportunity to increase capacity utilization and sales.

15.3.6 Mediated Threshold Discounting

In the most popular implementations of threshold discounting, the service provider offers threshold discounts through an intermediary (such as Groupon), which features the deal on its website in exchange for a commission. The main advantage of going through a third party is to reach a larger number of customers: in this case, threshold discounting can generate word-of-mouth effects, as Jing and Xie (2011) analyze. From the operational perspective of our study, however, an intermediary provides no clear advantage to the firm, though the need for an intermediary may still arise as a way for a firm to obtain the necessary visibility and reach his customers, possibly because customers are not aware of the firm's website. If threshold discounting is offered through an intermediary, decision rights are a key consideration: who decides on the characteristics of the deal (the activation threshold and the discounted price), the service provider or the intermediary? If it

is the service provider, then the intermediary is simply an extra cost, and threshold discounting is preferable to the traditional approach only insofar as the advantages outweigh the intermediation costs. In this case our analysis above applies with the cost of intermediation subtracted from the service provider's profit.

In practice, however, the intermediary has a large role in shaping the characteristics of the deal, because of the inexperience of the service provider, for example, or because of its high bargaining power – the intermediary may be a local monopolist, as Groupon was before the emergence of competition. In these cases, it is imperative to learn how the incentives of the intermediary differ from the those of the service provider. Based on our interactions with Groupon management, the contract arrangement most often used in practice, possibly due to its simplicity, observability and objectivity, is such that the intermediary earns a percentage of the revenues from those customers that subscribed to the offer through its website. The profit functions for the service provider under mediated threshold discounting (superscript *med*) and the profit of the intermediary (superscript *in*) as a function of the deal parameters, for any positive intermediation fee η , are then given by

$$\begin{aligned} \Pi_t^{\text{med}}(r_t, n_t | \eta) &= (r_h - c) \int_0^{n_t \bar{\alpha}_{t,\gamma}^{-1}} [\min(k, x) - c_F] dG(x) \\ &+ \int_{n_t \bar{\alpha}_{t,\gamma}^{-1}}^{+\infty} [\min(k, \alpha_{t,\gamma} x) (r_h - c) + \min(k, \bar{\alpha}_{t,\gamma} x) (\bar{\eta} r_t - c) - 2c_F] dG(x) \\ \Pi_t^{\text{in}}(r_t, n_t | \eta) &= \eta r_t \int_{n_t \bar{\alpha}_{t,\gamma}^{-1}}^{+\infty} \min(k, \bar{\alpha}_{t,\gamma} x) dG(x). \end{aligned}$$

with $\bar{\eta} = 1 - \eta$, $\alpha_{t,\gamma}$ and $\bar{\alpha}_{t,\gamma}$ defined as in Eq. 15.7 and their dependence from (r_t, n_t) being omitted for brevity.

Theorem 5 *Under mediated threshold discounting*

- the intermediary chooses a lower slow period price than the firm would, for any given activation threshold n ; formally $r_t^{\text{in}}(n) < r_t^{\text{med}}(n) \forall n > 0$;
- the intermediary chooses a lower activation threshold than the firm would, for any given slow period price r , if the fraction of strategic customers is not excessive; formally $\forall r_t \exists \bar{\gamma}(r_t) \in (0, 1]$: $\gamma \leq \bar{\gamma}(r_t) \Rightarrow n_t^{\text{in}}(r_t) < n_t(r_t)$; and
- the service provider earns a lower profit, even when the intermediation fee is negligible, i.e., when $\eta \rightarrow 0^+$,

where

$$\begin{aligned} r_t^{\text{in}}(n) &= \arg \max_{r_t} \Pi_t^{\text{in}}(r_t, n | \eta) \quad \text{s.t. } r_t \leq r_h, \\ n_t^{\text{in}}(r_t) &= \arg \max_{n_t} \Pi_t^{\text{in}}(r_t, n_t | \eta) \quad \text{s.t. } n_t > 0, \\ r_t^{\text{med}}(n_t) &= \arg \max_{r_t} \Pi_t^{\text{med}}(r_t, n_t) \quad \text{s.t. } r_t \leq r_h, \quad \text{and} \\ n_t(r_t) &= \arg \max_{n_t} \Pi_t(r_t, n_t) \quad \text{s.t. } n_t > 0. \end{aligned}$$

The profit of the intermediary differs from the profit of the service provider in three important ways. First, the intermediary earns profit only on customers who purchase during the slow period; second, the intermediary does not incur any additional fixed cost if the service provider opens also in the slow period; and third, the intermediary earns profit only when the deal is active. The first two differences provide strong incentives for the intermediary to charge a *lower slow-period price* than the service provider would. One reason is that the intermediary has much higher incentives to shift demand to the slow period – for he earns nothing when customers purchase on the hot day – and this is best achieved by lowering the price. Another reason is that the intermediary is willing to open during the slow period as long as this brings one cent more in revenues, while the service provider is wary of the fixed costs that such decision brings along.

The second and third differences imply instead that, compared to the service provider, the intermediary prefers a deal that is much more likely to be active, meaning a *lower activation threshold*. The reason is that the intermediary takes all the benefits of an active deal – higher revenues during the slow period – without getting most of the costs associated with it – costs of opening, since these are incurred by the service provider, and cost due to the cannibalization of the hot period sales by the slow period, since the intermediary gains nothing from selling during the hot period. The only cost for the intermediary in lowering the activation threshold comes from reducing the strategic scarcity effect – a lower threshold sends a weaker signal to strategic customers upon deal activation – which results in lower sales in the slow period.¹⁴ However, this cost is often negligible. To see why, one must consider the interaction between the two effects in the Theorem 5: once the intermediary lowers the price during the slow period, demand will further shift to the hot period; this demand shift weakens the strategic scarcity effect, which is based on the difference in availability between the two periods, and in how signaling a high market size via the deal activation makes such a difference more prominent in the eyes of the customers. Once the additional price reduction favored by the intermediary has weakened the strategic scarcity effect, lowering the threshold is going to have little consequences on the slow period sales.

In summary, the intermediary is better off with a lower activation threshold and a lower price, both of which undermine the advantages of threshold discounting for the service provider. In this case responsive duality is severely diminished, both because the deal would be activated in market states in which it would be best not to activate the deal, and because an excessive fraction of demand would be redirected to the slow period, reducing the operational benefit of price discriminating between periods. Further, the strategic scarcity effect would also be reduced on account of the lower threshold. This logic indicates that the deal preferred by the intermediary, one

¹⁴When the hot period is busier than the slow period; otherwise, reducing the strategic effect increases sales in the slow period, and the intermediary always prefers a lower threshold than the service provider does.

with a very low – if not zero – threshold and a deep discount, could substantially reduce the profit of the service provider.¹⁵ We therefore conjecture that threshold discounts work well when administered directly by the service provider, but not when administered by an intermediary.

15.3.7 Design Considerations in Threshold Discounting Offers

The above analysis examined one particular design of threshold discounting, that in which the firm pre-commits to the activation threshold for the deal, if the deal is active or not is announced before the beginning of both time periods, and the discount can be used only during the slow period. In practice, we encountered numerous variations of this basic setup and, at different points of time, Groupon experimented with other arrangements. Hence, in this section we examine alternative designs and compare them with the original design in Sect. 15.3.3, henceforth referred to as *classic* threshold discounting. We consider a population that comprises only of strategic customers since this simplifies the exposition and it has no impact on the results.

15.3.7.1 Opaque Activation Rule

In classic threshold discounting, the firm commits to a discounted price and an activation threshold before customers make the decision to subscribe or not: this commitment ties the service provider's hands, forcing him to abide by a specific activation rule. A potentially better design is one in which the provider does not publicly commit to a decision rule for activation, and instead makes the activation decision after he observes how many customers have subscribed: it is in fact well-known that postponing a decision to a later time is beneficial if this allows the acquisition of new information that is relevant for that decision – as in this case, where subscriptions contain new information on the market size, which is relevant to making the activation decision. With such a design (Fig. 15.4) the service provider announces the discount price r_o before the customers' subscription decision, yet does not commit to any activation rule. After customers have subscribed, the provider observes the number of subscribers, and only then announces whether or not the deal is active. Designs of this kind are quite common in Customer Voting

¹⁵In an extensive numerical study conducted using ranges of plausible values for all parameters and that simulates more than 800 different scenarios, we find that the intermediary prefers an activation threshold equal to zero, i.e., prefers a deal that is always active, in over 99% of the scenarios. Results from the numerical study have been omitted to keep the exposition short and available from the authors upon request.

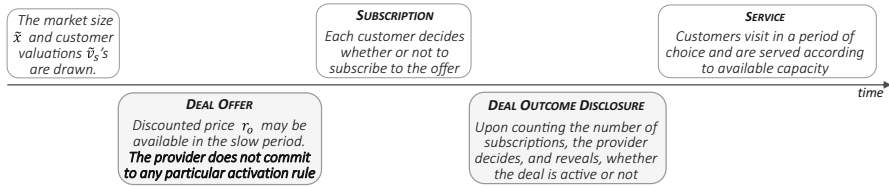


Fig. 15.4 Timeline of threshold discounting with opaque activation rule (*Shaded boxes are used to highlight the main differences relative to classic threshold discounting*)

Systems, whereby customers may be asked to vote for new product designs that could be developed by the firm in the near future, but the firm does not commit to any specific development rule in advance.¹⁶

Theorem 6 *Offering threshold discounts with a committed threshold-activation rule, as in the classic threshold discounting, is better for the firm than offering discounts with an opaque activation rule.*

Postponing the activation decision to a later time is detrimental to the firm, for two reasons. First, by not committing to a specific threshold in advance the firm loses the strategic advantage of being able to use the activation threshold to “steer” customers towards the desired equilibrium. In addition, postponing the activation decision does not provide the firm with any informational advantage, despite the fact that this allows the firm to acquire new information that is relevant to making the deal activation decision. The reason for this unexpected result is that, though relevant, the information contained in the subscriptions always leads to the optimal activation rule being a threshold decision, which the firm can determine already with the information available before the customer subscription stage. Hence, committing to a threshold activation rule upfront provides strategic benefits and no informational disadvantage, and a classic threshold discounting outperforms one with an opaque activation rule.¹⁷

15.3.7.2 Time When the Outcome of the Deal Is Announced

In classic threshold discounting, the service provider releases information about the outcome of the deal, i.e., whether or not it is active, before both time periods begin, allowing customers to make a consumption decision knowing if the threshold has

¹⁶See Marinesi and Girotra (2013).

¹⁷There can be cases in which some additional relevant information is exogenously revealed between the time the deal is announced and the time subscriptions are closed, as uncertainty over weather conditions in the case of an outdoor performance: in these cases, postponing the activation rule may lead to an informational advantage, and which design is better depends on the relative strength of the benefits of commitment and those of postponement.

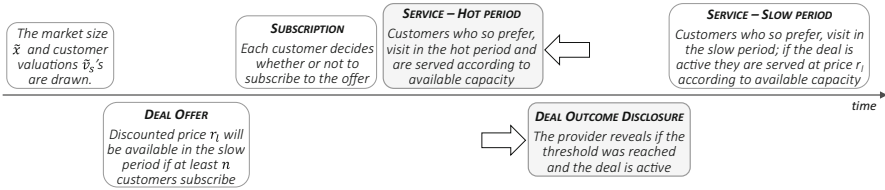


Fig. 15.5 Timeline of threshold discounting with late disclosure (*Shaded boxes are used to highlight the main differences relative to classic threshold discounting*)

been reached. However, in cases in which there is enough time between the hot and slow periods, the service provider could decide to disclose such information after the hot period is over but before the slow period begins.¹⁸ Strategic customers are responsive to price reductions, but also to changes in perceived availability. Liu and van Ryzin (2008) and Yin et al. (2009) have shown how a firm dealing with strategic customers can benefit from increasing the rationing risk they perceive. It is therefore important to study the impact of postponing the deal outcome revelation to customers, since doing so increases the uncertainty – hence the risk – of their subsequent visit decisions, and could therefore lead to a similar effect. The sequence of decisions and information revelation is described in Fig. 15.5. As in classic threshold discounting, the terms of the deal – the discount and the activation threshold – are announced upfront. The only difference is that the outcome disclosure stage now follows the hot period, whereas in the original model it preceded both the hot and slow periods.

Under late disclosure, the profit of the service provider takes the form

$$\Pi_l^* = \max_{r_l, n_l} \left[(r_h - c) \int_0^{n_l/\bar{\alpha}_l} (\min(k, \alpha_l x) - c_F) dG(x) + \int_{n_l/\bar{\alpha}_l}^{+\infty} (\min(k, \alpha_l x)(r_h - c) + \min(k, \bar{\alpha}_l x)(r_l - c) - 2c_F) dG(x) \right],$$

where $\bar{\alpha}_l$ and α_l are the fractions of subscribers and non-subscribers, and are a function of the deal discount, r_l , and activation threshold, n_l .

Theorem 7 *Classic threshold discounting, i.e. with early disclosure, achieves a higher profit for the firm than threshold discounting with late disclosure.*

Unlike in other similar settings, inducing a rationing risk on strategic customers – by postponing the disclosure decision – turns out to be unwise. Late disclosure of the deal outcome has two main implications for the firm. First, it impairs the intertemporal substitutability of demand. Specifically, in the event that the deal is not active, the firm loses sales to those customers who subscribed to the deal and did not

¹⁸If the slow period comes before the hot period, this scheme is obviously not viable.

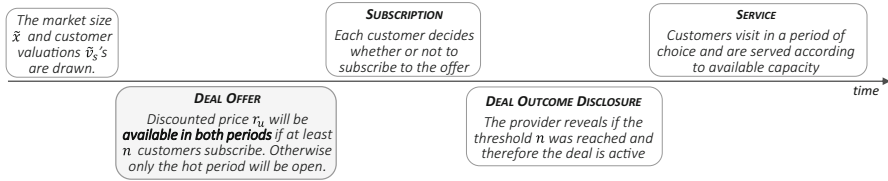


Fig. 15.6 Timeline for unrestricted threshold discounting (Shaded boxes are used to highlight the main differences relative to classic threshold discounting)

visit the service provider during the hot period because they intended to visit during the slow period. The second implication is a consequence of the first, and it is of a strategic nature. Given that subscribing to the deal and waiting for the slow period does not guarantee that the provider will be open at that time, strategic customers are less willing to visit during the slow time period than in the case of early disclosure, i.e. $\alpha_l(r, n) > \alpha_l(r, n)$ for every price $r \geq \bar{r}$ and every activation threshold $n > 0$. This has negative implications for profit, because the service provider needs to offer customers a higher discount for them to visit during the slow time period, further reducing margins. Basically, this strategic implication works in the opposite way of the strategic scarcity effect described in the discussion of Theorem 1, reducing the effectiveness of discounts as inter-temporal demand-balancing devices.

Taken together, the previous results show that providing customers with a transparent activation rule and full and timely information on the activation of the deal makes threshold discounting schemes most potent, or put differently, the less the uncertainty on the customer side, the more effective threshold discounting becomes at increasing capacity utilization and profit.

15.3.7.3 Time Restricted Discounts

While classic threshold discounting restricts the use of the discount to slow periods, discounted offers featured by Groupon and its numerous copycats often place no constraints on the time period of service, i.e., if activated, the discount can be used during hot and slow periods alike. The timeline for these type of deals, henceforth named unrestricted threshold discounting (subscript u), is otherwise the same as for classical threshold discounting, and it is described in Fig. 15.6.

Theorem 8 *Classic threshold discounting achieves a strictly higher profit for the firm than unrestricted threshold discounting.*

Classic threshold discounting is strictly better than unrestricted threshold discounting: by allowing customers to enjoy a reduced price in any period of their choice, the service provider cripples the main advantage of price reductions, that is, the ability to price discriminate between the hot and slow periods in order to redirect some demand to the latter and improve capacity utilization. Despite charging the same price in both periods, unrestricted discounting can still redirect some demand

to the slow period; in fact, a price reduction increases the service surplus in both periods, increasing the surplus loss for a customer from not obtaining a unit of service, thus making customers more willing to visit in the slow period where availability is higher. However, the magnitude of this demand-balancing effect is small compared to what can be achieved using regular discounting. Moreover, the cost associated with a price reduction under unrestricted threshold discounting is much higher than under classic threshold discounting, as the service provider reduces his margin in both time periods. As a result, under unrestricted threshold discounting, price reductions come at a higher cost and yield a smaller operational benefit than classic threshold discounting.

While the overall benefit of unrestricted threshold discounting will ultimately depend on the sum of many effects (see for example Edelman et al. 2016), from a purely operational point of view this design has severely unattractive features, and in many cases a service provider would be better off simply using the traditional approach.¹⁹ This may help explain the oft-repeated assertion that Groupon-like deals were worse for many businesses than just following the traditional approach to managing demand and capacity.²⁰ Perhaps the wide use of unrestricted discounting has been a consequence of the incentive misalignment caused by the commonly employed deal revenue contracts illustrated in Sect. 15.3.6: in fact, under unrestricted threshold discounting, the amount of revenues earned on subscribers is substantially higher compared to classic threshold discounting – and so is the commission earned by the intermediary.

15.3.7.4 Focused Threshold Discounting

One way to potentially improve threshold discounting is to observe that not all customers need to be incentivized to visit during the slow period. Customers with a high enough slow-period valuation value the hot period almost as much as the slow period and prefer to visit the firm during the slow period even when no discount is offered due to higher service availability. If we let n_t^* be the equilibrium activation threshold in classic threshold discounting, then this is true for $v_s \geq \hat{v}_t(r_h, n_t^*)$. This means that classic threshold discounting is inefficient, in that it ends up providing unnecessary monetary incentives to these customers, a source of inefficiency that could be remedied by focusing the incentives only on those customers who actually need them. Next, we explain the intuition behind focused threshold discounting, and then study it formally.

¹⁹Results from an extensive numerical study not reported in this study show that even the traditional approach is better than unrestricted threshold discounting in 90% of our scenarios, resulting on average in a 2% higher profit.

²⁰See for example, “Groupon Was The Single Worst Decision I Have Ever Made As A Business Owner,” Jun 9, 2011, <http://tcrn.ch/2xHeQ3G>

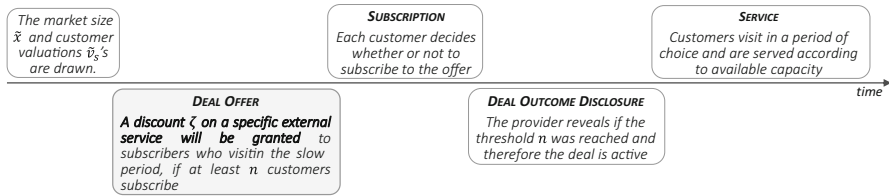


Fig. 15.7 Timeline of a focused threshold discounting (*Shaded boxes are used to highlight the main differences relative to classic threshold discounting*)

Consider an opera house performing *Rigoletto* on Saturday and Sunday nights. Potential customers are comprised of active workers, who prefer Saturday over Sunday – albeit with different degrees of preference – and retired workers, who don’t care about time and therefore prefer Sunday evening due to higher availability. Consider a service desired by active workers but not desired by retired workers, such as baby-sitting. Then a focused threshold discount that offers free baby-sitting service to subscribers for the Sunday night show could redirect the desired number of customers to Sunday without offering unnecessary discounts to retired workers, thus improving profits.

To formalize this intuition, suppose that customers can be divided into two segments, one characterized by strong time preferences ($v_s < \bar{v}_s$) that attach to the external service a positive value $V > 0$, and another with weak time preferences ($v_s \geq \bar{v}_s$) that find no value in the service, and that the frontier valuation \bar{v}_s is high enough. Specifically, let v_e be the value attached by a customer to an external service, where

$$v_e = \begin{cases} V & \text{if } v_s < \bar{v}_s, \\ 0 & \text{otherwise,} \end{cases} \quad \text{with } \bar{v}_s \in [\hat{v}_t(r_h, n_t^*), v_h - \epsilon], \forall \epsilon > 0. \quad (15.8)$$

Focused threshold discounting consists in promising subscribers not a discounted price, but rather a discount on the external service if they visit during the slow period (Fig. 15.7).

Theorem 9 *Under the conditions in Eq. 15.8, focused threshold discounting improves profit for the firm compared to classic threshold discounting.*

Indeed, when customers can be segmented in a way that links their time preferences to their interest for some other service, the firm can employ focused threshold discounting to improve the efficiency of his incentive system, effectively offering an incentive only to those customers who need them, and thus achieving a higher profit.

15.4 Discussion

This paper studies the operational advantages of threshold discounting schemes when used by a capacity-constrained service provider that offers two vertically differentiated services to a random-sized population of strategically-acting customers. We show that threshold discounting outperforms traditional approaches on account of two phenomena: its responsive duality, which allows a firm to match its pricing and closing decisions to different market states, and a strategic scarcity effect, which improves the operational effectiveness of price reductions by signaling lower hot-period availability in high market states to strategic customers. Atypically, the presence of strategic customers increases firm profits in our context. We find that when offered through an intermediary, threshold discounting can lose its effectiveness if the deal specifications are chosen by the intermediary, due to incentive misalignment caused by commonly used contracts. We further expand the understanding of the design of threshold discounting schemes by showing that the optimal design involves a transparent threshold, early deal disclosure, as well as restricted discounting, and we suggest an idea for improving threshold discounting by providing focused incentives to specific consumer segments.

Our model includes assumptions to avoid unnecessary complications to the analysis. We consider customer heterogeneity only with respect to valuation for the slow period, which is rich enough to both model customer preferences as heterogeneous and create vertical differentiation between the two service periods. Nevertheless, our results continue to hold with a more sophisticated bivariate distribution that accounts for heterogeneity of valuations over both time periods. Another assumption we make is that the firm takes pricing and closing decisions only with respect to the slow period. Our results are robust to endogenizing the closing and pricing decision in the hot period; however, making the pricing decision non-trivial requires the bivariate model of customer preferences mentioned above which, coupled with endogenous prices, substantially complicates the exposition of our results. In the classic threshold discounting studied in Sect. 15.3.3, the firm conditions both pricing and opening decisions during the slow period on the number of subscribers. A weaker form of threshold discounting is such that the number of subscribers merely affects pricing, and the firm is always open during both time periods. The results in this case are very similar, because when no discount is offered most customers shun the slow period; specifically, threshold discounting still grants the beneficial effects described in Sect. 15.3.3 and always outperforms regular discounting, but is less effective at managing fixed costs, so that when these are high enough seasonal closure becomes a better choice.

In our analysis, we assume that all customers prefer the hot-period service to the slow-period service; that is, that service periods are vertically differentiated. This may not always be the case, as some consumers may have preferences that differ from the majority. Our results continue to hold under more general preference functions, where each period is preferred by a fraction of customers, except for

the special case in which each period is preferred by exactly half of the consumer population, since in this case there is no need to re-balance demand through the use of discounts.

The sudden rise and (partial) fall of threshold discounting offers that motivated this study exemplify an important lesson for firms who intend to engage with their customers via new technologies: that any approach that leverages customer information to and from the firm and acts upon it, as threshold discounts do, holds the potentials to boost operational performance and profit. But also that, if not understood in all its implications, any such approach can backfire and ultimately cripple the same operational performance that it was meant to improve.

References

- Anand KS, Aron R (2003) Group buying on the web: a comparison of price-discovery mechanisms. *Manag Sci* 49(11):1546–1562
- Aviv Y, Pazgal A (2008) Optimal pricing of seasonal products in the presence of forward-looking consumers. *Manuf Serv Oper Manag* 10(3):339–359
- Boyacı T, Özer Ö (2010) Information acquisition for capacity planning via pricing and advance selling: when to stop and act? *Oper Res* 58(5):1328–1349
- Cachon GP, Swinney R (2009) Purchasing, pricing, and quick response in the presence of strategic consumers. *Manag Sci* 55(3):497–511
- Cachon GP, Swinney R (2011) The value of fast fashion: quick response, enhanced design, and strategic consumer behavior. *Manag Sci* 57(4):778–795
- Chen Y, Zhang T (2014) Interpersonal bundling. *Manag Sci* 61(6):1456–1471
- Chen J, Chen X, Song X (2007) Comparison of the group-buying auction and the fixed pricing mechanism. *Decis Support Syst* 43(2):445–459
- Chun SY, Ovchinnikov A (2017) Strategic consumers, revenue management, and the design of loyalty programs. Research paper No. 2606791, Georgetown McDonough School of Business
- Coase RH (1972) Durability and monopoly. *J Law Econ* 15:143
- Crew MA, Fernando CS, Kleindorfer PR (1995) The theory of peak-load pricing: a survey. *J Regul Econ* 8(3):215–248
- Dana JD Jr (1999) Equilibrium price dispersion under demand uncertainty: the roles of costly capacity and market structure. *RAND J Econ* 30(4):632–660
- Deneckere R, Peck J (1995) Competition over price and service rate when demand is stochastic: a strategic analysis. *RAND J Econ* 26(1):148–162
- Edelman B, Jaffe S, Kominers SD (2016) To Groupon or not to Groupon: the profitability of deep discounts. *Mark Lett* 27(1):39–53
- Fudenberg D, Tirole J (1991) *Game theory*. MIT Press, Cambridge
- Hu M, Shi M, Wu J (2013) Simultaneous vs. sequential group-buying mechanisms. *Manag Sci* 59(12):2805–2822
- Jing X, Xie J (2011) Group buying: a new mechanism for selling through social interactions. *Manag Sci* 57(8):1354–1372
- Li X, Wu L (2018, Forthcoming) Herding and social media word-of-mouth: evidence from Groupon. *MIS Q*
- Li J, Granados N, Netessine S (2014) Are consumers strategic? Structural estimation from the air-travel industry. *Manag Sci* 60(9):2114–2137
- Liu Q, van Ryzin GJ (2008) Strategic capacity rationing to induce early purchases. *Manag Sci* 54(6):1115–1131

- Lus B, Muriel A (2009) Measuring the impact of increased product substitution on pricing and capacity decisions under linear demand models. *Prod Oper Manag* 18(1):95–113
- Marinesi S, Girotra K (2013) Information acquisition through customer voting systems. Working paper No. 2013/99/TOM, INSEAD
- Marinesi S, Girotra K, Netessine S (2017) The operational advantages of threshold discounting offers. *Manag Sci. Articles in Advance*. <https://doi.org/10.1287/mnsc.2017.2740>
- Netessine S, Tang CS (2009) Consumer-driven demand and operations management models: a systematic study of information-technology-enabled sales mechanisms, vol 131. Springer, New York
- Schelling T (1960) *The strategy of conflict*. Harvard University, Cambridge
- Su X (2007) Intertemporal pricing with strategic customer behavior. *Manag Sci* 53(5):726–741
- Su X, Zhang F (2008) Strategic customer behavior, commitment, and supply chain performance. *Manag Sci* 54(10):1759–1773
- Talluri KT, Van Ryzin GJ (2006) *The theory and practice of revenue management*, vol 68. Springer, New York
- Wu J, Shi M, Hu M (2014) Threshold effects in online group buying. *Manag Sci* 61(9):2025–2040
- Yin R, Aviv Y, Pazgal A, Tang CS (2009) Optimal markdown pricing: implications of inventory display formats in the presence of strategic customers. *Manag Sci* 55(8):1391–1408

Chapter 16

Innovation and Crowdsourcing Contests



Laurence Ales, Soo-Haeng Cho, and Ersin Körpeoğlu

Abstract In an innovation contest, an organizer seeks solutions to an innovation-related problem from a group of independent agents. Agents, who can be heterogeneous in their ability levels, exert efforts to improve their solutions, and their solution qualities are uncertain due to the innovation and evaluation processes. In this chapter, we present a general model framework that captures main features of a contest, and encompasses several existing models in the literature. Using this framework, we analyze two important decisions of the organizer: a set of awards that will be distributed to agents and whether to restrict entry to a contest or to run an open contest. We provide a taxonomy of contest literature, and discuss past and current research on innovation contests as well as a set of exciting future research directions.

16.1 Introduction

Everybody has a creative potential and from the moment you can express this creative potential, you can start changing the world. — Paulo Coelho

Our best ideas come from clerks and stockboys. — Sam Walton

Many organizations today look beyond their boundaries to elicit innovation. With advances in information technology and global access to skilled individuals, contests (also known as tournaments) have emerged as a popular and cost-effective tool to elicit innovative solutions to challenging problems. A contest usually starts when a contest organizer announces a problem along with contest rules such as a set of awards (called “award scheme”) and whether the contest is open to the public or

L. Ales (✉) · S.-H. Cho
Tepper School of Business, Carnegie Mellon University, Pittsburgh, PA, USA
e-mail: ales@cmu.edu; soohaeng@andrew.cmu.edu

E. Körpeoğlu
School of Management, University College London, London, UK
e-mail: e.korpeoglu@ucl.ac.uk

not. Then, agents who are interested in the contest make efforts to develop solutions to the problem, and submit them to the organizer. Finally, the organizer evaluates these solutions, and awards the best one(s) according to the announced rule.

This chapter focuses on two popular types of contests: innovation and crowdsourcing contests. In an innovation contest, an organizer seeks solutions to an innovation-related problem from a (not necessarily large) group of agents, and in a crowdsourcing contest, the organizer seeks (not necessarily innovative) solutions from a large group of agents. While pointing out this subtlety, we will refer to both types of contests as “innovation contests” throughout the chapter as these contests mostly overlap in practice. To illustrate how these contests work in practice, consider the following example from Ales et al. (2017a) (see their introduction for detailed examples). Since 2012, Samsung has organized several innovation contests, called Samsung Smart App Challenge, seeking innovative apps for its products. The contest started with Samsung’s announcement of contest rules. For example, Samsung Smart App Challenge 2013 for Galaxy S4 was open to anyone who wished to participate, and distributed a total of \$800,000 prizes for top ten apps. The judging criteria were uniqueness, commercial potential, functionality, usability, and design.

Innovation contests are utilized for a broad set of topics ranging from mining solutions (e.g., Goldcorp Challenge which seeks proposals for identifying potential gold mining targets) to design (e.g., a logo design contest for FIFA World Cup) and software development (e.g., Samsung Smart App Challenge). While some organizations run their own contests, others employ contest platforms such as Challenge.gov, Ennomotive, InnoCentive, Inocrowd, and TopCoder that intermediate contests on behalf of their clients. For example, InnoCentive crowdsources innovation on behalf of a diverse group of clients such as AARP Foundation, Eli Lilly, NASA, and P&G (InnoCentive 2017). InnoCentive organizes ideation, theoretical, and reduction-to-practice (RTP) challenges (in which agents develop ideas, theoretical solutions, and prototypes, respectively) in subject areas such as chemistry and social innovation. Agents of different background compete in these free-entry open innovation contests for awards ranging from \$5,000 to \$1 million. As another example, TopCoder crowdsources software solutions on behalf of a large client base including Best Buy, Comcast, HP, and IBM (TopCoder 2017). Agents around the world compete in various software development contests that are open to public, and winners are awarded cash prizes around \$10,000, and the performance of all participants is converted into a continually updated TopCoder rating.

Through open innovation processes, an organizer can tap into a large number of experts outside of its firm boundary, and can select the most promising solution from many submitted solutions. Despite this benefit of having a large number of participants, the organizer does not need to pay every participant, the organizer may pay only one agent, the “winner,” who has submitted the best solution. In such a winner-take-all contest, all agents except the winner bear all the costs of their efforts. Yet, with many contest participants, agents expect their individual chance of winning a contest to be low, and hence may not have sufficient incentives to exert their best efforts. Thus, the contest organizer should carefully choose the right award scheme, and determine whether to restrict the number of participants to increase the probability of winning for individual agents.

The objective of this chapter is to present a general model framework for innovation contests, and provide insights into two of the organizer's decisions that are essential from both practical and theoretical points of view. The first decision we study is a set of awards through which an organizer incentivizes agents to participate in a contest and make costly efforts. From a practical point of view, it is important for an organizer to assess when to adopt a winner-take-all award scheme and when to offer multiple prizes. From a theoretical point of view, the winner-take-all award scheme is almost a standard assumption in the contest literature, and it is important to determine when this assumption is justified. The second decision we analyze in this chapter is whether it is optimal for a contest organizer to hold an open contest without restricting entries to the contest. From a practical point of view, this analysis addresses when open innovation initiatives that rely on the "wisdom of crowds" are desirable. From a theoretical point of view, several papers in the contest literature implicitly assume that an organizer does not impose any entry restrictions to agents. Thus, it would be interesting to examine when it is indeed optimal for an organizer to hold an open contest. Because this decision of an organizer is closely related to agents' incentives, we will also discuss how participating agents change their level of effort when additional agents enter the contest.

While discussing award scheme and entry restriction decisions based on Ales et al. (2017a,b) and Körpeoğlu and Cho (2017), we present a general model framework that encompasses several models that have been studied in the literature. In particular, our framework captures the following three main features of typical innovation contests in practice:

- When an organizer seeks the best K solutions, where K is a positive integer between one and the total number of participants, Ales et al. (2017a) say that there are K "contributors" among participants in a contest. In some contests, an organizer is interested in only the best solution, so K equals one (cf. Taylor 1995). For example, in a logo design contest for FIFA World Cup, the organizer was interested in finding the best logo to adopt. In other contests, the organizer seeks several good solutions; for example, Samsung sought many useful apps in Smart App Challenge.
- Following Ales et al. (2017b), we consider two sources of uncertainty that agents face. The first source of uncertainty is referred to as "technical uncertainty," and this stochastic element is often modeled as a search process for the best solution from a number of trials (e.g., Dahan and Mendelson 2001). For example, a logo designer may experiment on several logo sketches, and s/he does not know the results of those experiments a priori. The second source of uncertainty, called "taste uncertainty," is due to the subjective or unknown taste of the organizer. For example, in Samsung Smart App Challenge, when submitting their apps, developers do not know how judges will evaluate their apps in subjective criteria such as uniqueness, usability, and design.
- We model agents' heterogeneity utilizing a "productivity-based" model introduced by Körpeoğlu and Cho (2017). In this model, agents are heterogeneous in their productivity levels so that one unit of effort from a high-productivity agent

creates higher value than that from a low-productivity agent. In practice, agents can feature heterogeneous productivity levels due to difference in experience, expertise, and overall ability. For instance, the TopCoder rating of an agent can indicate his/her ability or experience level because high rating indicates successful past performance.

The remainder of this chapter is organized as follows. In Sect. 16.2, we present our general model framework and discuss how this model framework encompasses existing models in the literature. In Sect. 16.3, we provide a brief taxonomy of the literature, and discuss several interesting work in the area of innovation contests. Then, we study the organizer's the award scheme and entry restriction decisions in Sects. 16.4 and 16.5. While we choose to focus on the above two decisions, we discuss some exciting open questions in Sect. 16.6.

16.2 A General Model Framework for Innovation Contests

In this section, we describe a fairly general environment that encompasses commonly used models when studying innovation contests. As discussed in Sect. 16.1, this general model essentially combines the models used in Ales et al. (2017a,b) and Körpeoğlu and Cho (2017). In what follows, we first present our model of agents, and then we present our model of the organizer. At the end of this section, we briefly discuss this model in comparison to other models in the literature.

Agents Suppose that there are N (≥ 2) agents who can participate in the contest. Each participating agent i ($\in \{1, 2, \dots, N\}$) develops a solution to the problem posed in the contest with solution quality (hereinafter “output”) $y_i \in \mathcal{Y} \subseteq \mathbb{R} \cup \{-\infty, \infty\}$. Following Ales et al. (2017b), we represent agent i 's output as a function of improvement effort q_i , a number of trials m_i , trial shocks $(\tilde{\epsilon}_{i1}, \dots, \tilde{\epsilon}_{im_i})$, and a taste shock $\tilde{\epsilon}_i$ as follows:

$$y(q_i, m_i, \tilde{\epsilon}_{i1}, \dots, \tilde{\epsilon}_{im_i}, \tilde{\epsilon}_i) = v(q_i) + \max\{\tilde{\epsilon}_{it}, t = 1, \dots, m_i\} + \tilde{\epsilon}_i. \quad (16.1)$$

This function combines the following three components. First, each agent i can exert “improvement effort” q_i , and this effort leads to a deterministic improvement $v(q_i)$ of agent's output, where v is an increasing and concave function of q_i . Second, each agent i may engage in a trial-and-error process by conducting several experiments, where the agent determines a number of trials m_i (hereinafter, “trial effort”). In each trial t ($= 1, 2, \dots, m_i$), the agent faces uncertainty in the outcome of a trial, which is modeled through a trial shock $\tilde{\epsilon}_{it}$ that follows a Gumbel distribution with $E[\tilde{\epsilon}_{it}] = 0$ and scale parameter μ . (Throughout the chapter, we use the tilde accent to represent random variables.) Each agent observes the outcome of these trials $(\tilde{\epsilon}_{i1}, \dots, \tilde{\epsilon}_{im_i})$ and submits the best one to the organizer. Third, each agent i 's output is subject to the taste of the organizer, which we model by a taste shock $\tilde{\epsilon}_i$. The taste shocks of agents, $\tilde{\epsilon}_i$'s, are independent and identically distributed (i.i.d.) random variables

with a general distribution and $E[\tilde{\varepsilon}_i] = 0$. Unlike trial shocks $\tilde{\varepsilon}_{it}$'s, each agent i is uncertain about the taste shock $\tilde{\varepsilon}_i$ even after the development process is over. For practical examples and details of these components, see Ales et al. (2017b).

We next define a general form for the utility of agent i , $U_a(q_i, m_i, x_i, c_i) : \mathbb{R}_+^4 \rightarrow \mathbb{R}$, which is defined over improvement effort q_i , trial effort m_i , monetary compensation x_i received from the organizer, and heterogenous cost coefficient c_i for exerting effort. The parameter c_i is a privately known cost coefficient for agent i , drawn from a continuous distribution Φ similar to Moldovanu and Sela (2001). The utility of the agent takes the following form:

$$U_a(q_i, m_i, x_i) = x_i - \psi(c_i(\tau_1 q_i + \tau_2 m_i)), \tag{16.2}$$

where $\tau_1 > 0$, $\tau_2 > 0$, and ψ is convex and increasing with $\psi(0) = 0$. This utility function is more general than Ales et al. (2017a,b) which consider identical agents, but it is similar to Körpeoğlu and Cho (2017). We define total effort as $e_i = \tau_1 q_i + \tau_2 m_i$. For example, for an agent with improvement effort q_i and trial effort m_i , total effort e_i may represent the total labor hours an agent spends, where τ_1 and τ_2 are time required for one unit of improvement and trial effort, respectively. Agent i 's cost of making effort e_i is $\psi(e_i) = \psi(c_i(\tau_1 q_i + \tau_2 m_i))$.

The following lemma shows that the output function y given in Eq. 16.1 can be simplified to a new function that depends only on the agent's total effort e_i and aggregate shock $\tilde{\xi}_i$. The first part of the lemma is shown by Ales et al. (2017b) and a special case of the second part of the lemma is shown by Körpeoğlu and Cho (2017) under linear cost of effort and no output uncertainty. We present the proof of Lemma 1(b) in "Appendix".

Lemma 1

- (a) (Lemma 1 of Ales et al. 2017b) *The output function in Eq. 16.1 can be simplified to $y(e_i, \tilde{\xi}_i) = r(e_i) + \tilde{\xi}_i$ in which e_i is the total effort, r is a concave and increasing function, and $\tilde{\xi}_i$ is a random shock that is independent of e_i . For example, if $v(q_i) = \kappa \log(q_i)$ for some $\kappa > 0$, then $r(e_i) = \gamma + \theta \log(e_i)$ where $\theta (> 0)$ and γ are constants.*
- (b) (Adapted from Körpeoğlu and Cho 2017) *The cost-based model in which agents are heterogeneous in their cost coefficients and the output function $y(e_i, \tilde{\xi}_i) = r(e_i) + \tilde{\xi}_i$ is equivalent to a productivity-based model with $y(a_i, e_i, \tilde{\xi}_i) = r(a_i e_i) + \tilde{\xi}_i$, where a_i is agent i 's heterogeneous productivity level drawn from distribution $G(a_i) = 1 - \Phi(1/a_i)$ with support $[\underline{a}, \bar{a}]$, and the cost of effort is $\psi(e_i)$.*

In the rest of this chapter, we use the simplified output function $y(a_i, e_i, \tilde{\xi}_i) = r(a_i e_i) + \tilde{\xi}_i$, and refer to a_i as agent i 's productivity, e_i as agent i 's effort, and $\tilde{\xi}_i$ as agent i 's output shock. This model adds uncertainty to the productivity-based model introduced by Körpeoğlu and Cho (2017), and adds heterogeneous productivity levels to the model of Ales et al. (2017a). The productivity level a_i is drawn from a general distribution G over the support $[\underline{a}, \bar{a}]$. Let $\tilde{a}_{(j)}^N$, $G_{(j)}^N$, and $g_{(j)}^N$ represent

the random variable, the distribution function, and the density function of the j -th highest productivity level among N agents, respectively. It is not difficult to verify that $g_{(j)}^N(a_i) = [N!/((j-1)!(N-j)!)](1-G(a_i))^{j-1}G(a_i)^{N-j}g(a_i)$. The output shock $\tilde{\xi}_i \in \mathcal{E}$ follows cumulative distribution H and density h with $E[\tilde{\xi}_i] = 0$ and $\mathcal{E} = [\underline{s}, \bar{s}]$ where $\underline{s} \in \mathbb{R} \cup \{-\infty\}$ and $\bar{s} \in \mathbb{R} \cup \{\infty\}$. Similarly, let $\tilde{\xi}_{(j)}^N$, $H_{(j)}^N$, and $h_{(j)}^N$ represent the random variable, the distribution function, and the density function for the j -th highest value among N output shocks, respectively.

The Organizer The profit of the organizer, $\hat{\Pi}(Y, X) : \mathcal{Y}^N \times \mathbb{R}^N \rightarrow \mathbb{R}$, is defined over the output vector Y and the compensation vector X . Following Ales et al. (2017a), we consider the case where the organizer benefits from the best $K \in \{1, \dots, N\}$ outputs, and refer to those agents who produce the best K outputs as “contributors.” Formally, we can extend the definition of a contributor in Ales et al. (2017a) (who assume that only the best output is awarded a fixed prize) by utilizing a general compensation vector as follows:

Definition 1 Let $Y^{(K)} = \{y_{(1)}[Y], \dots, y_{(K)}[Y]\}$ where $y_{(j)}[Y]$ represents the j -th highest output in Y – for ease of notation, we use $y_{(j)}$ in short. The organizer’s profit has K contributors if for all $Y \in \mathcal{Y}^N$, $X \in \mathbb{R}_+^N$,

- (i) There exists a continuously differentiable function Π so that $\hat{\Pi}(Y, X) = \Pi(Y^{(K)}, X)$;
- (ii) For all $j = 1, 2, \dots, K$, $\partial \Pi(Y^{(K)}, X) / \partial y_{(j)} > 0$.

A compensation rule $\phi : \mathcal{Y}^N \rightarrow \mathbb{R}^N$ maps the output vector $Y = (y_1, \dots, y_N)$ to a vector of compensations the organizer pays to agents, $X = (x_1, \dots, x_N)$. As in many contests in practice, we restrict attention to the relative (also called ranked-order) compensation rule which compensates each agent based on the agent’s relative rank of the output. Formally, a compensation rule is called the relative compensation rule when there exists some constant $A_{(j)}$ such that $\phi_i(y_{(j)}[Y]) = A_{(j)}$ for all $i \in \mathcal{N}$, $j = \{1, \dots, N\}$ and $Y \in \mathcal{Y}^N$. Thus, the relative compensation rule consists of a vector of N prizes (awards), denoted by $(A_{(1)}, \dots, A_{(N)})$, such that the agent who produces the j -th best output receives a prize of $A_{(j)}$. We refer to this vector of prizes as “award scheme.” Furthermore, we refer to the agent who produces the best output as the winner, and to the award scheme that awards only the winner as the winner-take-all (hereinafter WTA) award scheme.

With K contributors, the organizer’s profit function under the relative compensation rule is:

$$\Pi(Y^{(K)}, (A_{(1)}, A_{(2)}, \dots, A_{(N)})) = \sum_{j=1}^K y_{(j)} - \sum_{j=1}^N A_{(j)}, \quad \forall Y \in \mathcal{Y}. \quad (16.3)$$

Whereas Ales et al. (2017a) consider a general utility function for the organizer that can allow risk aversion and other complex functional forms (see Sect. 5 in Ales et al. 2017a), in this chapter, we restrict attention to a risk-neutral organizer who maximizes profit as in Ales et al. (2017b) and Körpeoğlu and Cho (2017).

We say that the organizer holds an “open contest” when all agents who wish to participate in a contest are allowed to do so. An open contest proceeds in the following sequence. First, the organizer announces the award scheme $(A_{(1)}, A_{(2)}, \dots, A_{(N)})$. Then, each agent $i \in \{1, 2, \dots, N\}$ privately learns a productivity level a_i , and then determines whether to participate in the contest and chooses an effort level e_i . An agent who chooses not to participate receives reservation utility 0. Each agent i who chooses to participate in the contest incurs the cost of effort $\psi(e_i)$. Next, each agent observes an output shock ξ_i , and produces an output $y_i = r(a_i e_i) + \xi_i$. Finally, the contest organizer collects solutions of all participating agents, and gives awards to agents based on the award scheme $(A_{(1)}, A_{(2)}, \dots, A_{(N)})$. We assume that all parameters except the productivity level a_i are common knowledge to both agents and the organizer, and we focus on symmetric pure-strategy Nash equilibria in which all agents with the same productivity level make the same effort.

We next present the agent’s and organizer’s problems. Because Ales et al. (2017a,b) and Körpeoğlu and Cho (2017) consider either agent uncertainty or heterogeneity, we will develop an original formulation that encompasses the formulations in those papers. Let $e^* : [\underline{a}, \bar{a}] \rightarrow \mathbb{R}_+$ denote the equilibrium effort function, where $e^*(a_i)$ corresponds to the equilibrium effort of an agent with productivity level a_i . We first derive $P_{(j)}^N[e_i|e^*, a_i]$, the probability that agent i with productivity a_i and effort e_i has the j -th highest output when all other $N - 1$ agents exert effort based on the equilibrium effort function e^* . Because agent i has no information about productivity levels of other agents, from agent i ’s perspective, another agent k has a random productivity level \tilde{a}_k and a random output shock $\tilde{\xi}_k$, and hence a random output $\tilde{y}_k = r(e^*(\tilde{a}_k)) + \tilde{\xi}_k$. Let F be the distribution function of \tilde{y}_i , and let f be the corresponding density function. It is not difficult to show that the support of \tilde{y}_i is $[\underline{s} + r(\underline{a}e^*(\underline{a})), \bar{s} + r(\bar{a}e^*(\bar{a}))]$ (because it can be shown that $a_i e^*(a_i)$ is increasing in a_i). We can calculate F as follows:

$$F(y_i) = P\{r(a_i e^*(\tilde{a}_i)) + \tilde{\xi}_i \leq y_i\} = \int_{[\underline{a}, \bar{a}]} H(y_i - r(ae^*(a))) g(a) da. \tag{16.4}$$

Let $\tilde{y}_{(j)}^N$ be a random variable with cumulative distribution $F_{(j)}^N$ and density $f_{(j)}^N$ that represents the j -th highest value among N outputs. Conditional on agent i having an output shock realization s , the probability that agent i outperforms agent k by exerting effort e_i is

$$\begin{aligned} P\{r(a_i e_i) + s \geq \tilde{y}_k\} &= P\{r(a_i e_i) + s \geq r(\tilde{a}_k e^*(\tilde{a}_k)) + \tilde{\xi}_k\} \\ &= \int_{[\underline{a}, \bar{a}]} H(r(a_i e_i) + s - r(ae^*(a))) g(a) da. \end{aligned}$$

Thus, we can write the unconditional probability that agent i with productivity level a_i has the j -th highest output among N agents as follows:

$$P_{(j)}^N[e_i|e^*, a_i] = \int_{s \in \Xi} \frac{(N-1)!}{(j-1)!(N-j)!} P\{r(a_i e_i) + s > \tilde{y}_k\}^{N-j} P\{r(a_i e_i) + s < \tilde{y}_k\}^{j-1} h(s) ds, \tag{16.5}$$

because $N - j$ agents are ranked lower than agent i , $j - 1$ agents are ranked higher than agent i , and they can be ordered in $(N - 1)! / ((j - 1)!(N - j)!)$ combinations. The organizer solves the following program:

$$\max_{N \geq K, (A_{(1)}, \dots, A_{(N)})} \Pi = \sum_{j=1}^K \int_{[\underline{s}+r(\underline{a}e^*(\underline{a})), \bar{s}+r(\bar{a}e^*(\bar{a}))]} y f_{(j)}^N(y) dy - \sum_{j=1}^N A_{(j)} \tag{16.6}$$

$$\text{s.t. } \sum_{j=1}^N A_{(j)} P_{(j)}^N[e^*(a_i)|e^*, a_i] - \psi(e^*(a_i)) \geq 0, \quad \forall a_i \in [\underline{a}, \bar{a}] \tag{16.7}$$

$$e^*(a_i) = \arg \max_{e_i \in \mathbb{R}_+} \sum_{j=1}^N P_{(j)}^N[e_i|e^*, a_i] A_{(j)} - \psi(e_i), \quad \forall a_i \in [\underline{a}, \bar{a}]. \tag{16.8}$$

The objective of the organizer given in Eq. 16.6 is to choose N ($\geq K$) and $(A_{(1)}, \dots, A_{(N)})$ that maximize his expected profit. Participation constraint Eq. 16.7 guarantees that each agent receives non-negative from the contest in equilibrium, and hence chooses to participate in the contest. Constraint Eq. 16.8 is the incentive compatibility constraint through which the organizer considers the agent’s utility maximization problem. In this problem, each agent i with productivity a_i chooses an effort e_i that maximizes the expected prize $\sum_{j=1}^N P_{(j)}^N[e_i|e^*, a_i] A_{(j)}$ less the expected cost $\psi(e_i)$, assuming that every other agent chooses an effort based on the function e^* in equilibrium.

Discussion The present model framework encompasses the main features of models used in the innovation contest literature as detailed in Table 16.1. This framework includes both heterogeneous agents and output shocks that affect agents’ outputs as well as the organizer’s payoff. Unfortunately, without making very restrictive assumptions, such a generic model has limited analytical tractability for two reasons. First, $f_{(j)}^N(y_i) = N! / ((j - 1)!(N - j)!) (1 - F(y_i))^{j-1} F(y_i)^{N-j} f(y_i)$ expression in the organizer’s objective Eq. 16.6 is highly complex because it contains the multiplication of N integrals stemming from of the distribution F in Eq. 16.4, and its density f . Second, the distribution F in Eq. 16.4 depends on the equilibrium effort e^* , so one needs to characterize the agent’s equilibrium effort e^* solving the agent’s problem in Eq. 16.8 before optimizing the organizer’s decisions. Yet, in Eq. 16.8, the agent’s probability of attaining rank j , $P_{(j)}^N[e_i|e^*, a_i]$, is highly complex, and so is the system of equations arising from agents’ first-order

Table 16.1 Review of innovation contest literature that use a variant of the present model framework

Paper	Model of uncertainty	Model of heterogeneity	Other features
Terwiesch and Xu (2008)	(i) Trial-and-error projects with no improvement effort and no taste shock; (ii) ideation projects with no trial effort and a Gumbel distributed taste shock	Expertise-based projects with heterogeneous expertise levels and no uncertainty	The organizer's payoff is in the weighted combination of the best output and the average of all outputs
Ales et al. (2017a)	A model that utilizes Lemma 1(a), and assumes a log-concave or increasing density for the output shock	Homogenous agents	A general utility function for the organizer that allows risk aversion and complementarity
Mihm and Schlapp (2017)	Ideation projects with no trial-and-error experiments and uniformly distributed taste shock	Expertise-based heterogeneity in the second period with feedback	A two-period model with two agents where feedback can be given to agents between periods
Nittala and Krishnan (2016)	Ideation projects as in Terwiesch and Xu (2008) with Gumbel distributed taste shock	Homogenous agents	Internal contests where the organizer incurs a cost from agents' efforts; and external contests where there is a risk for linkage of intellectual property
Körpeoğlu and Cho (2017)	No output shock	Productivity-based projects that encompass cost-projects and expertise-based projects	The same profit function for the organizer as T&X; fixed cost of entry when driving equilibrium
Ales et al. (2017b)	Same model of uncertainty with the present chapter	Homogenous agents	No other features
Hu and Wang (2017)	A model that utilizes Lemma 1(a), and assumes a symmetric log-concave density for the output shock	Two agents model where each agent has high ability in exactly one of two attributes	Two attribute model with the option of running one contests per each attribute or a single contest for both attributes
Körpeoğlu et al. (2017)	A model that utilizes Lemma 1(a) and assumes a log-concave density for the output shock	Homogenous agents	Multiple contest organizers and a more general cost function that allows economies of scope across different contests
Stouras et al. (2017)	A common taste shock that does not affect agents' relative rank but impacts their absolute outputs	Heterogeneous expertise levels	Fixed cost of entry leading to uncertain number of participants

conditions. Due to these technical complications, most papers in the literature have chosen one of two pathways: either focus on agents' uncertainty by suppressing their heterogeneous ability levels or focus on agents' heterogeneity by suppressing the uncertainty they face. Accordingly, we use this separation while discussing the literature in the following section, and we analyze these two analytically tractable special cases separately in Sects. 16.4 and 16.5.

16.3 A Brief Taxonomy of Contest Literature

In this section, we briefly discuss contest literature in general, and then discuss the distinguishing factors of innovation contests. Although there is a stream of empirical studies on contests, we restrict attention to theoretical work.

The research on contests is not new. Following the pioneering works of Tullock (1967, 1980) and Lazear and Rosen (1981), contests have been used in various settings such as labor tournaments (e.g., Green and Stokey 1983; Nalebuff and Stiglitz 1983) in which employers aim to incentivize employees to exert more effort, and sales contests (e.g., Kalra and Shi 2001) in which firms elicit effort from salespeople. Several topics have been explored such as the optimal set of awards that an organizer should distribute (e.g., Moldovanu and Sela 2001; Kalra and Shi 2001), the risk-taking behavior of agents in a contest (Hvide 2002), having multiple rounds or a single round in a contest (Moldovanu and Sela 2006), and designing auction-based mechanisms in which heterogeneous agents have different costs (Che and Gale 2003; Siegel 2009). Vojnović (2015) provides a detailed overview of such contests. Different from these classical contests, innovation contests possess two important distinct features: (i) an organizer is interested in only the best solution(s) rather than all solutions (i.e., $K < N$) and (ii) agents' uncertainty impacts an organizer's profit from a contest, and hence the organizer considers agents' uncertainty as well as their effort while determining contest rules.

As discussed in Sect. 16.2, due to tractability issues, the literature on innovation contests has been divided into two streams. The first stream focuses primarily on innovation contests in which agents exert effort or conduct random trials when their outcomes are uncertain, while suppressing agent heterogeneity. Terwiesch and Xu (2008) show that agents' efforts always decrease with more participants but an open contest is always optimal when considering agents' Gumbel distributed shocks. Ales et al. (2017a) show that more agents may lead to increased or decreased effort from agents depending on the distribution H of the output shock, and further show that an open contest is optimal for a general distribution only when the output shock distribution is sufficiently spread out or the organizer seeks many diverse solutions. Meanwhile, Ales et al. (2017b) characterize the optimal set of awards in this environment, and prove that when agents' uncertainty has a log-concave or increasing density function, the winner-take-all award scheme is optimal. Mihm and Schlapp (2017) compare different types of feedback (e.g., public, private, or no feedback) that can be used to improve the contest outcome. Nittala and Krishnan

(2016) compare internal innovation contests within firms, in which the organizer incurs a cost from agents' efforts (as they are employees), with external contests where the organizer utilizes independent agents. Hu and Wang (2017) study a case where the organizer seeks two attributes, and compare running a single contest for both attributes with running two contests – one for each attribute. Körpeoğlu et al. (2017) study multiple contests tackled by the same set of agents, and show that when organizers seek innovative solutions rather than low-novelty tasks, it may be better for organizers to allow agents to freely participate in multiple contests rather than to restrict them to a single contest. They further characterize the optimal number of parallel contests, and show that this optimal number increases with the novelty of the solutions organizers seek.

A second stream of the literature studies contests in which heterogeneous agents compete but with no uncertainty in agents' outputs. These papers build on prior research in economics such as Moldovanu and Sela (2001), who analyze the optimal set of awards in a cost-based model where agents are heterogeneous in their cost of effort. Liu et al. (2007) use a similar model to Moldovanu and Sela (2001) to study prize structure, segmentation, and handicapping in a consumer contest where the organizer aims to stimulate consumption of a good. Terwiesch and Xu (2008) analyze an expertise-based model in which agents are heterogeneous in their initial expertise, and show that an open contest can be optimal under certain conditions. Körpeoğlu and Cho (2017) propose an alternative productivity-based model to unify cost-based model of Moldovanu and Sela (2001) and expertise-based model of Terwiesch and Xu (2008). They show that an agent's equilibrium effort can increase with more participants, and offer a precise explanation to this result by detailing two opposing drivers. Körpeoğlu and Cho (2017) further show that an open contest is more likely to be optimal than what prior studies asserted. Recently, Stouras et al. (2017) analyze how an organizer can promote agents' participation and effort when only a random number of agents participate in the contest because agents incur large fixed costs of entry, which discourage some agents from participating.

Besides the two streams of research on innovation contests discussed above, there are some papers that use a different, more tailored modelling framework to study special types of innovation contests. Taylor (1995) considers a contest among a pool of identical agents, in which each agent conducts random trials until the best output of those trials reaches a pre-determined quality level. Fullerton and McAfee (1999) analyze a contest in which an organizer auctions entry into a contest. Both of these papers show that more agents in a contest leads to a lower equilibrium effort for every agent, but unlike Terwiesch and Xu (2008) and Ales et al. (2017a), these papers conclude that the organizer should restrict entry to the contest. Erat and Krishnan (2012) study design contests in which each agent selects one design approach among a finite set of approaches. Bimpikis et al. (2016) study information extraction and disclosure strategies that keep agents active in dynamic contests.

16.4 Contests with Uncertainty

In this section, we analyze innovation contests where the output uncertainty plays a larger role than the heterogeneity of agents. To implement this, we suppress the agent heterogeneity by setting $\bar{a} = \underline{a} = 1$. With this assumption, in a symmetric equilibrium, each agent exerts the same equilibrium effort $e^*(a_i) = e^*$. In this case, we can simplify agent i 's probability of producing the j -th highest output in Eq. 16.5 as follows:

$$P_{(j)}^N[e_i|e^*] = \int_{s \in \mathcal{E}} \frac{(N-1)!}{(j-1)!(N-j)!} H(s+r(e_i)-r(e^*))^{N-j} \times (1-H(s+r(e_i)-r(e^*)))^{j-1} h(s) ds. \tag{16.9}$$

Then, we can rewrite the organizer problem Eqs. 16.6, 16.7, and 16.8 as follows:

$$\max_{N \geq K, (A_{(1)}, \dots, A_{(N)})} \Pi = Kr(e^*) + E \left[\sum_{j=1}^K \tilde{\xi}_{(j)}^N \right] - \sum_{j=1}^N A_{(j)} \tag{16.10}$$

$$\text{s.t. } \frac{1}{N} \sum_{j=1}^N A_{(j)} \geq \psi(e^*) \tag{16.11}$$

$$e^* = \arg \max_{e_i \in \mathbb{R}_+} \sum_{j=1}^N P_{(j)}^N[e_i|e^*] A_{(j)} - \psi(e_i). \tag{16.12}$$

In Sect. 16.4.1, we analyze the optimal award scheme, and in Sect. 16.4.2 we study the decision of the organizer to restrict entry or not.

16.4.1 Optimal Award Scheme

This section discusses the optimal award scheme based on Sect. 3 of Ales et al. (2017b). As discussed in Sect. 16.2, a tournament organizer determines an award scheme by choosing a set of prizes $(A_{(1)}, A_{(2)}, \dots, A_{(N)})$ for each ranked agent. It is common in the literature to focus on environments where the WTA is used. However, the WTA scheme may not always be optimal. To examine when the WTA scheme is justified, Ales et al. (2017b) derive a necessary and sufficient condition in their Proposition 1 under which the WTA scheme is optimal. Specifically, they link the optimality of the WTA scheme to (i) the distribution of the output shock and (ii) whether the participation constraint Eq. 16.11 is satisfied. Without going into details about this condition, we will discuss when it is violated and when it is satisfied.

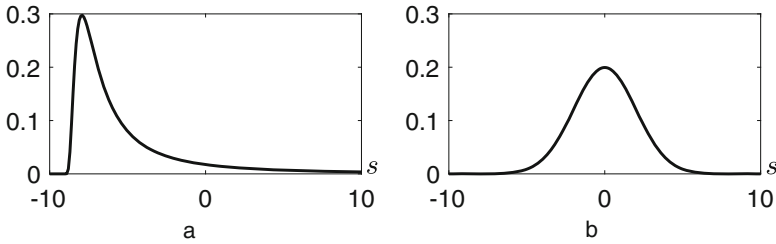


Fig. 16.1 (a) Frechet density with mean 0, shape parameter $\beta = 1.2$, and scale parameter $\mu = 1.5$ (under which the WTA scheme is suboptimal) and (b) Gumbel density with mean 0 and $\mu = 1$, which is log-concave (under which the WTA scheme is optimal)

Proposition 1 (Propositions 2 of Ales et al. 2017b) *For any given A , the winner-takes-all (WTA) award scheme is suboptimal when one of the following conditions is satisfied:*

- (i) $\lim_{s \rightarrow \bar{s}} h(s) = 0$, $\lim_{s \rightarrow \bar{s}} |h'(s)/h(s)| < \infty$, and

$$\int_{s \in \mathcal{E}} [H_{(j)}^N(s) - H_{(1)}^N(s)] \left(\frac{h'(s)}{h(s)} \right)' ds > 0, \tag{16.13}$$

where Eq. 16.13 holds if $h(s)$ is strictly log-convex (i.e., $d^2 \log h(s)/ds^2 > 0, \forall s$).

- (ii) $\frac{A}{N} - \psi \left(\left(\frac{\psi'}{r'} \right)^{-1} \left(A \int_{s \in \mathcal{E}} (N - 1) H(s)^{N-2} h(s)^2 ds \right) \right) < 0$.

We first discuss condition (i), using an example that satisfies this condition. Observe that the density h in Fig. 16.1a features a large highly convex and decreasing region between its peak point and its fat right tail. In this example, an agent’s effort may be more effective in increasing the agent’s probability of attaining some rank $j (> 1)$ than that of becoming the winner. Thus, reducing the winner prize $A_{(1)}$ and increasing award $A_{(j)}$ corresponding to this rank j makes the agent’s effort more effective to win a prize, and hence the agent finds it optimal to increase effort. In practice, this may occur when it is likely that most agents generate low outputs while a few agents generate very high outputs in the contest; for example, when agents’ outputs are evaluated based on popularity among consumers and only few solutions are expected to be extremely popular (e.g., evaluation based on download counts for apps in the 2012 Samsung Smart App Challenge).

Condition (ii) in Proposition 1 specifies when agents do not find it beneficial to participate a contest under the WTA scheme. In this case, the participation constraint Eq. 16.11 is violated under the solution to the agent’s problem in Eq. 16.12, because each agent’s effort in equilibrium is too high in a WTA contest to be justified by the expected winner prize. In this case, the WTA scheme cannot be optimal to the organizer because there is no equilibrium under the WTA. Thus, the organizer may offer multiple awards to strategically reduce agents’ effort in order to guarantee their

participation. Ales et al. (2017b) show that the condition in Proposition 1(ii) holds when agents' uncertainty is sufficiently low and/or their cost function ψ has low convexity. This suggests that, all else being equal, the WTA scheme is more likely to be optimal in a contest that seeks highly technical or innovative solutions that demand more substantial increase in agents' marginal costs of effort (e.g., an RTP challenge at InnoCentive).

Next, we discuss sufficient conditions for $\tilde{\xi}_i$ under which the WTA scheme is optimal.

Proposition 2 (Proposition 3 of Ales et al. 2017b) *Suppose that Eq. 16.11 holds under the WTA award scheme. Then, for any given A , the WTA award scheme is optimal when the density $h(s)$ of the output shock $\tilde{\xi}_i$ is log-concave (i.e., $d^2 \log h(s)/ds^2 \leq 0, \forall s$) or increasing in s .*

According to Proposition 2, the WTA award scheme is optimal when the output shock density is log-concave or increasing. When a density function is log-concave or increasing, no portion of its support is highly convex and decreasing (e.g., see Fig. 16.1b), and hence such a density violates Eq. 16.13 in Proposition 1. In fact, the class of distributions proposed in Proposition 2 is fairly large because many of the commonly-used distributions are either log-concave (e.g., Gumbel, exponential, normal, uniform, and logistic distributions, and Weibull distribution with a shape parameter greater than 1) or increasing (e.g., Weibull distribution with a shape parameter less than or equal to 1). Thus, in practice, the WTA scheme may be appropriate in contests where homogenous agents expect that their outputs will be evenly distributed rather than a few agents generating very high outputs.

16.4.2 Open Innovation and Agents' Incentives

In this section, we build our discussion on Sect. 4 of Ales et al. (2017a) that shows when the organizer should hold an open contest that allows entry of all agents who wish to participate in the contest. The number of participants N directly impacts the organizer's profit $\Pi = Kr(e^*) + E[\sum_{j=1}^K \tilde{\xi}_{(j)}^N] - A$ in two ways. First, N affects the agent's equilibrium effort e^* and hence $Kr(e^*)$, since K is fixed and $r(\cdot)$ is increasing. Second, N affects $E[\sum_{j=1}^K \tilde{\xi}_{(j)}^N]$, which represents the expected value of the best K outcomes from N random shocks. It is easy to see that this term increases with N ($\geq K$) for any K because a more diverse set of solutions increases the expected value of the best K outputs. Thus, for a given award A , depending on how e^* changes with N , Π can increase or decrease with N . When Π is increasing with N , it is optimal for the organizer to choose an open contest. In the remainder of this section, we first study how the agent's equilibrium effort e^* changes with N , and then when the organizer should choose an open contest.

As the number of participants N increases, one may expect that agents would decrease their effort e^* because their individual chance of becoming the winner decreases. Yet, Ales et al. (2017a) show, counter-intuitively, that more participants

do not always induce lower efforts from agents. To discuss this finding, we can derive the equilibrium effort e^* using the condition $\psi'(e^*)/r'(e^*) = AI_N$, where $I_N \equiv \int_{s \in \Xi} (N - 1) H(s)^{N-2} h(s)^2 ds$. Because ψ'/r' is increasing, the effort e^* is increasing (resp., decreasing) in N whenever I_N is increasing (resp., decreasing) in N ; see the following example for illustration.

Example 1

- (i) When $\tilde{\xi}_i$ follows a Weibull distribution with mean 0, shape parameter $\beta = 1$, and scale parameter μ , we have $I_N = (N - 1)/(\mu N)$ increasing in N . Thus, e^* increases with N as well.
- (ii) When $\tilde{\xi}_i$ follows a Gumbel distribution with mean 0 and scale parameter μ , we have $I_N = (N - 1)/(\mu N^2)$. In this case, I_N is decreasing in N , and so is e^* .

Ales et al. (2017a) explain the intuition for why more participants can increase the equilibrium effort e^* by analyzing I_N as follows. From Eq. 16.12, the agent’s marginal benefit of increasing effort is $A(P_{(1)}^N)'[e^*|e^*] = Ar'(e^*)I_N$, and it increases with $(P_{(1)}^N)'[e^*|e^*] = r'(e^*)I_N$, which represents a marginal change of the winning probability with additional effort. Thus, how e^* changes with N depends not on the winning probability but on the marginal impact of additional effort on the winning probability. When $I_{N+1} > I_N$, more intense competition due to a larger number of agents increases the marginal benefit of an agent’s additional effort on the probability of winning. In this case, agents increase effort when faced with more intense competition.

Building on this observation, we next presents a necessary and sufficient condition on the output shock $\tilde{\xi}_i$ under which the equilibrium effort e^* decreases with the number agents N , and presents sufficient conditions under which e^* increases with N .

Proposition 3 (Proposition 1 in Ales et al. 2017a)

- (a) *The equilibrium effort e^* is non-increasing for any $N \geq 2$ if and only if the density $h(s)$ of the output shock $\tilde{\xi}_i$ satisfies*

$$\int_{s \in \Xi} (1 - H(s))H(s)h'(s) ds \leq 0. \tag{16.14}$$

- (b) *Suppose $h(s)$ is increasing in s or the symmetric function of h with respect to y -axis, i.e., $h_r(s) \equiv h(-s)$ for all s , satisfies Eq. 16.14 strictly. Then, e^* is increasing up to some N^* (where $N^* = \infty$ for increasing h).*

Condition Eq. 16.14 in Proposition 3(a) ensures that the density h is sufficiently right-skewed as in Example 1(ii), and this condition is satisfied by any symmetric log-concave density (e.g., normal, logistic) as well as Gumbel and exponential densities. This implies that when agents believe that a bad outcome is at least as likely as a good outcome, they tend to decrease effort with more participants. On the other hand, whenever the necessary and sufficient condition given in Eq. 16.14 is violated, Proposition 3(b) shows that the equilibrium effort e^* is increasing in

N up to some N^* . For example, this condition is violated by a reversed Gumbel distribution or a Weibull distribution. This implies that when agents expect good outcomes with high likelihood, they tend to increase effort with more participants in the contest. This finding is supported from experimental results of List et al. (2014), which demonstrate that in contests with small size, when agents know that they have a high chance of getting favorable outcomes, increasing the number of participants may have positive impact on agents' efforts (see Ales et al. (2017a) for detailed discussion).

We next discuss the findings of Ales et al. (2017a) about when the organizer should hold an open contest that allows all agents who wish to participate in the contest to do so. When the equilibrium effort e^* is increasing in the number of agents N , the organizer's profit increases with N because more participants in the contest also provide a more diverse set of solutions to the organizer (i.e., increases $\sum_{j=1}^K E[\tilde{\xi}_{(j)}^N]$ as discussed above). Thus, it is optimal for the organizer to hold an open contest.

When the equilibrium effort e^* is decreasing in the number of agents N , the organizer's profit may increase or decrease with N , depending on whether the benefit of having a diverse set of solutions outweighs the agents' reduced effort. To quantify the benefit of having a more diverse set of solutions for a general output shock distribution $H(s)$, the notion of a scale transformation is used. When the output shock $\tilde{\xi}_i$ is transformed with scale parameter α , the transformed output shock (i.e., $\hat{\xi}_i = \alpha\tilde{\xi}_i$) has the same mean as $\tilde{\xi}_i$ at 0, and its variance is α^2 times the variance of $\tilde{\xi}_i$. When $\alpha > 1$, $\hat{\xi}_i$ has a larger variance and its density is more spread out. The following proposition of Ales et al. (2017a) shows that when the output shock density h is sufficiently spread out, an open contest is optimal.

Proposition 4 (Proposition 2 of Ales et al. 2017a) *For any distribution H of the output shock $\tilde{\xi}_i$, there exist $\bar{\alpha}$ such that under a scale transformation of $\tilde{\xi}_i$ with $\alpha \geq \bar{\alpha}$, an open contest with unrestricted entry is optimal for any number of contributors K .*

Proposition 4 shows that when the agents' output uncertainty is sufficiently large, an open contest is optimal. In practice, agents can face large uncertainty when the organizer seeks innovative solutions (e.g., writing a software that matches 3D objects with 2D images) rather than low-novelty tasks (e.g., findings bugs in a software). Similarly, how broadly the organizer's problem is defined or how objective the evaluation criteria are can play a role in agents' uncertainty. Overall, Proposition 4 shows that open innovation initiatives are justified when the organizer seeks innovative solutions for broadly defined problems and/or with subjective judging criteria.

Ales et al. (2017a) further show in their Proposition 3 that the threshold scale parameter $\bar{\alpha}$, which is the minimum α required for an open contest, decreases with the number of contributors K . This suggests that an open contest is more likely to be optimal when there are more contributors. This result, in conjunction with Proposition 4, generates insights that are consistent with practice. For example, Samsung Smart App Challenge and Goldcorp Challenge are open contests, probably

because agents face large uncertainty, and anticipate a large number of contributors. On the other hand, in the design contest for the official emblem of the 2014 FIFA World Cup, participating agencies were restricted to 25 (James 2014). Although this contest also involves uncertainty, the restricted entry may be because there is a single contributor.

16.5 Contests with Heterogenous Agents

In this section, we go back to our general model, and analyze contests where agents feature heterogenous productivity levels, while suppressing agents' uncertainty. This model may be suitable for contests in which agents engage in low-novelty tasks, and their ability levels are highly heterogeneous. For ease of illustration, we focus on a case with a single contributor (i.e., $K = 1$) and a linear cost of effort $\psi(e_i) = ce_i$. This model corresponds to a special case of Körpeoğlu and Cho (2017) by assuming that the organizer is interested in only the best solution.

In a symmetric equilibrium, an agent with productivity level a_i chooses an effort level according to the equilibrium effort function $e^*(a_i)$, and creates an output $y^*(a_i)$. In this case, each agent can decide on an output level y_i by choosing an appropriate effort $e_i = r^{-1}(y_i)/a_i$ because $y_i = r(a_i e_i)$. Since agent i does not know other agents' ability levels, the equilibrium output $\tilde{y}^* = y^*(\tilde{a})$ is uncertain, where \tilde{a} is a random variable that represents another agent's unknown productivity level. Assuming that y^* is an increasing function of a productivity level (verified later), we can write the probability that agent i is better than another agent as $P(y_i \geq \tilde{y}^*) = G((y^*)^{-1}(y_i))$. Thus, each agent i 's problem in Eq. 16.8 can be rewritten as

$$\max_{y_i} \left\{ \sum_{j=1}^N A_{(j)} \frac{(N-1)!}{(j-1)!(N-j)!} G((y^*)^{-1}(y_i))^{N-j} (1 - G((y^*)^{-1}(y_i)))^{j-1} - \frac{cr^{-1}(y_i)}{a_i} \right\}. \tag{16.15}$$

In equilibrium, $y_i = y^*(a_i)$ for all agents with productivity a_i . Thus, for agent i to participate, the utility from the contest must be non-negative; i.e.,

$$\sum_{j=1}^N A_{(j)} \frac{(N-1)!}{(j-1)!(N-j)!} G(a_i)^{N-j} (1 - G(a_i))^{j-1} - \frac{cr^{-1}(y^*(a_i))}{a_i} \geq 0. \tag{16.16}$$

Lastly, given the equilibrium effort $e^*(a_i) = r^{-1}(y^*(a_i))/a_i$, the organizer's profit in Eq. 16.6 becomes:

$$\Pi = \int_a^{\bar{a}} r(a_i e^*(a_i)) g_{(1)}^N(a_i) da_i - A. \tag{16.17}$$

In Sect. 16.5.1 we study the optimal award scheme, and in Sect. 16.5.2 we analyze the decision of the organizer to hold an open contest or restrict entry to the contest.

16.5.1 Optimal Award Scheme

In this section, as in Sect. 16.4.1, we discuss when the WTA award scheme is optimal. The result of this section is new, so its proof is presented in “Appendix”. Suppose that the organizer distributes two prizes to the winner and the runner-up with a total prize of A . (The analysis can easily be generalized to multiple prizes.) Let $\alpha \in [0, 0.5]$ be a proportion of the total prize that is awarded to the runner-up. Then., the winner prize $A_{(1)} = (1 - \alpha)A$ and the runner-up prize $A_{(2)} = \alpha A$. To investigate when the WTA (i.e., $\alpha = 0$) is optimal, we use a specific functional form for the effort function $r(e) = \theta(e^{1-b} - 1)/(1 - b)$ (where $b \geq 0$), which is a Constant Relative Risk Aversion (CRRA) function. The CRRA effort function collapses to the linear effort function of Moldovanu and Sela (2001) and Mihm and Schlapp (2017) when $b = 0$ (i.e., $\lim_{b \rightarrow 0} \theta(e^{1-b} - 1)/(1 - b) = \theta e$), and to the logarithmic effort function of Terwiesch and Xu (2008) when $b = 1$ (i.e., $\lim_{b \rightarrow 1} \theta(e^{1-b} - 1)/(1 - b) = \theta \log e$).

Proposition 5 *Let $r(e) = (e^{1-b} - 1)/(1 - b)$. There exists b_0 such that for all $b \leq b_0$, it is optimal for the organizer to set $\alpha = 0$. In contrast, it is optimal for the organizer to set $\alpha > 0$ when*

$$\int_a^{\bar{a}} a_i (a_i e^*(a_i))^{-b} \frac{\partial e^*(a_i)}{\partial \alpha} g_{(1)}^N(a_i) da_i > 0.$$

Proposition 5 shows that the WTA scheme is optimal when the concavity of the effort function (captured in parameter b) is small; this result is also illustrated in Fig. 16.2a, b. In contrast, as Fig. 16.2c depicts, when b is large, the organizer’s profit improves by increasing the weight on the second prize, so the WTA scheme is suboptimal. To understand the intuition behind this result, we need to analyze the derivative of the organizer’s profit with respect to the weight on the second prize α :

$$\frac{\partial \Pi}{\partial \alpha} = \theta \int_a^{\bar{a}} a_i (a_i e^*(a_i))^{-b} \frac{\partial e^*(a_i)}{\partial \alpha} g_{(1)}^N(a_i) da_i. \tag{16.18}$$

As the organizer increases the weight on the second prize, agents with low productivity increase effort (i.e., $\partial e^*(a_i)/\partial \alpha > 0$ for small a_i) and agents with high productivity reduce effort (i.e., $\partial e^*(a_i)/\partial \alpha < 0$ for large a_i). There are two forces that determine whether the former effect or the latter effect dominates. On the one hand, because the organizer is interested in the best output, the organizer has larger weight on the effort of the high-productivity agents than low-productivity agents. On the other hand, because the equilibrium output $y^*(e_i)$ is increasing in

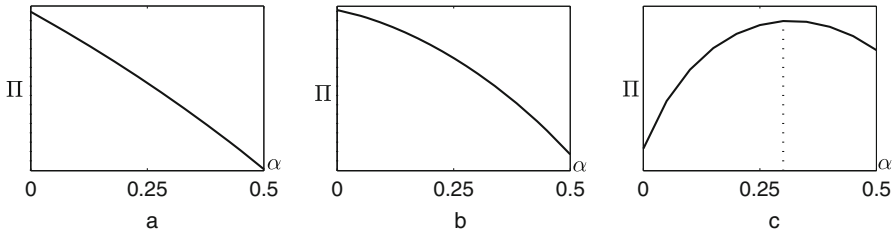


Fig. 16.2 The organizer’s profit Π as a function of the weight on the runner-up prize α when $\tilde{a} \sim \text{Uniform}(0, 1)$, $N = 10$, $r(e) = e^{1-b}/(1 - b)$, $A = 1$, and $c = 0.1$. (a) $b = 0.5$ (b) $b = 1$ (c) $b = 1.5$

a_i , and the effort function r is concave, additional effort by low-productivity agents leads to larger increase in their outputs (i.e., $(a_i e^*(a_i))^{-b}$ is decreasing in a_i). When the effort function is linear or close to linear (i.e., b is small), the negative effect of the second prize (i.e., αA) on the equilibrium effort of high-productivity agents outweighs its positive effect on the equilibrium effort of low-productivity agents, so the WTA scheme is optimal. When the effort function is highly concave (i.e., b is large), additional effort by low-productivity agents leads to significant increase in their outputs so it is optimal for the organizer to offer a second prize.

Proposition 5 is similar to Propositions 2 and 5 of Moldovanu and Sela (2001), who study a cost-based heterogeneity model with $r(e_i) = e_i$, and all contributors (i.e., $K = N$). Their Proposition 2 shows that the WTA scheme is optimal when the cost of effort ψ is linear. Proposition 5 in Moldovanu and Sela (2001) assumes convex ψ , and proposes a necessary and sufficient condition for the WTA scheme to be suboptimal. Our Proposition 5 extends their results to the productivity-based model where the organizer is interested in the best solution (i.e., $K = 1$), and shows that the concavity of the effort function r is another factor that affects the optimality of the WTA scheme.

16.5.2 Open Innovation and Agents’ Incentives

In this section, we discuss when the organizer should allow the entry of all agents who wish to participate in a winner-take-all contest. We build on Sect. 3 of Körpeoğlu and Cho (2017). As in Sect. 16.4.2, we first discuss how the equilibrium output changes with the number of agents in the contest. Then, we present an original result regarding the impact of the number of agents on the organizer. Before describing how the equilibrium effort and output change with the number of agents N , we present the following result from Körpeoğlu and Cho (2017). The lemma characterizes the equilibrium effort and output under the WTA scheme, while generalizing this lemma to the case with multiple awards in Lemma 3 of “Appendix”.

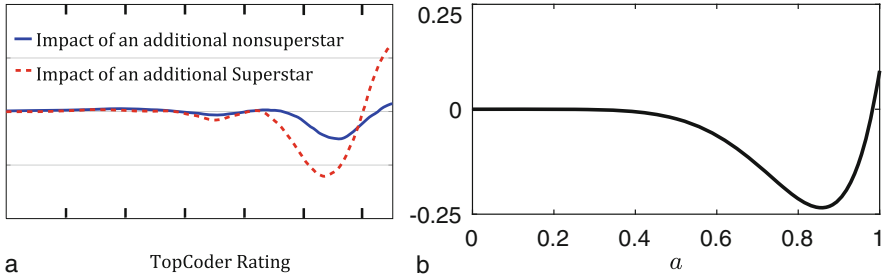


Fig. 16.3 The impact of an additional agent on the agent’s output: (a) empirical observation in Boudreau et al. (2012), and (b) our theoretical prediction of $y^{*,N+1} - y^{*,N}$ when $G \sim \text{Beta}$ with parameters 1 and 0.75; $N = 10$, $r(e) = e^{0.9}/0.9$, $A = 1$, and $c = 0.1$

Lemma 2 (Lemma 1 of Körpeoğlu and Cho 2017) *In a productivity-based project with a general productivity distribution G and a general effort function r , an agent with productivity a_i has equilibrium effort $e^*(a_i)$ and equilibrium output $y^*(a_i)$, where*

$$e^*(a_i) = \frac{A}{ca_i} \int_a^{a_i} ag_{(1)}^{N-1}(a) da \quad \text{and} \quad y^*(a_i) = r\left(\frac{A}{c} \int_a^{a_i} ag_{(1)}^{N-1}(a) da\right).$$

We discuss how the agent’s output y^* and effort e^* change with the number of agents N , by contrasting the implications of our model with empirical observations. Let $y^{*,N}$ and $e^{*,N}$ denote the agent’s output and effort, respectively, when there are N agents in the contest. Figure 16.3a, adapted from Fig. 7 of Boudreau et al. (2012), depicts how the agent’s output changes with an additional high-ability “superstar” (dotted curve) or an additional lower-ability “non-superstar” (normal curve) in software development contests organized by TopCoder. For both cases, an additional agent has a *minimal* effect on low-ability agents with TopCoder rating less than 2000, whereas it has a *negative* effect on moderate-ability agents with TopCoder rating between 2000 and 2400, and it has a *positive* effect on high-ability agents with TopCoder rating over 2400. To compare such empirical observation with our theoretical prediction, we illustrate the impact of an additional agent on the output of agents with different productivity levels in Fig. 16.3b by plotting $y^{*,N+1}(a_i) - y^{*,N}(a_i)$ over a_i . One can clearly see that the patterns in Fig. 16.3a are strikingly similar to those in Fig. 16.3b.

In order to identify the factors that derive the patterns in Fig. 16.3a, b, we utilize the findings of Körpeoğlu and Cho (2017). Specifically, substituting the expression of $e^{*,N}$ in their Eq. 6 into $y_i^{*,N} = r(a_i e^{*,N})$, we can write $y^{*,N}$ as follows:

$$y^{*,N}(a_i) = r\left(\frac{A}{c} G_{(1)}^{N-1}(a_i) E[\tilde{a}_{(1)}^{N-1} | \tilde{a}_{(1)}^{N-1} < a_i]\right). \tag{16.19}$$

In Eq. 16.19, there are two opposing forces that influence agent i 's equilibrium output with an increase of N . First, a higher N reduces agent i 's probability of winning the contest, which corresponds to the probability of having a higher productivity than all other agents; i.e., $P(\tilde{a}_{(1)}^{N-1} < a_i) = G_{(1)}^{N-1}(a_i)$, decreases with N . Second, Körpeoğlu and Cho (2017) show in their Proposition 1 that a larger N raises the expected productivity of the runner-up, given that agent i is the winner, $E[\tilde{a}_{(1)}^{N-1} | \tilde{a}_{(1)}^{N-1} < a_i]$. This second effect creates positive incentives for some agents to exert higher effort and improve output in order to win the contest. Depending on which of these two opposing forces dominates, agent i may generate a better or worse output $y^{*,N}(a_i)$. Low-ability agents are hardly affected by increased competition because they already exert minimal effort due to low chances of winning. Moderate-ability agents tend to have lower effort and hence worse outputs because the impact of increased competition on their winning probability (i.e., $G_{(1)}^N(a_i)$) is dominant. High-ability superstars, who have higher winning probabilities, tend to increase effort and hence improve outputs because the incentives for exerting higher efforts to win the contest are stronger for them (i.e., an increase of $E[\tilde{a}_{(1)}^{N-1} | \tilde{a}_{(1)}^{N-1} < a_i]$ outweighs a decrease of $G_{(1)}^N(a_i)$).

Finally, we discuss when an open contest is optimal. An open contest is optimal when the organizer's profit increases with additional agents in the contest. Due to agents' heterogeneous response to additional agents in the contest, the organizer faces a trade-off when determining whether to allow more agents in the contest. More agents in the contest induce higher efforts from high-productivity agents, but reduce efforts of moderate-productivity agents. Because the organizer knows only the distribution of productivity levels of agents but does not know their exact productivity levels a priori (for example, it is possible that all agents have moderate productivity), it is not clear whether the organizer should hold an open contest. To illustrate when an open contest is optimal, we consider a special case with a CRRA effort function $r(e) = \theta(e^{1-b} - 1)/(1 - b)$ (with $b \in [0, 1]$) and a generalized beta distribution that encompasses a beta distribution (when $\bar{a} = 1$) with parameters d and 1 including uniform as shown in Fig. 16.4a.

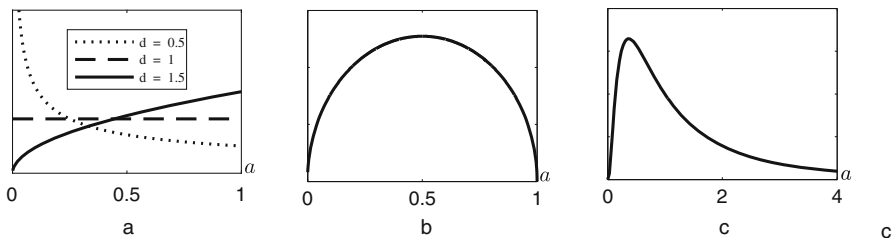


Fig. 16.4 The density $g(a)$ for (a) beta distribution with parameters d and 1 ($\tilde{a} \sim \text{Beta}(d, 1)$), (b) beta distribution with parameters 1.5 and 1.5 ($\tilde{a} \sim \text{Beta}(1.5, 1.5)$), and (c) log-normal distribution with log-scale parameter 0 and shape parameter 1 ($\tilde{a} \sim \log N(0, 1)$)

Proposition 6 *When $r(e) = \theta(e^{1-b} - 1)/(1 - b)$, and $\tilde{a} \in [0, \bar{a}]$ follows $G(a_i) = a_i^d / \bar{a}^d$ (where $b \in [0, 1]$ and $d > 0$), an open contest is optimal.*

To build intuition for Proposition 6, we rewrite the organizer’s profit Π as follows:

$$\Pi = \int_{\underline{a}}^{\bar{a}} y^{*,N}(a_i) g_{(1)}^N(a_i) da_i - A = E[r(\tilde{a}_{(1)}^N e^{*,N}(\tilde{a}_{(1)}^N))] - A. \tag{16.20}$$

The number of agents N has three effects on the organizer’s profit Π . First, an increase in N reduces the equilibrium effort $e^{*,N}$ of moderate-productivity agents (i.e., $e^{*,N}(a_i)$ is decreasing in N for moderate values of a_i). Second, a higher N raises $e^{*,N}$ for high-productivity agents. Third, the productivity level of the highest-productivity agent in the contest, $\tilde{a}_{(1)}^N$, stochastically increases with N (i.e., $\tilde{a}_{(1)}^{N+1}$ first-order stochastically dominates $\tilde{a}_{(1)}^N$). Proposition 6 indicates that the second and third effects outweigh the decreased effort from moderate-productivity agents. Thus, the organizer’s profit increases with the number of agents N , so an open contest is optimal. While Proposition 6 shows the optimality of open contests for the generalized beta distribution, our supplementary numerical analysis verifies that this is also true for various other distributions such as symmetric beta distribution (Fig. 16.4b) or a log-normal distribution (Fig. 16.4c).

16.6 Conclusion and Future Research

Innovation contests are becoming an ever more popular instrument for research and development. This transformation makes the research in the optimal design of contests of first-order importance. This chapter contributes to this research agenda by proposing a general model framework that encompasses commonly used models in the literature, and discussing two of the organizer’s important decisions: How to award agents and whether to allow unrestricted entry to the contest. Our hope is that this chapter can serve as a building block for future contest research, and insights we provide can help both theorists and practitioners.

Research in innovation contests is still relatively young, and there are many interesting open questions. First, prior literature as well as this chapter has assumed a fixed contest duration, but the duration of a contest is also a strategic decision that organizers make in practice. The exploration of the optimal contest duration is an important future research direction. Second, we have adopted a relative compensation rule in awarding agents. Comparison of this compensation rule with other possible compensation rules may shed some light on why the relative compensation rule is so popular in practice. Third, we have considered a case in which the organizer is interested in a fixed number of solutions, and an interesting future research direction is to analyze a case in which the number of solutions organizer utilizes is endogenous to agents’ solution qualities and the cost of implementing those solutions. Finally, characterizing equilibrium under both agent heterogeneity and uncertainty in a general form is an important research to pursue.

Overall, pioneering work in innovation contests has demonstrated that even the questions that have already been studied by prior economics literature can have completely different answers when considering the unique properties of innovation contests such as the impact of agents' uncertainty on the organizer's profit and the fact that the organizer is interested in only the best solution(s). Furthermore, the rapid growth of contest platforms such as InnoCentive poses new questions that were not relevant before. With abundant potential for interesting, practically relevant, and important research questions, innovation contests are an exciting area for future research.

Appendix

Proof of Lemma 1(b) We use superscript P to denote productivity-based model and superscript C to denote cost-based model. Consider the productivity-based model with the output function $y(a_i, e_i, \xi_i) = r(a_i e_i) + \xi_i$, where a_i is a heterogeneous productivity level. Let $v_i = a_i e_i$. Let $v^*(a_i) = a_i e^*(a_i)$ be a best-response function for agent i with productivity a_i , where e^* is the best-response effort. In this model, from agent i 's perspective, another agent's output is a random variable $\tilde{y}^{*,P} \equiv y^{*,P}(\tilde{a}_k) = r(v^{*,P}(\tilde{a}_k)) + \tilde{\xi}_k$. Thus, in a productivity-based model, an agent i solves:

$$\max_{v_i} \sum_{j=1}^N P_{(j)}^N [v_i, v^{*,P}] A_{(j)} - \psi(v_i/a_i), \tag{16.21}$$

where:

$$P_{(j)}^N [v_i, v^{*,P}] = \int_{s \in \mathcal{E}} \frac{(N-1)!}{(j-1)!(N-j)!} P\{r(v_i) + s > \tilde{y}^{*,P}\}^{N-j} P\{r(v_i) + s < \tilde{y}^{*,P}\}^{j-1} h(s) ds.$$

In a cost-based model, all agents except agent i have $v^{*,C}(c_i)$. We will construct a bijective mapping $\eta : R_+ \rightarrow R_+$ from an agent's cost c_i to a productivity a_i (i.e., $\eta(c_i) = a_i$) such that given that all other agents have $v^{*,P}(a_i) = v^{*,C}(c_i)$, agent i will have the same best-response v . Define agent i 's productivity as $a_i = \eta(c_i) = 1/c_i$. Given $v^{*,P}(a_i) = v^{*,C}(c_i)$, another agent's output is the following random variable: $\tilde{y}^{*,C} \equiv r(v^{*,C}(\tilde{c}_j)) + \tilde{\xi}_j = r(v^{*,P}(\tilde{a}_j)) + \tilde{\xi}_j = \tilde{y}^{*,P}$. Then, in a cost-based model,

$$P_{(j)}^N [v_i, v^{*,C}] = \int_{s \in \mathcal{E}} \frac{(N-1)!}{(j-1)!(N-j)!} P\{r(v_i) + s > \tilde{y}^{*,C}\}^{N-j} \times P\{r(v_i) + s < \tilde{y}^{*,C}\}^{j-1} h(s) ds,$$

and hence

$$\begin{aligned} & \arg \max_{v_i} \sum_{j=1}^N P_{(j)}^N [v_i, v^{*,C}] A_{(j)} - \psi(c_i v_i) \\ & = \arg \max_{v_i} \sum_{j=1}^N P_{(j)}^N [v_i, v^{*,P}] A_{(j)} - \psi(v_i/a_i), \end{aligned} \tag{16.22}$$

where the equality follows because $\tilde{v}^{*,C} = \tilde{v}^{*,P}$ and $c_i = \eta^{-1}(a_i) = 1/a_i$. Thus, the agent’s problem in a cost-based model is equivalent to the agent’s problem in a productivity-based model. As a result, given that all other agents have output $v^{*,P}(a_i) = v^{*,C}(c_i)$, by using the mapping η , we obtain the same best response for agent i under both models. Thus, in equilibrium, cost-based and productivity-based models satisfy $v^{*,P}(a_i) = v^{*,P}(\eta(c_i)) = v^{*,C}(c_i) = v^{*,C}(1/a_i)$. Finally, using $\tilde{a} = \eta(\tilde{c}) = 1/\tilde{c}$, we obtain

$$G(a_i) = P(\tilde{a} \leq a_i) = P(1/\tilde{c} \leq a_i) = P(1/a_i \leq \tilde{c}) = 1 - \Phi(1/a_i).$$

□

Lemma 3 *In a productivity-based project with two prizes, a general productivity distribution G and effort function r , an agent with productivity a_i has the following equilibrium effort:*

$$e^*(a_i) = \frac{1}{a_i} \int_a^{a_i} \frac{a}{c} [A_{(1)} g_{(1)}^{N-1}(a) + A_{(2)} (N-1) (g_{(1)}^{N-2}(a) - g_{(1)}^{N-1}(a))] da. \tag{16.23}$$

Proof of Lemma 3 First, suppose that all agents except agent i have output based on the best-response output function $y^*(a_i)$, which is assumed to be continuously differentiable and increasing in the productivity level a_i . We can write the best-response effort as $e^*(a_i) = r^{-1}(y^*(a_i))/a_i$. Output y_i of agent i with productivity level a_i is determined by the following problem:

$$\begin{aligned} & \max_{y_i} \left\{ A_{(1)} G_{(1)}^{N-1}((y^*)^{-1}(y_i)) \right. \\ & \quad \left. + A_{(2)} (N-1) [G_{(1)}^{N-2}((y^*)^{-1}(y_i)) - G_{(1)}^{N-1}((y^*)^{-1}(y_i))] - cr^{-1}(y_i)/a_i \right\}. \end{aligned}$$

The first-order condition when evaluated at $y_i = y^*(a_i)$ gives (note that $y^*(a_i) = r(a_i e^*(a_i))$)

$$\left[A_{(1)} g_{(1)}^{N-1}(a_i) + A_{(2)} (N-1) (g_{(1)}^{N-2}(a_i) - g_{(1)}^{N-1}(a_i)) \right] \frac{1}{(y^*)'(a_i)} - \frac{c}{a_i r'(r^{-1}(y^*(a_i)))}$$

$$\begin{aligned}
 &= \frac{[A_{(1)}g_{(1)}^{N-1}(a_i) + A_{(2)}(N-1)(g_{(1)}^{N-2}(a_i) - g_{(1)}^{N-1}(a_i))]}{r'(a_i e^*(a_i))[a_i(e^*)'(a_i) + e^*(a_i)]} \\
 &\quad - \frac{c}{a_i r'(a_i e^*(a_i))} = 0.
 \end{aligned} \tag{16.24}$$

Multiplying both sides of Eq. 16.24 with $a_i r'(a_i e^*(a_i))[a_i(e^*)'(a_i) + e^*(a_i)]/c$, we obtain

$$\frac{a_i}{c}[A_{(1)}g_{(1)}^{N-1}(a_i) + A_{(2)}(N-1)(g_{(1)}^{N-2}(a_i) - g_{(1)}^{N-1}(a_i))] - [a_i(e^*)'(a_i) + e^*(a_i)] = 0. \tag{16.25}$$

Since $y^*(a_i)$ is increasing, in a contest with $N > 2$, the least productive agent cannot win $A_{(1)}$ or $A_{(2)}$, so exerts zero effort (i.e., $e^*(a) = 0$). Thus,

$$e^*(a_i) = \frac{1}{a_i} \int_a^{a_i} \frac{a}{c} [A_{(1)}g_{(1)}^{N-1}(a) + A_{(2)}(N-1)(g_{(1)}^{N-2}(a) - g_{(1)}^{N-1}(a))] da$$

is the solution to the solution of Eq. 16.25. Therefore, the equilibrium output function $y^*(a_i)$ is

$$y^*(a_i) = r \left(\int_a^{a_i} \frac{a}{c} [A_{(1)}g_{(1)}^{N-1}(a) + A_{(2)}(N-1)(g_{(1)}^{N-2}(a) - g_{(1)}^{N-1}(a))] da \right). \tag{16.26}$$

Finally, we verify that the equilibrium output function $y^*(a_i)$ is continuously differentiable and increasing in a_i . Since all of the terms inside the integral in Eq. 16.26 are continuously differentiable in a_i , and r is continuously differentiable, so is y^* . Taking the derivative of $y^*(a_i)$ with respect to a_i , we obtain $(y^*)'(a_i) = r'(\int_a^{a_i} \phi(a) da) \times \phi(a_i)$, where $\phi(a_i) \equiv (a_i/c)[A_{(1)}g_{(1)}^{N-1}(a_i) + A_{(2)}(N-1)(g_{(1)}^{N-2}(a_i) - g_{(1)}^{N-1}(a_i))]$. Thus, y^* is increasing because $r' > 0$, and $A_{(1)} \geq A_{(2)}$ implies

$$\begin{aligned}
 \phi(a_i) &\geq \frac{a_i}{c} [A_{(2)}g_{(1)}^{N-1}(a_i) + A_{(2)}(N-1)(g_{(1)}^{N-2}(a_i) - g_{(1)}^{N-1}(a_i))] \\
 &= \frac{a_i A_{(2)}}{c} (N-2)(N-1)G(a_i)^{N-3}g(a_i)[1 - G(a_i)] > 0.
 \end{aligned}$$

Then, $y^*(a_i) = r(a_i e^*(a_i))$ is the agent's equilibrium output proposed in the lemma. □

Proof of Proposition 5 The derivative of V with respect to α is

$$\frac{\partial V}{\partial \alpha} = \int_a^{\bar{a}} a_i r'(a_i e^*(a_i)) \frac{\partial e^*(a_i)}{\partial \alpha} g_{(1)}^N(a_i) da_i. \tag{16.27}$$

To evaluate Eq. 16.27, we need the equilibrium effort $e^*(a_i)$ and its derivative with respect to α . If we substitute $A_{(1)} = (1 - \alpha)A$, $A_{(2)} = \alpha A$ into Eq. 16.23, and take derivative of $e^*(a_i)$ with respect to α , we obtain

$$\frac{\partial e^*(a_i)}{\partial \alpha} = \frac{A}{a_i} \int_{\underline{a}}^{a_i} \frac{a}{c} [(N - 1)g_{(1)}^{N-2}(a) - Ng_{(1)}^{N-1}(a)] da.$$

Under CRRA function, noting that $r'(e) = \theta e^{-b}$, Eq. 16.27 becomes

$$\frac{\partial V}{\partial \alpha} = \theta \int_{\underline{a}}^{\bar{a}} a_i (a_i e^*(a_i))^{-b} \frac{\partial e^*(a_i)}{\partial \alpha} g_{(1)}^N(a_i) da_i. \tag{16.28}$$

If $\partial V / \partial \alpha > 0$, it is optimal for the organizer to set $\alpha > 0$ which proves the second part of the proposition.

When $b = 0$ (i.e., $r(e) = \theta e$), Eq. 16.28 becomes

$$\begin{aligned} \frac{\partial V}{\partial \alpha} &= \int_{\underline{a}}^{\bar{a}} \int_{\underline{a}}^{a_i} \frac{A\theta a}{c} [(N - 1)g_{(1)}^{N-2}(a) - Ng_{(1)}^{N-1}(a)] g_{(1)}^N(a_i) da da_i \\ &= \int_{\underline{a}}^{\bar{a}} \frac{A\theta a}{c} [g_{(2)}^{N-1}(a) - g_{(2)}^N(a)] \frac{1 - G_{(1)}^N(a)}{1 - G(a)} da, \end{aligned} \tag{16.29}$$

because

$$\begin{aligned} Ng_{(1)}^{N-1}(a) &= N(N - 1)G(a)^{N-2}g(a) \\ &= \frac{N(N - 1)(1 - G(a))G(a)^{N-2}g(a)}{1 - G(a)} = \frac{g_{(2)}^N(a)}{1 - G(a)}. \end{aligned}$$

Thus,

$$\frac{\partial V}{\partial \alpha} = \frac{A\theta}{c} \left(E \left[\tilde{a}_{(2)}^{N-1} \left(\sum_{j=1}^{N-1} G(\tilde{a}_{(2)}^{N-1})^j \right) \right] - E \left[\tilde{a}_{(2)}^N \left(\sum_{j=1}^{N-1} G(\tilde{a}_{(2)}^N)^j \right) \right] \right) < 0,$$

where the inequality follows because $a(\sum_{j=0}^{N-1} G(a)^j)$ is an increasing function of a , and $\tilde{a}_{(2)}^N$ is larger than $\tilde{a}_{(2)}^{N-1}$ in the sense of first-order stochastic dominance (cf. Theorem 1.A.8 of Shaked and Shanthikumar 2007). Thus, it is optimal for the organizer to set $\alpha = 0$. When $b > 0$, it is not difficult to verify that $\partial V / \partial \alpha$ is continuous in b because all terms in Eq. 16.28 are continuous in b . Then, for sufficiently small b , we have $\partial V / \partial \alpha < 0$. Therefore, there exists $b_0 > 0$ such that for all $b < b_0$, it is optimal for the organizer to set $A_{(1)} = A$ and $A_{(2)} = 0$. \square

Proof of Proposition 6 Suppose that $\tilde{a} \in [0, \bar{a}]$ follows $G(a_i) = a_i^d / \bar{a}^d$. Substituting the effort

$$e^*(a_i) = \frac{Ad(N-1)}{c(d(N-1)+1)} \left(\frac{a_i}{\bar{a}}\right)^{d(N-1)} \quad \text{and} \quad r(e) = \theta \frac{e^{1-b} - 1}{1-b}$$

in $y^*(a_i) = r(a_i e^*(a_i))$ yields

$$y^*(a_i) = \frac{\theta}{1-b} \left(\frac{a_i Ad(N-1)}{c(d(N-1)+1)} \left(\frac{a_i}{\bar{a}}\right)^{d(N-1)} \right)^{1-b} - \frac{1}{1-b}.$$

Noting that $g_{(1)}^N(a_i) = nG(a_i)^{N-1}g(a_i) = (Nd/\bar{a})(a_i/\bar{a})^{d(N-1)+d-1}$, we can express the organizer’s profit as

$$\begin{aligned} \Pi &= \int_0^{\bar{a}} \left[\frac{\bar{a}Ad(N-1)}{c(d(N-1)+1)} \right]^{1-b} \frac{(a_i/\bar{a})^{[d(N-1)+1](1-b)} - 1}{1-b} \\ &\quad \times \left[\frac{Nd}{\bar{a}} \left(\frac{a_i}{\bar{a}}\right)^{d(N-1)+d-1} \right] da_i - A \\ &= \left[\frac{\bar{a}Ad(N-1)}{c(d(N-1)+1)} \right]^{1-b} \frac{d}{1-b} \left[\frac{N}{[d(N-1)+1](1-b)+Nd} \right] - \frac{1}{1-b} - A. \end{aligned}$$

Let

$$W(N) = \frac{\bar{a}d(N-1)}{c(d(N-1)+1)} \left(\frac{dN}{[d(N-1)+1](1-b)+Nd} \right)^{\frac{1}{1-b}}.$$

Noting that $\Pi = (A^{1-b}W(N)^{1-b})/(1-b) - 1/(1-b) - A$ is concave in A , the optimal winner prize is $A^* = W(N)^{(1-b)/b}$. Substituting A^* back to Π , we get

$$\begin{aligned} \Pi &= \frac{W(N)^{(1-b)((1-b)/b+1)}}{1-b} - \frac{1}{1-b} - W(N)^{(1-b)/b} \\ &= W(N)^{(1-b)/b} \left[\frac{1}{1-b} - 1 \right] - \frac{1}{1-b} = \frac{W(N)^{(1-b)/b}}{1-b} b - \frac{1}{1-b}. \end{aligned}$$

Because $W(N)$ is increasing with N , so is Π . □

References

Ales L, Cho S, Körpeoğlu E (2017a) Innovation tournaments with multiple contributors. Working paper, Carnegie Mellon University
 Ales L, Cho S, Körpeoğlu E (2017b) Optimal award scheme in innovation tournaments. *Oper Res* 65(3):693–702
 Bimpikis K, Ehsani S, Mostagir M (2016, Forthcoming) Designing dynamic contests. *Oper Res*

- Boudreau KJ, Helfat CE, Lakhani KR, Menietti M (2012) Field evidence on individual behavior and performance in rank-order tournaments. Working paper
- Che YK, Gale I (2003) Optimal design of research contests. *Am Econ Rev* 93(3):646–671
- Dahan E, Mendelson H (2001) An extreme-value model of concept testing. *Manag Sci* 47(1):102–116
- Erat S, Krishnan V (2012) Managing delegated search over design spaces. *Manag Sci* 58(3):606–623
- Fullerton RL, McAfee RP (1999) Auctioning entry into tournaments. *J Polit Econ* 107(3):573–605
- Green JR, Stokey NL (1983) A comparison of tournaments and contracts. *J Polit Econ* 91(3):349–364
- Hu M, Wang L (2017) Joint vs. separate crowdsourcing contests. Working paper, University of Toronto
- Hvide HK (2002) Tournament rewards and risk taking. *J Labor Econ* 20(4):877–898
- InnoCentive (2017) Corporate info. <http://www.innocentive.com/about-us>. Accessed on 1 May 2017
- James H (2014) Cultural representation in FIFA logos: 1990–2014. <http://www.logocontestreviews.com/cultural-representation-in-fifa-logos/>
- Kalra A, Shi M (2001) Designing optimal sales contests: a theoretical perspective. *Market Sci* 20(2):170–193
- Körpeoğlu E, Cho S (2017) Incentives in contests with heterogeneous solvers. *Manag Sci. Articles in Advance*. <https://doi.org/10.1287/mnsc.2017.2738>
- Körpeoğlu E, Körpeoğlu ÇG, Hafalır IE (2017) Contest among contest organizers. Working paper, University College London
- Lazear EP, Rosen S (1981) Rank-order tournaments as optimum labor contracts. *J Polit Econ* 89(5):841–864
- List JA, Van Soest D, Stoop J, Zhou H (2014) On the role of group size in tournaments: theory and evidence from lab and field experiments. Working paper, NBER
- Liu D, Geng X, Whinston AB (2007) Optimal design of consumer contests. *J Market* 71(4):140–155
- Mihm J, Schlapp J (2017, Forthcoming) Sourcing innovation: on feedback in contests. *Manag Sci*
- Moldovanu B, Sela A (2001) The optimal allocation of prizes in contests. *Am Econ Rev* 91(3):542–558
- Moldovanu B, Sela A (2006) Contest architecture. *J Econ Theory* 126:70–96
- Nalebuff BJ, Stiglitz JE (1983) Prizes and incentives: towards a general theory of compensation and competition. *Bell J Econ* 14(1):21–43
- Nittala L, Krishnan V (2016) Designing internal innovation contests. Working paper, University of California San Diego
- Shaked M, Shanthikumar JG (2007) *Stochastic orders*. Springer, New York
- Siegel R (2009) All-pay contests. *Econometrica* 77(1):71–92
- Stouras KI, Hutchison-Krupat J, Chao RO (2017) Motivating participation and effort in innovation contests. Working paper, University of Virginia
- Taylor CR (1995) Digging for golden carrots: an analysis of research tournaments. *Am Econ Rev* 85(4):872–890
- Terwiesch C, Xu Y (2008) Innovation contests, open innovation, and multiagent problem solving. *Manag Sci* 54(9):1529–1543
- TopCoder (2017) Press releases. <http://www.topcoder.com/press/>. Accessed 19 May 2017
- Tullock G (1967) The welfare costs of tariffs, monopolies, and theft. *Econ Inq* 5(3):224–232
- Tullock G (1980) Efficient rent seeking. In: Buchanan JM, Tollison RD, Tullock G (eds) *Toward a theory of the rent-seeking society*. Texas A&M University Press, College Station
- Vojnović M (2015) *Contest theory*. Cambridge University Press, New York

Part IV
Context-Based Operational Problems in
Sharing Economy

Chapter 17

Models for Effective Deployment and Redistribution of Shared Bicycles with Location Choices



Mabel C. Chou, Qizhang Liu, Chung-Piaw Teo, and Deanna Yeo

Abstract We develop practical OR models to support decision making in the design and management of public car-sharing or bicycle-sharing systems. We develop a network flow model with proportionality constraints to estimate the flow of bicycles within the network, and to estimate the number of trips supported by the system given an initial allocation of bicycles at each station. Furthermore, the number of docks needed at each station, to support the flow, can also be estimated. We also examine the impact of periodic redistribution of bicycles in the network to support more flows, and the location choices of bicycle stations. We conduct our numerical analysis using transit data from the train and bus operators in Singapore. Given that a substantial proportion of the passengers in the train system commute short distance – more than 16% of the passengers alight within 2 stops from the start station – this forms a latent segment of demand for the bicycle-sharing program. We argue that for the bicycle-sharing system to be most effective for this customer segment, the system must deploy the right number of bicycles at the right place, as this affects the utilization rate of the bicycles, how the bicycles circulate within the system, and also the effectiveness of any redistribution strategy. The same approach can be extended to incorporate the issue of station location choices, by incorporating the proportional flow constraints into the MIP formulation. Using a set of bus transit data, we implemented this approach to identify the ideal locations for the bicycle stations in a new town in Singapore, to support the movement of passengers from residential areas to the train station.

M. C. Chou (✉) · Q. Liu · C.-P. Teo
Department of Analytics and Operations, NUS Business School, National University
of Singapore, Singapore, Singapore
e-mail: bizchoum@nus.edu.sg; bizlqz@nus.edu.sg; bizteocp@nus.edu.sg

D. Yeo
Department of Analytics and Operations, NUS Business School, National University
of Singapore, Singapore, Singapore

GE Healthcare, Singapore, Singapore

17.1 Introduction

In recent years, the sharing economy has morphed into a major part of the global economy, impacting various aspects of consumers' lives. Among many other prominent examples, bike sharing which has been around since 1965 has developed dramatically in the 2000s with the introduction of new information technology. The number of bicycle-sharing systems (BSSs) around the world increased from around 200 in 2010 (MetroBike LLC 2011) to approximately 1200 in operation (MetroBike LLC 2017) and more than 350 under construction or being planned for the near future (Russell and DeMaio 2017) at the end of 2016.

With heightened concerns about global oil prices, carbon emissions, and traffic congestion, governments around the world are exploring ways to “nudge” urban residents to commute using public transport instead of private automobiles. Many cities have set up public BSSs to facilitate short trips within the city. As of August 2017, there were approximately 172,700 bikes shared worldwide in 132 cities.

17.1.1 *Review of the Bicycle-Sharing Systems*

Bicycle-Sharing System (BSS) is perceived as a green, healthy and sustainable mode of public transport, which helps decrease greenhouse gas emissions through reducing road congestion and fuel consumption, and improves the first/last mile connection to other modes of transport. The citizens of Amsterdam in the Netherlands initiated arguably the world's first generation of bicycle-sharing program with “White Bicycles” on July 28, 1965. This system operated on a rather ad hoc basis, i.e., one could locate a white bicycle on the street, ride it to the destination, and leave it there for the next trip. Unfortunately, due to theft and abuse, the program only survived for several days. The second generation of the BSS came with improved product and system design features. For example, the bicycle was specifically designed for urban use, its components were not usable on other types of bicycles, and public bicycle-sharing stations were equipped with coin-deposit machines. The citizens of Copenhagen in Denmark launched such a system in 1995. However, this system still faced the problem of bicycle theft since the system was not able to track the identity of the users. This gave rise to the impetus to develop the third generation of BSSs with user tracking abilities – Smart Bicycles, equipped with electronically-locking racks, telecommunication systems, magnetic stripe cards or smart cards, and mobile phone access. The first third generation BSS took off in 2005, with the launch of the Vélo'v in Lyon. In a typical third generation BSS, base stations are located around the city with pre-determined number of bicycles at the beginning of each day. Due to the unbalanced usage of bikes in the system, there is either congestion (dock unavailability) or starvation (bicycle unavailability) of bicycles at the stations each day, which results in a lot of unmet demands (Ghosh et al. 2017). Through analyzing social media data for the BSSs in Spain, Serna et al.

(2017) concluded that 16.4% of demand would be lost due to system unavailability that was caused by uneven usage. As such, managing the redistribution activities to increase system availability is crucial to any successful implementation of such BSSs (The Economist 2011). Given that potential dock unavailability is a major issue, recently a tech-on-bike/dockless system, in which the locking and rental technology is located on the bike itself, was developed.

Bikes are now the new frontier to on-demand transportation. The number of BSSs is expected to continue to grow as new technology is utilized to improve the current systems. Currently, there are a few main types of Bicycle-Sharing Systems (BSSs) in use, namely Ad-hoc, Kiosk-based, Tech-on and Managed Fleet system.

Ad-hoc Bicycle-Sharing System An ad-hoc system involves the operator purchasing and distributing marked bicycles across the community without any locking technology or bike stations. Such systems are usually informal by design and rely on the integrity of individuals to use the bikes in an appropriate manner. The ad-hoc bicycle-sharing system typically works in closed campuses, such as universities, or companies located in Silicon Valley, where the area tends to be large, and the usage of public transport is minimal.

Kiosk-Based System In a kiosk system, bikes are secured to and rented from tech-enabled docking stations. These stations range in sophistication from simple bike racks with key lockboxes to digital automatic locking kiosks with integrated rental systems. Kiosk systems are more common in the public BSSs around the world. Public BSSs can either be introduced under the government, or under joint efforts made by a company and the government. Such examples include Taiwan's YouBike (collaboration between the government and a local company) and Korea's Seoul Bike (managed by the city government).

Tech-on Bike/Dockless System Technology advancement in recent years have enabled tech-on-bike systems, in which the locking and rental technology is located on the bike itself. Riders check out bikes using a smartphone app which allows them to release the lock that secures the bike to the rack. Several Asian countries such as China and Singapore (of0 and Mobike) (Forbes 2017) have adopted this system. A challenge in this sharing system is that it usually results in a large number of discarded bikes. Development of new technology which enables better tracking of the whereabouts of the bikes is needed to handle this problem.

Managed Fleet System In a managed fleet system, bikes are stored in a central location and managed by an employee. Typically, such bikes are equipped with locks, and bike users would have to obtain the keys from the employee at a centralized location, such as a student center on a campus, to unlock the bike for usage.

While Kiosk-Based Systems are still the most commonly used systems around the world at present, future technology development which enables better tracking capability may increase the adoption rate for the Tech-on Bike/Dockless Systems. However, dockless or not, managing the redistribution activities to increase system availability is still crucial since a system without dock unavailability problems is likely to still face bike unavailability problems.

17.1.2 *Research Issues and Structure of the Chapter*

Most of the earlier work on bike redistribution focused on the operational problem of moving bikes using special vehicles deployed for this purpose. Benchimol et al. (2011) studied the static bike repositioning problem (SBRP) and developed a station balancing technique based on the traveling salesman problem (TSP) to reposition bikes by a single vehicle. Raviv et al. (2013) developed four mixed integer linear program formulations to solve large instances with the objective of minimizing the user dissatisfaction in face of stochastic demand. Angeloudis et al. (2014) introduced a novel strategic repositioning algorithm to tackle SBRP, addressing both routing and assignment problems.

Li et al. (2016) considered multiple classes of bikes in SBRP, whereas Kloimüller and Raidl (2017) considered only fully loaded vehicles for movement among the rental stations in SBRP. Schuijbroek et al. (2017) solved SBRP by combining inventory and vehicle routing issues in the BSSs, and proposed a new cluster-first route-second heuristic to search for the optimal solution.

Although the aforementioned techniques are effective in reducing the repositioning cost to some extent, these solutions could not incorporate the real-time demand of the users in their approach. To tackle this issue, Shu et al. (2013) proposed a model on bicycle deployment and flow in BSSs and used the model to address various pertinent issues in managing bicycle-sharing networks. In the rest of this chapter, we address the following questions by introducing the work in Shu et al. (2013) and extending the model and discussion to incorporate location decisions in Sect. 17.2.

- Given the location of the stations, what is the appropriate number of bicycles to deploy in the network? The availability of bicycles affects the number of bicycle trips made and also the bicycle utilization rate. The former measures how much of existing demand can be captured, whereas the latter affects the economic viability of the system. Given the demand pattern, we need an optimal number of bicycles, appropriately located, to make effective use of the resources available to meet demand.
- Impact of redistribution: The flow of the bicycles is dictated by the travel patterns of the commuters. To deal with flow imbalances and to improve bicycle utilization, we may have to do periodic redistribution of bicycles within the system. However, if the bicycles are already heavily (or under) utilized, the periodic redistribution strategy may have only a limited impact on the performance of the system (measured by the additional number of bicycle trips supported). Given the high operational cost associated with redistribution in the BSS, it is thus crucial to estimate the improvement in performance prior to its adoption in actual operation.
- The number of bicycle docks to be installed at each station: To make the bicycle-sharing program implementable, we need to consider how many bicycle docks to install at each station so that commuters can return their bicycles upon arrival at the destination station. Clearly, the number of docks needed at each station

depends on the utilization rate of the bicycles and how the flows are supported in the system, and whether periodic redistribution is used to match supply with demand. In fact, redistribution has the potential to reduce the number of docks needed at each station. The dock design of the bicycle-sharing network is thus intimately tied to the operational decision on bicycle redistribution. It will be useful to have a system to help estimate the number of docks to be installed at each station to support the flow in the system.

In this chapter, we introduce the work in Shu et al. (2013) about a simple proportional network flow model to help to address the above issues. We discuss the theory and intuition for this model in the next section. To validate the findings from the model, a set of commuter data in a Singapore mass rapid transit system to develop the demand model is used for the BSS in Sect. 17.3. By focusing on this segment of the market, and through comparison with extensive simulation results, we demonstrate that the proposed model can be used to approximate the flow of bicycles in the system to a reasonable level of accuracy. In Sect. 17.2.2, we describe how the model can be adopted for general bicycle-sharing network design to incorporate the selection choices of bicycle stations in the network. This leads to a mixed integer proportional network flow model with 0–1 decision variables to denote station choices. We validate the findings in Sect. 17.4 using transit data on bus trips in a new town in Singapore. Finally, we conclude the chapter in Sect. 17.5.

17.2 The Stochastic Network Flow Model

We assume that there are an initial allotment of bicycles at each train station. For each time period, passengers arrive randomly at the station to use the bicycles to travel to their destinations. Data from existing BBSs shows that bike trips are normally within short distance. Our data also shows that in Singapore more than 16% of the train system commuters alight within two stops from the start station. Therefore our study focus on origin-destination demands within two stops in the transit network. The goal is to analyze and estimate the number of such trips that can be supported and substituted by the public bicycle-sharing system, based on the initial allotment of bicycles and the passenger arrival process. This is a technically challenging problem.

Formally, let \mathbf{S} denote the set of stations in the network. In each time period t , the number of passengers who arrive with plan to travel from station i to station j follows a Poisson process, with rate $r_{ij}(t)$. The total number of passengers arriving to use bicycles at station i is thus given by a Poisson process with rate $\sum_{j:j \neq i} r_{ij}(t)$. Within each time period, let $D_{ij}(t)$ and $D_i(t)$ denote the number of arrivals traveling on each link and into each station, respectively. We assume that all rides can be completed within a single time period. Note that bicycles are allocated to the passengers on a first-come-first-serve basis, so that whenever the initial stock of bicycles at a station is depleted, the late comers will not be able to ride to their

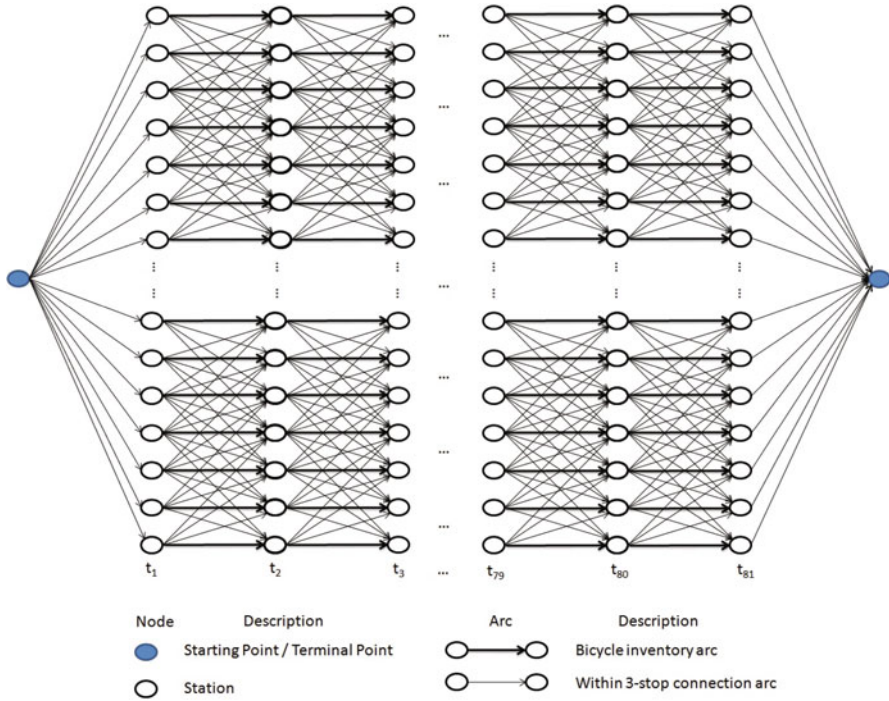


Fig. 17.1 Bicycle flow network: time expanded graph

destinations using bicycles and such demands are considered lost. Figure 17.1 shows the time expanded view of the entire network, where the flow on each arc depends on the realization of the number of passengers arriving in each period to each station, the order of arrivals, and also the number of bicycles available at the station.

To gain better insights into this problem, we consider an initial allotment of bicycles $x_i(t)$ at station i in time period t . The number of bicycle trips that will materialize in time period t will be $\min(x_i(t), D_i(t))$. However, the number of bicycles flowing from i to j will depend on the order of arrivals of passengers at station i , and is more complicated to track. For $0 < p < 1$, let $D_i(t)[p]$ denote the number of tagged passengers where each passenger is tagged with probability p upon arrival. More formally, let $\{\eta_i(p)\}$ denote a sequence of independent Bernoulli r.v.s with mean p , then

$$D_i(t)[p] = \sum_{k=1}^{D_i(t)} \eta_k(p).$$

By the well known Poisson Thinning Lemma, $D_i(t)[p]$ is Poisson with rate $p \times (\sum_{j:j \neq i} r_{ij}(t))$. Let $p_{ij}(t) = r_{ij}(t) / \sum_{k:k \neq i} r_{ik}(t)$. Hence

$$D_{ij}(t) \sim D_i(t)[p_{ij}(t)].$$

By a slight abuse of notation, for some number of bicycles $x_i(t)$, let

$$\min(x_i(t), D_i(t))[p] = \sum_{k=1}^{\min(x_i(t), D_i(t))} \eta_k(p).$$

If there are $x_i(t)$ bicycles at station i , the number of bicycles leaving station i at time t is clearly $\min(x_i(t), D_i(t))$. The number of bicycles traveling from i to j however depends on the order of arrival of the customers traveling to different destinations. In particular, the number of passengers traveling to j follows the distribution of

$$\min(D_i(t), x_i(t)) [p_{ij}(t)].$$

The number of bicycles at station i at the end of the time period is given by

$$\begin{aligned} x_i(t+1) &= x_i(t) - \underbrace{\min(D_i(t), x_i(t))}_{\text{total departures}} + \underbrace{\sum_{j:j \neq i} (\min(D_j(t), x_j(t)) [p_{ji}(t)])}_{\text{total arrivals}} \\ &= x_i(t) - \sum_{j:j \neq i} (\min(D_i(t), x_i(t)) [p_{ij}(t)]) \\ &\quad + \sum_{j:j \neq i} (\min(D_j(t), x_j(t)) [p_{ji}(t)]) \end{aligned} \quad (17.1)$$

The expected number of trips traversed using bicycles is given by

$$\sum_{t=0}^N \sum_{i \in S} \sum_{j:j \neq i} E(\min(D_i(t), x_i(t)) [p_{ij}(t)]).$$

Let $y_i(t) = E(x_i(t))$, and

$$\begin{cases} y_{ij}(t) = E(\min(D_i(t), x_i(t)) [p_{ij}(t)]), \\ y_{ii}(t) = y_i(t) - \sum_{j:j \neq i} y_{ij}(t). \end{cases}$$

By the above definition, $y_{ij}(t)$ stands for the expected number of bicycles traveling from station i to station j during time period t . We next describe some simple structural properties of $y_{ij}(t)$.

Lemma 1 $y_{ij}(t) : y_{il}(t) = r_{ij}(t) : r_{il}(t)$.

Lemma 2 $y_{ij}(t) \leq r_{ij}(t)$.

Lemma 3 $y_i(t+1) = y_i(t) - \sum_{j:j \neq i} y_{ij}(t) + \sum_{j:j \neq i} y_{ji}(t)$.

Let Z^* denote the optimal objective value to the following linear programming problem:

$$\begin{aligned}
 Z^* &= \max \sum_{t=0}^N \sum_{i \in S} \sum_{j:j \neq i} y_{ij}(t) \\
 \text{subject to} \quad &y_i(t+1) = y_i(t) - \sum_{j:j \neq i} y_{ij}(t) + \sum_{j:j \neq i} y_{ji}(t), \quad \forall i, t; \\
 &y_i(t) = y_{ii}(t) + \sum_{j:j \neq i} y_{ij}(t), \quad \forall i, t; \\
 &\frac{y_{ij}(t)}{y_{il}(t)} = \frac{r_{ij}(t)}{r_{il}(t)}, \quad \forall i, j, l, t; \\
 &y_i(0) = x_i(0), \quad \forall i; \\
 &0 \leq y_{ij}(t) \leq r_{ij}(t), \quad \forall t, i \neq j.
 \end{aligned}$$

The second constraint depicts, for each station i , the number of bicycles which are available at the beginning of period t equals to the number of bicycles which remain at station i and the number of bicycles which leave station i during period t . Given an initial allotment of bicycles at station i denoted by $x_i(0)$, the mean number of bicycle trips supported in the BSS on each link is a feasible solution to the above LP. Hence we have:

Theorem 1 Z^* denotes an upper bound to the expected number of bicycle trips in the system when the initial allotment of bicycles to station i is given by $x_i(0)$.

The above LP is surprisingly effective in providing a simple estimate on the performance (based on the number of bicycle trips that the system can support) of the BSS with an initial bicycle inventory position $x_i(0)$. We will use this model extensively in the next section to examine the issues of bicycle utilization and the value of bicycle redistribution.

Example 1 To see that the above LP is not exact, consider a 3-station example where there are 2 bicycles at station 3 initially, and none at the other 2 stations. Suppose $r_{31}(0) = r_{32}(0) = 1, r_{23}(t) = r_{32}(t) = 1$ for all $t > 1$, and $r_{ij}(t) = 0$ otherwise. In this case, to support the maximum number of flow in the network, the optimal LP solution suppress the flow of bicycles from station 3 to station 1 and 2 in the first period, so that 2 bicycles will remain in station 3 from period 1 onwards to serve the flow between station 2 and station 3, without parking any bicycle at station 1. This LP solution dominates the expected number of trips in the stochastic network flow model.

17.2.1 Equilibrium State in Time Invariant System

In the rest of this section, we further analyze the properties of this formulation (simple network flow with proportionality constraints) to gain insight on the problem.

Suppose the Poisson arrival in each time period is stationary with rate r_{ij} . Is there a way to characterize the number of bicycles in the equilibrium state of the bicycle-sharing network? We modify the LP to provide a glimpse to the answer to this problem.

In the equilibrium state, we expect $y_i(t+1) = y_i(t)$ as $t \rightarrow \infty$. Let

$$y_{ij} = \lim_{t \rightarrow \infty} y_{ij}(t).$$

The total number of bicycles in the system is denoted by N . Let y_{ij}^* denote the optimal solution to the following LP.

$$\begin{aligned} Z^*(\infty) &= \max \sum_{i,j \in \mathbf{S}: j \neq i} y_{ij} \\ \text{subject to} \quad & \sum_{j:j \neq i} y_{ij} = \sum_{j:j \neq i} y_{ji}, \quad \forall i; \\ & \frac{y_{ij}}{y_{il}} = \frac{r_{ij}}{r_{il}}, \quad \forall i, j, l; \\ & 0 \leq y_{ij} \leq r_{ij}, \quad \forall i, j; \\ & \sum_i \left(y_{ii} + \sum_{j:j \neq i} y_{ij} \right) = N. \end{aligned}$$

It can be seen easily that there exists i^* such that $y_{i^*j} = r_{i^*j}$ for all $j \neq i^*$, otherwise we could scale the solution to improve the objective value. We call such nodes the *sink* stations. Furthermore, if there exists i such that $y_{ii}^* > 0$ but $y_{ij}^* < r_{ij}$ for all $j \neq i$, then we could modify the solution by shifting y_{ii}^* to the station i^* , without affecting the feasibility and quality of the solution. i.e.,

$$y_{ii}^* \leftarrow 0, \quad y_{i^*i^*}^* \leftarrow y_{i^*i^*}^* + y_{ii}^*.$$

We call such nodes where $y_{ij}^* < r_{ij}$ the *transient* stations. Note that WLOG, we can assume that $y_{ii}^* = 0$ when i is transient.

Let $z_i^* = \sum_{j:j \neq i} y_{ij}^*$. By the proportionality constraints, it is easy to see that

$$y_{ij}^* = \frac{r_{ij}}{\sum_{k:k \neq i} r_{ik}} z_i^*.$$

Note that z_i^* is a solution to the following system of linear equations:

$$z_i = \sum_{j:j \neq i} \frac{r_{ji}}{\sum_{k:k \neq j} r_{jk}} z_j, \quad i = 1, \dots, n. \tag{17.2}$$

If the transition probability matrix constructed using $r_{ji}/\sum_{k:k \neq j} r_{jk}$ is irreducible, then the above system of equations has essentially a unique solution scaled to a constant. Note that $z_i^* \leq \sum_{k:k \neq i} r_{ik}$ since $y_{ij}^* \leq r_{ij}$, and $\sum_i z_i^* \leq N$. Since our objective is to maximize $\sum_i z_i^*$, the solution to the linear system (17.2) is scaled in such a way that either (i) $\exists \mathbf{S}$ such that $z_i^* = \sum_{k:k \neq i} r_{ik}$ for all $i \in \mathbf{S}$, and $z_i^* < \sum_{k:k \neq i} r_{ik}$ otherwise; or (ii) $\sum_i z_i^* = N$ and $z_i^* < \sum_{k:k \neq i} r_{ik}$ for all i . \mathbf{S} corresponds to the set of sink nodes in the system. In case (i), the surplus $N - \sum_i z_i^*$ can be distributed to any of the y_{ii}^* variables for $i \in \mathbf{S}$ without affecting the optimality of the solution.

Theorem 2 *The linear program $Z^*(\infty)$ may have multiple optimal solutions, but the flow solution $y_{ij}^*, i \neq j$, is uniquely determined by the rates r_{ij} , if the transition probability matrix is irreducible. The “surplus” denoted by y_{ii}^* for the sink nodes are however non-determined and can be distributed across different sink nodes.*

Since the surplus y_{ii}^* have zero weights in the objective function, having large surplus does not help improve the quality of the solution. This result indicates that given the rates r_{ij} 's, there is a limit N^* such that any number of bicycles beyond this limit N^* will not help to improve the performance of the system.

Example 2 As shown in Fig. 17.2, we have three stations which are connected to each other. The number beside each direct arc (i, j) stands for the arrival rate r_{ij} . Station 1 has a net outflow of three passengers per unit, whereas stations 2 and 3 have net inflow of two and one passenger, respectively. Naively, we expect the average number of bicycles at station 1 to drain down to 0 quickly, with the bulk of bicycles building up at stations 2 and 3. However, note that once the bicycle at station 1 drains down to 0, stations 2 and 3 immediately receive less inflow and in fact station 3 will now have a net outflow of 2 passengers per unit.

Fig. 17.2 Numerical example with 3 stations

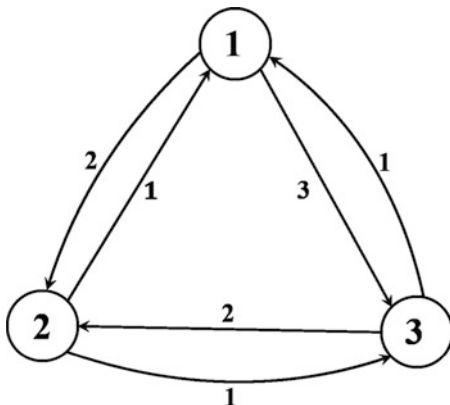


Table 17.1 Computational results of the numerical example with three stations

	Simulation		Deterministic	Gap %
	$t = 0$	$t = 50$		
Avg no. of bicycles at station 1	5	1.656600000	1.65666651648	0.004015241
Avg no. of bicycles at station 2	2	5.506000000	5.50643297152	0.007863631
Avg no. of bicycles at station 3	3	2.837400000	2.83690051200	-0.017603722

We use the outputs from the simulation model to plot the time average level of bicycles at each station over 2,000 periods and summarize the computational results in Table 17.1. In particular, the gap is calculated as $100\% \times (\text{output of the deterministic model} - \text{output of the simulation model (average of 2,000 simulations)}) / \text{output of the simulation model}$. We observe that the time average number of bicycles at each station stabilizes after 10 time periods, and the LP model gives very accurate prediction to the time average level of bicycles in the stochastic system.

17.2.2 Bicycle-Sharing System Design with Location Choice

The previous model assumes that the stations are fixed. This is reasonable when the location choices are obvious – like the train stations in the C-Bike system. It is however inevitable, just like the C-bike system in Kaohsiung, for the network to extend its reach into hot spots and residential areas to capture new passengers who would otherwise not use the public train system for their transport needs. It is therefore essential that we incorporate the location decision of the bicycle-sharing stations as one of the key decisions in our model.

The difficulty in the modeling approach is to incorporate the proportionality constraints into the formulation. We have seen in the earlier section that this class of constraints is crucial for the LP model to approximate the performance of the stochastic network flow model. To see that this is hard to incorporate into the formulation, suppose station i , j , and k have been set up, but l has been omitted in the model. Then we would need to incorporate the proportionality constraint on the flow from i into j and k , i.e., $y_{ij}(t)/y_{ik}(t) = r_{ij}(t)/r_{ik}(t)$, but not for flow from i to l , since the flow $r_{il}(t)$ would be lost as l has not been selected as a node in the bicycle-sharing network.

To deal with complications in the location choice formulation, we need to introduce 0–1 decision variables into our formulation. With a slight abuse of notation, we can redefine \mathbf{S} as the set of potential bicycle dock stations which includes MRT stations and neighborhood locations. Let f_i be the setup cost of installing bicycle docks at location i , q_{ij} the environmental benefit/amount charged by the operator of trip ij , and N the number of planning periods.

We also let z_i be a binary variable which represents the presence of bicycle station at location i . We can modify the Linear Program developed in the earlier section to

account for location choice as follows: Let $u_i(t)$ be a decision variable to denote the effective rate of demand substitution from station i at time t . All positive flows out of i , to any other station j at time t , normalized over the demand rate $r_{ij}(t)$, must be identical to this ratio $u_i(t)$. Then, the model of the BSS design with location choice can be formulated as:

$$Z_L^*(\beta) = \max_{x_i(0), y_{ij}(t)} \left(\sum_{t=0}^N \sum_{i \in S} \sum_{j: j \neq i} q_{ij} y_{ij}(t) - \sum_{i \in S} f_i z_i \right)$$

subject to

$$y_i(t+1) = y_i(t) - \sum_{j: j \neq i} y_{ij}(t) + \sum_{j: j \neq i} y_{ji}(t), \quad \forall i, t; \quad (17.3)$$

$$y_i(t) = y_{ii}(t) + \sum_{j: j \neq i} y_{ij}(t), \quad \forall i, t; \quad (17.4)$$

$$u_i(t) - (1 - z_j) \leq \frac{y_{ij}(t)}{r_{ij}(t)} \leq u_i(t) + (1 - z_j), \quad \forall i, j, t; \quad (17.5)$$

$$y_i(0) = x_i(0), \quad \forall i; \quad (17.6)$$

$$y_{ij}(t) \leq r_{ij}(t) z_i, \quad \forall i, j, t; \quad (17.7)$$

$$y_{ij}(t) \leq r_{ij}(t) z_j, \quad \forall i, j, t; \quad (17.8)$$

$$0 \leq u_i(t) \leq 1, \quad \forall i, t; \quad (17.9)$$

$$z_i \in \{0, 1\}, \quad \forall i. \quad (17.10)$$

Constraint (17.5) is crucial for this formulation – when there is a station setup at location j , then $z_j = 1$ and $y_{ij}(t) = u_i(t)r_{ij}(t)$. This forces all flow from i to other locations with stations setup to follow the proportionality constraint and the effective rate of demand substitution for all trips leaving i will equal to $u_i(t)$. Otherwise, the flow from i to j is zero. Constraints (17.7) and (17.8) model the fact that if there are bicycle stations at both location i and j , the flow between the stations will be no more than the demand rate. If either location i or j does not have a station, the flow between i and j will be 0. Constraint (17.9) forces the effective rate of demand substitution to be between 0 and 1. All the other constraints follow from the models derived in the previous sections.

17.3 Bicycle Sharing as Substitute for Train Rides

The Mass Rapid Transit (MRT) system of Singapore operates around 5.00 a.m. to 01.00 a.m. each day, with morning peak hour traffic occurring at around 7.30 a.m. to 9.30 a.m., and evening peak at around 5.30 p.m. to 7.30 p.m. To construct a numerical example for our model, we use a one-week sample of train service passenger-flow data (covering more than 10 million trips) to construct our demand

model. Interestingly, we found that about 16% of the trips are short trips, i.e., with passengers leaving the train system within two stops from their starting stations. The longest trip can take up to 33 stops, but the average number of stops traversed is only around 7.7 stops. Note that except for a handful of stations, commute time between neighboring stations are around two to three minutes. The statistics thus show that a significant proportion of passengers (around 16%) commute up to at most six minutes on the train on a daily basis. An alternate public transport system such as a public bicycle-sharing service, located at the MRT stations, is an attractive alternate for such commuters, especially during the morning and evening peak hours. The challenge however is to determine the right level of bicycles to deploy at each station, and how the utilization rates are affected by the demand pattern.

We compare next the proposed proportional network flow model with a simulation model to identify the operational characteristics of the bicycle-sharing service.

17.3.1 *Bicycle Deployment and Utilization*

We split the horizon into 15-min intervals, starting from 05:00 am, to collect passengers data on those alighting within two stations. There are 80 time intervals for each day and 560 time intervals for a week. We use a directed time-expanded network to model each MRT station at each time interval on each day. Let \mathcal{A} denote the arc set in the time-expanded network. There are two types of arcs in \mathcal{A} . The first is the one which links station i in time t to station j in time $t + 1$, for all t in which station j is within two stops away from station i . The other arcs are inventory arcs, joining the same station across two consecutive time periods.

We also adopt the following notations:

- We define the system bicycle utilization rate $\alpha(t)$ for each time period t as follows:

$$\alpha(t) \equiv \frac{\sum_{i,j:i \neq j} y_{ij}(t)}{\sum_i x_i(0)},$$

where $\sum_i x_i(0)$ represents the total number of bicycles positioned at all the stations at the beginning of the planning horizon. Hence $\alpha(t)$ is the proportion of bicycles in use at time t .

- Similarly, since the number of bicycles in the system is a constant,

$$\beta = \sum_t \alpha(t)$$

measures the total number of rides in the system, divided by the total number of bicycles available, i.e., the (average) number of times each bicycle is being used.

To support a larger number of trips on bicycles, we might have to deploy more bicycles in the system, but the average bicycle utilization rate might decrease in this case. On the other hand, if we want to enhance the bicycle utilization rate, we could try to deploy relatively fewer bicycles in the system. Therefore, there is a one-to-one relationship between the number of bicycles (optimally) deployed and the utilization rate of each bicycle. To design the BSS, we opt to first determine the desired level of β in the system. Note that β determines the economic viability of the BSS – the bicycle needs to be used more than a threshold value within a stipulated number of years to justify the initial investment in the bicycle.

With the above defined notations, we can modify the linear program developed in the earlier section to account for the utilization rate:

$$Z^*(\beta) = \max_{x_i(0), y_{ij}(t)} \sum_{t=0}^N \sum_{i \in S} \sum_{j: j \neq i} y_{ij}(t)$$

$$\text{subject to } y_i(t+1) = y_i(t) - \sum_{j: j \neq i} y_{ij}(t) + \sum_{j: j \neq i} y_{ji}(t), \quad \forall i, t; \quad (17.11)$$

$$\sum_{t=0}^N \sum_{i \in S} \sum_{j: j \neq i} y_{ij}(t) \geq \beta \sum_i x_i(0); \quad (17.12)$$

$$y_i(t) = y_{ii}(t) + \sum_{j: j \neq i} y_{ij}(t), \quad \forall i, t; \quad (17.13)$$

$$\frac{y_{ij}(t)}{y_{il}(t)} = \frac{r_{ij}(t)}{r_{il}(t)}, \quad \forall i, j, l, t; \quad (17.14)$$

$$y_i(0) = x_i(0), \quad \forall i; \quad (17.15)$$

$$0 \leq y_{ij}(t) \leq r_{ij}(t), \quad \forall t, i \neq j. \quad (17.16)$$

Note that constraint Eq. 17.12 requires the weekly bicycle utilization rate to be at least β . The above LP determines the total number of bicycles and their deployment at the beginning (i.e., $x_i(0)$) of the planning horizon, to attain the desired utilization rate of β for the system. We solve the above model using the CPLEX LP solver. We solve the above program to obtain the maximum number of substituted trips using bicycles, the number of bicycles positioned at each station initially, and bicycle utilization rate $\alpha(t)$ at each time period.

We also compare the solutions obtained from the deterministic model with a simulation model. The detailed steps in implementing the simulation are given as follows.

- We fix a β , solve the deterministic model outlined in this section, and obtain the optimal $x_i(0)$ to be deployed at each station i .

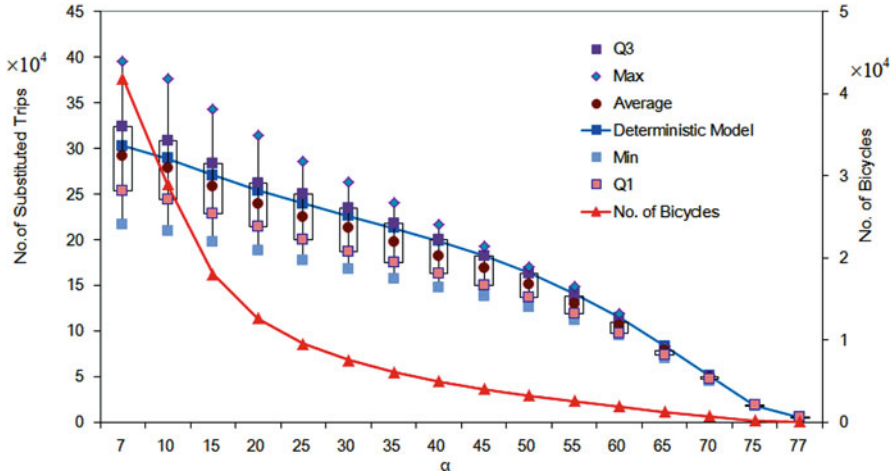


Fig. 17.3 Short trip substitution boxplot

- We use $x_i(0)$ as the input to run the simulation model for stochastic network flow system with Poisson demand at each arc in the network. We run the simulation 100 times for each β to obtain the sample average of the system performance.
- In each simulation, we use the direct time-expanded network and assume the number of passengers arriving at each station during each 15 min time interval follows a Poisson process. In particular, the mean of the inter-arrival time for passengers arriving at Station i with destination Station j at time index t equals to $15/r_{ij}(t)$. We then sort the passengers at each node according to their arrival time at node i and discard those arrivals after 15 min. The bicycles at station i are used by the passengers arriving on a first-come-first-serve basis. We run this for a whole week to obtain the number of trips on bicycles and the utilization rate.

Figure 17.3 shows the performance of the bicycle-sharing network when short trips (within 2 stops) can be completely substituted. In the figure, the x -axis corresponds to the average daily utilization rate (denoted by α , where $\alpha = \beta/7$). The y -axis on the left shows the number of trips using bicycles, and the y -axis on the right shows the number of bicycles deployed in the system. The box plots (obtained via simulation) show that the variations in the number of bicycle trips increase when the daily utilization rate decreases. More amazingly, the numerical results show that the deterministic LP model yields very tight estimate (upper bound) to the average number of trips on bicycles in the stochastic network flow model.

The relationship between the number of bicycle trips supported by the system and the daily utilization rate appears to be almost linear – the number of trips decreases linearly as the targeted utilization rate of the system increases. However, the number of bicycles needed is inversely proportional to the targeted daily utilization rate, in the optimal configuration.

These trade-offs have important implications – it appears that an appropriate targeted utilization rate to operate is around the region $\alpha = 30\text{--}40$ in this test case – at the rate above this level, we need to deploy a significantly many bicycles to support a small increase in the number of supported trips. However, in this range, the service level will not be high as a significant portion of demand for rides cannot be supported. The average total demand within the system is around 308,000 trips, whereas the system operating at $\alpha = 40$ can only support on average 50,000 trips, under the assumption that all short trips can be substituted by bicycle rides if possible.

17.3.2 Number of Bicycle Docks Needed

Technically, we need to set up enough number of bicycle docks at each station so that passengers have space to return the bicycles when they reach their destinations. We calculate the number of docks needed for each station as the maximum bicycle quantities at each station across all time periods. Figures 17.4, 17.5, and 17.6 give the maximum bicycle quantities at each station among all time periods for $\alpha = 10, 40, 70$, for the network flow and the simulation model.

Interestingly, the data extracted from the network flow model is pretty close to the actual peak inventory level at each station obtained from the simulation model. Furthermore, the number of docks needed to support storage of peak inventory decreases with increased utilization (i.e. less number of bicycles deployed).¹

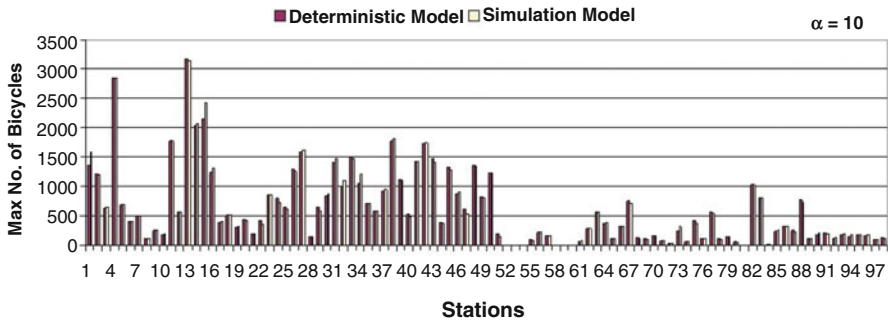


Fig. 17.4 Two-stop no. of docks: deterministic model vs. simulation model

¹Note that we have assumed all passengers will use bicycles to substitute their short distance MRT trips (within 2-Stop), upon the availability of the bicycles. We have thus actually obtained a gross over-estimate on the total volume of trips that can be substituted by bicycles. In reality, only a small percentage of the short distance passengers captured in the data will choose to use bicycles, say 10%. Therefore, all our numbers must be scaled down by a factor of 10 accordingly. In this case, we can see that for $\alpha = 40$, the maximum number of bicycle docks we need to setup among all stations is no more than 80 for our system.

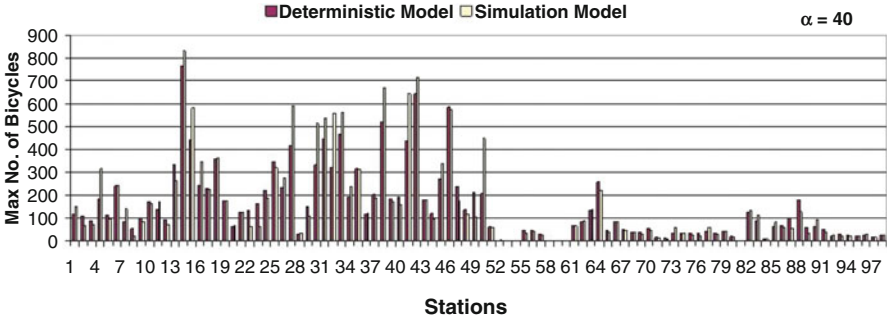


Fig. 17.5 Two-stop no. of docks: deterministic model vs. simulation model

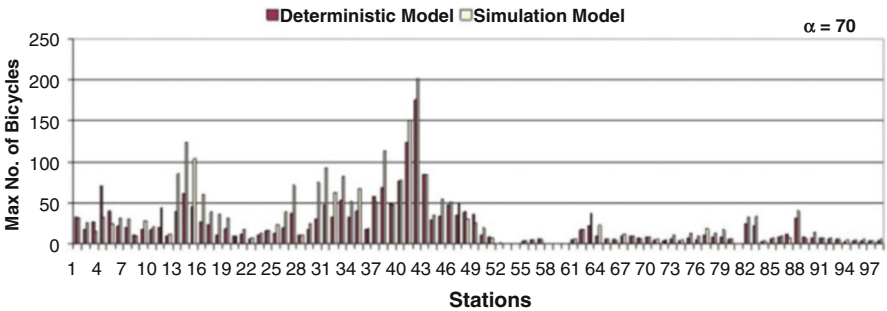


Fig. 17.6 Two-stop no. of docks: deterministic model vs. simulation model

The computational results also show interestingly that with a smaller number of bicycles deployed ($\alpha = 70$), the system should deploy more bicycles near and around the stations in the central business district (Station 30–50 in the chart), leading to a relatively higher number of docks at these stations. However, with more bicycles available ($\alpha = 40$ or $\alpha = 10$), the deployment of the additional bicycles should move towards other congested areas such as the stations near the interchange in the East (station 1–20), leading to a surge in the number of docks there. This suggests that the operators should focus first in the central business district area with a small number of bicycles and stations, to capture the maximum number of trips supported, before branching out into major residential areas as the scale of the system grows.

17.3.3 Effectiveness of Bicycle Redistribution

With a slight abuse of notation, we redefine the time-expanded network to model the passengers flow for each day k . Let N_k denote the time index in the network on day k . We conduct the experiments as follows. We first solve the deterministic model $Z^*(\beta)$ proposed earlier based on the one-week data to obtain the number of

bicycles deployed (denoted by C_β). We then use these as input to run the following program ($Z_k^*(\beta)$) for each day k :

$$\begin{aligned}
 Z_k^*(\beta) &= \max_{x_i(0), y_{ij}(t)} \sum_{l \in N_k} \sum_{i \in S} \sum_{j: j \neq i} y_{ij}(t) \\
 \text{subject to} \quad &y_i(t + 1) = y_i(t) - \sum_{j: j \neq i} y_{ij}(t) + \sum_{j: j \neq i} y_{ji}(t), \quad \forall i, t; \\
 &\sum_i y_i(0) = C_\beta; \\
 &y_i(t) = y_{ii}(t) + \sum_{j: j \neq i} y_{ij}(t), \quad \forall i, t; \\
 &\frac{y_{ij}(t)}{y_{il}(t)} = \frac{r_{ij}(t)}{r_{il}(t)}, \quad \forall i, j, l, t; \\
 &0 \leq y_{ij}(t) \leq r_{ij}(t), \quad \forall i, j, t.
 \end{aligned}$$

The above linear program computes the optimal way to locate the C_β bicycles in the system, given the travel patterns of the day. Note that we solve an LP for each β .

The redistribution strategy has impact on the performance of the BSS. Figure 17.7 shows that this strategy prevents surplus bicycles from building up at stations, and thus reduces the need to build large number of docks at each station. For $\alpha = 40$, it reduces the peak docking stations needed from 800 to around 700.²

Although redistribution strategy can enhance the system performance in terms of the number of substituted trips supported in the system and the number of docks needed in each station, it is a very time consuming and expensive task. The concern is how often shall we conduct the redistribution in the system? In the rest of this

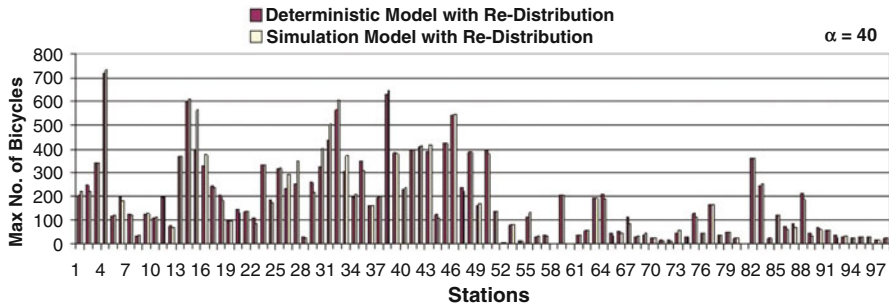


Fig. 17.7 Two-stop no. of docks: deterministic model vs. simulation model with redistribution

²If we assume that the take-up rate for bicycle trip is only 10% of the full demand, then the corresponding number of docks needed will be reduced by 90%, i.e., from 700 to 70 docks.

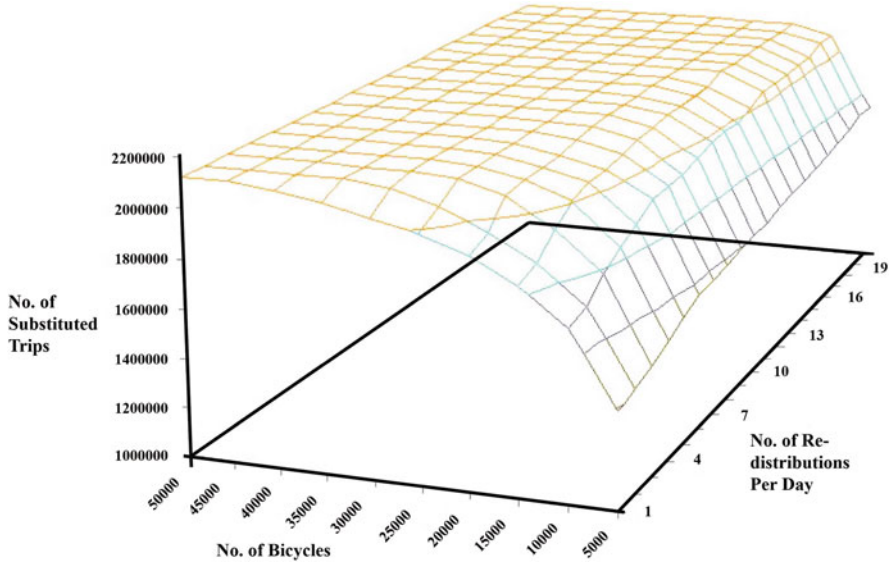


Fig. 17.8 3-D illustration of periodic redistribution

section, we discuss the value of periodic redistribution. Figure 17.8 shows the number of substituted trips supported in the system under given combinations of the total number of bicycles and the number of periodic redistributions per day. In this set of experiments, we subdivide the time horizon evenly into 80 smaller time intervals in a day, and perform periodic redistribution at equal time interval over a day. For certain cases, when 80 time intervals is not divisible by the number of redistributions per day, we keep the remainder in the last time interval of the day. Figure 17.8 shows the end result: when the total number of bicycles invested into the system is more than 30,000, frequent periodic redistribution does not add much to the number of bicycle trips supported by the system. Furthermore, a small number of daily redistribution (says 2–4) suffices, since more frequent redistribution will not add much to the total supported bicycle trips.

17.4 Case Study on Bicycle Sharing with Location Decisions

Punggol is a neighborhood in the northeastern region of Singapore. Initially an area populated by farms, Punggol has been developed into a residential new town. Currently, the district is home to 17,980 HDB³ flats and has an estimated residential

³Housing and Development Board – a statutory board of the Singapore Government responsible for public housing.

population of 59,200. There are plans to develop Punggol as Singapore’s first Eco-Town to enhance the living environment in its estates and encourage residents to do their part for the environment (Housing Development Board, 2010). A BSS would be an appropriate addition to the Punggol landscape as it actively promotes different forms of sustainable transportation.

The Punggol district is served by one Mass Rapid Transit (MRT) station, 29 bus stops, and a Light Rail Transit (LRT) network of eight stations. A total of eight bus services serve the area. The LRT was set up in 2005 as an alternative transportation mode or feeder service within the neighborhood. At present, private bicycles are already being used as a mode of transportation within the Punggol area. Residents were spotted riding their personal bicycles both on the roads as well as on the footpaths. Also, bicycles were seen parked at LRT stations and there is even a designated Bike Park area at the Punggol MRT Station cum Bus Interchange for residents to park their bicycles.

We use a set of commuters data on bus and LRT services to design the bike-sharing network. For this reason, it would be most appropriate for the candidate locations to be at the bus stops, LRT stations, and MRT station. The 16 candidate locations are as shown in Fig. 17.9.

There are two peaks in the travel pattern in this town. Figure 17.10 shows the number of trips during each time period, where each time period corresponds to a 15-min period. Period 1 starts at 05:15 a.m. while Period 79 ends at 01:00 a.m. It can also be concluded that throughout the day, the LRT is the more utilized transportation mode.

The most traversed route is from Location 11 to 5, while the route from Location 5 to 11 is ranked second. Location 5 corresponds to the MRT Station while Location 11 corresponds to the area around a LRT Station named Meridian. It would be



Fig. 17.9 Candidate locations

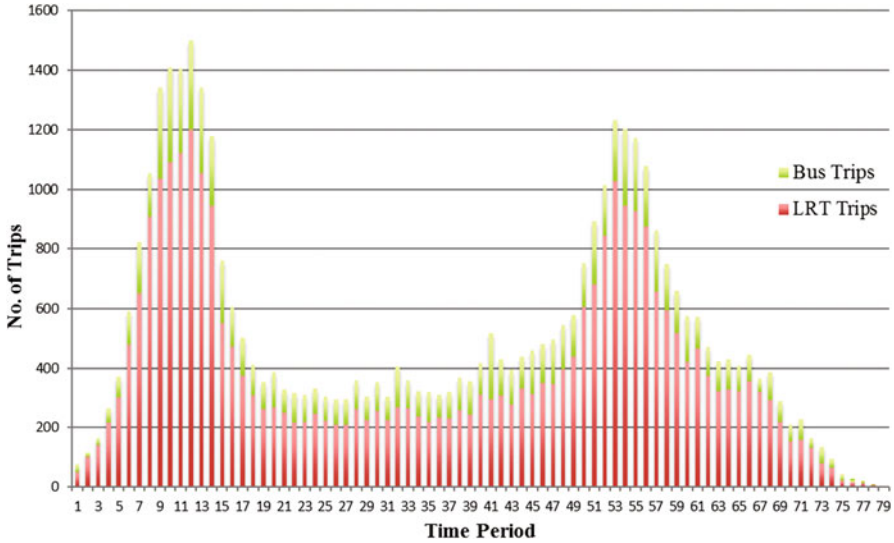


Fig. 17.10 Demand pattern

expected that these 2 routes are the most heavily traversed as Location 11 has the greatest number of residential blocks surrounding it. It should also be noted that some routes see no trips the entire day, hence suggesting that there would be little or no demand for a BSS to cover these routes. Based on the distribution of trips across locations, the frequency at which Location 5 is involved is significantly higher than any of the other locations. This suggests that Location 5 would be a strong candidate location to be chosen for a bike station. However, note that the routes involving Location 5 also suffer more from trip imbalance throughout the day. In the morning, there is a huge demand for bikes to go to the MRT station while in the evening the demand transfers to bike leaving the MRT station. For the rest of the day, the flow of passengers to and from the MRT station is much less, which might result in bikes being stranded at the MRT station in the middle of the day. This would affect the overall utilization of the bikes. Therefore, in solving the model, a case where Location 5 is excluded will be solved in order to see if this results in higher utilization of the bikes.

However, we need to account for the fact that there would not be 100% uptake of the BSS. Based on an informal survey conducted, as well as reference to the average uptake rates predicted for overseas systems (Dector-Vega et al. 2008), it is predicted that the average uptake rate in Singapore should lie between 4 to 6%. Furthermore, this uptake rate is unlikely to be constant throughout the day. It is expected that the uptake would be higher in the mornings and evenings when the weather is cooler. Fewer people are likely to switch to cycling in the afternoon when the weather is hotter. This is taken into consideration by allocating a predicted uptake rate of 6% for the periods between 05:15 a.m.–11:45 a.m. and 05:00 p.m.–01:00 a.m. and an uptake rate of 4% for the period between 11:45 a.m.–05.00 p.m.

Table 17.2 Optimization results

Case	Uptake rate	β	No of locations selected	locations selected	Initial no of bikes at location	Final no of bikes at location
Without accounting for redistribution costs						
1	6% for periods 1–26, 48–79 4% for periods 27–47	12	5	5	0	0
				9	4	1.6
				11	0	4.7
				12	8	11.6
				15	18	12.1
2	4% throughout the day	10	5	5	0	0
				9	6	3.3
				11	0	6.1
				12	10	10.1
				15	14	10.5
3	3% throughout the day	8	5	5	0	0
				9	8	4.4
				11	0	6.6
				12	9	10.1
				15	13	8.9
Accounting for redistribution costs						
4	6% for periods 1–26, 48–79 4% for periods 27–47	12	5	5	0	0
				9	2	2
				11	4	4
				12	12	12
				15	12	12
5	4% throughout the day	10	5	5	0	0
				9	4	4
				11	4	4
				12	11	11
				15	11	11

We use the model introduced in Sect. 17.2.2 to design an optimal Bike-Sharing network given 30 bikes and around four to five bike stations. The model was coded and solved using the CPLEX MIP solver in General Algebraic Modelling System (GAMS). We also account for the inventory imbalance at the start and end of the day, and use that to penalize for redistribution cost. Table 17.2 summarizes the results obtained for several cases where uptake rates were varied. The optimal number of bike stations, optimal locations for the bike stations, and optimal number of bikes to locate at each station initially were determined and tabulated for each case.

The model was also run with an additional term in the objective function to penalize for redistribution costs at the end of the day, in order to determine the effects this had on the solution.

By inspecting the output of our location model, we obtained the following interesting observations:

- Regardless of the uptake rate, it is found that the same number of stations is selected. However, the utilization rate that can be achieved is reduced according to the uptake rate. For an uptake rate of 3% throughout the day, the maximum utilization that can be achieved is 8, compared to 10 for an uptake rate of 4% throughout the day.
- It is interesting to note that regardless of the uptake and utilization rates, and whether redistribution is accounted for or not, the same 5 locations are always chosen. The location decision is thus insensitive to the accuracy of the uptake and utilization rates.
- Even though the same locations for the bike stations are chosen for each case, there are slight differences in the initial number of bikes that should be located at each station. The common feature in all cases is that no bikes should be located at Location 5 initially.
- For cases 1–3 where redistribution costs are not accounted for in the model, it is found that the number of bikes found at each station at the end of the day varies pretty significantly from the initial number of bikes. On average, 7.5 bikes would have to be redistributed at the end of each day.

Take for instance the scenario corresponding to case 4 in our experiment. The model proposed that bike stations should be installed at locations 5, 9, 11, 12, and 15. Also, at the start of the day, there should be 2 bikes at Location 9, 4 bikes at Location 11, 12 bikes at Location 12, 12 bikes at Location 15, and none at Location 5. The model predicted that with this configuration, the bikes should circulate throughout the day such that at the end of each day, the number of bikes at each station will be back to the number at the beginning of the day. We compare the results obtained with a simulation output. The optimal number of bikes that should be located at each station initially, which was determined by the MIP model, was used as input in the simulation model. Again, we assumed that the number of passengers arriving at each station i with destination j during each time period follows a Poisson process. The simulation was run for 100 days. Table 17.3 confirms that the stochastic model behaves more or less as predicted, with utilization rate of 11.22, which is close to the rate of 12 predicted by the MIP model. Based on the maximum number of bikes present at each of the locations over all the time periods, the number of bike docks that should be installed at each station can also be determined.

Table 17.3 Comparison with simulation output

	Station	Initial number	Maximum number	Utilization
Simulation model	5	0	30	
	9	2	4	
	11	4	6	11.22
	12	12	14	
	15	12	17	
MIP model	5	0	25	
	9	2	8	
	11	4	8	12
	12	12	13	
	15	12	13	

17.5 Concluding Remarks

Despite the many problems and success stories of the third generation BSS, there seems to be further advancement to the fourth generation, where more emphasis will be placed on improved efficiency, sustainability, and usability (cf. DeMaio 2009). This can be achieved by focusing on improving the deployment and tracking of bicycles, improving the installation and powering of bicycle stations, creating new business models, and building both intra- and inter-transport system integration (Forbes 2017; MetroBike LLC 2012).

In recent years, BSS vendors have emerged and created their own systems which they sell to local operators. Also, startups like CityRide are converting bike rides into carbon offset that can be sold on the carbon market. The evolution in business strategies and pricing strategies allows the different BSSs to seek out a business model that would be profitable, thus ensure that new BSSs will continue to be set up all around the world, regardless of the goals or scale.

Nevertheless, the fundamental issue of deployment remains a challenge. The deployment can be improved by balancing the supply and demand at each of the bicycle stations, and providing relevant incentives in order to steer demand towards the less popular bicycle stations or routes. In particular, operations research tools can be used to design a network with an improved deployment of bicycles.

In this chapter, we review a novel bicycle-sharing model proposed in Shu et al. (2013) in which passengers use bicycles to substitute their short distance trips. We use a deterministic LP model to approximate the system performance of the stochastic system, and show that the deterministic model can imitate the actual system performance very closely based on actual Singapore MRT ridership data. We use extensive numerical experiments to discuss the important issues such as the bicycle utilization rate, the value of redistribution of the bicycles, and the number of bicycle docks that should be set up at each station.

Our model can be extended to incorporate the scenario of using bicycles to transport between MRT stations and neighborhoods. We implemented our model

using a set of bus transit data in a new town in Singapore, and identified the ideal locations to set up the bicycle stations for the bicycle-sharing network. Our numerical results suggest that the optimal location choices are robust to input errors – for various demand scenarios, the same set of locations are identified as optimal.

Our approach is general enough to incorporate various other features in practice. For example, when the passengers are not able to reach their destination station using bicycles within 15 min time period, we only need to slightly modify the arcs in the time-expanded network defined in the earlier section to allow them to extend across multiple time periods. The same LP based approach can be used to model the flow of commuters in the network. Of course, in the most general case, we need to use queueing network based approach⁴ to model the flow of bicycles in this system. However, the associated optimization problem becomes intractable using this approach, due to the time varying nature of the travel patterns.

The performance of the LP model can also be further enhanced, exploiting recent advances in stochastic optimization (cf. Natarajan et al. 2009, 2011). In particular, a promising direction is to enhance the model further using the constraint that

$$x_i(t) - \underbrace{\min(D_i(t), x_i(t))}_{\text{total departures}} > 0$$

if and only if all passengers arriving in time t to station i can find a bicycle. Thus

$$\left(x_i(t) - \sum_{j:j \neq i} \left(\min(D_i(t), x_i(t)) [p_{ij}(t)] \right) \right) \times \left(D_{ik}(t) - \left(\min(D_i(t), x_i(t)) [p_{ik}(t)] \right) \right) = 0$$

for every OD pair i, k in all realization of the stochastic system. This can be handled by lifting the problem into a higher dimension, and using the copositive cone approach in Natarajan et al. (2011) to deal with the quadratic constraints.

Another interesting direction of research is to explore the usage of incentive schemes to balance the flow. Our approach hinges crucially on the fact that system parameters $r_{ij}(t)$ are given as input. When they are endogenous to the model, i.e., that promotional activities can be used to influence the flow rate between i and j , then the problem is still unsolved.

We leave these and other issues to future research.

Acknowledgements We thank Singapore Mass Rapid Transit and Land Transport Authority for providing the data used in this research. This research was supported in part by NUS Academic Research Fund R-314-000-078-112.

⁴We thank Prof Gideon Weiss for pointing this out.

References

- Angeloudis P, Hu J, Bell MG (2014) A strategic repositioning algorithm for bicycle-sharing schemes. *Transportmetrica A* 10(8):759–774
- Benchimol M, Benchimol P, Chappert B, De La Taille A, Laroche F, Meunier F, Robinet L (2011) Balancing the stations of a self service “bike hire” system. *RAIRO-Oper Res* 45(1):37–61
- DeMaio P (2009) Bicycle-sharing: history, impacts, models of provision, and future. *J Public Transp* 12(4):41–56
- Dector-Vega G, Snead C, Phillips A (2008) Feasibility study for a central london cycle hire scheme. Technical report, Transport for London
- Forbes (2017). <https://www.forbes.com/sites/ywang/2017/06/20/worth-1-billion-but-whats-really-driving-chinas-bike-sharingboom/#608d7e69427e>
- Ghosh S, Varakantham P, Adulyasak Y, Jaillet P (2017) Dynamic repositioning to reduce lost demand in bike sharing systems. *J Artif Intell Res* 58:387–430
- Kloimüllner C, Raidl GR (2017) Full-load route planning for balancing bike shaing systems by logic-based Benders decomposition. *Networks* 69(3):270–289
- Li Y, Szeto WY, Long J, Shui CS (2016) A multiple type bike repositioning problem. *Transp Res Part B Methodol* 90:263–278
- MetroBike LLC (2011) The bike sharing blog. <http://bike-sharing.blogspot.com/>. Accessed 1 Oct 2011
- MetroBike LLC (2012) Have card, will travel. <http://bike-sharing.blogspot.com/>. Accessed 17 Jan 2012
- MetroBike LLC (2017) The bike sharing blog. <http://bike-sharing.blogspot.com/>. Accessed 19 Aug 2017
- Natarajan K, Song M, Teo CP (2009) Persistency model and its applications in choice modeling. *Manag Sci* 55(3):453–469
- Natarajan K, Teo CP, Zheng Z (2011) Mixed zero-one linear programs under objective uncertainty: a completely positive representation. *Oper Res* 59(3):713–728
- Raviv T, Tzur M, Forma IA (2013) Static repositioning in a bike-sharing system: models and solution approaches. *EURO J Transp Logist* 2(3):187–229
- Russell M, DeMaio P (2017) The bike sharing world map. <http://bike-sharing.blogspot.com/>
- Schuijbroek J, Hampshire RC, van Hoeve WJ (2017) Inventory rebalancing and vehicle routing in bike sharing systems. *Eur J Oper Res* 257(3):992–1004
- Serna A, Gerrickagoitia JK, Bernabe U, Ruiz T (2017) A method to assess sustainable mobility for sustainable tourism: the case of the public bike systems. In: *Information and communication technologies in tourism 2017*. Springer, Cham, pp 727–739
- Shu J, Chou MC, Liu Q, Teo CP, Wang IL (2013) Models for effective deployment and redistribution of bicycles within public bicycle-sharing systems. *Oper Res* 61(6):1346–1359
- The Economist (2011) Why a Boris bike can be an existential hell. <http://www.economist.com/blogs/gulliver/2011/04/londonscycle-hirescheme/>

Chapter 18

Bike Sharing



Daniel Freund, Shane G. Henderson, and David B. Shmoys

Abstract We discuss planning methods for bike-sharing systems that operate a set of stations consisting of docks. Specific questions include decisions related to the number of docks to allocate to each station, how to rebalance the system by moving bikes to match demand, and expansion planning. We describe linear integer programming models, specially tailored optimization algorithms, and simulation methods. All of these methods rely on careful statistical analysis of bike-sharing data, which we also briefly review. Our discussion of the issues is informed by our 4-year collaboration with Citi Bike in New York City, and its parent company Motivate.

18.1 Introduction

Bicycle-sharing programs are now ubiquitous. These programs allow a user to borrow a bike at one location and return it to another. Such programs enable both bicycle commutes and tourism use. Since users employ a bike episodically, bikes are shared across many users. Pricing schemes for such programs vary, with a common model being subscription based, with the first 30–45 min of each use being free.

Bike-sharing systems vary in their design. At Citi Bike in New York City, with which we have been working since May of 2013 when they first began operations, there are fixed station locations around the city. Each station consists of a number of docks (also known as racks) and riders must pick up and return bikes from stations. Another model, also having fixed stations, does not use docks, but instead uses geo-fencing, whereby bikes need only be returned to the general proximity of a central kiosk. Yet another model does not use stations at all. Instead, bikes are simply left at any convenient location; the global positioning system (GPS) is used through a smartphone app to locate a nearby bike. Indeed, almost all systems provide

D. Freund (✉) · S. G. Henderson · D. B. Shmoys
Cornell University, Ithaca, NY, USA
e-mail: df365@cornell.edu; sgh9@cornell.edu; david.shmoys@cornell.edu

smartphone apps to aid usage, and it is conceivable that bike-sharing programs are enabled by such technology and would be far less popular without it.

In this chapter we focus on the design used by Citi Bike, with stations and docks. In doing so, we in no way imply that the other designs are less important. Indeed, some of the very largest systems in the world are based on those other designs. However, the Citi Bike design is very common, and is the one we know well. Moreover, some of the methods we describe partially extend to other designs, particularly those based on geofencing.

The key questions in systems like that used by Citi Bike relate to capacity planning and rebalancing. In capacity planning, one tries to determine how to size stations, i.e., determine the number of docks at each station. Key objectives in such planning relate to the user experience. One wants to avoid situations where a user cannot find a bike, or cannot find an empty dock to which they can return their bike. With geo-fencing designs, the latter issue is less of a concern. Capacity planning operates on a quarterly or longer timescale, because dock repositioning is non-trivial. It is, however, quite feasible in New York, where docks typically come in sets of 3 and are somewhat portable. In rebalancing, one attempts to move bikes between stations to improve the user experience. Rebalancing is expensive, so it is important to do so judiciously. Citi Bike uses both motorized and non-motorized rebalancing. Box trucks are particularly effective in moving large quantities of bikes over long distances, especially overnight, when streets are less congested. During the day, box trucks are complemented by non-motorized means, including bikes towing trailers that can hold up to 18 bikes. Optimizing these operations requires an appropriate statistical analysis of past data to forecast future demand patterns.

Most research on rebalancing has focused on the optimization of truck routes. A particularly important paper in this context is by Raviv and Kolka (2013) who define a user dissatisfaction function to measure the number of out-of-stock events at an individual station. Different ways of computing this cost function have been suggested by Schuijbroek et al. (2017), O'Mahony (2015), and Parikh and Ukkusuri (2015). Subsequent work by Raviv et al. (2013) defined a routing problem based on the user dissatisfaction function; such routing problems, and attempts to solve them to optimality for larger and larger instances, were further investigated by Forma et al. (2015), Ho and Szeto (2014), and Szeto et al. (2016), among others. Similarly, a line of work, starting with Rainer-Harbach et al. (2013) and followed by Raidl et al. (2013) and Kloimüller et al. (2014) investigated greedy strategies for the rebalancing problem, though they considered a slight variation (i.e., a fluid version) of the user dissatisfaction function. The work by Kloimüller et al. (2014) stands out in that regard in that it also applies to the dynamic case, in which unsatisfied demand also occurs during the rebalancing process. An orthogonal approach to rebalancing has been taken by Shu et al. (2013), O'Mahony (2015), and Jian and Henderson (2015); all of these papers aim to find the optimal configuration of bikes at the beginning of some period. Shu et al. (2013) assume complete knowledge of the future and solve a flow problem; O'Mahony (2015) employs the user dissatisfaction function; Jian and Henderson (2015) use a simulation-optimization based approach to capture network effects. In these three versions, limited means for rebalancing are

disregarded since the focus is solely on the optimal allocation of bikes. Contardo et al. (2012), Vogel et al. (2014), and Nair et al. (2013) are similar to Shu et al. (2013) in that they solve particular flow problems rather than routing problems. Nair et al. (2013) aim to obtain certain service levels with at least some probability. Vogel et al. (2014) presents an NP-hard flow model that also takes into consideration a rebalancing cost. All of these assume that not only the rate of rentals and returns at each station is known, but also the routing probability of each customer, i.e., the probability of a customer at a given station having a particular destination. An approach similar to that of Jian and Henderson (2015) was pursued by Datner et al. (2017), in which they also account for the cost of longer travel times due to out-of-stock events rather than minimizing only the number of out-of-stock events.

A disjoint line of work has focused on minimizing the length of the route of a single capacitated truck, or the combined length of routes for a fleet of such trucks, that needs to visit nodes with demand and supply. The paper by Benchimol et al. (2011) is an early example of such work. They give an approximation algorithm, a hardness result, and a polynomial-time algorithm for instances, wherein the underlying graph is a tree. The same problem has been studied by Chemla et al. (2013) and Dell'Amico et al. (2014) from a mixed-integer programming perspective. Further works in the same spirit have been pursued by Erdoğan et al. (2014), Erdoğan et al. (2015), and Bulhões et al. (2018). The last of these introduces multiple visits, i.e., provides an IP formulation that allows vehicles to visit the same location repeatedly. Interestingly, Casazza et al. (2017) prove that conditions exist that guarantee that multiple visits are not needed. Liu et al. (2016) use weather-data and trips that have started already to predict demand for bikes and docks online, obtain targets for stations, and then solve a routing problem minimizing travel time. Interestingly, Di Gaspero et al. (2013), in a sense, combine the approaches of maximizing impact and minimizing travel time: given fixed targets for each station, the authors aim to minimize a weighted combination of travel time and absolute value distance (summed over all stations) between the targeted bike allocation and the one resulting from rebalancing.

Some recent papers have taken different approaches based on robust optimization. Ghosh et al. (2016) study a repositioning approach based on an iterative two-player game, in which the environment generates a demand scenario out of feasible demand scenarios; they apply this approach to small systems with 20 stations. They also develop a simulation model, which Lowalekar et al. (2017) use to demonstrate the benefit of multi-stage stochastic optimization. Ghosh et al. (2017) make explicit the distinction between *routing* and *repositioning* with the former being about minimizing travel time and the latter being about finding the best obtainable allocation.

In contrast to the work outlined on rebalancing with trucks, O'Mahony and Shmoys (2015) also investigate the use of trailers in bikesharing systems; later work by Freund et al. (2016) also considers so-called corrals.

Separate from the literature on rebalancing, there is also a long line of literature related to forecasting. Most of the forecasting relates to prediction of demand based on historical data; examples include Li et al. (2015), Rudloff and Lackner (2014),

Salaken et al. (2015), O'Mahony and Shmoys (2015), and Riquelme et al. (2017). Several other forecasting questions have been studied as well: Kaspi et al. (2016) try to detect which bikes in a system are broken given the usage data at each station, a question relevant for routing problems such as the budgeted prize-collecting traveling salesman problem studied by Paul et al. (2017), Hsu et al. (2016) use a discrete choice model to study the behavior of users when faced with out-of-stock events; Zhang et al. (2016) predict the destination and destination time of customers given the origin, the time and personal information about the user (gender/age); an approach to predicting pairwise demand, rather than incoming/outgoing demand at individual stations, can be found in Singhvi et al. (2015) and Chen et al. (2016) dynamically cluster stations to predict which stations will run out of available bikes/docks.

Finally, there is a line of work on the design of such systems. Kabra et al. (2015) apply techniques from econometrics to study the density with which stations should be placed. O'Mahony (2015) defines an integer program to investigate what allocation of docks, given a budget of docks, to existing stations minimizes out-of-stock events (using the local user dissatisfaction function at each station); Freund et al. (2017) extend this question in various ways and provide an efficient algorithm to solve it. Jian et al. (2016), using simulation optimization in the same manner as Jian and Henderson (2015), aim to find the optimal allocation of docks, though it allows for network effects that cause non-convexities. All of these papers are based not on rebalancing but on the question of what the result of optimal rebalancing would look like. A similar approach is used by Saltzman and Bradford (2016), who investigate the augmentation problem, that is, the problem of (optimally) adding docks to existing stations. While Saltzman and Bradford (2016) use simulation, this question can also be approached using the methodology of O'Mahony (2015) and Freund et al. (2017).

In this chapter we first describe optimization-based methods for overnight and real-time rebalancing. The optimization models developed for rebalancing extend to methods for deciding how to allocate, or reallocate, docks across the stations in a city, which is important when the initial allocations of docks need to be refined, as is invariably the case, or where usage patterns are evolving. Such evolution is a fact of life in New York City, where the system is constantly growing. We also discuss methods for expansion planning, whereby additional stations are to be added to a system. Finally, we mention incentive schemes, such as that investigated by Singla et al. (2015), to help make a system "self balancing." Such schemes are important, given the expense of manually relocating bikes.

What do we not do? Perhaps the key question we do not address in this chapter is that of pricing. How should a bike-sharing system charge its customers to ensure profitability, or at least self-financing, while ensuring that the resulting system-level behavior is desirable from various standpoints? This question, which is related to, but more complex than, our discussion of incentive schemes, is still far from settled. We refer the reader to the literature on queueing-theoretic models of such systems, including George et al. (2012), Waserhole and Jost (2016), and Banerjee et al. (2017).

18.2 Data and Statistical Challenges

Bike-sharing systems generate a great deal of data. Most systems will log every transaction, namely where and when a bike is checked out, where and when it is returned, and by whom. This data is easily anonymized by removing the user identity, or by creating user indices that are not linked to user names. Bikes in some systems are also equipped with Global Positioning System (GPS) units that allow them to be tracked as they are ridden.

We use this data to fit a stochastic model of a bike-sharing system, which then forms the foundation for optimization models that are discussed later. This stochastic model captures the arrival processes of bikers wanting to pick up bikes at stations, the selection of a destination station, and the time spent biking from the origin station to the destination station. More precisely, we assume, as is typical in modeling bike-sharing systems, that potential bikers arrive to station i according to a non-homogeneous Poisson process with (instantaneous) arrival rate function $(\mu_i(t) : t \geq 0)$, and that the Poisson arrival processes at each station are mutually independent. Why a Poisson arrival process? The Palm-Khintchine theorem, e.g., p. 221 of Karlin and Taylor (1975), Çinlar (1972), and p. 107 of Nelson (2013), tells us, roughly speaking, that the process that results when many potential users each have a small chance of arriving in each time step is very well approximated by a Poisson process, and the result extends to time-varying rates as in our case. Moreover, the theorem extends to *spatial* arrival processes as in our case. One might argue that the arrival processes at adjacent stations might exhibit some dependence; when potential bikers at one station find it empty, they might look to borrow a bike from an adjacent station. However, in our model we assume a loss model wherein bikers not finding a bike simply leave the system, presumably finding another mode of transportation. Assuming that a biker *does* find a bike at station i , our model assumes that the biker then proceeds to station j with probability P_{ij} independent of all else. The matrix of routing probabilities is thus stochastic (i.e., is non-negative with row sums equal to 1), and will typically also be time-dependent, so we denote the matrix process by $(P(t) : t \geq 0)$. Finally, we assume that the successive biking times between stations i and j are independent and identically distributed with distribution function G_{ij} , and that all biking times, irrespective of origin and destination, are independent of all else.

Fitting this model to data is challenging. First we must deal with *censoring*. Censoring arises when a potential biker comes to a station, sees that there are no bikes, and so leaves. Our data does not record a transaction, so we have no direct observation of this event, despite the fact that it should be represented in the arrival process. We term the process that includes such customers the *nominal* process; the *realized* process does not include them. The nominal process is the one we want to fit. A relatively simple approach can be used, at least reasonably, to deal with this censoring. We assume that the normal rate function $(\mu_i(t) : t \geq 0)$

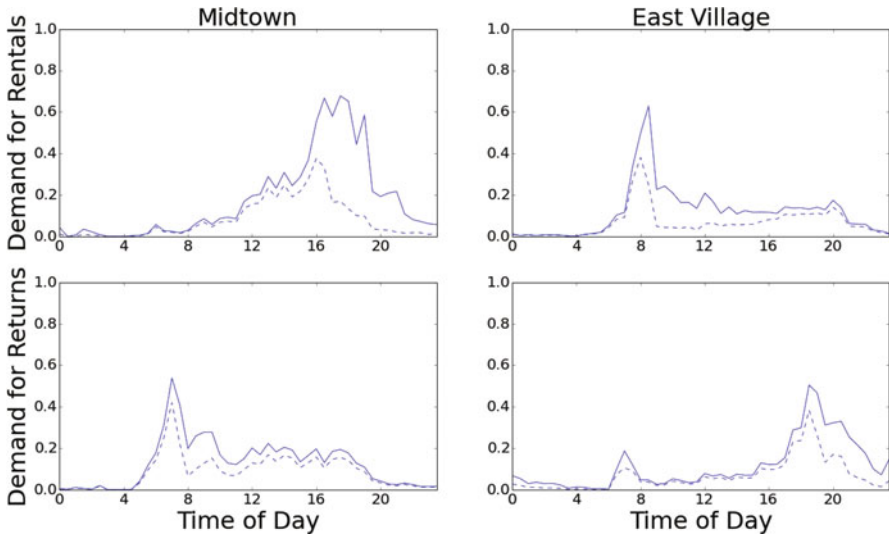


Fig. 18.1 Per-minute demand rates, censored (dashed) and decensored (solid), for a station in Midtown and a station in the East Village (NYC). Notice that censoring of demand for rentals/returns has a stronger impact at times when the stations are more likely to be empty/full

of the arrival process at station i is piecewise constant over intervals of perhaps 30 min. If μ_{ik} is the arrival rate at station i in the k th 30-min period, then we can estimate μ_{ik} by

$$\hat{\mu}_{ik} = \frac{N_{ik}}{\tau_{ik}}.$$

Here, $N_{ik}(\tau_{ik})$ is the cumulative number of bikes that were checked out (cumulative amount of time that bikes were available) from station i in the k th period over multiple instances of the arrival process. For example, if we are estimating the arrival rate for Monday mornings from 9 a.m.–9:30 a.m. and have observed 12 weeks of data, then N_{ik} is the total number of bikes checked out from station i from 9 a.m.–9:30 a.m. over the 12 Mondays. The quantity τ_{ik} is the cumulative amount of time within the interval 9 a.m.–9:30 a.m. over those 12 weeks during which at least one bike was available. So $\tau_{ik} \leq 6$ h, with strict inequality arising if station i was ever empty in that timeframe. Figure 18.1 gives a sample of the fitted arrival rate functions for two stations and a typical business day.

There is a second form of censoring that we cannot handle so readily. Users learn, over time, the availability patterns of bikes at stations. If a station rarely has bikes within a certain period, then users who would otherwise attempt to borrow a bike at that time may learn to not even check for availability. We do not know how to handle this issue, apart from periodically re-learning the nominal rates after system changes that (hopefully) increase the availability of bikes.

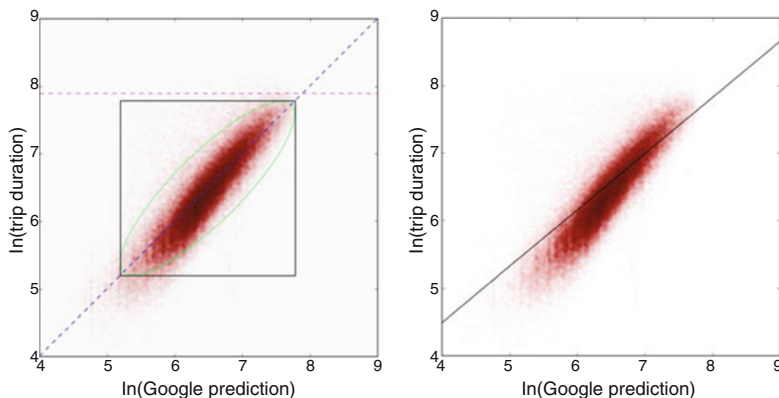


Fig. 18.2 The regression line is $\ln(\text{observed}) = 0.83 \ln(\text{Google prediction}) + 1.16 + \epsilon$, where ϵ is normally distributed with mean 0 and variance 0.168. The R^2 value of the fit is 0.82. Durations are measured in seconds. The blue line is the identity, the magenta line indicates 45 min, that is, the maximum time commuters may ride without surcharge (Based on Figure 2 in Jian et al. 2016)

Fitting the transition matrix $P_{ij}(t)$ is relatively straightforward. Again we assume that $P_{ij}(\cdot)$ is piecewise constant. We then simply estimate it from the empirical distribution for the destinations of bikes borrowed from station i in the appropriate 30-min period. This is the usual maximum-likelihood estimator of the transition matrix of a Markov chain. However, it also suffers from another form of censoring. When a biker's destination station is full, the biker must return the bike to another station, or wait until a rack becomes free. We do not observe, in our data, such events. Elegant methods for handling this destination censoring would be welcome.

Finally, we want to model the biking times. These are again relatively straightforward to estimate from data. We use linear regression, predicting the log of the trip durations (denoted $\ln(T_{ij})$, where i and j are station indices) seen in data to the log of the predicted cycling durations obtained from Google Maps (denoted $\ln(D_{ij})$). In this notation we suppress the fact that many station pairs have a very large number of individual trip durations corresponding to the station pair, and we use a single regression model to fit biking durations for *all* station pairs. The resulting distribution G_{ij} of trip durations between Stations i and j is lognormal, with parameters obtained from the regression. Figure 18.2 shows a scatter plot in log scale with the fitted regression line (right plot) based on 85% of the data that lies inside the central ellipse (left plot). We used this method to attempt to ensure that the fit is not unduly influenced by extreme data points that may represent data errors. This approach works for all station pairs where $i \neq j$. When $i = j$, so that the biker returns the bike to the same station, we again use lognormal trip durations, but this time with parameters that are specific to station i .

All of these fitting methods must be used with caution. In particular, Citi Bike operations are heavily seasonal, and have been in a constant state of growth since operations began in 2013, not only in the sense of an increasing footprint of stations,

but also in terms of the number of rides within an existing footprint. Accordingly, we use data from an appropriate window of time that matches, to some degree, the time period for which we are fitting rates. For example, when fitting parameters that are used to model operations during the summer months, we do not use winter data, and we do not use summer data from more than one year ago. We will say more about fitting rate parameters in the context of expansion planning in Sect. 18.6.

There is one more challenge that we have yet to overcome. Weather certainly has a pronounced effect on ridership. When we use our models to make decisions for, e.g., the next day, we could exploit the very reliable weather forecasts that are available for that day. This is important because, e.g., if a storm is expected in the morning, then ridership will be down in the morning and bikes will not be transported in great numbers from regions where apartments are concentrated, like the East Village in New York, to regions associated with the workplace, like the Financial District in New York. Therefore, if the weather is also likely to be fine in the afternoon, then great shortages of bikes will likely ensue in the afternoon. In principle, this is easily handled by fitting parameters that are weather specific. Nevertheless, we are yet to overcome the challenges associated with doing so in our work with Citi Bike.

18.3 Motorized Rebalancing

We begin this section by first defining the user dissatisfaction function introduced by Raviv and Kolka (2013). Based upon this, we then present an optimization problem introduced by O'Mahony (2015) that finds the optimal placement of bikes across the system at a given time. The optimization problem and the solution thereof do not depend on the actual number of bikes at each station or the rebalancing means available; instead, it focuses on finding the allocation that could be attained with unlimited rebalancing resources. Thereafter, we focus on what can actually be attained in practice by sketching the various integer programming formulations that route resources in an attempt to achieve near-optimal allocations.

18.3.1 User Dissatisfaction Function

The user dissatisfaction function of Raviv and Kolka (2013) estimates the number of out-of-stock events at a station over a finite time interval as a function of the number of bikes present at the beginning of said interval. To do so, it defines a birth-death process that is bounded above by the capacity of the station, i.e., the number of docks there, to track the number of bikes over the course of the interval. Mathematically, this corresponds to an $M/M/1/K$ queue where K denotes the number of docks. In the birth-death process, bike rentals correspond to deaths and bike returns correspond to births. Whenever there are no bikes present, i.e., the birth-death

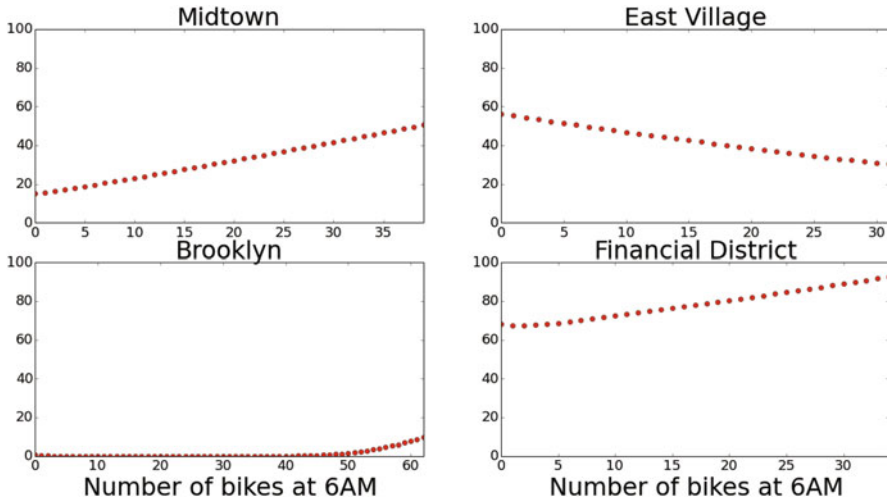


Fig. 18.3 User dissatisfaction functions for various stations in New York

process is at 0, a rental experiences an out-of-stock event and the objective increases by 1. Similarly, whenever the birth-death process is at the capacity of the station, meaning that all docks are occupied by bikes, an attempted bike return experiences an out-of-stock event, with respect to docks, and the objective increases by 1. Using, for example, the decensoring techniques described in the last section, the birth and death rates can be estimated from historical data. The work of Schuijbroek et al. (2017) and O'Mahony (2015) then provides different ways of using time-invariant rates to compute the expected number of out-of-stock events in an interval as a function of the initial number of bikes. Parikh and Ukkusuri (2015) observe that such intervals can be stitched together through standard stochastic recursion techniques, thus generalizing the techniques to piecewise constant rates. The resulting functions are convex as was first observed by Raviv and Kolka (2013) (cf. Fig. 18.3).

18.3.2 Optimal Allocation Before the Rush

Given the user dissatisfaction functions, it is natural to ask what the optimal allocation of bikes across the system looks like at the beginning of a time interval (e.g., at 6 a.m. on a weekday). We introduce the following notation in order to present the resulting optimization problem: for each station i , we denote its capacity by K_i . Thus, a station that holds b_i bikes at the beginning of the interval, has $K_i - b_i$ empty docks at that time. The user dissatisfaction function at such a station i is then written as $c_i(K_i - b_i, b_i)$. Writing c_i as a function of two variables, rather than just

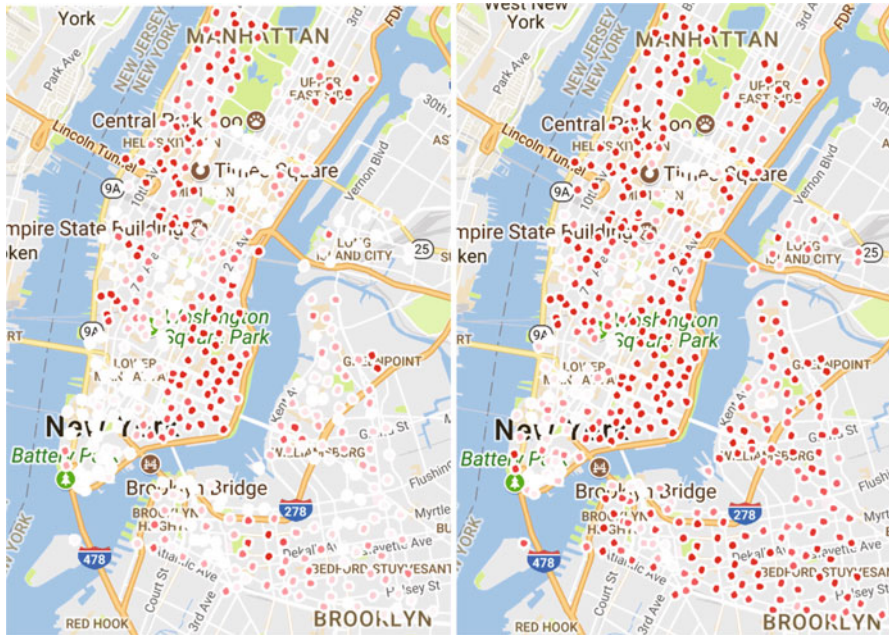


Fig. 18.4 Optimal allocation for $B \in \{5,000, 10,000\}$ at 6 a.m. in July 2017. The more red a station is, the larger the fraction of docks that is occupied in the optimal solution. (Map data: Google)

as a function of b_i , facilitates the use of consistent notation in this section and the next. Given a total number of bikes B , the optimal allocation of bikes across the system is then given by the solution to

$$\begin{aligned}
 & \underset{\mathbf{b}}{\text{minimize}} && \sum_i c_i(K_i - b_i, b_i) && (18.1) \\
 & \text{s.t.} && \sum_i b_i \leq B, \\
 & && 0 \leq b_i \leq K_i, \\
 & && b_i \in \mathbb{Z}.
 \end{aligned}$$

As one would expect, the resulting solution (cf. Fig. 18.4) suggests a great number of bikes should be placed in residential areas, e.g., the East Village and the Upper West Side, whereas very few bikes should be placed in commercial areas such as the Financial District or Midtown.

In preparation for the afternoon, this image is reversed; stations that should be full/empty in the morning should be empty/full in the afternoon (cf. Fig. 18.5).

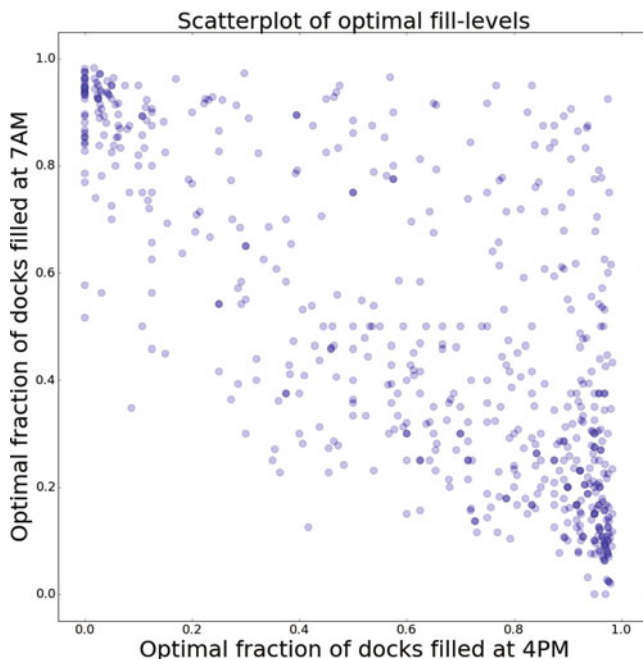


Fig. 18.5 Each point corresponds to a station, the coordinates correspond to the fraction of the station's docks that are filled in the optimal solution ($B = 12,000$) at 7 a.m. and at 4 p.m.

While the described optimization problem helps understand the needs of the system, the optimal configuration is very much unattainable. We now turn our attention to routing problems that combine the minimization of out-of-stock events with the constraints induced by limited rebalancing resources.

18.3.3 Resulting Routing Problems

The cleanest routing problem for motorized rebalancing in bike-sharing systems deals with rebalancing overnight. Since systems experience very little demand during the late evening and early morning hours (cf. Fig. 18.6), the hours in between give operators a large time window in which to prepare for the morning rush.

Just as in the previous optimization, the integer program suggested by Raviv et al. (2013) aims to minimize the expected number of out-of-stock events over the course of the rush. However, rather than having a budget constraint on the number of bikes in stations, we now have a fleet of capacitated trucks, each of which has an initial location, a final destination, and a number of bikes that are initially loaded. Given a bound on the length of the route of each truck to ensure that the routes have been completed by the time the rush begins, the integer program aims to find routes

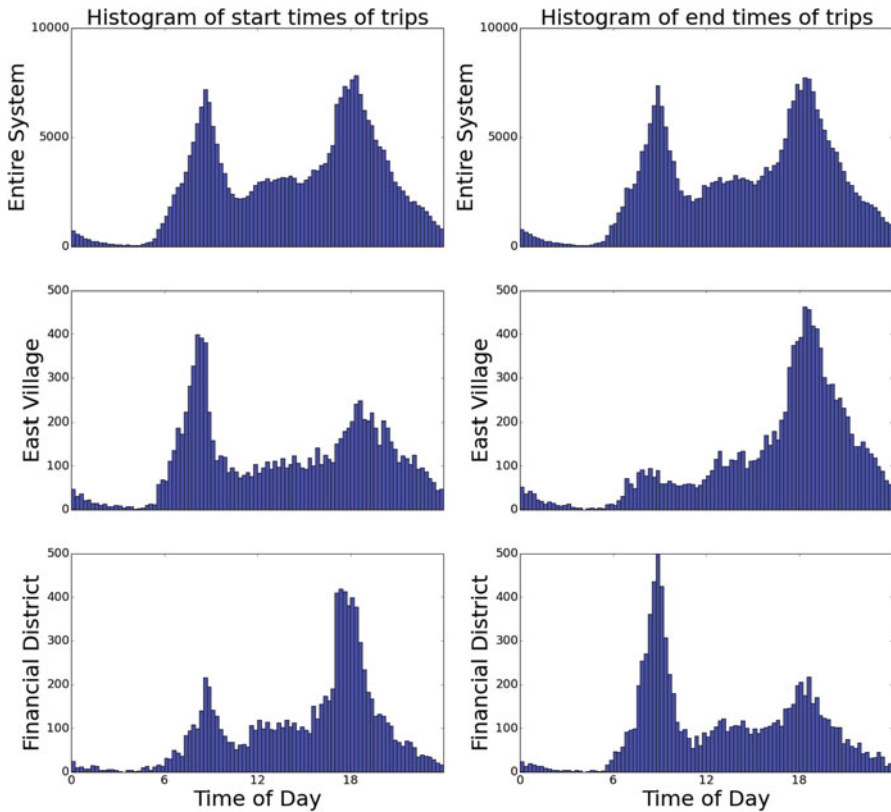


Fig. 18.6 Average number of trips starting on weekdays in the week of July 10th 2017 partitioned into 15-min buckets. The left (right)-hand side shows trips that start (end) at the given location

so that the resulting allocation minimizes the objective subject to no route being longer than the bound, and each truck starting and ending its route at the specified locations. Further, the number of bikes on each truck is bounded between 0 and its capacity at all times, and can only change when the truck is at a station – in that case, the change in number of bikes aboard the truck equals the negative of the change in number of bikes in the station. Subsequent to Raviv et al. (2013), several heuristics have been suggested by Forma et al. (2015), Ho and Szeto (2014), and Szeto et al. (2016) to solve larger instances.

Variations of the integer program have been studied as well, e.g., by Freund et al. (2016) who introduce a trade-off between travel-time and the number of stops to account for the time required for the trucks to park.

Rebalancing during rush hours is often more difficult than before the rush hour, which is one reason the focus of the rebalancing literature has been on overnight rebalancing. The difficulties of rebalancing during rush hours are mainly due to two complications. First, traffic congestion slows down trucks and thus reduces their

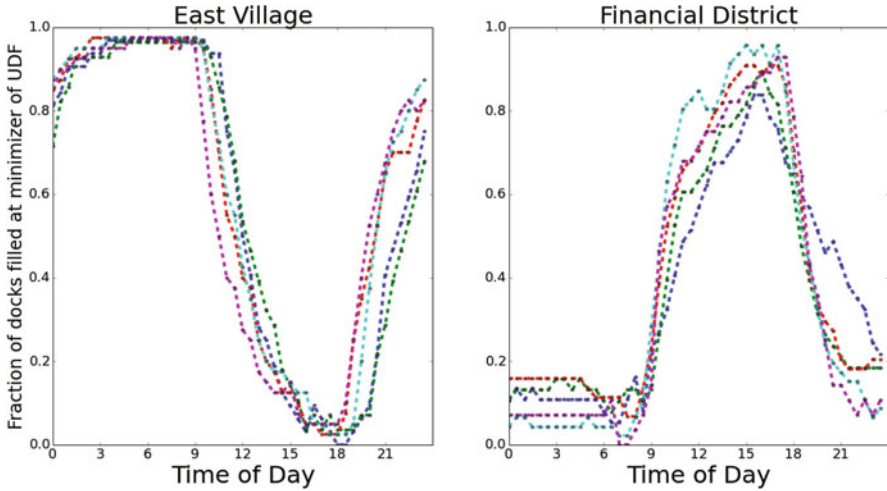


Fig. 18.7 Minimizer of the user dissatisfaction function at each half-hour interval over the course of the day for five stations in the East Village and five stations in the Financial District

efficiency. Second, since the system is dynamic, it is more difficult to plan ahead. For example, a route during rush hour may involve picking up bikes at a station that, by the time the truck arrives at that station, is already out of bikes. Moreover, towards the end of the rush, the question arises whether or not rebalancing should be aimed at the current rush or the subsequent period. To make this more precise, consider a station in a residential area like the East Village. During the morning (cf. Fig. 18.6), the station experiences much greater demand for rentals than for returns, so it is easy to see that before the rush, say at 6 a.m., rebalancing should ensure the station is close to being full (cf. Fig. 18.7). A little while into the rush, say at 8 a.m., this is still the case. However, by 10 a.m., rebalancing decisions with respect to routing are made that only take effect when rebalancing is actually performed some time later – by that time, say at 10:30 or 10:45 a.m., the demand has drastically changed, as can be observed in the sharp drop in the value of the minimizer of the user dissatisfaction function in Fig. 18.7. As the survey by de Chardon et al. (2016) points out, this issue can lead to rebalancing actions performed at the end of the morning rush hour being reversed by further rebalancing in the afternoon (or vice versa).

It is thus evident that it is easier to derive appropriate rebalancing formulations for rebalancing between peak traffic times than it is for rebalancing during peak traffic times. However, there is a trade-off between (1) rebalancing at the beginning of the peak when it is easy to identify where to add/take bikes, but congestion slows down rebalancing, (2) in between peaks when it is easy to identify where to add/take bikes, but the system is less imbalanced, (3) towards the end of peaks when the system is imbalanced, but it is much harder to identify where bikes should be taken/added. System operators rebalance during peak times and

require technological support for decision-making at those times. To that end, most academic approaches adapt the existing approaches for in-between rush-hour rebalancing in natural ways. For example, Kloimüller et al. (2014) adapt the approach of Rainer-Harbach et al. (2013) and Raidl et al. (2013) under the assumption that demand in any time interval matches its expectation.

18.4 Allocating Capacity

Forecasting and optimal motorized rebalancing, especially with a focus on overnight activity and minimizing out-of-stock events, have received a great deal of academic attention, but related system-design questions have not. In this section we focus on one of these questions, namely the question how many docks should be allocated to each station.

Given the user dissatisfaction function explored in the previous section, it is natural to focus on station capacities; after all, the user dissatisfaction function is a function of the number of docks at the station (cf. Fig. 18.8).

In this section, we first outline the integer program suggested by O’Mahony (2015) to find the optimal allocation of bikes and docks, for given user dissatisfaction functions. We then describe the results of Freund et al. (2017) extending the integer program and review structural results that allow us to efficiently solve the integer program to optimality. We next explore an extension of the user dissatisfaction, described by Freund et al. (2017), that captures the advantages of adding docks to *self-balancing* stations. At the end of the section we use real system-data to indicate what improvements might be possible through reallocation of dock capacity and discuss methods to evaluate the effect after implementation.

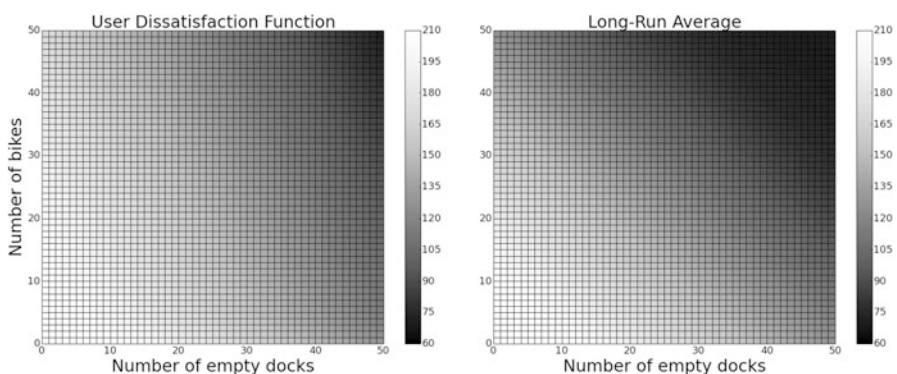


Fig. 18.8 User dissatisfaction function as a function of capacity and number of bikes and its long-run average, only a function of its capacity. Since the cost of the long-run average depends on K_i , but not on the initial number of bikes, it is constant on every diagonal $y = c - x$. Further, since its value on each diagonal is a convex combination of the values of the user dissatisfaction function on the diagonal, it is bounded between the minimum and the maximum along that diagonal

We assume that the demand at each station is given exogenously and can be estimated as in Sect. 18.2. In particular, we assume that our demand estimates for bike usage already capture the latent demand that could be served through additional capacity. Further, we assume that the demand for bike returns is not affected by added capacity elsewhere within the system. While the latter assumption seems strong, it is necessary to obtain a tractable problem. To justify this assumption, we refer to the work of Jian et al. (2016), which uses heuristic methods from simulation optimization whilst making only the first of the two assumptions; the results of Jian et al. (2016) indicate that the second assumption does not have a strong effect on the solutions.

18.4.1 Model formulation

Similar to the integer programs described in Sect. 18.3, we again aim to minimize the system-wide expected number of out-of-stock events. However, in contrast to the previous integer program, we now treat the number of docks K_i allocated to station i as a decision-variable with an associated budget K on the total number of docks.

$$\begin{aligned}
 & \underset{\mathbf{K}, \mathbf{b}}{\text{minimize}} && \sum_i c_i(K_i - b_i, b_i) && (18.2) \\
 & \text{s.t.} && \sum_i K_i \leq K, \\
 & && \sum_i b_i \leq B, \\
 & && \forall i \in [n] : 0 \leq b_i \leq K_i, \\
 & && \forall i \in [n] : l_i \leq K_i \leq u_i.
 \end{aligned}$$

It has been observed by Freund et al. (2017), as well as by Kaspi et al. (2017), that c_i is a multimodular function, meaning that it fulfills particular convexity/diminishing return properties. We refer the reader to Murota (2003) and the references therein for an overview of the literature related to multimodularity. It is known that multimodular functions can be efficiently minimized.

A practically relevant extension asks to minimize the objective whilst moving a limited number of docks, that is, given the current capacities \hat{K}_i for each station i and a reallocation budget \mathcal{R} it aims to minimize with the additional constraint $\sum_i |K_i - \hat{K}_i| \leq \mathcal{R}$. Freund et al. (2017) prove that by first solving the integer program with $l_i = \hat{K}_i = u_i$, i.e., without allowing capacities to change, and then running \mathcal{R} gradient-descent steps, with an appropriately defined notion of gradient, this extension can be solved optimally as well.

18.4.2 Long-Run Average

Part of the motivation to study the capacity allocation stems from the notion that limited resources to rebalance thwart attempts to attain the optimal allocation of bikes, given a fixed allocation of docks. A legitimate criticism of the integer program above thus questions whether the new allocation of docks, given by the optimal K_i s from the integer program, would give fewer out-of-stock events on a typical day, on which the stations are not optimally rebalanced (as in Fig. 18.4), even under the assumptions on the demand estimates mentioned before.

We address these concerns by explaining how Freund et al. (2017) extend the user dissatisfaction function to study a long-run average of the expected number of daily out-of-stock events in a regime in which no rebalancing occurs and the number of bikes at the beginning of a day is thus solely a function of the number of bikes at the beginning of the previous day and the realized demand over the course of the previous day. Mathematically, this extension is equivalent to computing the user dissatisfaction function over (infinitely) many days rather than just one. Computationally, Freund et al. (2017) compute transition probabilities between the possible numbers of bikes in a station, i.e., $0, 1, 2, \dots, K_i$, at the beginning of successive days. Obtaining the steady-state probability of the discrete Markov chain induced by these transition probabilities, the long-run average of the user dissatisfaction function can be calculated as a convex combination of the cost-values with $0, 1, 2, \dots, K_i$ bikes.

Given the extreme contrast between the two regimes (*perfect rebalancing* versus *no rebalancing*), it would not be surprising if the optimal dock allocations to the two were very different. In fact, one can easily construct examples of demand estimates for which the optimal solution for each of the two regimes is unboundedly far away from the optimal solution in the other regime. Such examples are based on stations with very asymmetric demand, that is, stations at which there is, over the course of a day, great demand to rent bikes, but none to return them (or vice-versa). At such a station, the perfect rebalancing regime can improve its objective through each additional dock by 1 as it is likely that such a dock, when initialized with a bike, will be used via a bike being rented. Since the long-run average reveals that the station would end up empty more often than not anyway, it does not reward an additional dock at such a station. Instead, the additional dock would be placed at a station that is more likely to balance itself over the course of a whole day and yet require the additional dock. In practice, however, such cases are rare; in the next paragraph we describe that operators can simultaneously improve both objectives.

In Table 18.1 we describe the results of the analysis for Boston (Hubway), Chicago (Divvy), and NYC (Citi Bike), based on Freund et al. (2017). The analysis is based on demand estimates from July 2016. The first three columns describe the size of the systems (in July 2016) in terms of number of stations, number of docks at all stations combined, and number of bikes. The next six columns provide the objectives under optimal allocation of bikes (c) and under long-run average (c^π) for three different allocations of docks: the current allocation of docks (Present), the

Table 18.1 Results of optimizing integer program (18.2) in New York, Chicago, and Boston with demand estimates stemming from July 2016 usage data. The columns headlined number summarize system statistics, the columns under c and c^π describe the objective with perfect rebalancing and long-run average respectively. The last column provides the number of docks that need to be moved to obtain the optimal allocation

City	Number			Present		1% of docks moved		OPT		Reallocated docks
	Stations	Docks	Bikes	c	c^π	c	c^π	c	c^π	
New York	449	14942	7500	4313	7275	4019	6985	2631	6231	2442
Chicago	582	9987	6000	1462	2340	1281	2165	0988	1978	0669
Boston	164	2861	1600	0614	0875	0567	0921	0479	0836	0213

allocation obtained by optimally moving 1% of the docks, and the optimal allocation of docks (OPT). Finally, the last column contains the number of docks that would need to be reallocated to get from the current allocation to the optimal one. Three interesting results arise from the analysis. One, moving even a small fraction of docks can significantly decrease the number of out-of-stock events. Two, moving to the optimal solution may require a large fraction of docks to be moved. Three, optimizing in either regime yields significant improvements in the other. Thus, bike-sharing systems can have their cake and eat it too: most of the improvement obtainable in the one regime is also obtained by optimizing for the other regime.

18.4.3 Measuring the Impact

The impact of reallocated capacity can be estimated in both an a priori and an a posteriori way. The *a priori analysis* is shown in Table 18.1: by comparing, for the current demand estimates, the objective of the continuous-time Markov chain model, we can estimate by how much reallocating docks would reduce out-of-stock events. Once docks have been reallocated, in the knowledge that demand patterns are likely affected, we can then estimate new demand rates, and compute an a posteriori estimate of the change in objective.

We can take the a posteriori analysis one step further by evaluating not only the objective of the user dissatisfaction function, based on a continuous-time Markov chain that uses demand estimates as its input, but rather use the observed, partially censored, demand explicitly. Here, we only sketch the idea of how this can be done.

Consider a station that used to have 40 docks, but had 10 docks taken from it. Suppose in the months after docks were taken away, the number of bikes in the station never gets to 30, i.e., even with the reduced capacity, no customer returning a bike ever experiences an out-of-stock event at the station. In that case, one can safely assume that taking the docks away did not cause any out-of-stock events. On the other hand, consider a station that used to have 30 docks and had 10 docks added to it. Suppose the station emptied out every morning, and then experienced

out-of-stock events, but then had 35 bikes returned in every afternoon rush-hour. In that case, one can estimate the reduction in out-of-stock events to be 10: in every afternoon rush-hour, 35 bikes are returned, 5 of which would not have been possible with the old capacity of 30. And in every morning rush-hour those 35 are all rented, only 30 of which would have been there if the capacity was only 30.

The idea outlined above to measure the impact a posteriori can be extended to estimate the number of out-of-stock events due to removed capacity, when that number is in fact greater 0. Implicit in this discussion is an assumption that no rebalancing occurs.

18.5 Beyond Motorized Rebalancing

While most operators of bike-sharing systems deploy trucks for rebalancing, some have developed targeted non-motorized techniques, which we briefly discuss here. We focus on Citi Bike's user incentive scheme, the routes of their trikes, and the placement of their corrals.

18.5.1 Incentives

Citi Bike ran a pilot of its incentive scheme, *Bike Angels*, in October 2015. In the pilot, the program targeted individual users with an incentive to drop off bikes at stations nearby to their usual trip destinations as well as to pick up bikes at stations nearby to their usual trip origins.

After running the pilot throughout the month of October, Citi Bike implemented a wider program that included many more stations. Eventually, it was also extended to all customers. In this new form, the program closely resembles the ideas outlined by O'Mahony (2015): stations are labeled to indicate whether a station needs additional rentals, additional returns, or is neutral. Originally, these labels were static in that stations kept their labels for every weekday morning rush hour (6 a.m.–12 p.m.) and each afternoon rush hour (4 a.m.–8 p.m.) over the course of two-week periods. This has since changed as the labeling of stations became dynamic in April 2017.

The decision to adopt dynamic labels involves a trade-off between two conflicting objectives: on the one hand, static labels provide a simpler user experience in which customers do not have to check the current state of the system to know whether they will be awarded points. On the other hand, static labels can lead to undesirable outcomes when either (i) customers are incentivized to drop off (pick up) bikes at stations that are full (empty) or (ii) customers are incentivized to drop off (pick up) bikes that lead to an increase in the expected number of out-of-stock events. The latter issue has been investigated, using Citi Bike's proprietary data, by Chung et al. (2018). In Fig. 18.9, we display the distribution of the impact on expected number of out-of-stock events due to rentals/returns for which points were awarded between

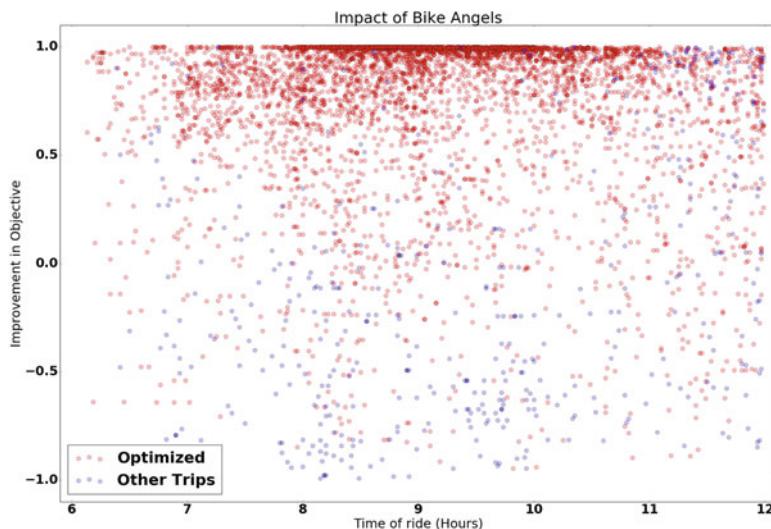


Fig. 18.9 Change in the user dissatisfaction function for each rental/return for which a point was awarded, mapped over the course of each morning rush hour, evaluated for October 15 through December 15, 2016

6 a.m. and noon in the last months of 2016. While the vast majority of such rides (awarded through a static policy) improved the state of the system, it is noticeable that some did not. On the one hand, a fully dynamic scheme would not award any points to rentals/returns with negative impact on the objective. On the other hand, the red points correspond to points that would have been awarded with a simple policy that determines an optimal time interval for each station over which rentals/returns are incentivized. Notably, the latter policy excludes many trips with negative impact and only few with positive impact.

18.5.2 Valets and Corrals

Valets ride trikes towing trailers that hold up to 18 bikes, which they move back and forth between stations – a picture of such a trailer can be found on page 2 of O’Mahony (2015). Due to the physical difficulties of moving the heavy trailers, the stations are not meant to be more than several blocks apart. Interestingly, as shown in Freund et al. (2016), under the right set of assumptions, one can extend the user dissatisfaction function to model the expected number of out-of-stock events for a pair of stations with a trailer moving bikes back and forth between them. Combining this with the expected number of out-of-stock events without a trailer, one can obtain, for each feasible pair of stations, the expected improvement due to adding a trailer. Finding the optimal trailer routes for k trailers then reduces to finding the

maximum-weight matching of cardinality k on the graph induced by stations as vertices and weighted edges between feasible pairs of stations, where the weight of each edge is the improvement due to the trailer between them.

Another non-motorized instrument Citi Bike and other operators take advantage of are so-called corrals. Corrals artificially extend the capacity of individual stations by having one employee store bikes in the space between racks and look after them over the course of the day. In the morning rush-hour, this allows additional bikes to be returned, which can then, in the afternoon rush-hour, be rented again. This approximately triples the capacity of the station. One could compute the expected change in out-of-stock events due to adding a corral using the framework of the user dissatisfaction functions; this would only require computing for each station the long-run average of the user dissatisfaction function with three times its current capacity. However, doing so would not be in line with the explicit purpose operators pursue with the use of corrals: rather than minimizing the number of out-of-stock events at one station, the aim is to minimize the system-wide number of return-related out-of-stock events. Indeed, a corral at a station i can reduce (or mitigate the effects of) out-of-stock events at nearby stations in addition to those at station i . In Freund et al. (2016), a maximum-coverage integer program is introduced that captures this objective.

18.6 Expansion Planning

Bikesharing programs are expanding worldwide, either through the introduction of a new program to a city, or through the expansion of an existing program.

When a new bike-sharing program is introduced to a new city, one must decide the number, location and capacity (number of docks) of stations, in addition to questions relating to the number of bikes that are required. In such settings, there is no existing data on bike-sharing demand. Indeed, we faced such a situation when Motivate asked us to predict the level of demand they might see in a bike-sharing network in San Francisco and the Bay Area. While there was an existing bike-sharing network in place, it was so small as to be of no help in predicting what might arise in a full-scale implementation. In this setting, perhaps the most viable tool, especially now that there are a large number of bike-sharing programs around the world, is to gauge one's target locale (in our study, San Francisco) against established programs in other cities. We used a regression model that predicted bike rides as a function of city population, demographic data, geography, and climate. We do not report on this study in detail here, instead simply noting that our resulting estimates of demand in San Francisco were very close to estimates produced internally, through a different method, at Motivate.

It is not possible to *perfectly* scale a new bike-sharing venture because of the difficulty in estimating usage rates. Fortunately, it is usually possible to “evolve” a bike-sharing design, expanding or contracting stations by adding or removing docks as demand patterns become clear. Indeed, in NYC, most of Citi Bike's docks come

in sets of 3 or 4, and are freestanding. The primary constraint is real estate, which is not owned by Citi Bike. Accordingly, Citi Bike must negotiate with the owners, often the city, to obtain clearance to change the capacity of bike stations.

The creation of the Citi Bike system in NYC has proceeded in a series of stages. In each stage, a new set of stations and bikes are installed in a geographical area not previously served, and some adjustments are made to the existing network to account for potential changes in flows. The initial installation in Manhattan and Brooklyn in May 2013 was complemented by major expansions in August 2015 and August 2016 of 140 stations each. For each of these expansions, the nominal flow (nominal demand) rates are unknown between new stations, but also between new and existing stations. Furthermore, the nominal flow rates between existing stations may change as a result of the expansion, although we typically expect these latter changes to be modest.

We now review the key ideas outlined by Singhvi et al. (2015) for exploiting existing ridership data in NYC, along with the models discussed above, to help Citi Bike in expansion planning.

Let S be the set of existing stations and S' be the set of new stations. Let $(\lambda_{ij}: i, j \in S)$ denote the flow rates between existing stations *before* expansion. Let $(v_{ij}: i, j \in S \cup S')$ denote the flow rates between all pairs of stations *after* expansion. (For simplicity we suppress the dependence of these rates on the hour of the week.) We pursued the following agenda.

1. Use demographic data, taxi usage data and potentially other data sources to estimate the parameters in a regression model that predicts the nominal flow rates between *existing* stations, $(\lambda_{ij}: i, j \in S)$, as a function of those data. Here we exploit the fact that these stations have been in operation for some time, so that rates are quite well known. This model is not very effective when used at the station level. However, when stations are aggregated into *neighborhoods*, which are simply collections of stations, and rates are estimated between neighborhoods, the predictions are more accurate.
2. Assume the nominal flow rates between existing stations do not change after expansion, i.e., take $(v_{ij}: i, j \in S) = (\lambda_{ij}: i, j \in S)$.
3. Use the estimated regression model to predict the remaining nominal flow rates after expansion, i.e., the flows from existing to new stations $(v_{ij}: i \in S, j \in S')$, the flows from new to existing stations $(v_{ij}: i \in S', j \in S)$, and the flows between new stations $(v_{ij}: i \in S', j \in S')$. In doing so, work first at the neighborhood level to predict inter-neighborhood trips, and then disaggregate these trips into station-specific trips.
4. Compute the user dissatisfaction curves for each number of bikes and each number of docks for the station-specific rates.
5. Solve the optimization problem (18.2) to determine how to allocate bikes and docks to new and existing stations. In doing so, one can fix the dock allocations to existing stations if modifying those allocations is deemed difficult or impossible.

The step in the above process that is perhaps most prone to error is that where we disaggregate the inter-neighborhood predictions to obtain station-specific predictions. This stage of the calculation could potentially benefit from further research.

Given that we are so uncertain about the inter-station rates after expansion, our approach of using point estimates for those quantities in the optimization problem 18.2 is a second point of concern. An intriguing potential area of future research is to explore versions of the optimization problem with robustness properties.

18.7 Conclusion

This chapter has introduced and discussed a range of logistical questions and strategic concerns related to bike-sharing programs that operate a collection of stations consisting of finite-capacity stations. The methods we have discussed have seen extensive use with Citi Bike in New York City. Beyond the work presented here, and still with the type of bike-sharing program operated by Citi Bike, there are at least two important directions that are yet to be fully explored. First, alternative pricing models might have the potential to better match supply and demand of both bikes and racks. One can also imagine a study of revenue management techniques for bike-sharing systems. A key challenge is that the per-ride cost of bike-sharing to a user is very small, so there is little room to adjust prices, and perhaps low sensitivity to price mechanisms. Second, bikes are often viewed as a solution to the “last mile” problem, where commuters use other modes of transport, such as the subway, in their daily commutes. Can one develop transport planning infrastructure that allows the planning of end-to-end travel, incorporating various modes of transport that might also include taxis and other modes?

Related planning methods for alternative bike-sharing designs are largely unexplored. Turnkey systems do not use docks, instead registering a bike as returned if it is within a certain distance of a fixed station. Such stations can be approximately modeled with the methods in this chapter if one simply assumes their capacity to be very large, or even infinite. While this should be sufficient for many locations, there would still be stations where demand is so great that this assumption is inappropriate. What then? Moreover, the chaos that can ensue when bikes are left in haphazard fashion may be untenable for certain locations, so this modeling approach may have to be augmented with other ideas. Still other bike-sharing programs employ GPS units on bikes, and bikes can be left in any location. Modeling approaches to support managerial decision making in such systems are just beginning to be explored.

Acknowledgements We thank our colleagues at Citi Bike, and its parent company Motivate, for our strong and ongoing collaboration. We also thank the many contributors to the work described herein, especially the students, both undergraduate and graduate, at Cornell University. This work was partially supported by National Science Foundation grants CCF-1526067, CMMI-1537394, CCF-1522054, and CCF-1740822, and Army Research Office grant W911NF-17-1-0094.

References

- Banerjee S, Freund D, Lykouris T (2017) Pricing and optimization in shared vehicle systems: an approximation framework. In: Proceedings of the 2017 ACM conference on economics and computation. ACM, p 517. arXiv preprint:1608.06819
- Benchimol M, Benchimol P, Chappert B, De La Taille A, Laroche F, Meunier F, Robinet L (2011) Balancing the stations of a self service “bike hire” system. *RAIRO-Oper Res* 45(1):37–61
- Bulhões T, Subramanian A, Erdoğan G, Laporte G (2018) The static bike relocation problem with multiple vehicles and visits. *Eur J Oper Res* 264:508–523
- Casazza M, Ceselli A, Calvo RW (2017) Inventory rebalancing in bike-sharing systems. In: Proceedings of the 15th cologne-twente workshop on graphs and combinatorial optimization, pp 35–38
- Chemla D, Meunier F, Calvo RW (2013) Bike sharing systems: solving the static rebalancing problem. *Discret Optim* 10(2):120–146
- Chen L, Zhang D, Wang L, Yang D, Ma X, Li S, Wu Z, Pan G, Thi-Mai-Trang Nguyen JJ (2016) Dynamic cluster-based over-demand prediction in bike sharing systems. In: Proceedings of the 2016 ACM international joint conference on pervasive and ubiquitous computing. ACM, pp 841–852
- Chung H, Freund D, Shmoys DB (2018) Bike angels: an analysis of Citi Bike’s incentive program. In: Proceedings of the 1st ACM SIGCAS conference on computing and sustainable societies. ACM, pp 5:1–5:9
- Çınlar E (1972) Superposition of point processes. In: Lewis PAW (ed) *Stochastic point processes: statistical analysis, theory, and applications*. Wiley Interscience, New York, pp 549–606
- Contardo C, Morency C, Rousseau L-M (2012) Balancing a dynamic public bike-sharing system. Technical report, CIRRELT, Sept 2012. <https://www.cirrelt.ca/Documents\penalty0Travail\CIRRELT2012-09.pdf>. Accessed in Mar 2018
- Datner S, Raviv T, Tzur M, Chemla D (2017) Setting inventory levels in a bike sharing network. *Transp Sci. Articles in Advance*. <https://doi.org/10.1287/trsc.2017.0790>
- de Chardon CM, Caruso G, Thomas I (2016) Bike-share rebalancing strategies, patterns, and purpose. *J Transp Geogr* 55:22–39
- Dell’Amico M, Hadjicostantinou E, Iori M, Novellani S (2014) The bike sharing rebalancing problem: mathematical formulations and benchmark instances. *Omega* 45:7–19
- Di Gaspero L, Rendl A, Urili T (2013) Constraint-based approaches for balancing bike sharing systems. In: Schulte C (ed) *Proceedings of the 19th international conference on principles and practice of constraint programming*. Springer, pp 758–773
- Erdoğan G, Laporte G, Calvo RW (2014) The static bicycle relocation problem with demand intervals. *Eur J Oper Res* 238(2):451–457
- Erdoğan G, Battarra M, Calvo RW (2015) An exact algorithm for the static rebalancing problem arising in bicycle sharing systems. *Eur J Oper Res* 245(3):667–679
- Forma IA, Raviv T, Tzur M (2015) A 3-step math heuristic for the static repositioning problem in bike-sharing systems. *Transp Res B Methodol* 71:230–247
- Freund D, Norouzi-Fard A, Paul A, Wang C, Henderson SG, Shmoys DB (2016) Data-driven rebalancing methods for bike-share systems. Working paper
- Freund D, Henderson SG, Shmoys DB (2017) Minimizing multimodular functions and allocating capacity in bike-sharing systems. In: Eisenbrand F, Koenemann J (eds) *Integer programming and combinatorial optimization proceedings*. Lecture notes in computer science, vol 10328. Springer, pp 186–198. arXiv preprint arXiv:1611.09304
- George DK, Xia CH, Squillante MS (2012) Exact-order asymptotic analysis for closed queueing networks. *J Appl Probab* 49(2):503–520
- Ghosh S, Trick M, Varakantham P (2016) Robust repositioning to counter unpredictable demand in bike sharing systems. In: Kambhampati S (ed) *Proceedings of the twenty-fifth international joint conference on artificial intelligence*. IJCAI/AAAI Press, pp 3096–3102

- Ghosh S, Varakantham P, Adulyasak Y, Jaillet P (2017) Dynamic repositioning to reduce lost demand in bike sharing systems. *J Artif Intell Res* 58:387–430
- Ho SC, Szeto W (2014) Solving a static repositioning problem in bike-sharing systems using iterated tabu search. *Transp Res E Logist Transp Rev* 69:180–198
- Hsu YT, Kang L, Wu YH (2016) User behavior of bikesharing systems under demand–supply imbalance. *Transp Res Rec J Transp Res Board* (2587):117–124
- Jian N, Henderson SG (2015) An introduction to simulation optimization. In: Yilmaz L, Chan WKV, Roeder TMK, Macal C, Rosetti M (eds) *Proceedings of the 2015 winter simulation conference*. IEEE, pp 1780–1794
- Jian N, Freund D, Wiberg H, Henderson SG (2016) Simulation optimization for a large-scale bike-sharing system. In: Roeder TMK, Frazier PI, Szechtman R, Zhou E (eds) *Proceedings of the 2016 winter simulation conference*. IEEE, Piscataway, pp 602–613
- Kabra A, Girotra K, Belavina E (2015) Bike-share systems: accessibility and availability. Working paper
- Karlin S, Taylor HM (1975) *A first course in stochastic processes*, 2nd edn. Academic, Boston
- Kaspi M, Raviv T, Tzur M (2016) Detection of unusable bicycles in bike-sharing systems. *Omega* 65:10–16
- Kaspi M, Raviv T, Tzur M (2017) Bike-sharing systems: user dissatisfaction in the presence of unusable bicycles. *IIEE Trans* 49(2):144–158
- Kloimüller C, Papazek P, Hu B, Raidl GR (2014) Balancing bicycle sharing systems: an approach for the dynamic case. In: Blum C, Ochoa G (eds) *European conference on evolutionary computation in combinatorial optimization*. Springer, pp 73–84
- Li Y, Zheng Y, Zhang H, Chen L (2015) Traffic prediction in a bike-sharing system. In: Bao J, Sengstock C, Ali ME, Huang Y, Gertz M, Renz M, Sankaranarayanan J (eds) *Proceedings of the 23rd SIGSPATIAL international conference on advances in geographic information systems*. ACM, pp 33:1–33:10
- Liu J, Sun L, Chen W, Xiong H (2016) Rebalancing bike sharing systems: a multi-source data smart optimization. In: Krishnapuram B, Shah M, Smola AJ, Aggarwal CC, Shen D, Rastogi R (eds) *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. ACM, pp 1005–1014
- Lowalekar M, Varakantham P, Ghosh S, Jena SD, Jaillet P (2017) Online repositioning in bike sharing systems. In: Barbulescu L, Frank J, Mausam, Smith SF (eds) *Proceedings of the 27th international conference on automated planning and scheduling (ICAPS)*. AAAI Press, pp 200–208
- Murota K (2003) *Discrete convex analysis*. SIAM monographs on discrete mathematics and applications. Society for Industrial and Applied Mathematics, Philadelphia
- Nair R, Miller-Hooks E, Hampshire RC, Bušić A (2013) Large-scale vehicle sharing systems: analysis of Vélib'. *Int J Sustain Transp* 7(1):85–106
- Nelson BL (2013) *Foundations and methods of stochastic simulation*. International series in operations research & management science, vol 187. Springer, New York
- O'Mahony E (2015) Smarter tools for (Citi) bike sharing. Ph.D. thesis, Cornell University, Ithaca
- O'Mahony E, Shmoys DB (2015) Data analysis and optimization for (Citi) bike sharing. In: *Twenty-ninth AAAI conference on artificial intelligence*, pp 687–694
- Parikh P, Ukkusuri S (2015) Estimation of optimal inventory levels at stations of a bicycle sharing system. In: *Transportation Research Board 94th annual meeting*. Transportation Research Board
- Paul A, Freund D, Ferber A, Shmoys DB, Williamson DP (2017) Prize-collecting TSP with a budget constraint. In: Pruhs K, Sohler C (eds) *25th annual European symposium on algorithms (ESA 2017)*, Schloss Dagstuhl–Leibniz-Zentrum für Informatik, Dagstuhl, Germany, Leibniz International Proceedings in Informatics (LIPIcs), vol 87, pp 62:1–62:14. <https://doi.org/10.4230/LIPIcs.ESA.2017.62>, <http://drops.dagstuhl.de/opus/volltexte/2017/7837>
- Raidl GR, Hu B, Rainer-Harbach M, Papazek P (2013) Balancing bicycle sharing systems: improving a VNS by efficiently determining optimal loading operations. In: Middendorf M, Blum C (eds) *International workshop on hybrid metaheuristics*. Springer, pp 130–143

- Rainer-Harbach M, Papazek P, Hu B, Raidl GR (2013) Balancing bicycle sharing systems: a variable neighborhood search approach. In: European conference on evolutionary computation in combinatorial optimization. Springer, pp 121–132
- Raviv T, Kolka O (2013) Optimal inventory management of a bike-sharing station. *IIE Trans* 45(10):1077–1093
- Raviv T, Tzur M, Forma IA (2013) Static repositioning in a bike-sharing system: models and solution approaches. *EURO J Transp Logist* 2(3):187–229
- Riquelme C, Johari R, Zhang B (2017) Online active linear regression via thresholding. In: Singh SP, Markovitch S (eds) Proceedings of the thirty-first AAAI conference on artificial intelligence. AAAI Press, pp 2506–2512
- Rudloff C, Lackner B (2014) Modeling demand for bikesharing systems: neighboring stations as source for demand and reason for structural breaks. *Transp Res Rec J Transp Res Board* 2430:1–11
- Salaken SM, Hosen MA, Khosravi A, Nahavandi S (2015) Forecasting bike sharing demand using fuzzy inference mechanism. In: Arik S, Huang T, Lai WK, Liu Q (eds) ICONIP 2015: proceedings of the 22nd international conference on neural information processing. Springer, pp 567–574
- Saltzman RM, Bradford RM (2016) Simulating a more efficient bike sharing system. *J Supply Chain Oper Manag* 14(2):36–47
- Schuijbroek J, Hampshire R, van Hoesel WJ (2017) Inventory rebalancing and vehicle routing in bike sharing systems. *Eur J Oper Res* 257(3):992–1004
- Shu J, Chou MC, Liu Q, Teo CP, Wang IL (2013) Models for effective deployment and redistribution of bicycles within public bicycle-sharing systems. *Oper Res* 61(6):1346–1359
- Singhvi D, Singhvi S, Frazier PI, Henderson SG, O’Mahony E, Shmoys DB, Woodard DB (2015) Predicting bike usage for New York City’s bike sharing system. In: Dilkina B, Ermon S, Hutchinson RA, Sheldon D (eds) AAAI workshop: computational sustainability. AAAI Press
- Singla A, Santoni M, Bartók G, Mukerji P, Meenen M, Krause A (2015) Incentivizing users for balancing bike sharing systems. In: Bonet B, Koenig S (eds) Proceedings of the twenty-ninth AAAI conference on artificial intelligence. AAAI Press, pp 723–729
- Szeto W, Liu Y, Ho SC (2016) Chemical reaction optimization for solving a static bike repositioning problem. *Transp Res D Transp Env* 47:104–135
- Vogel P, Saavedra BAN, Mattfeld DC (2014) A hybrid metaheuristic to solve the resource allocation problem in bike sharing systems. In: Blesa MJ, Blum C, Voß S (eds) International workshop on hybrid metaheuristics. Springer, pp 16–29
- Waserhole A, Jost V (2016) Pricing in vehicle sharing systems: optimization in queuing networks with product forms. *EURO J Transp Logist* 5(3):293–320
- Zhang J, Pan X, Li M, Philip SY (2016) Bicycle-sharing system analysis and trip prediction. In: Chow C, Jayaraman PP, Wu W (eds) 2016 17th IEEE international conference on mobile data management (MDM), vol 1. IEEE, pp 174–179

Chapter 19

Operations Management of Vehicle Sharing Systems



Long He, Ho-Yin Mak, and Ying Rong

Abstract The emerging sharing economy has encouraged the rapid rise of vehicle sharing businesses. Much of this growth is due to the innovation of the free-float model, which allows users to start and end rentals at any location within a defined service region. Compared with conventional models of vehicle sharing, the free-float model offers its users the flexibility to make one-way, two-way and multi-stop trips, and as a result offer a more viable alternative to individual vehicle ownership. On the other hand, the flexibility of free-float model leads to a number of operations management challenges that must be overcome for such vehicle sharing systems to be economically sustainable. In this chapter, we review several operations management problems in vehicle sharing including system design, vehicle repositioning, fleet sizing, dynamic pricing and reservation policy. In particular, we discuss the optimization models for service region design and fleet repositioning in details.

19.1 Introduction

Recent novel business models emerging along with the sharing economy aim to improve resource utilization. Because of passenger cars' low utilization rates (on average idle 92% of the time Time Magazine 2012), high fixed costs to own (on average \$6,500 per year) and relatively low variable costs to operate, it becomes a

L. He (✉)

NUS Business School, National University of Singapore, Singapore, Singapore
e-mail: longhe@nus.edu.sg

H.-Y. Mak

Saïd Business School, University of Oxford, Oxford, UK
e-mail: ho-yin.mak@sbs.ox.ac.uk

Y. Rong

Antai College of Economics and Management, Shanghai Jiao Tong University, Shanghai, China
e-mail: yrong@sjtu.edu.cn

prime candidate for sharing business. The worldwide market size of car sharing is forecast to grow from 2.3 million users in 2013 to 12 million users in 2020 (Navigant Research 2013).

Different from the conventional station-based model, e.g., Zipcar, where customers are required to return vehicles to the same stations where the rental trip started, firms like car2go adopt a new operating model that allows one-way and multi-stop trips in their free-float car sharing systems. For instance, customers are able to check out and return cars anywhere within the service region at any on-street parking space. In free-float systems, customers can use smartphone apps to rent, reserve, and drop off vehicles on demand, wherever they choose. This innovative model in car sharing has also been extend to bike sharing systems, e.g, Mobike and ofo.

Meanwhile, electric vehicles (EVs) have been widely adopted in various vehicle sharing business models. For instance, car2go operates EV fleets in Amsterdam, Madrid and Stuttgart. Drivenow/Reachnow also operates with partially electric fleets in various cities in Europe and North America. In 2017, Telepod started its e-scooter sharing business with 7 stations in Singapore. EVs typically have higher fixed (purchase) costs but lower variable (fuel and maintenance) costs than their gasoline counterparts. The sharing business model allows the higher fixed costs to be shared among multiple users, and takes advantage of the lower variable costs through increased vehicle utilization. These innovative operations models allow more customers to experience EVs on a daily basis without long-term commitments, which possibly helps overcome the major barriers to EV adoption, such as the range and resale anxiety discussed in Lim et al. (2015).

While the free-float model and the use of EVs offers new opportunities in the vehicle sharing business, the resulting operations management can be challenging. In this chapter, we discuss several operations problems, including service region design, fleet sizing, fleet repositioning, dynamic pricing and reservation policy. We classify these into *strategic planning* and *operational decision* problems in Table 19.1. Due to the different decision time frames, the modeling approaches and data requirement for analyzing them differ significantly. In particular, while static (mostly mathematical programming) models are generally more suitable for strategic planning, the operational decision making usually calls for dynamic models (e.g., dynamic programming). We provide a brief discussion for each problem below.

Table 19.1 Operations management problems in vehicle sharing systems

Problems	Strategic decision	Operational decision
Service region design	X	
Fleet sizing	X	
Fleet repositioning		X
Dynamic pricing		X
Reservation management		X

Service Region Design A key strategic decision in free-float vehicle sharing systems is designing the service region. On the one hand, expanding geographical coverage attracts more service adoption and thus higher revenue. On the other hand, doing so brings operational challenges, such as the fleet repositioning to ensure availability under imbalanced demand and recharging operations in the case of EV fleets. The optimal service region design requires carefully modeling the trade-off between these factors.

Fleet Sizing To ensure high vehicle availability (i.e., service level of the system) and cost effectiveness, the operator must deploy an optimal fleet size given the service region (in free-float systems) or network of stations (in station-based systems). The model for fleet sizing should capture the relationship between vehicle availability, utilization and fleet size (i.e., cost). Depending on the flexibility of fleet deployment or redeployment, fleet sizing can be strategic or tactical decision. With the close connection with the geography of the service region, it is often appropriate to jointly optimize fleet sizing and service region design in strategic planning.

Fleet Repositioning Fleet management is one of the main operational challenges in vehicle sharing systems, especially for free-float systems. Under time-varying and spatially-unbalanced travel patterns of customers, the vehicles must be repositioned throughout the day to ensure availability at locations where customers need them. This requires analyzing and solving a dynamic optimization problem. One of the major challenges is to overcome the curse of dimensionality as the problem size grows with the spatial (geographical location) and temporal (planning horizon) dimensions. Such problem is common among various vehicle sharing systems, including both car sharing and bike sharing.

Dynamic Pricing The idea of dynamic pricing not only plays a big role in ride sharing (e.g., surge pricing by Uber), but also in vehicle sharing systems as well. Particularly in free-float systems, dynamic pricing can also serve as an instrument to encourage customers to move vehicles and rebalance the fleet. Drivenow offers discounts to customers based on the location and time of the vehicle being rented. Optimizing these discount offerings dynamically could incentivize customers to help reposition vehicles from less desirable locations to where they are needed.

Reservation Management While vehicle sharing systems generally provide on-demand service, it is often important to allow customer reservations to enhance operational efficiency. In station-based sharing systems, the number of parking spaces is limited in each station. To better manage capacity, operators design parking reservation policies to better utilize the parking lots and improve the system's flexibility by allowing one-way rental where parking availability at the destination becomes a key constraint. On the other hand, in free-float vehicle sharing systems, allowing for advanced reservation can enhance customer experience by eliminating the uncertainty on availability. In this case, the operator faces the problem of dynamically allocating capacity between advanced reservations and real-time demand arrivals. Designing effective policies, e.g., limiting the number and duration of reservations, are important to balance customer satisfaction and operational efficiency of the system.

In the next sections, we cover these problems in more details by reviewing some of the existing work and discussing some potential research directions.

19.2 Service Region Design

Location design for vehicle sharing systems is critical in attracting customer demand as well as managing operational costs. We first review the strategic decision of service region design in free-float vehicle sharing systems. In particular, we consider a sharing system equipped with EVs where recharging scheduling is also involved, based on a recent paper He et al. (2017).

The strategic planning problem of service region design for free-float EV sharing systems entails two major challenges. First, the travel pattern and adoption behavior of potential customers are highly uncertain at the planning stage. Moreover, the firm may not possess accurate data to describe the demand uncertainty before its operations for a sufficient period, that further exacerbates such challenge. Because strategic commitments, e.g., the acquisition of land for stations and charging outlets, are often made in conjunction with service region design, a robust planning methodology is imperative. Second, the operational details of EV sharing, such as repositioning and recharging of EVs, depend on the demands from the service region. Hence, the firm must also conscientiously account for operational cost drivers when determining the service region with only limited data available. Taking these factors into account, He et al. (2017) propose an integrated service region design model that considers customers' satisficing behavior in service adoption together with various operational characteristics of a free-float EV sharing system.

19.2.1 Basic Model

He et al. (2017) consider an urban area consisting of a set N of non-overlapping geographical locations (e.g., districts). The operator chooses a subset of N to be its service region. We use binary decision variables x_i to denote whether location i is covered in the service region ($x_i = 1$) or not ($x_i = 0$). Throughout this chapter, we use boldface letters to denote matrices or vectors consisting of scalar parameters or variables denoted by the same letter. Therefore, \mathbf{x} is the vector that consists of x_i for $i \in N$. We also define the inner product $\langle \mathbf{A}, \mathbf{B} \rangle = \text{trace}(\mathbf{BA})$, where \mathbf{A} and \mathbf{B} are two matrices.

We assume that customers exhibit *satisficing* behavior in their service adoption. Suppose a_{ij} is a customer's utility of being able to travel from origin i to destination j , which depends on both the trip frequency and the trip value to the customer. Therefore, once the origin i is in the service region, the customer has a potential utility of $\sum_{j \in N} a_{ij}x_j$ by adopting the service. Under the satisficing behavior framework in Simon (1957), the customer would adopt the service when the total

utility from the service exceeds a certain threshold b , known as the aspirational level. Hence, if the origin i is served, the customer adoption decision can be modeled as the indicator function below:

$$\mathbf{1}\left(\sum_{j \in N} a_{ij}x_j \geq b\right) = \begin{cases} 1, & \text{if } \sum_{j \in N} a_{ij}x_j \geq b \\ 0, & \text{otherwise.} \end{cases}$$

The utility values a_{ij} are heterogeneous among individual customers, even if they reside in the same region. From the operator’s perspective, the values of a_{ij} can be viewed as independent realizations of a random variable. Thus, we define the adoption rate at the aggregate level as the probability of customer adoption or the proportion of customers who adopt. The adoption rate q_i can be written as the expectation over a_{ij} in the indicator function:

$$q_i = \mathbb{E}\left[\mathbf{1}\left(\sum_{j \in N} a_{ij}x_j \geq b\right)\right] = \text{Prob}\left(\sum_{j \in N} a_{ij}x_j \geq b\right).$$

The operator’s profit consists of membership revenue and operational profit. Each customer pays an fixed annual membership fee f when signing up for the service, and pays for usage based on the duration of each rental. There is also a fixed cost g_i to the operator to serve region i . This may include infrastructure investment, e.g., for EV charging equipment, and parking permit costs. The service region design problem is then formulated as follows:

$$\max_{\mathbf{q}, \mathbf{x}} \left\{ \sum_{i \in N} f Q_i q_i - \sum_{i \in N} g_i x_i + \Theta(\mathbf{q}, \mathbf{x}) \right\} \tag{19.1}$$

$$\text{s.t. } q_i \leq \text{Prob}\left(\sum_{j \in N} a_{ij}x_j \geq b\right), \quad \forall i \in N \tag{19.2}$$

$$q_i \leq x_i, \quad \forall i \in N \tag{19.3}$$

$$x_i \in \{0, 1\}, \quad \forall i \in N.$$

The objective is to maximize the expected net profit in Eq.(19.1), i.e., total revenue less fixed and operational costs from charging, repositioning and fleet investment. For notational brevity, we temporarily denote the operational profit by $\Theta(\mathbf{q}, \mathbf{x})$, which is a function if the service region \mathbf{x} and associated adoption rates \mathbf{q} . The expected revenue from membership fees is given by $\sum_{i \in N} \sum_{k \in K} f Q_i q_i$, where Q_i is the market size in region i . Constraint (19.2) defines the adoption rate as the probability that a customer’s potential utility exceeds the aspirational level. Furthermore, constraint (19.3) stipulates that a location must be covered for its residents to adopt the service.

19.2.2 Customer Adoption

We derive a tractable formulation for the probability constraint (19.2). To evaluate the adoption rate exactly, complete information on the joint distribution of \mathbf{a} is required. In practice, however, perfect information is often unavailable when strategic decisions, such as service region design, are made. Specifically, as the firm only possesses limited operations data (e.g., from pilot studies or surveys) in the planning stage, fitting the joint distribution of travel patterns with confidence is generally not an easy task. Furthermore, in terms of computational tractability, the total utility $\sum_{j \in N} a_{ij}x_j$ is difficult to evaluate in general due to the need for convolutions, even if one assumes that the distributions are known and the components of \mathbf{a} are independent. Therefore, He et al. (2017) adopt a pragmatic approach that features both distributional robustness and computational tractability under limited information in their optimization model.

Relaxing the data requirement, He et al. (2017) assume that the operator only has knowledge of certain descriptive statistics of \mathbf{a} , in particular, their means and covariance matrices. They study a robust model that provides the worst-case adoption rate, i.e., the lowest adoption rate among all possible distributions \mathcal{P} of the utility parameters \mathbf{a} with the given means and covariance matrices:

$$q_i \leq \inf_{p \in \mathcal{P}} \mathbb{E}_p \left[\mathbf{1} \left(\sum_{j \in N} a_{ij}x_j \geq b \right) \right]. \tag{19.4}$$

The utility parameter a_{ij} is a nonnegative random variable, as utility would not be reduced as more locations are covered in the service region. In the following, we use $\mathbf{a}_i = (a_{ij})_{j \in I}$ to denote the random vector of random utility parameters associated with the origin i . We assume that the mean vector $\bar{\mathbf{a}}_i = (\bar{a}_{ij})$ and covariance matrix $\Gamma_i = [\text{cov}(a_{ij_1}, a_{ij_2})]$ of \mathbf{a}_i are known for each location $i \in I$. Then, the second moment matrix Σ_i is given by:

$$\Sigma_i := \mathbb{E} \begin{bmatrix} \mathbf{a}_i \\ 1 \end{bmatrix} \begin{bmatrix} \mathbf{a}_i \\ 1 \end{bmatrix}^T = \begin{bmatrix} \mathbf{S}_i & \bar{\mathbf{a}}_i \\ \bar{\mathbf{a}}_i^T & 1 \end{bmatrix}, \quad \text{where } \mathbf{S}_i := \Gamma_i + \bar{\mathbf{a}}_i \bar{\mathbf{a}}_i^T.$$

We assume that the covariance matrix is positive semidefinite (PSD), i.e., $\Gamma_i \geq 0$, which implies that the second moment matrix is also PSD, i.e., $\Sigma_i \geq 0$. With a given \mathbf{x} , the worst-case adoption rate can be obtained by solving a convex optimization problem with PSD constraints. In other words, certain constraints in the formulation take the form of requiring certain decision variable matrices to be in the cone of PSD matrices. Note that a symmetric matrix \mathbf{M} is said to be PSD (denoted $\mathbf{M} \geq 0$) if it satisfies $\mathbf{v}^T \mathbf{M} \mathbf{v} \geq 0, \forall \mathbf{v} \in \mathbb{R}^n$.

Constraint (19.2) in Problem (19.1) is approximated by constraint (19.4) that is further transformed into second-order cone constraints. A computationally tractable formulation for the worst-case adoption rate in Eq. (19.4) is provided in Proposition 1.

Proposition 1 (Proposition 1 in He et al. 2017) *Adoption rate q_i satisfies the following set of inequalities, with second-order cone constraints (19.5), if and only if it satisfies constraint (19.4).*

$$\left\| \begin{array}{c} 1 - q_{ik} - v_{ik} \\ 2\Gamma_i^{1/2} \mathbf{x} \end{array} \right\|_2 \leq 1 - q_i + v_i, \quad \forall i \in N \quad (19.5)$$

$$v_i = b^2 + \sum_{(j_1, j_2) \in I \times I} (\bar{a}_{ij_1} \bar{a}_{ij_2} + \sigma_{ij_1 j_2}) z_{j_1 j_2} - 2b \sum_{j \in N} \bar{a}_{ij} x_j, \quad \forall i \in N \quad (19.6)$$

$$(\mathbf{x}, q_i) \in \mathcal{X}_i, \quad \forall i \in N, \forall k \in K \quad (19.7)$$

$$(z_{j_1 j_2}, x_{j_1}, x_{j_2}) \in \mathcal{Z}, \quad \forall j_1, j_2 \in I \quad (19.8)$$

where v_i and $z_{j_1 j_2}$ are auxiliary decision variables, and \mathcal{Z} and \mathcal{X}_i are feasible regions characterized by linear constraints provided in He et al. (2017).

19.2.3 Operational Profit

It is essential to maintain service level in free-float sharing systems. Let the service level α be the probability that customers will find available EVs at their origins when they intend to travel. For the ease of discussion, we consider homogeneous service level α across regions the discussion below. To consider location-specific service level, one can replace α with α_i for region i .

We describe the system dynamics using the closed queueing network in Fig. 19.1a. Different from the typical models of call centers, in this queueing network, the EVs, rather than the customers, are the considered entities. In a closed network, the total number of EVs is fixed. There are four possible states of an EV at any given time instant: (i) being idle and awaiting the next customer rental, (ii) in transit from one idle node (i.e., the region) to another with a customer, (iii) in transit from one idle node to another due to repositioning, or (iv) recharging at a charging station. Thus, we define four types of queues for these four types of activities. In what follows, we use the terms “queue” and “node” interchangeably.

We first discuss the basic queueing dynamics corresponding to activities (i) and (ii) as if no repositioning or recharging was involved, and then further discuss how activities (iii) and (iv) can be incorporated.

First, each idle node i , that represents region i , is considered as a queue where EVs in state (i) discussed above remain idle until rented by a customer. Customer requests for trips are assumed to follow Poisson processes with origin- and destination-specific rates. Therefore, whenever there is an EV in the idle queue, the time until the next rental is exponentially distributed. Suppose that EVs are picked up in a first-in-first-out order, the dynamics of EV movements are analogous to a $\cdot/M/1$ queue in which EVs queue to enter “service” (i.e., wait to be picked up

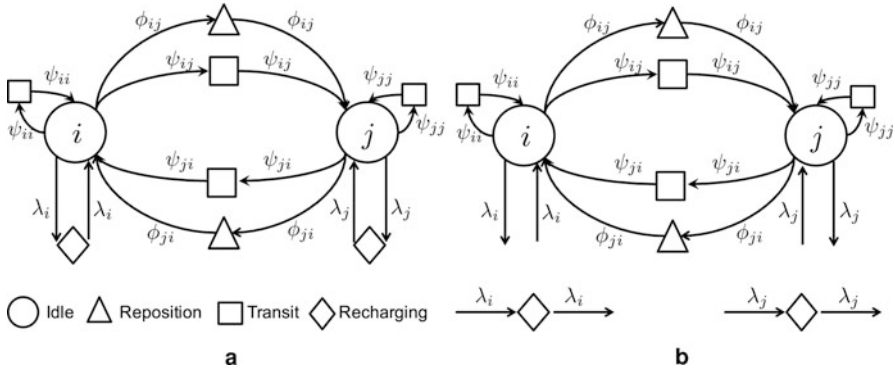


Fig. 19.1 EV sharing operations as queueing networks. (a) Closed queueing network, (b) open queueing network

by a customer). Since the EVs are all identical and interchangeable, and we are only interested in the probability that there are no EVs in the queue, it is no restrictive to impose the first-in-first-out assumption here in our analysis. Any customer requests (which still follow the same Poisson process), that find no EV idle in the queue, will be lost. The service level is then equivalent to the probability that the queue is not empty.

Once an EV is picked up from idle queue i by a customer traveling to destination j , it enters the transit node ($i-j$). The duration it stays in the transit node is the travel time from i to j . Assuming it follows a general distribution, the transit node is then a $\cdot/G/\infty$ queue with infinite capacity. At the departure from the transit node, the EV enters the idle queue j at destination.

The flows of EVs induced by customer trips among idle nodes are usually not balanced, due to the nature of travel patterns. Thus, it is important to conduct repositioning (iii) to maintain availability of EVs at all locations. Similar to the transit node in modeling customer-driven EV flows, the repositioning flows by the operator are modeled using the reposition nodes in Fig. 19.1a, which are also $\cdot/G/\infty$ queues, with different (faster, as repositioning is to be conducted economically) “service” rates than the corresponding transit queues.

The additional challenge in free-float EV sharing is that the EVs must be recharged when their battery levels are low. In this model, we consider a simplified recharging operations without tracking the battery level for individual EVs. We assume that an EV departure from a transit queue, i.e., arrives at its destination, there is P_c probability, as specified by the operator, to be re-routed to a recharging queue. Again, assuming sufficient charging capacity and generally distributed charging time, the recharging queue is modeled as an $\cdot/G/\infty$ queue.

Now, we are ready to formulate the operational profit function $\Theta(\cdot)$ mathematically. The first step is to calculate the rates of EV flows among three types of nodes: transit flows, repositioning flows and recharging flows. In the idle queue i , denote the parameter μ_i as the maximum outbound trip demand rate when all customer

adopts the service. We use the travel distribution P_{ij} to describe the proportion of trips with destination j from region i , if all regions are served. Thus, by definition, $\sum_{j \in N} P_{ij} = 1$ holds. When the customer adoption rate is q_i , the trip rate Ψ_{ij} from i to j is:

$$\Psi_{ij} = P_{ij}\mu_i q_i, \quad \forall i, j \in N. \quad (19.9)$$

However, not all demand for trips in Ψ_{ij} is fully satisfied for the following two reasons. First, those trips whose intended destinations are not within the service region will not be allowed. Second, due to the availability of EVs, some demand that find no EV available will be lost. We further introduce an auxiliary decision variable ψ_{ij} to be the realized transit flow. Therefore, we let

$$\psi_{ij} = \alpha \Psi_{ij} x_j, \quad \forall i, j \in N. \quad (19.10)$$

Note that the constraint (19.10) is nonlinear. It can be linearized as follows:

$$\psi_{ij} \leq \alpha \Psi_{ij}, \quad \forall i, j \in N \quad (19.11)$$

$$\psi_{ij} \leq P_{ij}\mu_i x_j, \quad \forall i, j \in N \quad (19.12)$$

$$\alpha \Psi_{ij} + P_{ij}\mu_i (x_j - 1) \leq \psi_{ij}, \quad \forall i, j \in N \quad (19.13)$$

We define the decision variable ϕ_{ij} as the rate of repositioning trips from i to j . For each idle queue i , the flows of EVs in and out must satisfy the flow balance constraint below:

$$\sum_{j \in N} \psi_{ij} + \sum_{j \in N} \phi_{ij} = \sum_{j \in N} \psi_{ji} + \sum_{j \in N} \phi_{ji}, \quad \forall i \in I \quad (19.14)$$

The next step is to decide on the fleet size required to ensure the desired service level. Recall that there are four possible states of EVs and that the total fleet size thus equals the sum of EVs at all these nodes in the closed queueing network. However, due to interdependence among flows of different nodes in the closed queueing network, the relationship between population (fleet) size and flow rates is not straightforward. To derive tractable formulation, we apply the fixed population mean (FPM) approximation (see, for example Whitt 2002). That is, we approximate the EV population in the closed queueing network by the steady state expected population in a closely-related *open* queueing network as illustrated in Fig. 19.1b. Specifically, the idle queues work as $M/M/1$ queues, and the transit, repositioning and recharging queues work as $M/G/\infty$ queues. By disconnecting the recharging queues from the rest of the network, we obtain an open queueing network.

The resulting open network approximation allows us to relate the fleet size to the flow rate for each node. In particular, the expected number of EVs awaiting in idle queue i is $\alpha/(1 - \alpha)$. Let t_{ij} and τ_{ij} be the expected trip durations from i to j for customer rental and repositioning trips respectively, and t_c be the expected

time to recharge an EV. Invoking Little's law for each queue, the expected fleet size constitutes the expected values of $\sum_{j \in N} \sum_{i \in N} t_{ij} \psi_{ij}$ EVs in transit nodes, $\sum_{i \in N} \sum_{j \in N} \tau_{ij} \phi_{ij}$ EVs in repositioning nodes, and $\sum_{i \in N} t_c \lambda_i$ EVs in recharging queues. Therefore, the fleet size N is no less than the following sum:

$$\sum_{i \in N} \frac{\alpha}{1 - \alpha} x_i + \sum_{j \in N} \sum_{i \in N} t_{ij} \psi_{ij} + \sum_{i \in N} \sum_{j \in N} \tau_{ij} \phi_{ij} + \sum_{i \in N} t_c \lambda_i \leq N. \quad (19.15)$$

Finally, we explicitly calculate the operational profit $\Theta(\mathbf{q}, \mathbf{x})$ that consists of operational revenue, repositioning cost, charging cost, and fleet investment. Let ξ be the scaling factor to unify the time unit, e.g., $\xi = 365$ to convert daily rates into yearly rates. The annual operational revenue from EV usage by customers is $\xi \sum_{j \in N} \sum_{i \in N} r t_{ij} \psi_{ij}$, where r is the per unit time usage price of an EV. Similarly, the annual repositioning cost is $\xi \sum_{i \in N} \sum_{j \in N} \eta \tau_{ij} \phi_{ij}$, where η is the repositioning cost per unit time. Let c be the average cost to fully recharge an EV, the total charging cost is then $\xi \sum_{i \in N} c \lambda_i$. Lastly, based on the price and typical life span in the EV sharing fleet, we use h as the annually amortized EV purchase cost. The annual operational profit is therefore summarized as:

$$\Theta(\mathbf{q}, \mathbf{x}) = \xi \left(\sum_{j \in N} \sum_{i \in N} r t_{ij} \psi_{ij} - \sum_{i \in N} \sum_{j \in N} \eta \tau_{ij} \phi_{ij} - \sum_{i \in N} c \lambda_i \right) - hN$$

By combining the results from the previous two sections, the integrated optimization model for the service region design problem is then formulated as a mixed integer second-order cone program (MISOCP):

$$\max_{\substack{x_i, q_i, N, \Psi_{ij} \\ \psi_{ij}, \phi_{ij}, \lambda_i}} \left\{ \sum_{i \in N} f Q_i q_i - \sum_{i \in N} g_i x_i + \xi \left(\sum_{j \in N} \sum_{i \in N} r t_{ij} \psi_{ij} - \sum_{i \in N} \sum_{j \in N} \eta \tau_{ij} \phi_{ij} - \sum_{i \in N} c \lambda_i \right) - hN \right\}$$

s.t. Constraints (19.3), (19.5), (19.6), (19.7), (19.8), (19.9), (19.11), (19.12), (19.13), (19.14) and (19.15)

$$q_i, N, \Psi_{ij}, \psi_{ij}, \phi_{ij}, \lambda_i \geq 0$$

$$x_i \in \{0, 1\}$$

The above formulation is readily solvable by optimization solvers, such as CPLEX and Gurobi, that can handle MISOCPs.

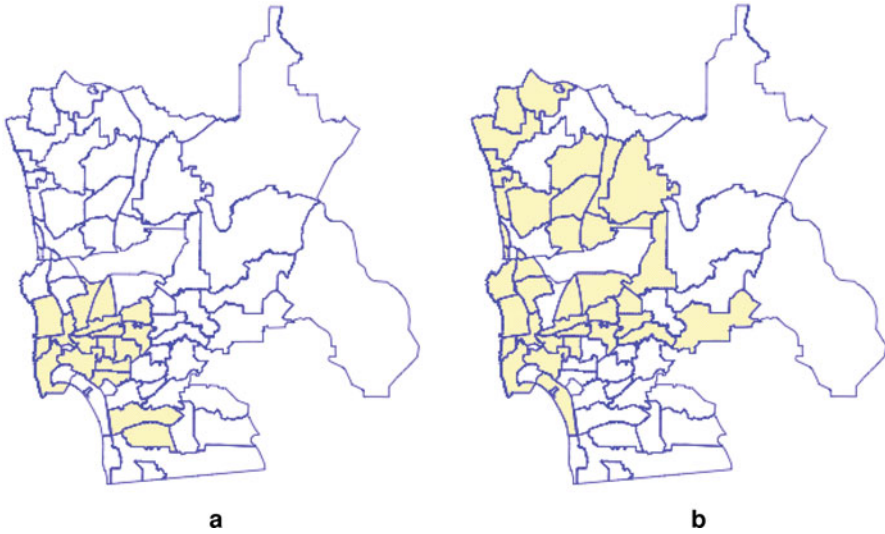


Fig. 19.2 Service region designs. (a) Service region of car2go, (b) optimal service region

19.2.4 Numerical Results

He et al. (2017) demonstrate their service region design optimization framework with a case study of car2go in San Diego, where the 61 zip codes in San Diego county are considered the candidate regions for service coverage. Besides the parameters estimated from public data sources, the study uses real operations data from car2go, data regarding travel characteristics from the California Household Travel Survey and EV charging station deployment data from the U.S. Department of Energy. In particular, the parameters related to travel patterns are estimated using the vehicle status data set of car2go in San Diego between March and April in 2014. Moreover, the Markovian assumption on customer arrivals is validated by the data set.

By solving the (MISOCP) proposed in the previous section, the optimal service region and fleet size are obtained. In Fig. 19.2a, b, we present the actual service region of car2go San Diego as the benchmark, and the optimal service region from solving the MISOCP. While both solutions agree in covering the downtown area of San Diego, where currently 49.88% of the recorded trips occurred, the MISOCP solution suggests more potential in the northern part. Furthermore, the optimal service region in Fig. 19.2b indicates an expansion of service region with a similar size of fleet.

Using the proposed MISOCP, He et al. (2017) observed a few insights from the computational results. We summarize the findings as follows.

1. EV sharing systems deliver higher environmental benefits, e.g., savings in CO₂ emissions, than replacing individually owned gasoline cars with EVs.

2. While faster charging technologies help improve profit and service coverage by improving fleet utilization, the benefits diminish as charging speed becomes higher. Thus, it is sufficient for the operator to deploy moderately fast, but not necessarily the fastest, charging technologies.
3. When customers' valuation of the availability of cars is lower or that of service coverage is higher, the optimal service region becomes larger.

19.3 Fleet Sizing

Fleet sizing is also an important planning problem for ensuring quality of service. In vehicle sharing systems, the burden of ownership of vehicles is on the firm. Therefore, fleet sizing involves relatively long-term investments and must be analyzed carefully at the planning stage. The service region design problem discussed in the last section provides the optimal fleet size in Eq. 19.15, which takes into account fleet repositioning, vehicle utilization and the spatial effect of the service region. However, as the focus of the model is to optimize the service region while considering fleet sizing as one of the cost drivers, a number of approximations were made regarding fleet utilization. Therefore, it is helpful to formulate fleet sizing models that captures these interactions in more detail, in the deployment phase once the service region has been fixed. For instance, planning for parking availability is one of the key factors influencing fleet sizing, as each vehicle needs a parking space either in garage or on street.

To capture the temporal and spatial characteristics of fleet management, approaches such as queueing networks are widely used modeling the fleet sizing problem, e.g., see George and Xia (2011) and Hu and Liu (2016). Using a simulation-based approach, Barrios and Godier (2014) analyze an agent-based model that optimizes fleet size together with fleet repositioning under the objective of maximizing the demand fulfillment. Lu et al. (2017) consider the planning problem of purchasing parking lots/permits and deploying an initial fleet in service regions under a budget constraint on total fleet size. We shall briefly review their two-stage stochastic optimization model below.

19.3.1 Two-Stage Stochastic Optimization Model

Lu et al. (2017) consider a hybrid vehicle sharing system that features both reservation-based and free-float demand. The operator of the system allocates a budget of S vehicles to a set of N regions, in order to maximize its profit together with quality of service (QoS) over a T -period horizon. In the first stage, the operator decides the number of parking spaces (w_i) to purchase for the reservation-based mode, the initial numbers of vehicles that require spaces in parking lots (for the reservation-based mode) and street parking permits (for the free-float mode), x_i^1 and

x_i^2 respectively, to allocate for each region $i \in N$. The associated cost parameters are c_i^{lot} and c_i^{loc} for acquiring one parking space and allocating a vehicle to the region in region i , and c^{ffp} for obtaining one free-float parking permit. Let $c_i^1 = c_i^{\text{loc}}$ and $c_i^2 = c_i^{\text{loc}} + c^{\text{ffp}}$ for each $i \in N$. The fleet allocation and parking planning problem is formulated as a two-stage stochastic optimization model as follows:

$$\min_{\mathbf{w}, \mathbf{x}^1, \mathbf{x}^2 \in \mathbb{Z}_+^N} \left\{ \sum_{i \in N} \left(c_i^{\text{lot}} w_i + c_i^1 x_i^1 + c_i^2 x_i^2 \right) + \Theta(\mathbf{w}, \mathbf{x}^1, \mathbf{x}^2) \right\} \quad (19.16)$$

$$\text{s.t.} \quad \sum_i \left(x_i^1 + x_i^2 \right) \leq S \quad (19.17)$$

$$x_i^1 \leq w_i, \quad \forall i \in N \quad (19.18)$$

where $\mathbf{w} = (w_i)$, $\mathbf{x}^1 = (x_i^1)$ and $\mathbf{x}^2 = (x_i^2)$.

In the first stage, the operator decides w_i (number of parking spaces), x_i^1 (allocation of vehicles without parking permits) and x_i^2 (vehicles with parking permits) in each region i , where all decision variables are required to be integer-valued. The total fleet allocation cannot exceed the fleet budget S , as required by constraint (19.17). Moreover, constraint (19.18) implies that vehicles without a parking permits only park at the purchased parking spaces in parking lots. The objective is to minimize the total cost of fleet allocation and parking space/permit purchase, together with the operational cost $\Theta(\mathbf{w}, \mathbf{x}^1, \mathbf{x}^2)$ in the second stage.

A spatial-temporal network is constructed for the second stage optimization to model the dynamic fleet repositioning, given the initial system state $(\mathbf{x}^1, \mathbf{x}^2)$ and parking space capacity \mathbf{w} . To evaluate $\Theta(\mathbf{w}, \mathbf{x}^1, \mathbf{x}^2)$, Lu et al. (2017) employ the Sample Average Approximation (SAA) method by generating i.i.d samples of random one-way and round-trip trip demand. Due to large number of nodes and arcs in the spatial-temporal network, the second stage optimization problem is computational challenging. To address such difficulty, Lu et al. (2017) generalize a branch-and-cut algorithm with a mixed-integer rounding subroutine to derive stronger cuts from Benders cuts.

19.3.2 Numerical Results

Using a set of Zipcar trip data from the Boston-Cambridge area, Lu et al. (2017) examine several issues in the operations of vehicle sharing. The data set contains a record the starting-ending times of trips and the origin and destination zip codes of each trip, from Oct 1 to Dec 1, 2014. Consistent with Zipcar's practice, the Boston-Cambridge area is divided into nine regions based on the travel patterns observed in the data set. In the numerical experiments, the period is set to one hour, and one-way trips are aggregated by the record quadruple (origin, destination, starting hour, ending hour) and round-trips are aggregated by the triple (origin, starting hour, ending hour).

In the extensive numerical experiments, the mean total rental hours is fixed to 1000 vehicle-hours in a 24-h period. Moreover, the second stage cost is evaluated with 1000 scenarios using the SAA method. The major insights from the numerical results are listed as below.

1. When the one-way demand is exogenous, higher one-way proportion can increase the systems profitability. If the one-way demand is endogenously driven by pricing and strategic customer behavior, higher one-way proportion could decrease profitability.
2. Effective fleet repositioning is important: the number of repositioning trips increases as one-way demand increases.
3. A larger fleet size improves the system's profitability and QoS, e.g., demand fulfillment, substantially.

19.4 Fleet Repositioning

Once strategic decisions such as service region design and fleet investment have been determined, it is important for the system operator to carefully and continuously model the operations of the system. In the strategic service region design problem discussed in the previous section, the repositioning activities are modeled as a queueing network where the stationary performance is evaluated. In the operational level, however, fleet repositioning needs to be considered in more detail. In this section, we study models designed for fleet repositioning operations to dynamically match vehicle supply and travel demand.

The repositioning problem in vehicle sharing has been investigated in the recent literature. Shu et al. (2013) develop a spatial-temporal network flow model and discuss the bicycle redistribution problem for bike sharing systems. Nair and Miller-Hooks (2011) use a similar stochastic model and formulate a mixed-integer program with joint chance constraints. O'Mahony and Shmoys (2015) use the operational data from New York's Citi Bike sharing to estimate the demand flows and solve a mixed-integer program for overnight repositioning. By assuming the demand to be deterministic or follow a Poisson process, optimization models are also developed in Boyacı et al. (2015), Febbraro et al. (2012), Kek et al. (2009) and Nourinejad et al. (2015) for fleet repositioning in one-way station-based systems. The repositioning problem is also closely connected to classical transportation problems such as rail-car distribution (Jordan and Turnquist 1983), empty container deployment (Crainic et al. 1993; Shu and Song 2013) and car rental logistics (Pachon et al. 2003), as well as the transshipment of inventories in supply chains, see, for example Rong et al. (2010), Robinson (1990) and Tagaras (1989).

Recent papers by Benjaafar et al. (2017) and He et al. (2018) consider stochastic dynamic programming formulations for the fleet repositioning problem, both of which identify the optimal reposition up-to and down-to policies in a 2-region system. While Benjaafar et al. (2017) focus on the structural properties of the

optimal repositioning policy in a general product rental network setting, He et al. (2018) develop a computationally tractable optimization framework to deal with spatial and temporal demand correlations, which are not captured in the dynamic program formulations. Specifically, Benjaafar et al. (2017) identify the no-repositioning region in the optimal policy, i.e., not to reposition any vehicles when the state (of vehicle distribution) lies in a certain region and to reposition to the boundary of said region when the state is outside of the region. By allowing spatial and temporal dependence of the demands, He et al. (2018) develop a distributionally robust optimization formulation which is shown to be computationally efficient in their numerical experiments.

19.4.1 Stochastic Dynamic Program Formulation

Different from conventional station-based systems, one-way trips are allowed in free-float systems where customers can pickup any available vehicle. Without having to inform the system about their destinations, the customers may return the vehicles anywhere in the service region at the end of their trips. As discussed in the service region design problem, fleet repositioning is critical in providing desired vehicle availability. In this section, we consider the stochastic dynamic formulation studied in Benjaafar et al. (2017) and He et al. (2018) for the fleet repositioning problem.

Similar to previous sections, the service region (which is exogenously fixed in this case) is partitioned into N regions as a network where customers can travel between any two regions in the network. With a slight abuse of notation, we also use N to denote the set of regions. The firm conducts repositioning regularly over T periods a day. For example, if the operator repositions in four time epochs in a day, we set $T = 4$. If the operator only performs overnight repositioning, then we have $T = 1$. Because the period length is usually not small, we assume that all trips, including those made by customers and repositioning, complete within a period. In the following, we use bold faced characters, e.g., $\mathbf{x} \in \mathbb{R}^N$ and $\mathbf{A} \in \mathbb{R}^{M \times N}$, to vectors and matrices, where x_i to the i th element of \mathbf{x} .

Below, we present the stochastic dynamic program for the fleet repositioning problem with finite horizon of T periods. The sequence of events is as follows. At the beginning of period t , the operator first observes system's state defined by the vector $\mathbf{x}_t = (x_{it})$, where x_{it} is the number of vehicles in region i . Before the customer demands are realized, the operator decides the repositioning quantities $\mathbf{r}_t = (r_{ijt})$, where r_{ijt} is the number of vehicles to be repositioned from i to j . The repositioning cost per vehicle from i to j is given by s_{ijt} . The outbound demand arrival d_{it} at region i is then realized. When the number of vehicles available at i is insufficient to meet the realized demand, the unmet demand is lost with a penalty c_{ijt} if the associated destination is j .

In a free-float system, e.g., Mobike for bike sharing, the operator is usually not informed by the customers about their destinations until they finish their trips.

Consequently, the operator is not able to ration the available vehicles based on customers' intended destinations. Similar to the setting in Sect. 19.2, the operator may obtain the travel distributions at the aggregate level based on historical trip data. That is, we assume that the operator knows that a customer picks up a vehicle from region i has probability P_{ijt} to travel to region j , where $\sum_{j \in N} P_{ijt} = 1$. Therefore, we can define the average penalty $\bar{c}_{it} = \sum_{j \in N} P_{ijt} c_{ijt}$ when an outbound demand from region i is lost. Let w_{it} be the total fulfilled customer trips from i . The average fulfilled customer trips from i to j can be written as $w_{ijt} = P_{ijt} w_{it}$.

To formulate the stochastic dynamic program, we assume the trip demand $\mathbf{d}_t = (d_{it})$ follows some joint probability distribution \mathbb{P} and is independent over time periods. The proposed stochastic dynamic program (DP) minimizes the expected total repositioning cost and lost sales penalty as below:

$$V_t(\mathbf{x}_t) = \min_{\substack{\mathbf{r}_t \geq 0 \\ 0 \leq \sum_{j \in N} r_{ijt} \leq x_{it}}} \left\{ \sum_{i,j \in N} s_{ijt} r_{ijt} + \mathbb{E}_{\mathbb{P}}[J_t(\mathbf{x}_t, \mathbf{r}_t, \mathbf{d}_t)] \right\}. \tag{19.19}$$

where

$$J_t(\mathbf{x}_t, \mathbf{r}_t, \mathbf{d}_t) = \sum_{i \in [N]} \bar{c}_{it} (d_{it} - w_{it}) + V_{t+1}(\mathbf{x}_{t+1}),$$

and

$$x_{i(t+1)} = x_{it} + \sum_{j \in [N]} r_{jit} - \sum_{j \in N} r_{ijt} + \sum_{j \in N} \alpha_{jit} w_{jt} - w_{it}, \quad \forall i \in N, t \in T$$

$$w_{it} = \min \left\{ d_{it}, x_{it} + \sum_{j \in N} r_{jit} - \sum_{j \in N} r_{ijt} \right\}, \quad \forall i \in N, t \in T.$$

and the terminal cost $V_{T+1}(\mathbf{x}_{T+1}) = 0$.

In Eq. (19.19), the constraints require that the number of repositioning trips must be nonnegative and that the total repositioning departures from i can not exceed the available vehicles x_{it} at the beginning of the period. In $J_t(\mathbf{x}_t, \mathbf{r}_t, \mathbf{d}_t)$, the first constraint updates the systems state after repositioning and demand fulfillment in period t . Moreover, the demand fulfillment at region i is the minimum of realized demand or the number of available vehicles.

Before discussing the optimal repositioning policy, we first characterize the properties of the value function $V_t(\mathbf{x}_t)$ under some mild regulatory condition regarding to the cost parameters. Lemma 1 provides an equivalent formulation for the constraint $w_{it} = \min \left\{ d_{it}, x_{it} + \sum_{j \in N} r_{jit} - \sum_{j \in N} r_{ijt} \right\}$ and shows the convexity of $V_t(\mathbf{x}_t)$.

Lemma 1 (Lemma 1 in He et al. 2018) Suppose $\bar{c}_{it} \geq \sum_{j \neq i} s_{ji(t+1)} P_{ijt}$ for any $i \in N$ and $t \in T$. Then,

$$\begin{aligned}
 J_t(\mathbf{x}_t, \mathbf{r}_t, \mathbf{d}_t) &= \min_{\mathbf{w}_t} \left\{ \sum_{i \in N} \bar{c}_{it} (d_{it} - w_{it}) + V_{t+1}(\mathbf{x}_{t+1}) \right\}, \\
 \text{s.t. } x_{i(t+1)} &= x_{it} + \sum_{j \in N} r_{jit} - \sum_{j \in N} r_{ijt} + \sum_{j \in N} \alpha_{jit} w_{jt} - w_{it}, \forall i \in N, \\
 w_{it} &\leq d_{it}, \forall i \in N, \\
 w_{it} &\leq x_{it} + \sum_{j \in N} r_{jit} - \sum_{j \in N} r_{ijt}, \forall i \in N,
 \end{aligned} \tag{19.20}$$

and $V_t(\mathbf{x}_t)$ is convex in \mathbf{x}_t for any $t \in T$.

The condition $\bar{c}_{it} \geq \sum_{j \neq i} s_{ji(t+1)} P_{ijt}$ above suggests that the average profit of a trip departing from region i exceeds the average cost of repositioning a vehicle back to i in the subsequent period. Furthermore, it holds when the system is stationary, e.g., $c_{ijt} = p_{ij}$, $s_{ijt} = s_{ij}$, and $p_{ij} \geq s_{ji}$. Under such condition, formulation Eq.(19.20) implies that even if the operator is not required to fully satisfy all demand, it still optimal to satisfy as much demand as possible in the current period, i.e., $w_{it} = \min \left\{ d_{it}, x_{it} + \sum_{j \in N} r_{jit} - \sum_{j \in N} r_{ijt} \right\}$.

While $V_t(\mathbf{x}_t)$ is convex under the given condition, the DP problem in Eq. (19.19) still suffers from the ‘‘curse of dimensionality’’. In this chapter, we illustrate the optimal repositioning policy for a system with 2 regions.

19.4.2 The 2-Region System

Suppose the operator partitions the system into regions 1 and 2. The repositioning decision can then be reduced to a single variable r_t for repositioning from 1 to 2 in period t , where $r_t > 0$ and < 0 represents repositioning from 1 to 2 and 2 to 1, respectively. Therefore, we can denote the repositioning amount from 1 to 2 as $r_t^+ = \max(r_t, 0)$ and that from 2 to 1 as $r_t^- = -\min(r_t, 0)$. We further assume that the cost parameters satisfies the condition in Lemma 1: $\bar{c}_{it} \geq s_{ji(t+1)} P_{ijt}$ for $i, j \in \{1, 2\}$ and $j \neq i$.

Suppose the fleet size C is given. In any period t , we have $x_{1t} + x_{2t} = C$. We are then able to reduce the state variable to simply the number of available vehicles at region 1. That is, by defining $x_t = x_{1t}$, the number of available vehicles at region 2 is given by $x_{2t} = C - x_t$. After the repositioning operations in period t , there are $y_t = x_t - r_t$ number of available vehicles at region 1. We simplify the stochastic dynamic program Eq. (19.19) into:

$$V_t(x_t) = \min_{x_t - C \leq r_t \leq x_t} \{s_{12t} r_t^+ + s_{21t} r_t^- + \mathbb{E}_{\mathbb{P}}[J_t(y_t, \mathbf{d}_t)]\}$$

where

$$\begin{aligned}
 J_t(y_t, \mathbf{d}_t) &= \min_{w_{1t}, w_{2t}} \{ \bar{p}_{1t}(d_{1t} - w_{1t}) + \bar{p}_{2t}(d_{2t} - w_{2t}) + V_{t+1}(x_{t+1}) \}, \\
 \text{s.t. } x_{t+1} &= y_t - \alpha_{12t}w_{1t} + \alpha_{21t}w_{2t}, \\
 w_{1t} &\leq \min(y_t, d_{1t}), \\
 w_{2t} &\leq \min((C - y_t), d_{2t}),
 \end{aligned}$$

and the terminal cost $V_{T+1}(x_{T+1}) = 0$.

In Proposition 2, the optimal policy is derived as the reposition up-to and down-to policy with thresholds specified as below.

Proposition 2 (Proposition 1 in He et al. 2018) *Suppose $\bar{c}_{it} \geq s_{ji(t+1)}P_{ijt}$ for $i, j \in \{1, 2\}$ and $j \neq i$. For each period t , there exist \underline{x}_t and \bar{x}_t such that*

$$r_t^*(x_t) = \begin{cases} x_t - \underline{x}_t, & x_t \in [0, \underline{x}_t), \\ 0, & x_t \in [\underline{x}_t, \bar{x}_t], \\ x_t - \bar{x}_t, & x_t \in (\bar{x}_t, C], \end{cases}$$

and

$$y_t^*(x_t) = \begin{cases} \underline{x}_t, & x_t \in [0, \underline{x}_t), \\ x_t, & x_t \in [\underline{x}_t, \bar{x}_t], \\ \bar{x}_t, & x_t \in (\bar{x}_t, C], \end{cases}$$

where \underline{x}_t and \bar{x}_t are the optimal reposition up-to and down-to levels respectively defined by the following two convex programs

$$\underline{x}_t = \arg \min_{0 \leq y \leq C} \{ s_{21t}y + \mathbb{E}_{\mathbb{P}}[J_t(y, \mathbf{d}_t)] \}, \quad \bar{x}_t = \arg \min_{0 \leq y \leq C} \{ -s_{12t}y + \mathbb{E}_{\mathbb{P}}[J_t(y, \mathbf{d}_t)] \}.$$

The optimal repositioning policy in Proposition 2 resembles the (s, S) policy in the literature of inventory management. Such policy smooths and balances the available vehicles between the 2 regions. There are two thresholds, \underline{x}_t and \bar{x}_t , that trigger repositioning. When there are insufficient vehicles in region 1, i.e., $x_t < \underline{x}_t$, it is optimal to increase the number of available vehicles to \underline{x}_t by repositioning from region 2. Similarly, when there are too many vehicles in region 1, i.e., $x_t > \bar{x}_t$, it is optimal to reduce the number of available vehicles to \bar{x}_t by repositioning some vehicles to region 2. If $x_t \in [\underline{x}_t, \bar{x}_t]$, it is optimal to not do any repositioning. In the special case when $s_{12t} = s_{21t} = 0$ as considered in Shu et al. (2013), the optimal policy reduces to $\bar{x}_t = \underline{x}_t$ which implies the operator to always conduct repositioning at no repositioning cost.

Generally, the no-repositioning interval $[\underline{x}_t, \bar{x}_t]$ is not a singleton (see Fig. 19.3 for example). Furthermore, it becomes larger when the repositioning cost becomes higher. That is, when it is more costly to reposition, there is higher chance for the operator to find it optimal not to conduct repositioning at all. We formally state the results in Corollary 1.

Corollary 1 (Corollary 1 in He et al. 2018) *Suppose $\bar{c}_{it} \geq s_{ji(t+1)}P_{ijt}$ for $i, j \in \{1, 2\}$ and $j \neq i$. For each period t , \underline{x}_t is decreasing in s_{21t} and \bar{x}_t is increasing in s_{12t} .*

19.4.3 The N -Region System

Due to the “curse of dimensionality”, solving the DP problem in Eq. (19.19) is generally challenging. Benjaafar et al. (2017) prove that the optimal policy can be characterized by a no-repositioning region in the state space. When the system state \mathbf{x}_t is inside the no-repositioning region, it is optimal for the operator not to conduct any repositioning. When \mathbf{x}_t is outside of the no-repositioning region, it is optimal to reposition the right amount of vehicles to the right place so that the new system state that after repositioning is on the boundary of the no-repositioning region. Nevertheless, the optimal repositioning quantities are not solved explicitly.

To obtain computational results that are implementable in practice, He et al. (2018) approximate solutions based on the distributionally robust optimization framework using the enhanced linear decision rules. The benefits of the proposed solutions are three folds. First, it addresses the issue of “curse of dimensionality” by using tractable approximations to the value functions. Second, instead of perfect knowledge of the joint distribution of the demands among all regions, the proposed solution only requires limited distributional information. Finally, it is easy to incorporate both the demand correlations across regions and time periods from the historical travel data directly.

19.5 Other Topics

In this section, we review recent literature on the topics of dynamic pricing and reservation policy in the context of vehicle sharing. Many of these problems are coupled with each other and with the ones discussed in the previous sections. For example, the operator may use dynamic pricing as an incentive tool to support fleet repositioning. Moreover, the repositioning decisions can be constrained by the parking capacity across stations, in a station-based system. We briefly summarize some recent developments in these problems.

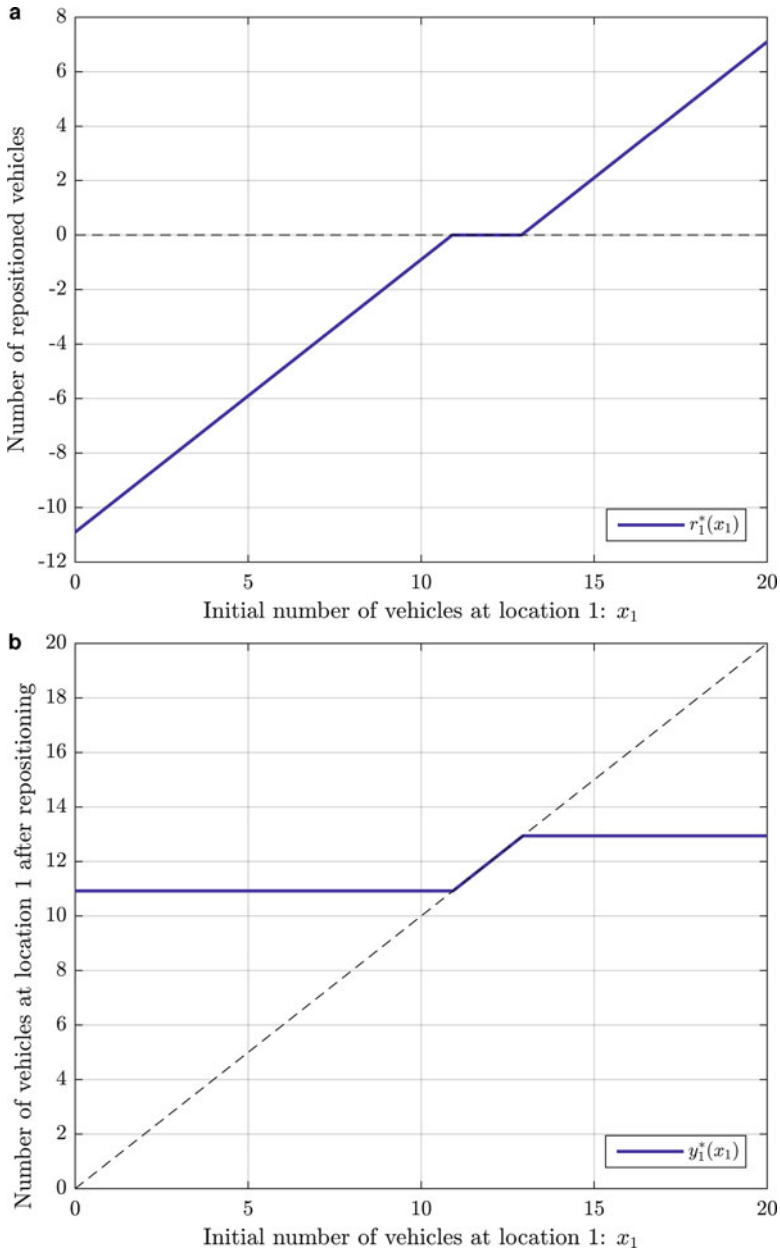


Fig. 19.3 Example of $r_1^*(x)$ and $y_1^*(x)$ when $C = 20, T = 5, s_{12t} = s_{21t} = 1, p_{11t} = 2, p_{12t} = 4, p_{21t} = 3, p_{22t} = 2, \alpha_{11t} = 0.6, \alpha_{12t} = 0.4, \alpha_{21t} = 0.7, \alpha_{22t} = 0.3,$ and $d_{1t} \sim N(12, 3), d_{2t} \sim N(8, 2),$ for $t = 1, \dots, 5.$ **a** $r_1^*(x_t),$ **b** $y_1^*(x_t)$

19.5.1 Dynamic Pricing

The fleet repositioning optimization model developed in Sect. 19.4 is referred to as the operator-based approach, where the operator is conduct the repositioning by itself. There is an alternative to reduce the burden of the operator by motivating customers to reposition the fleet, i.e., the user-based approach. Naturally, dynamic pricing offers the customers incentive to reposition the fleet to the right places.

Febbraro et al. (2012) consider an user-based approach to encourage customers to end their trips at a destination close to the zone with a shortage of vehicles. To model the complex system dynamics in a stochastic environment, a discrete event system (DES) is modeled and a relocation method was proposed based on a linear integer programming formulation. In the one-way vehicle sharing system they consider, customers are required to disclose their trip destinations such that the operator is able to predict the vehicle locations in the near future and optimize the fleet repositioning accordingly. Based on the DES, Febbraro et al. (2012) propose a 2-phase algorithm where the optimal repositioning quantities are determined in the first phase and fare discount to the customers for changing their destinations in the second phase.

Based on a closed queueing network with finite buffer and time-dependent service times, Waserhole et al. (2013) compare several heuristics for optimal pricing: a scenario approach, a fluid approximation, simplified stochastic models and asymptotic approximations. Pfrommer et al. (2014) also study real-time price incentives as a means to shape demand and reduce the need for excessive repositioning. Using simulations in computational experiments, Pfrommer et al. (2014) show that paying customers to reposition may be more cost efficient than hiring staff to reposition bicycles, when the objective is to minimize operating costs under a given desired service level. Focusing on a one-way station-based vehicle sharing system, Jorge et al. (2015) develop a mixed integer nonlinear programming that sets prices for trips to maximize profit. From the computational results in a case study of a network of 75 stations in Lisbon (Portugal), it is demonstrated that the trip pricing strategy can increase profit by inducing a more balanced system.

19.5.2 Reservation Management

Reservations offer an effective mechanism to coordinate the matching between supply and demand. Not only it can control the demand arrivals, it also provides the operator advanced demand information for operations planning by reducing uncertainty. In the case of station-based vehicle sharing, reservations applies to both the vehicles and parking spaces.

A complete parking reservation (CPR) policy for parking spaces at destinations, studied in Kaspi et al. (2014), can be summarized as follows: when a customer starts a trip, she also reveals her destination station and reserves a vacant parking

space in that station. Such policy ensures that there is a parking space available at the destination for vehicle return. On the other hand, it may lower the utilization of parking spaces as other customers may not be able to access them during the reserved period. This complicates the previous discussion on satisfying demands for vehicles at the time of rental, as the operator needs to also consider the demand for vacant parking spaces upon return. Based on a Markovian model, Kaspi et al. (2014) compare the above policy with a no-reservation policy by measuring the total excess time customers spend due to unfulfilled demand requests or delays in returning the vehicles, where the excess time is defined as the difference between the actual journey time and the shortest possible travel time. Through both analytical and numerical studies, it is shown that implementing parking reservations in the proposed policy generally improves the performance of one-way vehicle sharing systems.

Kaspi et al. (2016) also extend the model by developing a mixed-integer linear programming models for designing parking reservation policies, where customer behavior is jointly considered. Based on the analysis of two case studies of real-world systems, the study find the following insights:

1. The CPR policy delivers a significant improvement over the no-reservation policy as shown theoretically.
2. The more reservation information is required of the customer, the better the performances of the proposed partial reservation policies.
3. Parking space overbooking is not likely to be beneficial.

In the context of free-float vehicle sharing where customers can use any available vehicles on street, customers are also allowed to reserve vehicles online up to certain time period before their trips. The reservation policy of vehicles for online customers also influences the availability of vehicles to on street customers. An interesting research direction, as pointed out in He et al. (2017), is to investigate efficient reservation policies, e.g., the optimal time window for vehicle reservation and possibly fees that depend on time and vehicle availability, that may help balance the customer trips and reduce repositioning efforts without jeopardizing the service level.

19.6 Discussion

Vehicle sharing is a growing business model in the sharing economy, in the form of short-term product rental. In this chapter, we discuss several operations management problems such as the strategic planning of service region design, the dynamic fleet management in daily operations, fleet sizing and allocation, dynamic pricing as well as reservation policy for vehicle and parking. There are further research directions to improve the sustainability and efficiency of vehicle sharing systems. For instance, to maintain fleet performance, especially in the case of EV sharing, supporting infrastructure needs to be deployed in the network. In particular, insufficient or

costly access to charging facilities discourages the use of EVs. Charging stations with sufficient number of chargers are required to recharge the fleet so that EVs have sufficient battery level to serve customer trips. In the selection of charging sites and chargers, the operators may choose between the centralized deployment with fewer charging stations and more chargers at each station, or decentralized deployment with more charging stations and fewer chargers at each station. The typical tradeoff is between the saving in chargers via risk pooling and the saving in repositioning via proximity to customers. Such tradeoff has been explored in Mak et al. (2013), where they study the location and inventory decisions for battery swapping stations in face of general travel demand using EVs. In the case of EV sharing, the operational characteristics of the vehicle sharing, e.g., fleet repositioning, and charging scheduling need to be explicitly modeled.

As an integral part of smart city development, the management of vehicle sharing operations will inevitably become more and more data driven. There is great potential in developing methodologies to integrate multiple data sources (e.g., real-time travel times and public transit data) to improve forecasting of both customer demand vehicle availability (e.g., in free-float systems, where will customers return their vehicles in a few hours?). Integrating these data-driven real-time forecasts, on the one hand, offers the potential to improve service efficiency and integration with the city's wider transportation system on the one hand, and calls for sophisticated dynamic optimization approaches on the other. The growing availability of vehicle sharing data is also making possible more sophisticated empirical studies to analyze various important questions in managing vehicle sharing systems and customers' behavior toward them. Using usage data from the bike sharing system in Paris, Kabra et al. (2016) estimates the impacts on ridership of accessibility and availability of the service. As more detailed operational data becomes available, more facets of these systems can be studied in more detail.

References

- Barrios J, Godier J (2014) Fleet sizing for flexible carsharing systems: simulation-based approach. *Transp Res Rec J Transp Res Board* 2416:1–9
- Benjaafar S, Li X, Li X (2017) Inventory repositioning in on-demand product rental networks. Working paper, University of Minnesota
- Boyacı B, Zografos KG, Geroliminis N (2015) An optimization framework for the development of efficient one-way car-sharing systems. *Eur J Oper Res* 240(3):718–733
- Crainic TG, Gendreau M, Dejax P (1993) Dynamic and stochastic models for the allocation of empty containers. *Oper Res* 41(1):102–126
- Febbraro A, Sacco N, Saeednia M (2012) One-way carsharing: solving the relocation problem. *Transp Res Rec J Transp Res Board* 2319:113–120
- George DK, Xia CH (2011) Fleet-sizing and service availability for a vehicle rental system via closed queueing networks. *Eur J Oper Res* 211(1):198–207
- He L, Hu Z, Zhang M (2018) Robust repositioning for vehicle sharing. *Manufacturing & Service Operations Management* (Forthcoming)

- He L, Mak HY, Rong Y, Shen ZJM (2017) Service region design for urban electric vehicle sharing systems. *Manuf Serv Oper Manag* 19(2):309–327
- Hu L, Liu Y (2016) Joint design of parking capacities and fleet size for one-way station-based carsharing systems with road congestion constraints. *Transp Res B Methodol* 93:268–299
- Jordan WC, Turnquist MA (1983) A stochastic, dynamic network model for railroad car distribution. *Transp Sci* 17(2):123–145
- Jorge D, Molnar G, de Almeida Correia GH (2015) Trip pricing of one-way station-based carsharing networks with zone and time of day price variations. *Transp Res B Methodol* 81:461–482
- Kabra A, Belavina E, Girotra K (2016) Bike-share systems: accessibility and availability. Working paper, INSEAD
- Kaspi M, Raviv T, Tzur M (2014) Parking reservation policies in one-way vehicle sharing systems. *Transp Res B Methodol* 62:35–50
- Kaspi M, Raviv T, Tzur M, Galili H (2016) Regulating vehicle sharing systems through parking reservation policies: analysis and performance bounds. *Eur J Oper Res* 251(3):969–987
- Kek AG, Cheu RL, Meng Q, Fung CH (2009) A decision support system for vehicle relocation operations in carsharing systems. *Transp Res E Logist Transp Rev* 45(1):149–158
- Lim MK, Mak HY, Rong Y (2015) Toward mass adoption of electric vehicles: impact of the range and resale anxieties. *Manuf Serv Oper Manag* 17(1):101–119
- Lu M, Chen Z, Shen S (2017) Optimizing the profitability and quality of service in carshare systems under demand uncertainty. *Manuf Serv Oper Manag*. Articles in Advance. Published Online: 16 Oct 2017, <https://doi.org/10.1287/msom.2017.0644>
- Mak HY, Rong Y, Shen ZJM (2013) Infrastructure planning for electric vehicles with battery swapping. *Manag Sci* 59(7):1557–1575
- Nair R, Miller-Hooks E (2011) Fleet management for vehicle sharing operations. *Transp Sci* 45(4):524–540
- Navigant Research (2013) Carsharing services will surpass 12 million members worldwide by 2020. <https://www.navigantresearch.com/newsroom/carsharing-services-will-surpass-12-million-members-worldwide-by-2020>. Accessed 22 Aug 2013
- Nourinejad M, Zhu S, Bahrami S, Roorda MJ (2015) Vehicle relocation and staff rebalancing in one-way carsharing systems. *Transp Res E Logist Transp Rev* 81:98–113
- O'Mahony E, Shmoys DB (2015) Data Analysis and Optimization for (Citi) Bike Sharing. In: Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, AAAI 15. AAAI Press, Austin, Texas, pp 687–694 ISBN:0-262-51129-0030. acmid:2887103
- Pachon JE, Iakovou E, Ip C, Aboudi R (2003) A synthesis of tactical fleet planning models for the car rental industry. *IIE Trans* 35(9):907–916
- Pfrommer J, Warrington J, Schildbach G, Morari M (2014) Dynamic vehicle redistribution and online price incentives in shared mobility systems. *IEEE Trans Intell Transp Syst* 15(4):1567–1578
- Robinson LW (1990) Optimal and approximate policies in multiperiod, multilocation inventory models with transshipments. *Oper Res* 38(2):278–295
- Rong Y, Snyder LV, Sun Y (2010) Inventory sharing under decentralized preventive transshipments. *Naval Res Logist (NRL)* 57(6):540–562
- Shu J, Song M (2013) Dynamic container deployment: two-stage robust model, complexity, and computational results. *INFORMS J Comput* 26(1):135–149
- Shu J, Chou MC, Liu Q, Teo CP, Wang IL (2013) Models for effective deployment and redistribution of bicycles within public bicycle-sharing systems. *Oper Res* 61(6):1346–1359
- Simon HA (1957) *Models of man: social and rational*. Wiley, Oxford
- Tagaras G (1989) Effects of pooling on the optimization and service levels of two-location inventory systems. *IIE Trans* 21(3):250–257
- Time Magazine (2012) Will car-sharing networks change the way we travel? <http://content.time.com/time/specials/packages/0,28757,2094921,00.html>. Accessed 7 Feb 2012
- Waserhole A, Jost V, Brauner N (2013) Pricing techniques for self regulation in vehicle sharing systems. *Electron Notes Discret Math* 41:149–156
- Whitt W (2002) *Stochastic-process limits: an introduction to stochastic-process limits and their application to queues*. Springer, New York

Chapter 20

Agent Pricing in the Sharing Economy: Evidence from Airbnb



Jun Li, Antonio Moreno, and Dennis J. Zhang

Abstract One of the major differences between markets that follow a “sharing economy” paradigm and traditional two-sided markets is that the supply side in the sharing economy often includes individual nonprofessional decision makers, in addition to firms and professional agents. Using a data set of prices and availability of listings on Airbnb, we find that there exist substantial differences in the operational and financial performance of professional and nonprofessional hosts. In particular, properties managed by professional hosts earn 16.9% more in daily revenue, have 15.5% higher occupancy rates, and are 13.6% less likely to exit the market compared with properties owned by nonprofessional hosts, while controlling for property and market characteristics. We demonstrate that these performance differences between professionals and nonprofessionals can be partly explained by pricing inefficiencies. Specifically, we provide empirical evidence that nonprofessional hosts are less likely to offer different rates across stay dates based on the underlying demand patterns, such as those created by major holidays and conventions.

20.1 Introduction

The widespread adoption of Internet infrastructure and smartphones has reduced the transaction costs associated with individuals sharing and trading their idle resources and capacity. This has enabled innovative business models that provide services

J. Li (✉)

Ross School of Business, University of Michigan, Ann Arbor, MI, USA
e-mail: junwli@umich.edu

A. Moreno

Harvard Business School, Boston, MA, USA
e-mail: amoreno@hbs.edu

D. J. Zhang

Olin Business School, Washington University in St. Louis, St. Louis, MO, USA
e-mail: denniszhang@wustl.edu

© Springer Nature Switzerland AG 2019

M. Hu (ed.), *Sharing Economy*, Springer Series in Supply Chain Management 6,
https://doi.org/10.1007/978-3-030-01863-4_20

485

using distributed capacity contributed by independent contractors. In some cases, the agents ultimately providing the service are nonprofessional individuals who share their spare resources, giving rise to the trend often referred to as “the sharing economy”, which revolutionized various industries in the past years.

Most of the sharing economy business models, such as Uber and Airbnb, are based on digital platforms (Parker et al. 2016) that connect individuals who possess excess resources with individuals who need resources, creating two-sided markets (Parker and Van Alstyne 2005; Eisenmann et al. 2006). On one side of the market, the platform “contracts” the service with the customers. On the other side, independent service providers deliver the service using their own assets. Frequently, the platform simply acts as an intermediary and does not directly employ the service providers nor has any ownership or control of the assets that are used to provide the service.

Without the need to invest on physical assets or maintain a large internal workforce, many of the sharing-economy platforms scale up quickly. In December 2014, Airbnb had a global portfolio of one million listings, exceeding the capacity of the largest hotel groups in the world—Hilton, InterContinental and Marriott.¹ On the other hand, platforms are limited in the tools they can use to manage their capacity. While Marriott can decide how many rooms are offered and at which prices in each market, Airbnb cannot make that type of decisions. The independent providers (hosts, in this case) decide whether they want to offer their properties to the market as well as the quantity and price. This represents a change of paradigm from traditional service models where such decisions are made within the boundaries of the firm by professional decision makers. This paper studies the implications of this change of paradigm represented by the sharing economy.

In particular, we focus on one of the critical differences between sharing economies and some of the traditional two-sided markets (e.g., credit card markets, software markets), which is that in the sharing economy the supply side often consists of both professional (experienced) players and nonprofessional (inexperienced) players. For example, on Airbnb.com, there are professional rental service providers as well as “amateurs” who rent out their apartments occasionally.² Studies in behavioral economics have found that nonprofessionals are more likely to suffer from behavioral biases such as loss aversion (Mayer 2001), limited attention (DellaVigna and Pollet 2009), and overconfidence (Malmendier and Tate 2008). These behavioral anomalies often change the prediction of traditional models based on complete rationality, as seen in recent operations management modeling literature (e.g., see Su 2008; Huang et al. 2013). If the paradigm of the sharing economy involves a shift towards services provided more and more by nonprofessionals, it is crucial to understand how their biases translate into market

¹ Airbnb will soon be booking more rooms than the world’s largest hotel chains. Quartz. January 20, 2015.

² Airbnb in the city. New York State Office of General Attorney. October, 2014.

outcomes, and what interventions may improve market efficiency. In this paper, we empirically study the performance and behavioral differences between professional and nonprofessional agents.

We developed a software procedure to scrape listing data from Airbnb.com for all stay dates in a four-month period from December 1, 2012 to March 31, 2013 in the Chicago area. We classify hosts according to the number of properties they list on the site. We call *nonprofessional hosts* those who only list one property through Airbnb, and *professional hosts* those who list multiple properties, which represent 18% of hosts in our sample. We then compare the performance of professional and nonprofessional hosts using performance metrics commonly used in the hospitality industry, including average daily revenue, occupancy rate and price. We find substantial discrepancies between professional and nonprofessional hosts. All else being equal, a property managed by a professional host earns 16.9% higher average daily revenue, and has a 15.5% higher occupancy rate, despite being offered for the same number of days per week at similar average price.

To understand the stability of the market, we ping the URLs of these listings one and half year later and find a high turnover rate: 49% of previously available listings have exited the market. In particular, properties managed by non-professional hosts are 13.6% more likely to exit the market, everything else being equal.

We then explore the source of these discrepancies. We show that they are in part explained by the pricing inefficiencies of nonprofessional hosts. While professional hosts are more likely to offer different prices across stay dates based on the underlying demand level, which results in higher occupancy rates and revenues, nonprofessional hosts fail to do so. We also find that both professional and nonprofessional hosts engage in minimal dynamic price adjustments across the booking horizon. That is, they almost never adjust prices upward nor downward even when the property is not rented out a few days prior to the stay date. Note that this is in contrast to the common practice in the hotel industry, where the prices for a given stay date often experience substantial changes along the booking horizon based on time left and changes in customer willingness-to-pay.

Overall, our findings suggest that peer-to-peer platforms could consider interventions that assist agents, in particular, non-professional ones, adjust their prices more efficiently, such as the price recommendation tool recently launched by Airbnb.

20.2 Literature Review and Hypothesis Development

20.2.1 Literature Review

Sharing economy business models are capturing an increasing attention from the academic community. Recent work has studied what drives owning vs. sharing (Benjaafar et al. 2015; Horton and Zeckhauser 2016), how to create successful matches in the sharing economy (Cullen and Farronato 2014), how to manage

distributed, self-scheduling capacity (Cachon et al. 2015; Gurvich et al. 2015), or how to design and operate urban bike sharing programs (Kabra et al. 2015).

Within in this line of research, some has explored the context created by Airbnb, one of the most prominent platforms in the sharing economy. For example, Zervas et al. (2014) study the effects of Airbnb on hotel revenues, Fradkin (2014) analyzes the consequences of search frictions using internal data from Airbnb, and Edelman et al. (2015) study racial discrimination using a field experiment on Airbnb. Our work contributes to this emerging stream of literature.

As a new form of two-sided market, sharing economy business models inherit important traits from the traditional two-sided markets: network externalities. That is, each side of the market benefits from the presence of the other (David 1985; Farrell and Saloner 1985; Katz and Shapiro 1985; Parker and Van Alstyne 2005). However, sharing economy markets can be less efficient due to the presence of nonprofessional service providers, who are more likely to be subject to behavioral constraints. The focus of our work is on understanding the differences in behavior of professional and nonprofessional service providers, and the consequences for the market.

Behavioral differences between amateurs and professional players have received a considerable amount of attention in the behavioral economics literature. For example, using observational data, Mayer (2001) show that investors outperform homeowners in the real estate market because homeowners exhibit larger loss aversion in pricing their properties. List (2003) demonstrates, with a series of field experiments, that professional players outperform amateurs in the card-trading market due to endowment effects. List (2004) later shows that this performance gap shrinks when nonprofessional players gain more experiences in the market. DellaVigna (2009) surveys the empirical literature on behavioral anomalies and the resulting performance discrepancies between professional and nonprofessional players. Our work considers similar discrepancies between professionals and nonprofessionals and their impacts, but in the context of the sharing economy. Given that the sharing economy represents a general paradigm shift towards nonprofessional service providers, it is particularly important to understand the implications of their behavioral differences on market outcomes.

Our work is also closely related to (1) service operations literature which considers human interactions in service contexts (e.g., Buell et al. 2015; Frei and Morriss 2012); (2) behavioral operations literature which studies bounded rationality and its implications on operational decisions (e.g., Su 2008; Huang et al. 2013); and (3) revenue management literature which studies theory and practice of demand- and capacity-based revenue management (Netessine and Shumsky 2005; Talluri and Van Ryzin 2006; Jerath et al. 2010), particularly in the hospitality industry (Anderson and Xie 2011; Bodea et al. 2009; Lederman et al. 2014). We find that there are substantial differences in the way revenue management tools are implemented by professional and nonprofessional agents, which translate into significant differences in market outcomes.

20.2.2 Hypotheses Development

As mentioned above, past research in behavioral economics shows that professional players have superior financial and operational performance compared to nonprofessionals in traditional markets. We hypothesize that nonprofessionals will have inferior financial and operational performance in the Airbnb market as well. In particular, we define our metrics as follows. Let Revenue_{it} represent the total revenue that property i collects for stay dates within time interval t . We can write Revenue_{it} as

$$\text{Revenue}_{it} = \text{NumDaysOffered}_{it} \times \text{DailyRevenue}_{it},$$

where $\text{NumDaysOffered}_{it}$ is the number of stay days that property i is offered during time interval t , often determined exogenously before pricing decisions are made. DailyRevenue_{it} measures average daily revenue *conditional on being offered*. We use DailyRevenue_{it} to measure property i 's host's financial performance, as opposed to total revenue, i.e., Revenue_{it} , because we do not want to penalize a host merely because he decides to offer the property for fewer days. Note that our definition of DailyRevenue_{it} is parallel to the Revenue Per Available Room (RevPAR), a performance metric commonly used by hotels. RevPAR is defined as total room revenue divided by the number of rooms available and the number of days available during the period under consideration. We hypothesize that:

Hypothesis 1 *A property managed by a professional host has higher average daily revenue than a property managed by a nonprofessional host, everything else being equal.*

If Hypothesis 1 is supported, we are also interested in identifying the main channel through which professionals earn higher daily revenue. It could be that professional hosts have higher occupancy rates, or that they can charge higher average rent prices (controlling for property and market characteristics), or both. We can rewrite the daily revenue as the combination of those channels, and test them independently:

$$\text{DailyRevenue}_{it} = \text{OccupancyRate}_{it} \times \text{AverageRentPrice}_{it}$$

where $\text{OccupancyRate}_{it}$ is the occupancy rate for property i in time interval t , calculated as the number of days occupied divided by the total number of days offered, and $\text{AverageRentPrice}_{it}$ is the average price at which property i is rented out during time interval t (which is calculated using the prices listed on the days in which the property was rented).

Several past studies have shown that one of the major differences between professional and nonprofessional agents in traditional markets is that professional agents are more likely to reach a deal (Mayer 2001; List 2003). This allows us to hypothesize as follows.

Hypothesis 2 *A property managed by a professional host has a higher occupancy rate than a property managed by a nonprofessional host, everything else being equal.*

Similarly, Hypothesis 1 can also be driven by the fact that professional hosts have a higher average rent price, i.e., average price when a property is rented out. This could be true, for example, if being a professional host signals better service quality. Consequently, we hypothesize that:

Hypothesis 3 *A property managed by a professional host has a higher average rented price than a property managed by a nonprofessional host, everything else being equal.*

Besides merely testing whether the direction established in Hypotheses 2 and 3 is supported by the data, we are interested in their relative magnitude so that we can identify the main driver of better revenue performance of professional hosts, if Hypothesis 1 is supported. The following equation sums up our three hypotheses:

$$\text{Revenue}_{it} = \text{NumDaysOffered}_{it} \times \underbrace{\text{OccupancyRate}_{it}}_{\text{Hypothesis 2}} \times \underbrace{\text{AverageRentPrice}_{it}}_{\text{Hypothesis 3}}.$$

Hypothesis 1

Finally, we are interested in not only the temporary operational and financial performance of different hosts, but also the consequences of such differences on market dynamics in the long term. As suggested by the economics literature (e.g., Ellison and Fudenberg 2003), one of the important long-term metrics of two-sided markets in defining market efficiency is the number of suppliers in the platform, which, in our case, is closely related to agents' exiting behavior. Since nonprofessional agents may suffer from behavioral anomalies and receive lower than expected revenues, they are probably more likely to exit the market, possibly in favor of other options, for instance, selling the property in the real estate market or renting the property in the long-term rental rather than short-term rental market.³ Therefore, we hypothesize that:

Hypothesis 4 *A property managed by a professional host is less likely to exit the market than a property managed by a nonprofessional host, everything else being equal.*

³We restrict our attention to properties offered as entire apartments or houses and exclude those properties where the hosts also reside, so that we focus on a relatively homogeneous group of hosts with similar levels of mobility.

20.3 Empirical Setting and Data

20.3.1 Empirical Setting: The Airbnb Platform

To study the differences in behavior between professionals and nonprofessionals, we use data from Airbnb. Airbnb is a sharing-economy platform that connects hosts with empty rooms to potential renters. Hosts on Airbnb list their spare rooms or apartments/houses and determine their own daily prices for rentals. Users visit the Airbnb website to search for desirable accommodations. Founded in 2008, the Airbnb's marketplace has experienced tremendous growth in the last few years. As of 2014, there are more than one million properties worldwide and 30 million guests who use the service. Like other traditional two-sided markets, Airbnb earns revenues from both sides. In particular, guests pay a 9–12% service fee on average for each reservation, depending on the length of stay and the location, while hosts pay a 3% service fee to cover the cost of processing payments by Airbnb. Currently, Airbnb's business model operates with little to no regulation in most locations. As a result, it becomes a major concern, for some local governments such as New York City, that professional rental businesses use Airbnb to avoid taxes, and this has been the subject of intense policy debates.⁴ The main focus of our study is not to contribute to the ongoing debate about regulation in Airbnb, but to use data from the platform as an example to study differences in behavior between professionals and nonprofessionals that can be relevant in other sharing-economy platforms as well.

We classify Airbnb hosts in two types: (1) inexperienced individuals who list their spare rooms or apartments/houses for rent, which we denote as *nonprofessional hosts*, and (2) professional agents who manage multiple properties at the same time, which we denoted as *professional hosts*. In this paper, we define professional hosts as those who offer two or more unique units on Airbnb. Our results do not change qualitatively if we follow the definition by New York State Attorney General's office and define hosts as professional hosts if they hold three or more unique. In our sample, among hosts who offer entire apartments for rent, 18% are professional hosts with at least two properties. The professionals who constitute these 18% hold 24% of all properties in our sample and account for 33% of all revenue in our sample period.

20.3.2 Airbnb Data: Listings and Transactions

To conduct this study, we developed a software procedure to scrape listings around the Chicago area on Airbnb.com for stay dates ranging from December 1, 2012 to March 31, 2013. This time horizon has the advantage that it is not affected by the presence of automatic pricing tools that have been developed more recently,

⁴“Airbnb, New York State Spar Over Legality Of Rentals.” NPR. October 16, 2014.

so it is adequate to study differences in agent behavior. The procedure works as follows: (1) the program logs on to Airbnb.com to search for available rooms in the Chicago area; (2) the program then follows the link to each listing and records the information about that listing, such as location, room type, number of bedrooms, number of bathrooms, guest reviews, identify of the host, etc.; (3) for each listing, the crawler searches for availability and price of all stay dates during the four-month travel period. To capture at least one month worth of availability and price history for each listing on each stay date, the program was run on a daily basis from November 1, 2012 until March 31, 2013. In order to study the entry and exit of Airbnb hosts, we re-scraped Airbnb.com 18 months later, in August 2014. Since Airbnb does not reuse the host ID, we can identify hosts who had delisted their properties and exited the market.

We restrict our attention to offerings of an entire house or apartment and exclude those offerings with just a part of a property. This is because hosts who provide just a room or a bed in their house or apartment tend to have different demographics, incur different costs of renting and sometimes rent their rooms out for different reasons (such as social reasons). We also focus only on listings targeting short-term stays rather than long-term stays (listings with minimum length of stay less than a week).

Documenting differences in listings between different types of hosts is informative in itself, but we also use calendar listings to impute bookings from dynamic changes in listing availability. Based on descriptions on Airbnb's website, when a property is unavailable for a stay date, either booked or not offered, the price is not displayed in the calendar. For example, if we observe on December 10th that a property is available at \$149 for the night of December 11th, it means that it has not been booked for the stay on December 11th and it is available as of December 10th. On the other hand, if a price was displayed on booking date December 9th for a stay on the 11th, but it is no longer displayed on December 10th, it implies that the property was booked for December 11th on December 10th.

Table 20.1 gives a summary of all offerings, where an offering is defined as the combination of property and stay date. Price is the last observed price along a 30-day booking horizon prior to the date of stay. Rented is equal to 1 if the property is rented out for the stay date. The table also displays observable property characteristics, including number of reviews, average ratings, number of bathrooms, and number of bedrooms.

Table 20.1 Summary statistics of listings

	N	Mean	St. dev.	Min	Max
Price (\$)	24,845	149.99	79.65	10	600
Rented	24,845	0.27	0.45	0	1
NumReviews	24,845	10.99	15.23	0	150
AvgRating	17,179	9.61	0.55	8	10
NumBathrooms	23,055	1.29	0.62	1	5
NumBedrooms	24,140	1.58	0.89	1	6

Inferring availability and transactions from the calendar data has some potential limitations and requires some assumptions. First, a property could become unavailable in the calendar and be classified as “booked” because the host no longer wants to offer the property for a particular night, and not because the property has been booked. Even though one cannot completely rule out such possibility, we believe that imputing transactions in this way offers a reasonable proxy for real bookings. Given that we focus on listings for an entire house or apartment rather than a single room or bed at a property, the chance that a property owner delists a property due to personal reasons is significantly reduced because the owner does not reside at the property. Moreover, given that we focus only on short-term rentals and Airbnb is the leading existing short-term rental marketplace for individual properties, the chance that a property is rented out through other channels is also greatly reduced.⁵ Second, a property could appear as “available” from the calendar but could actually be unavailable. This could happen, for example, when the host has not updated the calendar to reflect the actual availability of the property, in what Fradkin (2014) refers to as a “stale vacancy.” Note that having access to internal data would not solve this problem.

Because the focus of this paper is to understand the differences between the behavior of professional and nonprofessional hosts, the aforementioned issues could be problematic if they affected professional and nonprofessional hosts differently. In the next subsection we present a comparison of professional and nonprofessional hosts and we report the results of two tests that suggest that these issues do not affect the two types of hosts differently.

Table 20.2 displays summary statistics at the weekly level for professional and nonprofessional hosts, respectively. The first part of the table simply shows the variables summarized in Table 20.1, for professional and nonprofessional hosts, aggregated at the weekly level. The second part of the table includes additional variables calculated at the weekly level. We do not observe any significant difference in the number of days offered per property per week between professional and nonprofessional hosts. It appears clear though, even before conducting any statistical analysis, that properties managed by professional hosts on average earn more per week, obtain a higher occupancy rate, and are less likely to exit the market. However, such discrepancies in performance can be driven by the fact that professional hosts offer more spacious properties (more bedrooms and more bathrooms), have more reviews (though they are not necessarily rated higher), and perhaps even are located in more popular districts. The rest of the paper studies this performance discrepancy systematically, introducing the relevant control variables in the analysis.

⁵In Sect. 20.5, we focus on the subset of hosts who make their properties available more than four days per week (50% of the time). Because of the high availability of their properties, it is less likely that these hosts will cancel availability for other reasons. We do not find any qualitative differences in our results by focusing on this subsample.

Table 20.2 Summary statistics for professional and nonprofessional hosts at weekly level

	Professional hosts		Nonprofessional hosts	
	Mean	St. dev.	Mean	St. dev.
Average price (\$)	167.87	83.35	146.94	99.5
Occupancy rate	0.31	0.38	0.27	0.35
NumReviews	12.63	15.47	11.39	17.26
AvgRating	9.36	0.63	9.69	0.51
NumBathrooms	1.45	0.88	1.24	0.52
NumBedrooms	2.02	1.19	1.47	0.78
Daily revenue (\$)	44.66	63.91	34.22	65.04
NumDays offered	5.89	1.94	5.79	2.02
NumWeekdays offered	4.22	1.56	4.12	1.64
NumWeekends offered	1.67	0.54	1.67	0.54
Exit [†]	0.28	0.45	0.57	0.50

[†]Exits are measured 18 months later from the original data collection period

20.4 Performance of Professional vs. Nonprofessional Hosts: Econometric Specifications and Results

In this section we describe the overall methodology of our analysis, the construction of our variables, and the corresponding empirical results.

20.4.1 Daily Revenue

In Hypothesis 1, we would like to understand how the daily revenue from renting out a property depends on whether the property is managed by a professional or a nonprofessional host. We focus on the daily revenue, DailyRevenue_{it} , which we average at the weekly level. There are two reasons why we construct our revenue measure in this way. First, as mentioned above, normalizing total revenue by the number of days each host makes their properties available allows us to tease out the effect of predetermined availability and focus on the performance metric driven by endogenous decisions such as pricing. Second, we choose the aggregation level at the weekly level because we want to construct our measure in a relative homogeneous period for each host and, at the same time, average out the day-of-week effect. We use a family of reduced-form specifications and model daily revenue as,

$$\log(\text{DailyRevenue}_{it}) = C_0 + \alpha_1 \text{Professional}_i + \beta X_{it} + v_m + v_t + \epsilon_{it}$$

where $\log(\text{DailyRevenue}_{it})$ is the natural log transformation of daily revenue for property i in week t , Professional_i denotes whether the property is owned by a

professional host. We control for various confounding factors that may potentially correlate with both the daily revenue (i.e., the dependent variable) and the host status (i.e., the treatment). First, we use zip-code-level and week-level fixed effects, v_m and v_t , to control for the possibility that certain markets are more attractive to travelers and meanwhile are also populated with more professional hosts.⁶ Second, we control for the characteristics of an offering, denoted by X_{it} , which includes the physical characteristics of the property (i.e., the number of bedrooms and bathrooms) and the quality of service (i.e., the number of guest reviews, the average review ratings, and the average response time of the host). Moreover, we control for the rank of an offering in the search result. If a host’s professional status (or factors correlated with it) is used as an input to the Airbnb’s search-engine ranking algorithm, then professional-host status can be correlated with performance through rank.

Table 20.3 shows the estimates obtained under different sets of control variables. Column 1 only has market and time fixed effects. Column 2 controls for the physical characteristics of the properties in addition to the fixed effects. Column 3 further

Table 20.3 Hypothesis 1 (Revenue)

	Dependent variable		
	LogDailyRevenue		
	(1)	(2)	(3)
Professional	0.233*** (0.072)	0.217*** (0.074)	0.169** (0.073)
NumWeekends	-0.232*** (0.061)	-0.236*** (0.061)	-0.147** (0.061)
NumBathrooms		0.070 (0.074)	0.209*** (0.074)
NumBedrooms		0.018 (0.051)	-0.028 (0.050)
NumReviews			0.020*** (0.002)
Rank			-0.001*** (0.0002)
Observations	4,297	4,297	4,297
R ²	0.142	0.143	0.185

Note: ResponseTime, week-level and zip-code-level dummy included
 *p<0.1; **p<0.05; ***p<0.01

⁶Ideally, we would like to use a fixed-effect model to control for a listing’s specific characteristics. However, since our independent variable of interest (i.e., whether a property is managed by a professional or a nonprofessional host) is time-invariant, including fixed effects in our model would absorb the effect of the variable of interest.

includes the quality of service in the control variables.⁷ All three columns show that α_1 is significantly greater than zero. Hence, our empirical evidence is consistent with Hypothesis 1. Properties managed by professional hosts on average earn higher daily revenue than properties managed by nonprofessional hosts, with the magnitude being 16.9% in the specification reported in Column 3.

20.4.2 Occupancy Rate and Average Rent Price

The fact that professional hosts earn higher daily revenue (Hypothesis 1) can be attributed to them having a higher occupancy rate (Hypothesis 2), or a higher rent price (Hypothesis 3), or both. In this section, we evaluate the two channels and discuss their relative importance.

We start by testing our second hypothesis, which suggests that a property managed by a professional host will have a higher weekly occupancy rate than one managed by a nonprofessional host. We employ the same model specifications as in the previous section:

$$\log(\text{Occupancy}_{it}) = C_0 + \alpha_2 \text{Professional}_i + \beta X_{it} + v_m + v_t + \epsilon_{it},$$

where $\log(\text{Occupancy}_{it})$ is the log transformation of the weekly occupancy rate of property i in week t , and all the other variables are defined as before.

Table 20.4 shows the estimates obtained under different sets of specifications. In all columns, α_2 is significantly greater than 0, which supports Hypothesis 2. Specifically, properties managed by professional hosts achieve 15.5% higher occupancy rates than properties managed by nonprofessional hosts, based on the estimates reported in Column 3.

Besides higher occupancy rate, do professional hosts also charge higher prices? We next evaluate the difference of rented price between professional and nonprofessional hosts. We adopt the same specification structure as above, changing the dependent variable to AvgRentPrice_{it} :

$$\log(\text{AvgRentPrice}_{it}) = C_0 + \alpha_3 \text{Professional}_i + \beta X_{it} + v_m + v_t + \epsilon_{it}.$$

Table 20.5 shows the estimates with different specifications. In Columns 2 and 3, α_3 is not significantly greater than 0, and the magnitude of the point estimates are small as well, therefore Hypothesis 3 is not supported. That is, professional hosts do not seem to charge a higher rental price on average. Since the dependent variable is calculated using the prices of properties that are rented, our results indicate

⁷We did not find a significant effect of average rating due to its lack of variation. Moreover, average rating is missing when there is no review available, which will limit the number of observations when included. Therefore, we decide to drop average rating in our analyses.

Table 20.4 Hypothesis 2
(Occupancy rate)

	Dependent variable		
	LogOccupancyRate		
	(1)	(2)	(3)
Professional	0.173** (0.068)	0.193*** (0.070)	0.155** (0.070)
NumWeekends	-0.283*** (0.058)	-0.279*** (0.058)	-0.191*** (0.058)
NumBathrooms		-0.048 (0.070)	0.078 (0.070)
NumBedrooms		-0.034 (0.048)	-0.076 (0.047)
NumReviews			0.020*** (0.002)
Rank			-0.0004** (0.0002)
Observations	4,297	4,297	4,297
R ²	0.143	0.144	0.186

Note: ResponseTime, week-level and zip-code-level dummy included

*p<0.1; **p<0.05; ***p<0.01

Table 20.5 Hypothesis 3
(Average rent price)

	Dependent variable		
	LogRentPrice		
	(1)	(2)	(3)
Professional	0.108*** (0.023)	0.032 (0.021)	0.022 (0.022)
NumWeekends	0.092*** (0.019)	0.082*** (0.017)	0.078*** (0.018)
NumBathrooms		0.240*** (0.022)	0.245*** (0.022)
NumBedrooms		0.115*** (0.015)	0.113*** (0.015)
NumReviews			-0.001** (0.001)
Rank			-0.0003*** (0.0001)
Observations	2,313	2,313	2,313
R ²	0.288	0.438	0.444

Note: ResponseTime, week-level and zip-code-level dummy included

*p<0.1; **p<0.05; ***p<0.01

that customers do not have a higher willingness-to-pay for properties managed by professional after controlling for the quality of the property and the service offered. This also alleviates potential concerns regarding to omitted variable biases.

Recall that Column 3 of Table 20.3 indicated that professional hosts on average earn 16.9% more revenue, controlling for property and market characteristics. We have shown that the additional revenue primarily comes from higher occupancy rates rather than higher rented prices.

Besides the main results shown above, we also conducted analyses to ensure the robustness of our results to various model specifications and omitted variable biases. We find our results are still consistent using propensity score matching estimators as well as under Rosenbaum bound sensitivity analysis (Rosenbaum 2002).

20.4.3 Exit Probability

Given that the performance of nonprofessional hosts is inferior to that of their professional counterparts, are properties managed by nonprofessional hosts also more likely to exit the market after a certain period of time? The literature on the two-sided markets suggests that the revenue of the platform and the social welfare depend critically on the size of the supply side (Armstrong 2006). Therefore, exit probability is an important measure to consider in analyzing the health and growth of any sharing-economy marketplace. We use a family of Logit specifications to test this hypothesis. In particular, we model the exit rate of a property as:

$$\text{Exit}_i = \text{logit}(C_0 + \alpha_4 \text{Professional}_i + \beta X_i + v_m + \epsilon_i),$$

where Exit_i indicates whether property i has exited Airbnb's market 18 months after the original sample period. All the other variables are as previously defined.

Table 20.6 shows the results. Under all three specifications the estimates of α_4 are statistically significantly negative, which supports Hypothesis 4. Computing the marginal effects from estimates in Column 3, we conclude that a property owned by a professional host is on average 13.6% less likely to exit the market after one and half years, measured by marginal effect at the means.

20.5 Understanding the Differences in Performance

In this section we examine the source of the performance differences among professional and non-professional hosts. We do so by analyzing the potential pricing inefficiencies of nonprofessional hosts.

Table 20.6 Hypothesis 4 (Exit)

	Dependent variable		
	Exit		
	(1)	(2)	(3)
Professional	-0.744*** (0.283)	-0.632** (0.290)	-0.599** (0.299)
NumBathrooms		0.026 (0.312)	-0.043 (0.319)
NumBedrooms		-0.330 (0.207)	-0.302 (0.211)
NumReviews			-0.019* (0.010)
Rank			-0.0002 (0.001)
Observations	317	317	317

Note: ResponseTime, zip-code-level dummy included
 *p<0.1; **p<0.05; ***p<0.01

The revenue management (RM) literature (Talluri and Van Ryzin 2006; Gallego and Van Ryzin 1994; Bitran and Caldentey 2003) has extensively documented the use and impact of various RM techniques in an array of industries providing perishable products or services, one of which being the hotel industry (Zhao and Zheng 2000). According to this literature, there are two RM tools widely adopted by hotels when setting room rates: (1) variable room rates across stay dates (Talluri and Van Ryzin 2006): hotels offer different room rates based on the day of the week, the season, or other observable factors affecting total demand; (2) variable room rates across booking dates (Su 2007): hotels offer different room rates based on time left to the stay date. We therefore define two measures of the intensity of pricing activity:

1. StayDateRateCnt_{it}, calculated as the total number of last observed price levels for all stay dates in week *t* minus 1. The last observed price equals to the rented price if a property is rented out eventually, or the listing price last observed along the booking horizon otherwise. StayDateRateCnt_{it} = 0, for example, means that there is no price variation across stay dates in week *t* because the last observed prices are constant.
2. BookingDateRateCnt_{it}. We first calculate the number of price levels along the 30-day booking horizon for each stay date in week *t* minus 1. We then take the sum of this measure over all stay dates in week *t*. BookingDateRateCnt_{it} = 0 indicates no price variation across booking dates for any stay date in week *t*, as listing prices are constant.

We first test in our context whether the use of such RM tools, as captured by the aforementioned variables, indeed leads to higher revenue. In particular, we

hypothesize that a property’s weekly revenue is higher if it has higher StayDateRatesCnt and BookingDateRatesCnt for that week. We test the hypothesis with the following model specifications:

$$\log(\text{DailyRevenue}_{it}) = C_0 + \theta_1 \text{StayDateRateCnt}_{it} + \theta_2 \text{BookingDateRateCnt}_{it} + \beta X_{it} + v_m + v_t + \epsilon_{it},$$

$$\log(\text{OccupancyRate}_{it}) = C_0 + \theta_3 \text{StayDateRateCnt}_{it} + \theta_4 \text{BookingDateRateCnt}_{it} + \beta X_{it} + v_m + v_t + \epsilon_{it},$$

Columns 1 and 2 of Table 20.7 show that a more intense pricing activity results in higher daily revenue and occupancy rates, controlling for property and market characteristics. The effect is mainly driven by the use of variable rates by stay date. We did not find significant revenue and occupancy effects of variable rates by booking date, which can be partially driven by lack of adoption of this practice—we observe that 75% of listings did not adjust their prices at all along the booking horizon and 95% of them adjust their prices at most once along the booking horizon.

Given that a more intense pricing activity leads to better performance, the next question is whether professional hosts indeed engage more in such practice and earn

Table 20.7 Impact and use of revenue management techniques

	Dependent variable			
	LogRevenue	LogOccupancyRate	StayDateRateCnt	BookingDateRateCnt
	<i>OLS</i>	<i>OLS</i>	<i>Poisson</i>	<i>Poisson</i>
	(1)	(2)	(3)	(4)
Professional			0.124** (0.053)	-0.217 (0.133)
StayDateRateCnt	0.328*** (0.048)	0.038*** (0.006)		
BookingDateRateCnt	0.002 (0.007)	-0.0002 (0.001)		
NumBathrooms	0.151** (0.069)	-0.001 (0.009)	0.178*** (0.063)	0.109 (0.143)
NumBedrooms	0.005 (0.045)	-0.002 (0.006)	-0.159*** (0.042)	-0.126 (0.095)
NumReviews	0.022*** (0.002)	0.003*** (0.0002)	0.005*** (0.002)	0.009** (0.004)
Rank	-0.0004** (0.0002)	-0.00000 (0.00002)	-0.0005*** (0.0001)	-0.002*** (0.0003)
Observations	4,743	4,743	4,743	4,743
R ²	0.134	0.127		

Note: Response Time, week-level and zip-code-level dummy included
 *p<0.1; **p<0.05; ***p<0.01

higher revenues in turn. In order to test this, we use Poisson model to analyze the levels of pricing sophistication:

$$\text{StayDateRatesCnt}_{it} = \text{Poisson}(C_0 + \theta_5 \text{Professional}_i + \beta X_{it} + v_m + v_t + \epsilon_{it}),$$

$$\text{BookingDateRatesCnt}_{it} = \text{Poisson}(C_0 + \theta_6 \text{Professional}_i + \beta X_{it} + v_m + v_t + \epsilon_{it}),$$

Column 3 of Table 20.7 shows that θ_5 is positive and significant, which indicates professional hosts vary property prices more often based on the date of stay. Computing the marginal effect, we find that properties managed by professional hosts vary prices 4.9% more frequently, calculated as marginal effect at mean. Since θ_6 is non-distinguishable from zero in Column 4, it indicates that professional hosts do not necessarily adjust their prices more often along booking dates, which may be driven by the lack of engagement in dynamic pricing across *booking dates* by all hosts.

Overall, we have shown that a more intense pricing activity results in higher occupancy rates and higher daily revenue. The fact that professional hosts are more likely to engage in intense pricing activity may partially explain their superior performances in this market. The evidence so far provides a mechanism (more intense price adjustments) through which Hypothesis 1 and 2 (higher daily revenue and occupancy rates for properties managed by professionals) may hold.

20.6 Conclusion

The sharing-economy business model comes with an increase in the use of nonprofessional labor. We have used Airbnb as the empirical setting to study the implications of this shift towards using nonprofessional service providers.

We have documented substantial discrepancies between professional and non-professional hosts. All else being equal, a property managed by a professional host earns more than a 16.9% higher average daily revenue, has a 15.5% higher occupancy rate. Moreover, properties managed by professional hosts are 13.6% less likely to exit the market compared with properties owned by nonprofessional hosts, controlling for property and market characteristics. We have shown that these discrepancies can be rationalized by the pricing inefficiencies of nonprofessional hosts. Our findings suggest platforms like Airbnb should try to assist nonprofessionals with their pricing and capacity-management decisions. An example of this is the pricing that Airbnb is currently providing to its hosts or the “heat maps” that Uber shows their drivers, to indicate areas where they are more likely to find a customer.

Although our empirical analysis has focused on Airbnb, we believe that our results provide meaningful insights that go beyond this specific setting. Other platforms such as Uber also use a combination of professionals (e.g., a full-time driver offering a “black car” service) and nonprofessionals (e.g., a student occasionally driving for Uber via their “UberX”). We expect that our findings, which

point to a lower efficiency of nonprofessionals, could play similarly in a service like Uber. Furthermore, as innovative business models are finding new ways of shifting risks to different parts of the value chain including final customers (Girotra and Netessine 2014), the inefficiencies that we observe arising from the use of nonprofessionals could become even more important.

References

- Anderson C, Xie X (2011) A choice-based dynamic programming approach for setting opaque prices. *Prod Oper Manag* 21(3):590–605
- Armstrong M (2006) Competition in two-sided markets. *RAND J Econ* 37(3):668–691
- Benjaafar S, Kong G, Li X, Courcoubetis C (2015) Modeling and analysis of collaborative consumption in peer-to-peer car sharing. Working paper, University of Minnesota
- Bitran G, Caldentey R (2003) An overview of pricing models for revenue management. *Manuf Serv Oper Manag* 5(3):203–229
- Bodea T, Ferguson M, Garrow L (2009) Choice-based revenue management: data from a major hotel chain. *Manuf Serv Oper Manag* 11(2):356–361
- Buell RW, Kim T, Tsay CJ (2015) Creating reciprocal value through operational transparency. Harvard Business School Technology & Operations Mgt Unit working paper (14–115)
- Cachon GP, Daniels KM, Lobel R (2015) The role of surge pricing on a service platform with self-scheduling capacity. Working paper
- Cullen Z, Farronato C (2014) Outsourcing tasks online: matching supply and demand on peer-to-peer internet platforms. Working paper
- David PA (1985) Clio and the economics of qwerty. *Am Econ Rev* 75:332–337
- DellaVigna S (2009) Psychology and economics: evidence from the field. *J Econ Lit* 47(2):315–372
- DellaVigna S, Pollet JM (2009) Investor inattention and Friday earnings announcements. *J Financ* 64(2):709–749
- Edelman BG, Luca M, Svirsky D (2015) Racial discrimination in the sharing economy: evidence from a field experiment. Harvard Business School NOM Unit working paper (16–069)
- Eisenmann T, Parker G, Van Alstyne MW (2006) Strategies for two-sided markets. *Harv Bus Rev* 84(10):92
- Ellison G, Fudenberg D (2003) Knife-edge or plateau: when do market models tip? *Q J Econ* 118:1249–1278
- Farrell J, Saloner G (1985) Standardization, compatibility, and innovation. *RAND J Econ* 16:70–83
- Fradkin A (2014) Search frictions and the design of online marketplaces. NBER working paper
- Frei F, Morriss A (2012) *Uncommon service*. Harvard Business Review Press, Boston
- Gallego G, Van Ryzin G (1994) Optimal dynamic pricing of inventories with stochastic demand over finite horizons. *Manag Sci* 40(8):999–1020
- Girotra K, Netessine S (2014) The risk-driven business model: four questions that will define your company. Harvard Business Press, Boston
- Gurvich I, Lariviere M, Moreno A (2015) Operations in the on-demand economy: staffing services with self-scheduling capacity. Working paper, Northwestern University
- Horton JJ, Zeckhauser RJ (2016) Owning, using and renting: some simple economics of the “sharing economy.” Technical report, National Bureau of Economic Research
- Huang T, Allon G, Bassamboo A (2013) Bounded rationality in service systems. *Manuf Serv Oper Manag* 15(2):263–279
- Jerath K, Netessine S, Veeraraghavan SK (2010) Revenue management with strategic customers: last-minute selling and opaque selling. *Manag Sci* 56(3):430–448

- Kabra A, Belavina E, Girotra K (2015) Bike-share systems: accessibility and availability. Chicago Booth Research Paper (15-04)
- Katz ML, Shapiro C (1985) Network externalities, competition, and compatibility. *Am Econ Rev* 75:424–440
- Lederman R, Olivares M, Ryzin GV (2014) Identifying competitors in markets with fixed product offerings. Columbia Business School research paper no. 14-10
- List JA (2003) Does market experience eliminate market anomalies? *Q J Econ* 118(1):41–72
- List JA (2004) Neoclassical theory versus prospect theory: evidence from the marketplace. *Econometrica* 72(2):615–625
- Malmendier U, Tate G (2008) Who makes acquisitions? CEO overconfidence and the market's reaction. *J Financ Econ* 89(1):20–43
- Mayer C (2001) Loss aversion and seller behavior: evidence from the housing market. *Q J Econ* 116:1233–1260
- Netessine S, Shumsky RA (2005) Revenue management games: horizontal and vertical competition. *Manag Sci* 51(5):813–831
- Parker GG, Van Alstyne MW (2005) Two-sided network effects: a theory of information product design. *Manag Sci* 51(10):1494–1504
- Parker G, Van Alstyne M, Choudary S (2016) Platform revolution. W. W. Norton & Company, New York
- Rosenbaum PR (2002) *Observational studies*. Springer, New York
- Su X (2007) Intertemporal pricing with strategic customer behavior. *Manag Sci* 53(5):726–741
- Su X (2008) Bounded rationality in newsvendor models. *Manuf Serv Oper Manag* 10(4):566–589
- Talluri KT, Van Ryzin GJ (2006) *The theory and practice of revenue management*, vol 68. Springer, New York
- Zervas G, Proserpio D, Byers J (2014) The rise of the sharing economy: estimating the impact of Airbnb on the hotel industry. Boston U School of Management research paper (2013-16)
- Zhao W, Zheng YS (2000) Optimal dynamic pricing for perishable assets with nonhomogeneous demand. *Manag Sci* 46(3):375–388

Chapter 21

Intermediation in Online Advertising



Santiago R. Balseiro, Ozan Candogan, and Huseyin Gurkan

Abstract In online advertising, impressions are sold via real-time auctions which are organized by central platforms referred to as ad exchanges. For technological or operational reasons, advertisers generally participate in the auctions run by exchanges through intermediaries which acquire impressions on their behalf. Intermediaries are specialized entities that provide targeted services for a particular segment of the market, and typically there are multiple stages of intermediation. Moreover, an advertiser may have private information, e.g., budget, targeting criterion or value attributed to an impression. The presence of intermediaries and this information asymmetry introduce several new research questions. In the first part of this chapter, we study the mechanism design problem of an intermediary who offers a contract to an advertiser with a private budget and a private targeting criterion. We characterize the optimal mechanism and establish that the presence of the intermediary results in simpler bidding policies. In the second part of this chapter, we study the strategic interaction among intermediaries organized in a chain network. We characterize a subgame perfect equilibrium of the resulting game among intermediaries and show that the most profitable position in the intermediation chain depends on the underlying value distribution of the advertiser.

S. R. Balseiro (✉)
Columbia University, New York, NY, USA
e-mail: srb2155@columbia.edu

O. Candogan
University of Chicago, Chicago, IL, USA
e-mail: ozan.candogan@chicagobooth.edu

H. Gurkan
Duke University, Durham, NC, USA
e-mail: hg67@duke.edu

21.1 Introduction

Online advertising is a rapidly growing market whose annual revenue exceeded 72.5 billion dollars in the United States in 2016 (Internet Advertising Bureau 2016). This growth has been accompanied by technological advancements which introduce novel tactical and operational challenges for both publishers (the supply side) and advertisers (the demand side) such as real-time bidding and sophisticated targeting. To overcome these challenges, publishers and advertisers increasingly work with intermediaries who have emerged to facilitate transactions between two sides by providing technological and managerial services.

Specifically, when a user visits a publisher's page, an advertising opportunity, referred to as an impression, is generated. This impression is supplied to an exchange, where they are auctioned. These auctions take place in milliseconds after the user's visit to the web page. Due to the real-time nature of these auctions, participants generally employ sophisticated algorithms for automatically targeting users based on specific metrics. Many advertisers work with advertising agencies, who often focus on serving similar clients, enabling them to become experts at campaign management for a specific industry (e.g., pharmaceutical, automotive). Advertising agencies are service-based organizations, and they serve as a managerial layer, typically on top of a licensed demand side platform. Demand side platforms (DSPs) aggregate demand from different market participants and provide real-time bidding service in the auctions run by the exchanges. In addition, there are other intermediaries who specialize in services such as re-targeting (tracking a particular impression in different websites), measurement, and analytics. These intermediaries support the entire trading infrastructure.

The presence of intermediaries in this industry introduces several new interesting questions. What kind of contracts should an intermediary offer to an advertiser? How should an intermediary bid on behalf of its customers? How does the presence of an intermediary affect the efficiency of the market? How does the structure of the intermediation network affect the profits of its participants? Do intermediaries prefer to be closer to the supply source or demand source? This chapter sheds light on these issues by reviewing two separate models which are studied by Balseiro and Candogan (2017) and Balseiro et al. (2017). In the first model (hereafter OCI), we characterize the optimal contract offered by an intermediary to an advertiser in a setting where the advertiser's budget and targeting criteria are private. Using our results, we show that the presence of the intermediary results in simpler bidding policies. In the second model (hereafter MSI), we focus on a multi-stage intermediation network, and analyze the relation between the intermediation network and the profits of the market participants. To do so, we provide a game theoretic model where intermediaries in a chain network between an exchange and an advertiser with private values, sequentially select their mechanisms from a practically relevant class of mechanisms. We characterize a subgame perfect equilibrium for this game, and show that the most profitable position in the intermediation chain depends on the underlying value distribution of the advertiser. In the remainder of this section, we detail our contributions and relate them to the existing literature.

21.1.1 Main Contributions

In the OCI model (see Sect. 21.2), we study the dynamic mechanism design problem of an intermediary who offers a contract to an advertiser with a private budget and targeting criteria. Since the private information of the advertiser is multi dimensional, this problem in general is hard to solve. Therefore, we develop a novel solution method which combines a performance space characterization technique and a duality-based approach. Specifically, we first characterize the performance space that consists of the expected cost and value achievable by any feasible (dynamic) bidding policy, and use our duality-based approach to reduce the optimal contract design problem to a tractable convex optimization problem. In this way, we obtain a crisp characterization of the intermediary's optimal bidding policy. The policy is stationary and has two notable features: (i) the policy bids a weighted average of the values associated with different types, and (ii) the bids are appropriately shaded. Here, bidding a weighted average of the values ensures truthful reporting of the advertiser while bid shading accounts for budget constraints. Using our results, we establish that the intermediary can profitably provide bidding service to a budget-constrained advertiser, and in some cases increase the overall market efficiency.

Differently from the OCI model, in the MSI model (see Sect. 21.3), we consider a setting where an advertiser seeks to acquire an impression from an exchange through a chain of intermediaries. Using a game theoretic model, we study the mechanisms offered by intermediaries when the advertiser's value is private. We characterize a subgame perfect equilibrium of the game between intermediaries within the class of second-price mechanisms which are commonly used in the display advertising market. We show that economic incentives are not necessarily aligned along the chain, i.e., profit-maximizing intermediaries have incentives to shade bids and not to allocate impressions, even when profitable for their downstream customers. Moreover, we establish that the position in the intermediation network has a significant impact on the profits of the intermediaries, and the most profitable position depends on the underlying value distribution of the advertiser.

The proofs of the claims in Sects. 21.2 and 21.3 can be found in Balseiro and Candogan (2017) and Balseiro et al. (2017), respectively.

21.1.2 Literature Review

The models considered in this chapter contribute to various streams of literature, namely to those of intermediary problems, online advertising, and mechanism design with budget constraints.

21.1.2.1 Intermediary Problems

Feldman et al. (2010) and Stavrogiannis et al. (2013) focus on settings where captive buyers bid through intermediaries to acquire impressions. In these papers, the value distribution of the buyer has bounded support, and intermediaries directly bid at the exchange. As opposed to the OCI model, these papers focus on captive advertisers who are not liquidity constrained. Moreover, in these studies there are no “multiple intermediation tiers” whereas the MSI model analyzes how multiple intermediaries share surplus under different value distributions. Additionally, in these papers, intermediaries are restricted to forwarding the highest bid they receive from the buyers. However, in our models, we show that at equilibrium strategic intermediaries shade their bids, as opposed to simply reporting upstream the highest downstream bid. Loertscher and Niedermayer (2007, 2012) and Niazadeh et al. (2014) consider fee setting mechanisms for an intermediary with the two-sided private information setting introduced by Myerson and Satterthwaite (1983). These papers consider a single intermediary (as opposed to the MSI model) with a captive seller and captive buyer without budget constraints (as opposed to the OCI model), and provide conditions under which an affine fee structure that is commonly used in practice is optimal. In another recent work, Manea (2018) focuses on a setting where intermediaries trade a single good in a network, by considering a complete information setting in which buyers’ values are deterministic and common knowledge in contrast to our models. We refer the reader to Condorelli and Galeotti (2016) for a review of the recent literature on intermediation network models.

There is also a stream of papers that study intermediary problems in the context of supply chains (see, e.g., Belavina and Girotra 2012; Wu 2004; Nguyen et al. 2016), but these papers do not feature private information settings, unlike our models. Moreover, there are papers which consider the problem of successive monopolies in manufacturer-retailer settings with posted pricing schemes (see, e.g., Bresnahan and Reiss 1985; Lariviere and Porteus 2001; Perakis and Roels 2007). In the MSI model, we consider a more general class of mechanisms than posted pricing that allows intermediaries to elicit the private values of downstream agents and acquire the impression from upstream only if the downstream agents signal interest. Standard posted pricing mechanisms do not allow for this “contingent sale” feature, which is a prevalent in online advertising.

21.1.2.2 Online Advertising and Ad Exchanges

Mansour et al. (2012) give a brief synopsis of the auction employed by Google’s exchange, and discusses the role of intermediaries in this auction. Although this auction is not optimal or incentive compatible, the authors argue that the intermediaries’ incentives to misreport valuations is small in these auctions. The study by Ghosh et al. (2009) considers the advertiser’s perspective. In particular, Ghosh et al. (2009) study the design of a bidding agent who first explores the competing bids in the market, and then bids according to the observed empirical distribution. Balseiro

et al. (2015) and Gummadi et al. (2011) study the dynamic interactions among budget-constrained advertisers bidding in exchanges by proposing an approximate equilibrium concept. Iyer et al. (2014) characterize the bidding equilibrium among advertisers who learn about their private value over time in a repeated auction setting by using a mean-field approximation. Jiang et al. (2014) provide a simple bidding strategy, which does not require any statistical knowledge, for a single bidder with an average budget constraint. However these papers assume that advertisers bid directly in the exchange and do not consider the presence of intermediaries.

21.1.2.3 Mechanism Design with Budget Constraints

Finally, the OCI model considered in this chapter is related to a stream of papers that study mechanism design with financially constrained bidders (Laffont and Robert 1996; Che and Gale 1998, 2000; Maskin 2000). In these papers, efficient and optimal auctions are studied by modeling valuations as private and budgets as either private or public information. When the budgets are private, standard auction formats such as first-price and second-price auction are suboptimal and not revenue equivalent. Moreover, the problem in fact becomes a multi-dimensional mechanism design problem which is hard to solve in general. Pai and Vohra (2014) study the problem of selling one item in a setting with multiple budget-constrained buyers by using linear programming, and show that the optimal mechanism is implementable via an all-pay auction. Additionally, Chawla et al. (2011) show that the problem of selling one item to unit-demand buyers who are budget constrained can be reduced to an unconstrained problem with a small loss in performance. In multi-unit auctions with budget-constrained bidders, Borgs et al. (2005) and Bhattacharya et al. (2010) study approximation algorithms to design revenue-optimal incentive-compatible mechanisms. Brusco and Lopomo (2008, 2009) also study multi-item settings with budget constrained bidders by focusing on ascending auction formats. Differently from these papers, here we consider an intermediary who has no inherent value for the items sold and does not own the items at the moment of contracting, but instead needs to specify a mechanism to procure impressions as they arrive at the exchange. In addition, the main objective of this chapter is to understand the presence and profitability of such intermediaries.

21.2 Optimal Contracts for Intermediaries in Online Advertising

In this section, we provide the OCI model and the optimal mechanism design characterization. We consider a setting where an advertiser acquires impressions from an exchange over a fixed horizon (e.g., a week or a month) by either bidding directly at the exchange, or contracting with an *intermediary* to acquire impressions

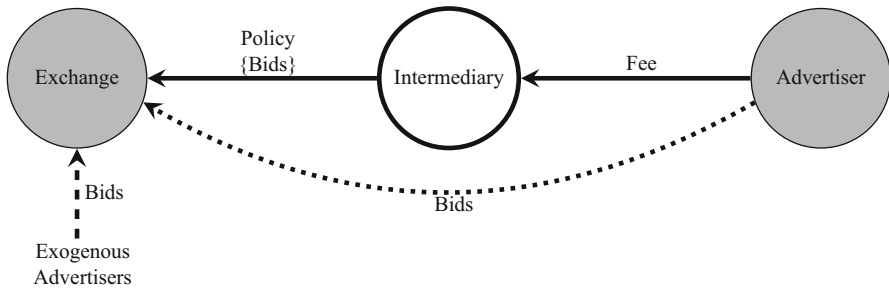


Fig. 21.1 The advertiser can submit bids directly to the exchange, or through the intermediary. The exchange has other (exogenous) bidders as well

on her behalf. This contracting advertiser is referred to as “the advertiser” when clear from the context, to distinguish her from other exogenous advertisers. We assume that a fixed number, n , of impressions arrive at the exchange over the horizon.¹ In the following, we denote vectors using boldface as in \mathbf{x} , and the transpose of this vector by \mathbf{x}^T .

Exchange The exchange sells impressions to the contracting advertiser and other exogenous advertisers using a second-price auction with no reserve (see Fig. 21.1). Rather than explicitly modeling the preferences of the exogenous advertisers, we denote the exogenous advertisers’ maximum competing bid by d_i for impression i . We assume that $(d_i)_{i=1}^n$ are i.i.d. and drawn from the cumulative distribution function $F_d(\cdot)$, with the strictly positive probability density function $f_d(\cdot) > 0$ over the compact support $[0, \bar{D}]$.

For each arriving impression i , the exchange announces some user information in the form of an *attribute* vector, which may affect the value of the impression perceived by the advertiser. We denote by $\alpha_i \in \mathcal{A}$ an attribute vector which contains relevant information for the advertiser’s targeting criterion (such as geographical location, age group of the viewer, tastes and interests obtained from her browsing history). The space of attributes \mathcal{A} is a compact subset of the Euclidean space. We assume that the random variables $(\alpha_i)_{i=1}^n$ are i.i.d. with cumulative distribution function $F_\alpha(\cdot)$, and strictly positive density $f_\alpha(\cdot) > 0$. We next describe how the attribute vector α_i impacts the advertiser’s valuation of the impressions.

Advertiser The advertiser has a budget and a targeting criterion which we denote by $b \in \mathbb{R}_+$ and $\theta \in \Theta$, respectively. Therefore, the type of the advertiser is given by the pair $t = (b, \theta)$ that belongs to a set $\mathcal{T} \subseteq \mathbb{R}_+ \times \Theta$. The type t is the private information of the advertiser. The set \mathcal{T} is assumed to be finite, and we denote by $T \triangleq |\mathcal{T}|$ its cardinality. The advertiser’s type is $t \in \mathcal{T}$ with probability $p_t > 0$. With some abuse of notation we denote by b_t and θ_t the budget and targeting

¹Random number of impressions can be accommodated in our model by considering dummy arrivals that are valued at zero by the advertiser.

criterion of type t , respectively. The value of a type t advertiser for an impression with attributes α is given by $v_t(\alpha)$. Here the function $v_t : \mathcal{A} \rightarrow \mathbb{R}_+$ is bounded and continuous in the impression attributes $\alpha \in \mathcal{A}$ for every type $t \in \mathcal{T}$.

Outside Option The outside option for the advertiser corresponds to running a campaign on her own, and thus the value of the outside option is the maximum surplus the advertiser can get by participating in the exchange directly. If an advertiser with type $t = (b, \theta)$ pursues the outside option, she is constrained to the set of feasible non-anticipative policies (i.e., policies that map the history to bids) that satisfy the budget constraint for every sample path. This set of policies is denoted by \mathcal{Z}_t . Since the exchange runs a second-price auction, a feasible policy $\zeta \in \mathcal{Z}_t$ should satisfy the inequality:

$$\sum_{i=1}^n \mathbf{1}\{z_i^\zeta \geq d_i\} d_i \leq b_t \text{ (almost surely),}$$

where $\mathbf{1}\{\cdot\}$ is the indicator function, and $z_i^\zeta \in [0, \bar{D}]$ corresponds to the bid from policy ζ for the i th impression. The optimal expected surplus of an advertiser with type t is denoted by V_t , and obtained by solving the following optimal control problem:

$$V_t \triangleq \sup_{\zeta \in \mathcal{Z}_t} \mathbb{E}_{\alpha, \mathbf{d}} \left[\sum_{i=1}^n \mathbf{1}\{z_i^\zeta \geq d_i\} (v_t(\alpha_i) - d_i) \right], \quad (21.1)$$

where the expectation is taken with respect to the vector of impression attributes $\alpha = (\alpha_i)_{i=1}^n$ and maximum competing bids $\mathbf{d} = (d_i)_{i=1}^n$.

21.2.1 Mechanism Design Problem

In this section we focus on the problem of the intermediary whose objective is to run a campaign on behalf of the advertiser that can alternatively pursue her own campaign by directly participating in the auctions of the exchange. Since the advertiser's type (her budget and targeting criterion) is her private information, the intermediary's problem can naturally be formulated as a mechanism design problem. By the Revelation Principle, without loss of optimality, we can focus on direct mechanisms where the advertiser reports her type t (possibly nontruthfully), and the intermediary responds to this report by choosing a required payment x_t , and a dynamic bidding policy $\zeta_t \in \mathcal{Z}$ he commits to running at the exchange on behalf of the advertiser. Here, we denote by \mathcal{Z} the set of all non-anticipative (and

potentially randomized) dynamic bidding policies that have no budget restrictions.² Specifically, a mechanism (\mathbf{x}, ζ) for the intermediary consists of a vector of payments $\mathbf{x} = (x_t)_{t \in \mathcal{T}}$, and a vector of non-anticipative dynamic bidding policies associated with different types $t \in \mathcal{T}$, $\zeta = (\zeta_t)_{t \in \mathcal{T}}$.

Note that the set \mathcal{Z} of all non-anticipative bidding policies is a high-dimensional set which includes all functions mapping every possible history to a bid. Therefore, instead of searching for the optimal mechanism by optimizing directly over this set (which may be computationally intractable), we provide an alternative technique which relies on first characterizing the performance that can be achieved by policies in \mathcal{Z} . The performance of a policy $\zeta \in \mathcal{Z}$ can be measured by two metrics: (i) the intermediary’s total expected cost for running this policy

$$\mathcal{C}(\zeta) \triangleq \mathbb{E}_{\alpha, \mathbf{d}} \left[\sum_{i=1}^n \mathbf{1}\{z_i^\zeta \geq d_i\} d_i \right],$$

and (ii) the total expected value

$$\mathcal{W}_t(\zeta) = \mathbb{E}_{\alpha, \mathbf{d}} \left[\sum_{i=1}^n \mathbf{1}\{z_i^\zeta \geq d_i\} v_t(\alpha_i) \right]$$

an advertiser of type $t \in \mathcal{T}$ derives from the impressions acquired by this policy. We consider an optimal mechanism design problem where any performance level in this achievable performance space can be chosen by the intermediary.

Definition 1 The *achievable performance space* \mathcal{P} is given by the set of points $(c, \mathbf{w}) \subseteq \mathbb{R} \times \mathbb{R}^T$ such that there exists a policy $\zeta \in \mathcal{Z}$ satisfying $\mathcal{W}_t(\zeta) = w_t$ for all $t \in \mathcal{T}$ and $\mathcal{C}(\zeta) \leq c$.

Following the performance space idea, the optimal mechanism can be derived by optimizing over the costs and total expected values associated with feasible policies (as opposed to optimizing over the policies themselves). In comparison to the set \mathcal{Z} , the performance space has a much smaller dimension because it does not scale with the time horizon, thereby yielding a more tractable formulation. In addition, after the optimal performance levels associated with the optimal mechanism are determined, an optimal policy that achieves the chosen performance level can be retrieved via a *synthesis* procedure (described in Sect. 21.2.2.3).

We next introduce some notation. We use $c_{t'}$ to represent the total expected cost incurred by the intermediary when the advertiser reports her type as t' , and the intermediary uses the policy $\zeta_{t'}$ to bid on behalf of her at the exchange. If the true type of this advertiser is t , we denote her total expected value for the

²Unlike the advertiser, the intermediary does not have stringent financial constraints, thus we do not restrict the intermediary’s policies $\zeta \in \mathcal{Z}$ to satisfy any budget constraints (unlike the set of policies \mathcal{Z}_t that can be employed by the advertiser of type t).

impressions acquired by this policy by $w_{t',t}$, and the matrix of total expected values by $\mathbf{W} = (w_{t',t})_{t',t \in \mathcal{T}} \in \mathbb{R}^{T \times T}$. Note that the policy run for type t' may yield different expected values for advertisers of different types, i.e., for types t_1, t_2 advertisers, we may have $w_{t',t_1} \neq w_{t',t_2}$ because their targeting criteria may differ. Using this notation, the mechanism design problem of the intermediary can be stated as follows:

$$\max_{\mathbf{x}, \mathbf{c} \in \mathbb{R}^T, \mathbf{W} \in \mathbb{R}^{T \times T}} \sum_{t \in \mathcal{T}} p_t (x_t - c_t) \quad (21.2a)$$

$$\text{(OPT)} \quad \text{s. t.} \quad w_{t,t} - x_t \geq w_{t',t} - x_{t'}, \quad \forall t, t' : b_{t'} \leq b_t, \quad (21.2b)$$

$$w_{t,t} - x_t \geq V_t, \quad \forall t, \quad (21.2c)$$

$$x_t \leq b_t, \quad \forall t, \quad (21.2d)$$

$$(c_{t'}, (w_{t',t})_{t \in \mathcal{T}}) \in \mathcal{P}, \quad \forall t'. \quad (21.2e)$$

In (OPT), we maximize the expected profit of the intermediary, which is given by the quantity $\sum_{t \in \mathcal{T}} p_t (x_t - c_t)$. Here, the variable x_t represents the payment of the advertiser whose report is t and c_t is the cost of running the campaign. The *incentive compatibility (IC)* constraint Eq. (21.2b) ensures that the advertiser maximizes her payoff by reporting her type truthfully. Note that the surplus of a type t advertiser whose report is t' is expressed by the quantity $w_{t',t} - x_{t'}$. For the IC constraint, we restrict attention to the cases where the advertiser does not report a budget larger than her true budget without loss of generality.³ Recall that the advertiser could alternatively run her own campaign. Therefore, we guarantee that the mechanism delivers a utility at least equal to the outside option to the advertiser with the *individual rationality (IR)* constraint Eq. (21.2c). Moreover, the payment collected by the mechanism does not exceed the advertiser's budget due to the budget constraint Eq. (21.2d). The constraint Eq. (21.2e) guarantees that given an optimal solution of (OPT), the intermediary can find a bidding policy ζ_t that delivers the performance (to all types) as required by the optimal solution. Specifically, this constraint ensures that the performance of the mechanism of the intermediary lies in the achievable performance space \mathcal{P} , thereby implying that the structure of \mathcal{P} plays a key role in the solution of (OPT). Hence, we conclude this section by emphasizing that the performance space \mathcal{P} is convex and closed, which implies that (OPT) is a convex optimization problem.

Lemma 1 *The performance space \mathcal{P} is convex and closed.*

³Specifically, the intermediary can prevent the advertiser from overstating her budget by requiring her to make an upfront payment equal to the reported budget, and returning the amount $b_t - x_t$ at the end of the advertising campaign (see, e.g., Che and Gale 2000).

21.2.2 Optimal Mechanism Characterization

Because the performance space \mathcal{P} does not have closed-form description, solving (OPT) directly is algorithmically challenging. To overcome this difficulty, we introduce a duality-based approach for the optimal mechanism design problem which consists of dualizing some constraints of (OPT). This dual problem is a tractable convex minimization problem that can be solved efficiently.

21.2.2.1 Dual Problem

In the dual approach we dualize the IC, IR and budget constraints of (OPT), while optimizing over the performance space. To provide a characterization of the dual problem, we employ the concept of support function of the performance space. The support function $\phi : \mathbb{R} \times \mathbb{R}^T \rightarrow \mathbb{R}$ associated with the performance space \mathcal{P} for a point $(\mu, \boldsymbol{\lambda}) \in \mathbb{R} \times \mathbb{R}^T$ is given by

$$\phi(\mu, \boldsymbol{\lambda}) \triangleq \sup_{(c, \mathbf{w}) \in \mathcal{P}} \boldsymbol{\lambda}^\top \mathbf{w} - \mu c. \quad (21.3)$$

The support function ϕ gives the intercept of the supporting hyperplane of the performance space \mathcal{P} with normal $(-\mu, \boldsymbol{\lambda})$. The support function ϕ is convex because it is obtained as the point-wise supremum of linear functions over the performance space.

We simplify the derivation of the dual problem by momentarily writing the IC, IR and budget constraints in matrix form. In particular, we define $\mathbf{w}_t \triangleq (w_{t,t'})_{t' \in \mathcal{T}}$ as the t th row vector of the matrix $\mathbf{W} \in \mathbb{R}^{T \times T}$, and rewrite (OPT) as follows:

$$\max_{\mathbf{x}, \mathbf{c}, \mathbf{W}} \sum_{t \in \mathcal{T}} p_t(x_t - c_t) \quad (21.4a)$$

$$\text{s. t.} \quad \sum_{t \in \mathcal{T}} \mathbf{d}_t x_t - \mathbf{A}_t \mathbf{w}_t \leq \mathbf{e} \quad (21.4b)$$

$$(c_t, \mathbf{w}_t) \in \mathcal{P}, \quad \forall t. \quad (21.4c)$$

where $\mathbf{d}_t \in \mathbb{R}^M$, $\mathbf{A}_t \in \mathbb{R}^{M \times T}$ and $\mathbf{e} \in \mathbb{R}^M$ capture the coefficients of the $M \leq T^2 + 2T$ linear inequalities corresponding to the IC, IR and budget constraints.

Let $\boldsymbol{\lambda} \geq 0$ in \mathbb{R}^M be the Lagrange multiplier of constraints Eq. (21.4b). Using these multipliers, we obtain the convex dual problem (D):

$$\min_{\boldsymbol{\lambda} \in \mathbb{R}^M} \boldsymbol{\lambda}^\top \mathbf{e} + \sum_{t \in \mathcal{T}} \phi(p_t, \boldsymbol{\lambda}^\top \mathbf{A}_t) \quad (21.5a)$$

$$\text{(D) s. t.} \quad \boldsymbol{\lambda}^\top \mathbf{d}_t = p_t, \quad \forall t \quad (21.5b)$$

$$\boldsymbol{\lambda} \geq 0. \quad (21.5c)$$

In the following theorem, we show that strong duality holds for (OPT) and (D). This result is established in Balseiro and Candogan (2017) by first constructing a feasible solution for which the performance level associated with each type t belongs to the relative interior of \mathcal{P} . Using this constraint qualification with the known results from duality theory, the result follows.

Theorem 1 (OPT) admits an optimal solution. Additionally, strong duality holds, that is, the optimal objective value of (OPT) and (D) coincide.

We next show that the support function ϕ can be efficiently evaluated, and then discuss how to construct the optimal mechanism based on an optimal dual solution.

21.2.2.2 Support Function Characterization

Note that it is possible to characterize the support function more explicitly, by restating the expression in Eq. (21.3) using Definition 1:

$$\phi(\mu, \lambda) = \sup_{\zeta \in \mathcal{Z}, c \in \mathbb{R}, \mathbf{w} \in \mathbb{R}^T} \{ \lambda^\top \mathbf{w} - \mu c \text{ s.t. } \mathcal{C}(\zeta) \leq c, w_t = \mathcal{W}_t(\zeta) \}. \quad (21.6)$$

We next provide a closed-form expression for the support function as well as the optimal solution of Eq. (21.6). Let $\zeta^\phi(\mu, \lambda)$ be a policy that bids $z_i^{\zeta^\phi(\mu, \lambda)} = \sum_t \gamma_t v_t(\alpha_i) / \mu$ for impression i with attributes α_i , and define $c^\phi(\mu, \lambda) \triangleq \mathcal{C}(\zeta^\phi(\mu, \lambda))$ as the total expected cost of the policy and $w_t^\phi(\mu, \lambda) \triangleq \mathcal{W}_t(\zeta^\phi(\mu, \lambda))$ as the total expected value type $t \in \mathcal{T}$ has for the impressions acquired by this policy. In the following proposition we denote by $x^+ = \max(x, 0)$ the positive part of $x \in \mathbb{R}$.

Proposition 1 Suppose that $\mu > 0$. In Eq. (21.6), an optimal ζ is $\zeta^\phi(\mu, \lambda)$, and the unique optimal c and \mathbf{w} are given by $c^\phi(\mu, \lambda)$ and $\mathbf{w}^\phi(\mu, \lambda)$, respectively. More explicitly, the support function is given by

$$\phi(\mu, \lambda) = n \mathbb{E}_{\alpha, d} \left[\left(\sum_{t \in \mathcal{T}} \gamma_t v_t(\alpha) - \mu d \right)^+ \right],$$

the optimal total expected cost is $c^\phi(\mu, \lambda) = n \mathbb{E}_{\alpha, d} [d \mathbf{1}\{\sum_{t \in \mathcal{T}} \gamma_t v_t(\alpha) \geq \mu d\}]$, and the optimal total expected value for type $t \in \mathcal{T}$ is given by

$$w_t^\phi(\mu, \lambda) = n \mathbb{E}_{\alpha, d} \left[v_t(\alpha) \mathbf{1} \left\{ \sum_{t \in \mathcal{T}} \gamma_t v_t(\alpha) \geq \mu d \right\} \right].$$

Note that the support function can be expressed as a simple expectation (of a piecewise-linear function) over the impression attributes and the competing bid (Proposition 1). Moreover, the convex dual problem has a compact representation

(recall that the dual problem has polynomial size). Hence, it follows that the dual problem is tractable and its optimal solution can be obtained by using standard convex optimization algorithms.

21.2.2.3 Synthesis

In this section, we provide a procedure in which the optimal solution of the dual is used to “synthesize” the optimal contract of the intermediary by relying on strong duality. In particular, given an optimal solution λ^* of the dual problem (D), we first construct the optimal policy of the intermediary ζ_t^* and then obtain the corresponding performance (c_t^*, \mathbf{w}_t^*) for all $t \in \mathcal{T}$. We then solve a linear feasibility problem that involves the constructed performance levels in order to determine the upfront fees $\{x_t^*\}_{t \in \mathcal{T}}$ associated with the optimal contract. More formally, the steps of our synthesis procedure can be given as follows:

- Step 1. Determine an optimal solution λ^* of the dual problem (D), and set $\lambda_t^* = (\gamma_{t,t'}^*)_{t' \in \mathcal{T}} = (\lambda^*)^\top \mathbf{A}_t$ and $\mu_t^* = p_t$ for $t \in \mathcal{T}$.
- Step 2. Set the policy of type $t \in \mathcal{T}$ to $\zeta_t^* \triangleq \zeta^\phi(\mu_t^*, \lambda_t^*)$. The total expected cost and value for this policy are respectively given as $c_t^* \triangleq c^\phi(\mu_t^*, \lambda_t^*)$, and $\mathbf{w}_t^* \triangleq \mathbf{w}^\phi(\mu_t^*, \lambda_t^*)$ (see Proposition 1).
- Step 3. Determine upfront fees $\mathbf{x}^* = (x_t^*)_{t \in \mathcal{T}}$ by solving

$$\sum_{t \in \mathcal{T}} \mathbf{d}_t x_t^* - \mathbf{A}_t \mathbf{w}_t^* \leq \mathbf{e} \quad \perp \quad \lambda^* \geq 0,$$

where \perp indicates that for each entry, at least one of these inequalities should hold with equality.

Theorem 2 *The synthesis procedure yields an optimal solution $(\mathbf{x}^*, \mathbf{c}^*, \mathbf{W}^*)$ for (OPT) and an optimal mechanism (\mathbf{x}^*, ζ^*) for the intermediary.*

We note that the mechanism (\mathbf{x}^*, ζ^*) provided in Theorem 2 is relatively easy to implement. In this mechanism, the intermediary posts a menu of contracts (associated with different payments and policies), and invites the advertiser to choose the one she prefers. We next turn to the optimal policies of the intermediary and obtain further insights on these bidding policies.

21.2.2.4 Optimal Bidding Policy

Theorem 2 implies that the optimal bidding strategy has a surprisingly simple structure. In particular, assume that an impression with an attribute vector α arrives at the exchange, and the advertiser is of type t . Proposition 1 suggests that under the optimal policy the intermediary bids:

$$z_t^*(\alpha) = \frac{1}{p_t} \sum_{t' \in \mathcal{T}} \gamma_{t,t'}^* v_{t'}(\alpha). \quad (21.7)$$

It is possible to obtain further intuition on the bidding strategy of the intermediary. In particular, let $\lambda^* = (\lambda^{\text{IC}}, \lambda^{\text{IR}}, \lambda^{\text{B}})$, where $\lambda^{\text{IC}} = (\lambda_{t',t}^{\text{IC}})_{t',t}$, $\lambda^{\text{IR}} = (\lambda_t^{\text{IR}})_t$, $\lambda^{\text{B}} = (\lambda_t^{\text{B}})_t$, and $\lambda_{t',t}^{\text{IC}} \geq 0$, $\lambda_t^{\text{IR}} \geq 0$, $\lambda_t^{\text{B}} \geq 0$ respectively denote the optimal Lagrange multipliers associated with constraints (21.2b), (21.2c), and (21.2d) of (OPT). The next result characterizes the optimal bidding policy of the intermediary in terms of these multipliers.

Corollary 1 *Let $\lambda^* = (\lambda^{\text{IC}}, \lambda^{\text{IR}}, \lambda^{\text{B}})$ be a dual optimal solution of (D). The optimal bidding policy of the intermediary ζ_t^* for type $t \in \mathcal{T}$ is such that for an impression with attribute vector α it bids:*

$$z_t^*(\alpha) = \left(1 - \frac{\lambda_t^{\text{B}}}{p_t}\right) v_t(\alpha) + \frac{1}{p_t} \sum_{t' \in \mathcal{T}} \lambda_{t,t'}^{\text{IC}} (v_t(\alpha) - v_{t'}(\alpha)). \quad (21.8)$$

This corollary reveals that if the IC and budget constraints are not active (and hence $\lambda^{\text{IC}} = \lambda^{\text{B}} = 0$), the intermediary simply reports the true value $v_t(\alpha)$. For example, this case occurs when the advertiser has a large budget, and the impression distribution is uniform (since in this case the IC constraints are not binding, see Balseiro and Candogan (2017) for more details). Similarly, if the budget constraint is binding, but the IC constraints are not, then the intermediary simply shades the true value $v_t(\alpha)$ by $1 - \lambda_t^{\text{B}}/p_t$ to account for the advertiser's budget constraint.

When the IC constraints are active, the bidding strategy takes into account that type t may have incentive to misreport her type as t' (the dual multiplier of the corresponding incentive compatibility constraint is $\lambda_{t,t'}^{\text{IC}} \geq 0$). For instance, suppose that types t and t' have the same budget but t' has lower average valuations for the impressions than t . Type t' typically pays a lower amount for the impressions she acquires, and hence effectively has a less stringent budget constraint. Thus, the intermediary needs to bid more aggressively to match the outside option of type t' , which may give the type with higher valuations an incentive to impersonate the lower type. The term $(1/p_t) \sum_{t' \in \mathcal{T}} \lambda_{t,t'}^{\text{IC}} (v_t(\alpha) - v_{t'}(\alpha))$ in Eq. (21.8) ensures that the intermediary bids more aggressively for impressions that type t highly values when compared to other types (i.e., $v_t(\alpha) > v_{t'}(\alpha)$), thereby eliminating the incentive of type t to misreport her type.

21.2.3 Economic Insights

The presence of intermediation affects the online advertising market in several ways. As an immediate example, the optimal bidding policy of an advertiser participating in the exchange on her own has a complex dynamic shading structure, which is

obtained by solving a dynamic program. On the other hand, the policy associated with the optimal contract (see Corollary 1) of the intermediary has a stationary structure. Thus, from an operational perspective, the presence of the intermediary results in simpler bidding policies in the exchange. In addition to this impact on the bidding policies, we shed light on the other economic insights derived from the optimal mechanism of the intermediary by considering the following aspects.

21.2.3.1 Intermediation Profit and the Advertiser Surplus

Recall that the advertiser has the option of bidding directly in the exchange. Therefore, the intermediary should guarantee that the advertiser has incentive to accept the contract. In other words, the contract of the intermediary should deliver surplus to the advertiser at least as high as the surplus of the outside option. Despite providing this surplus to the advertiser, interestingly the intermediary still manages to profit. The reason behind this observation can be explained by the capability of the intermediary to deliver the surplus of the advertiser's outside option at a lower expected cost than the advertiser could achieve on her own. Specifically, the advertiser's bidding policies, when bidding directly in the exchange, have to satisfy her budget at every realization (of impression attributes and competing bids) whereas the intermediary has no budget constraints. Therefore, the advertiser's bidding policies might need to significantly shade her bid in each auction while the intermediary can implement bidding policies that exceed the upfront fee in some realizations. In other words, the absence of these financial constraints equips the intermediary with a richer set of policies, and allows him to deliver the same surplus to the advertiser at a lower cost; therefore making it profitable to intermediate.

21.2.3.2 Market Efficiency

Note that the optimal mechanism of the intermediary and the optimal bidding policy are provided for general valuation structures. Balseiro and Candogan (2017) provide further insights by focusing on a special case where advertisers' targeting criteria exhibit symmetry, e.g., uniformly distributed impression attributes. Under this assumption, the optimal mechanism can be characterized in closed form. Balseiro and Candogan (2017) also establish that the presence of the intermediary in the market not only allows him to profit, but also increases the surplus of the advertiser and market efficiency as given by the sum of the revenue of the exchange, the intermediary's profit, the contracting advertiser's surplus, and the exogenous bidders' surpluses.

21.3 Multi-stage Intermediation in Display Advertising

In this section, we consider a setting with multiple intermediaries positioned between an advertiser and an exchange (see Fig. 21.2). The exchange sells a unique indivisible impression via a second-price auction with no reserve in which exogenous advertisers might participate. Since we focus on the trade that occurs through the intermediation chain, as before, we model the highest bid of those exogenous advertisers by a random variable d with support $[0, \bar{D}]$. The advertiser seeks to purchase the impression from the exchange, and her value for the impression is captured by a random variable denoted by v .⁴ The realizations of v are drawn from the cumulative distribution function $G_v(\cdot)$, with the strictly positive probability density function $g_v(\cdot)$ over the support $\mathcal{V} \subseteq [0, \infty)$. We assume that the distribution of v is common knowledge, and its realizations are the private information of the advertiser. We also assume that the expected value of v is finite, i.e., $\mathbb{E}[v] < \infty$.

As opposed to the OCI model, the advertiser has no budget constraint, however, she is constrained to purchase impression from the exchange through a chain of $m \geq 1$ intermediaries. In Fig. 21.2, the intermediation chain is illustrated. We refer to the intermediaries closer to the advertiser as *downstream* intermediaries, and the intermediaries closer to the exchange as *upstream* intermediaries. The intermediaries have no value for the impression, thus they only profit in case of purchasing the impression from upstream and reselling to downstream. Therefore, a mechanism for an intermediary should (i) map reports from the downstream intermediary to an upstream bid to purchase the impression from upstream, and (ii) decide on an allocation and a payment to resell the impression to downstream.

For the set of available mechanisms, we focus on the set of second-price mechanisms \mathcal{M} where a mechanism $(r, Y) \in \mathcal{M}$ consists of a reserve price $r \in \mathbb{R}_+$ and a nondecreasing reporting function $Y : \mathbb{R}_+ \rightarrow \mathbb{R}_+$. After receiving a downstream report x , intermediary j with mechanism (r_j, Y_j) reports $Y_j(x)$ if $x \geq r_j$, otherwise reports zero. The intermediary allocates the impression only if she wins it from upstream. The payment is determined as the minimum amount which guarantees winning of the downstream agent. Formally, given the exogenous bid $d \in [0, \bar{D}]$ at the exchange, we denote by \mathcal{W}_j the set of bids of intermediary I_j that guarantees winning, i.e.,

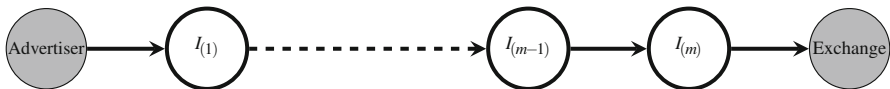


Fig. 21.2 A chain of intermediaries

⁴Note that in the OCI model we denote by $v_t(\alpha)$ the value of a type t advertiser for an impression with attributes α , thus the notation v_t represents a function. However, in the MSI model, the notation v represents a random variable which captures the value of the advertiser.

$$\mathscr{W}_j \triangleq \{x \geq 0 \mid x \geq r_{j+1}, Y_{j+1}(x) \geq r_{j+1}, \dots, Y_m \circ \dots \circ Y_{j+1}(x) \geq d\},$$

where Y_i and r_i respectively denote the reporting function and reserve of intermediary $I_{(i)}$ for $i = j + 1, \dots, m$, and \circ denotes the composition operator. The intermediary $I_{(j)}$ makes a payment to $I_{(j+1)}$ only in case of winning, and the payment amount is given by the smallest element of the winning reports \mathscr{W}_j , i.e., $\inf(\mathscr{W}_j)$. This payment rule is a natural extension of the second-price auction, which is commonly used in the context of display advertising. Therefore, the set of second-price mechanisms \mathscr{M} extends second-price auctions to a setting with intermediaries. We next provide some important properties of the set \mathscr{M} by the following lemma.

Lemma 2 *Suppose that intermediary $I_{(j)}$ selects her mechanism $(r_j, Y_j) \in \mathscr{M}$ for $j = 1, \dots, m$, and the exchange runs a second-price auction. Then,*

(a) *The composition of upstream mechanisms faced by the agent $I_{(j)}$ for $j = 0, \dots, m - 1$ (where $I_{(0)}$ represents the advertiser) is equivalent to a second-price auction with an exogenous random bid d^* and reserve price r^* given by*

$$d^* = Y_{j+1}^{-1} \circ \dots \circ Y_m^{-1}(d)$$

$$r^* = \max(r_{j+1}, Y_{j+1}^{-1}(r_{j+2}), Y_{j+1}^{-1} \circ Y_{j+2}^{-1}(r_{j+3}), \dots, Y_{j+1}^{-1} \circ \dots \circ Y_m^{-1}(0))$$

where we define the inverse of a reporting function Y as $Y^{-1}(x) = \inf\{\tilde{x} \geq 0 : Y(\tilde{x}) \geq x\}$.

(b) *Truthful bidding is an optimal strategy for the advertiser.*

The second item in this lemma implies that the truthful bidding is a best response for the advertiser regardless the mechanisms of the intermediaries. Therefore, we exclude the strategic behavior of the advertiser and focus on the game among intermediaries. We model the corresponding mechanism design game among intermediaries as a Stackelberg game where intermediaries move sequentially from upstream to downstream. In particular, the timing of the events is as follows:

1. The advertiser privately draws her value.
2. Intermediary $I_{(j)}$ determines a mechanism $(r_j, Y_j) \in \mathscr{M}$ after observing the mechanism of $I_{(j+1)}$ for $j = m, \dots, 1$.
3. The advertiser bids truthfully.
4. Intermediary $I_{(j)}$ learns the bid of $I_{(j-1)}$ and submits her own bid to $I_{(j+1)}$ according to her own mechanism, for $j = m, \dots, 1$.
5. The exchange $I_{(m+1)}$ runs a second-price auction and the impression is allocated either to exogenous bidders or to the advertiser through the intermediation chain.
6. Intermediary $I_{(j)}$ learns her payment to $I_{(j+1)}$ and charges a payment to $I_{(j-1)}$ for $j = m, \dots, 1$.

21.3.1 Equilibrium Characterization

We study the outcome of the strategic interaction among intermediaries by focusing on the subgame perfect equilibria (SPE) of the induced game as a solution concept and provide an SPE. Before stating our results, we provide some definitions that would be useful to characterize an SPE of this game. We first introduce the virtual value function and its inverse for a generic random variable X with a finite mean $\mathbb{E}[X] < \infty$, and c.d.f. $G_X(\cdot)$ and strictly positive p.d.f. $g_X(\cdot)$ over a nonnegative support $\mathcal{X} \subseteq [0, \infty)$.⁵ The virtual value function of X is given by

$$\phi_X(x) = x - \frac{1 - G_X(x)}{g_X(x)}$$

for $x \in \mathcal{X}$. Moreover, the inverse of the virtual value function is defined as $\phi_X^{-1}(x) \triangleq \inf\{\tilde{x} \in \mathcal{X} \mid \phi_X(\tilde{x}) \geq x\}$. We next define the *projected virtual value* function for random variables with strictly increasing virtual functions. This function projects the virtual value function to nonnegative reals, while extending its domain to \mathbb{R} .

Definition 2 Suppose X is an absolutely continuous random variable with a strictly increasing virtual value function $\phi_X(\cdot)$ and support \mathcal{X} . The projected virtual value function of X is given by

$$\psi_X(x) \begin{cases} \sup \mathcal{X} & x \geq \sup \mathcal{X} , \\ \phi_X(x) & z_X \leq x < \sup \mathcal{X} , \\ 0 & \text{otherwise,} \end{cases} \quad (21.9)$$

where the projection point is given by $z_X = \phi_X^{-1}(0)$. If the random variable X has an atom at zero and is absolutely continuous elsewhere in its support \mathcal{X} , then we define $\psi_X(\cdot)$ and z_X similarly by replacing $\phi_X(\cdot)$ in Eq. (21.9) with the virtual value $\phi_{X|X>0}(\cdot)$ of the strictly positive part of X , denoted by $X|X > 0$.⁶

In order to characterize an SPE of the game in the chain of intermediaries, we first focus on the mechanism design problem of a single intermediary positioned between the advertiser and the exchange, i.e., $m = 1$. The following lemma provides an optimal mechanism for this case.

⁵The advertiser's value v satisfies these requirements.

⁶Note that $\phi_X(x) = \phi_{X|X>0}(x)$ for $x \in \mathcal{X} \setminus \{0\}$. This can be seen by using the definition of the virtual value function and noting that the conditional random variable $X|X > 0$ has c.d.f. $G_{X|X>0}(x) = (G_X(x) - G_X(0))/(1 - G_X(0))$, and p.d.f. $g_{X|X>0}(x) = g_X(x)/(1 - G_X(0))$. Thus, focusing on the virtual value of $X|X > 0$ as opposed to X , excludes the atom at zero, without impacting the (projected) virtual values elsewhere.

Lemma 3 *Suppose that there exists a single intermediary between the advertiser and the exchange, i.e., $m = 1$. Then, an optimal mechanism $(r, Y) \in \mathcal{M}$ for the intermediary is given by*

$$Y(x) = \psi_v(x),$$

$$r = z_v.$$

When choosing her bidding strategy for the upstream auction, an intermediary needs to optimally tradeoff the probability of winning the impression in the auction of the exchange with the cost incurred when winning the impression (both of which increase with the bid of the intermediary). This lemma reveals that the intermediary's optimal bidding strategy takes a simple structure: the intermediary first determines the advertiser's virtual value function, and then bids at the exchange the projected virtual value of the report from downstream.

In the case of multiple intermediaries, there are two factors which can affect the strategy of an intermediary: (i) the mechanisms chosen by upstream intermediaries, and (ii) the reaction of downstream intermediaries. For the first factor, Lemma 2 suggests that in settings with multiple intermediaries, the upstream mechanism that an intermediary faces can equivalently be represented by a second-price auction. Therefore, in light of Lemma 3, an intermediary along the chain is not influenced by the upstream decisions and, in turn, her actions do not influence downstream mechanisms. For the second factor, note that the reports observed by an intermediary $I_{(j)}$ can be obtained by composing the reporting functions of intermediaries $I_{(j-1)}, \dots, I_{(1)}$ with the report of the advertiser. Specifically, each intermediary in the chain can be thought as a single intermediary which connects a "downstream agent" with "anticipated reports" induced by the optimal downstream mechanisms along the equilibrium path to an upstream second-price auction. This observation allows for characterizing the optimal mechanisms of intermediaries recursively via Lemma 3, starting with the downstream intermediaries whenever the downstream reports satisfy the regularity conditions as the advertiser's value does. Therefore, we first formally define the anticipated reports and provide the regularity assumption, and characterize an SPE of the game between multiple intermediaries.

Definition 3 The anticipated report of the advertiser is $W_0 = v$, and the anticipated report of intermediary $I_{(j)}$ is

$$W_j = \psi_{W_{j-1}}(W_{j-1}).$$

Note that the anticipated report W_j coincide with the report of intermediary $I_{(j)}$ to the upstream mechanism if all intermediaries use the projected virtual value functions of the downstream bids as reporting functions.

Assumption 3 *The anticipated report W_j of intermediary $I_{(j)}$ is well-defined and has finite expected values for $j = 1, \dots, m$. Moreover, the anticipated reports have strictly increasing virtual values, i.e., $\phi_{W_j}(\cdot)$ (or $\phi_{W_j|W_j>0}(\cdot)$ if W_j has an atom at zero) is strictly increasing for $j = 1, \dots, m$.*

The following theorem provides an SPE under Assumption 3.⁷

Theorem 4 *Suppose that Assumption 3 holds, and consider the mechanisms given by*

$$Y_j^*(x) = \psi_{W_{j-1}}(x),$$

$$r_j^* = z_{W_{j-1}},$$

for $j = 1, \dots, m$. Then, the strategy profile where the mechanism (r_j^*, Y_j^*) is chosen by intermediary $I_{(j)}$ for $j = 1, \dots, m$ constitutes an SPE.

21.3.2 Economic Insights

Using the SPE provided in Theorem 4, we numerically explore the impact of the advertiser's value distribution on the reports and profits of intermediaries in different positions. Finally, we compare our results with the double-marginalization literature.

Intermediaries' Bids Since the virtual value function lies below the 45 degree line, Theorem 4 shows that the intermediaries always shade their bids (i.e., the report functions satisfy $Y_j^*(x) \leq x$). Although bid shading is always present, our numerical studies illustrate that it takes a different structure depending on the advertiser's value distribution. Figure 21.3 indicates that as a result of bid shading in long intermediation chains, unless the value of the buyer is significantly large, the bid submitted by the intermediaries to the auction of the exchange can be equal to zero. In such cases, even when it is profitable for the buyer, the intermediation chain does not allocate the impression to the advertiser, thereby causing inefficiency.

Intermediaries' Profits We next study the impact of an intermediaries position in the chain on her profits. On one hand, downstream intermediaries closer to the advertiser receive higher bids. On the other hand, upstream intermediaries closer to the exchange incur lower payments to acquire the impression. The total contribution of these opposing effects is indeterminate, and the impact of the position on profits depends on the distribution of the advertiser's value.

⁷Balseiro et al. (2017) show that this assumption holds for Generalized Pareto Distributions, a large family of distributions including uniform, exponential and Pareto distributions.

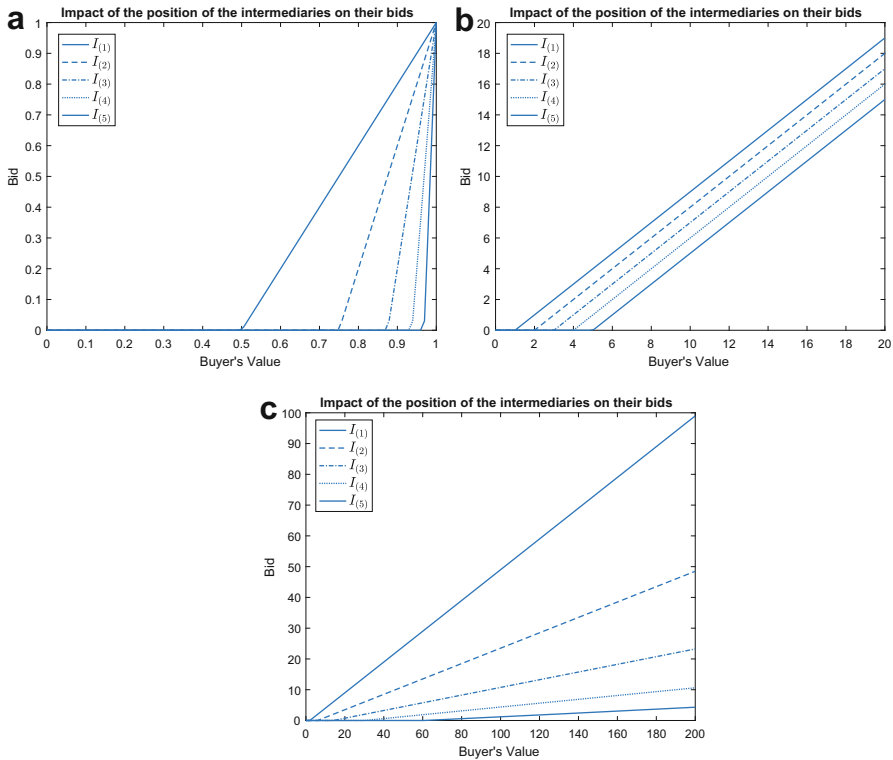


Fig. 21.3 This figure plots the bids of the intermediaries for different realization of the advertiser’s value v in a market with $m = 5$ intermediaries when the value distribution is uniform, exponential and shifted Pareto ((a) Uniform distribution ($\mathcal{V} = [0, 1]$, $\mathbb{E}[v] = 1/2$), (b) Exponential distribution ($\mathcal{V} = [0, \infty)$, $\mathbb{E}[v] = 1$), (c) S. Pareto distribution ($\mathcal{V} = [0, \infty)$, $\mathbb{E}[v] = 2/3$))

For example, when values are exponentially distributed the expected profit of each intermediary is the same, while for the uniform and shifted Pareto distributions the profits depend on the position of the intermediary in the chain. In particular, Fig. 21.4 plots the impact of the intermediary’s position on her profits conditional on winning the impression in a market with $m = 5$ intermediaries when the distribution of values is uniform, exponential and shifted Pareto. This figure shows that for uniform and Pareto distributions, profits as a function of an intermediary’s position in the chain, exhibit different trends. When the advertiser has a heavy-tailed value distribution, such as the shifted Pareto distribution, the downstream intermediaries who are closer to the advertiser have higher profits. Intuitively, this result stems from the fact that for such distributions, with significant probability the value of the advertiser for the impression is large, and hence the intermediaries that are closer to the advertiser can claim significant profits. Conversely, when the advertisers value distribution has a light-tail or a bounded support, the intermediaries find it more profitable to be closer to the exchange, due to lower costs of acquiring impressions.

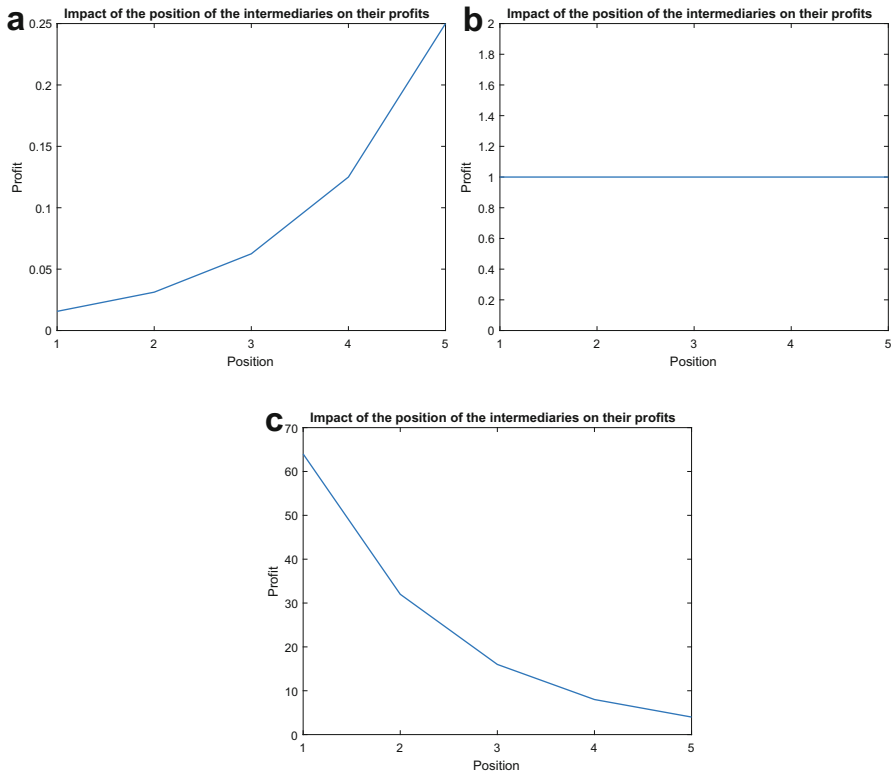


Fig. 21.4 This figure plots the profits of intermediaries for different positions in a market with $m = 5$ intermediaries when the value distribution is uniform, exponential and shifted Pareto ((**a**) Uniform distribution ($\mathcal{V} = [0, 1]$, $\mathbb{E}[v] = 1/2$), (**b**) Exponential distribution ($\mathcal{V} = [0, \infty)$, $\mathbb{E}[v] = 1$), (**c**) S. Pareto distribution ($\mathcal{V} = [0, \infty)$, $\mathbb{E}[v] = 2/3$))

This result suggests that depending on the advertiser’s value distribution for the impressions, intermediaries may prefer to participate in different stages of the intermediation process.

Comparison with Double-Marginalization Literature Our work is closely related to the double-marginalization literature (see, e.g., Tirole 1988), and some of the insights from this literature translate to our settings. In the basic double-marginalization framework, a manufacturer supplies a good to a single downstream retailer, who resells the good as a monopolist. Compared to the vertically integrated industry, the theory predicts that in the case of double-marginalization the price paid by the consumers is higher, industry profits, and the overall market efficiency are lower.

In our setting the first-best corresponds to the case when the advertiser bids directly in the exchange’s auction. Compared to the first-best, it is not hard to see that as the number of intermediaries increases, the expected surplus of the advertiser,

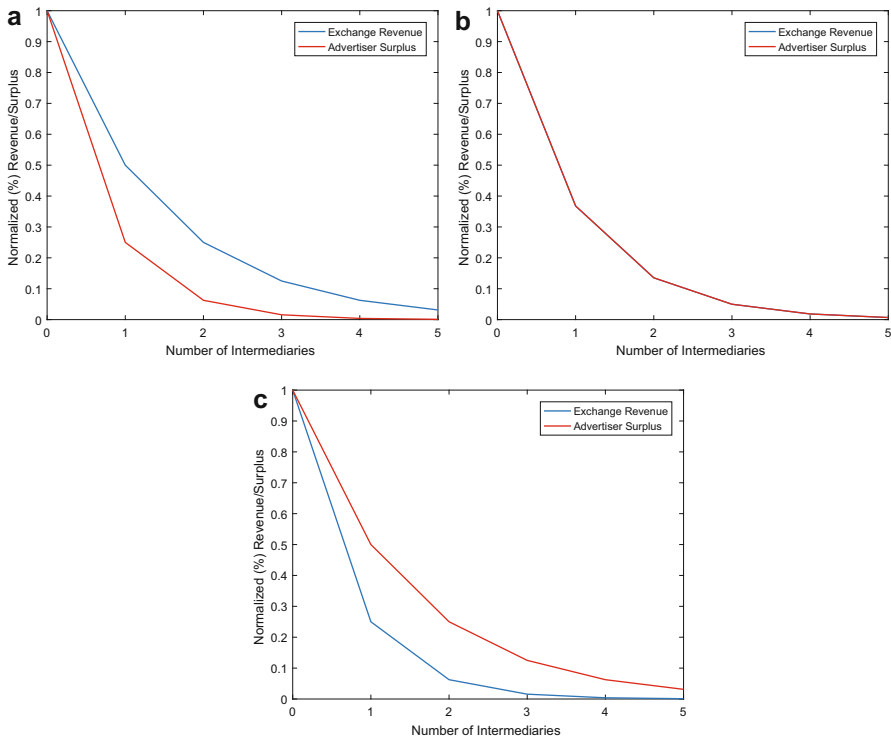


Fig. 21.5 Impact of the total number of intermediaries (m) on the expected surplus of the advertiser ($I_{(0)}$), and the expected revenue of the exchange ($I_{(m+1)}$). All results are relative to the first-best, which corresponds to the case when the advertiser bids directly in the exchange’s mechanism ($m = 0$). In these plots, the value of the exogenous bid d is constant and equal to $\mathbb{E}[v]$ ((a) Uniform distribution ($\mathcal{V} = [0, 1]$, $\mathbb{E}[v] = 1/2$), (b) Exponential distribution ($\mathcal{V} = [0, \infty)$, $\mathbb{E}[v] = 1$), (c) S. Pareto distribution ($\mathcal{V} = [0, \infty)$, $\mathbb{E}[v] = 2/3$))

revenue of the exchange, are affected negatively. This is a direct consequence of the bid shading behavior exhibited by the intermediaries. Figure 21.5 plots the expected surplus of the advertisers and the expected revenue of the exchange relative to that of first-best. This figure suggests that even one intermediary can significantly decrease the exchange’s revenue and the surplus of the advertiser, while multiple intermediaries can quickly decrease these quantities to almost zero. This result is aligned with the double marginalization literature. Our analysis extends these double-marginalization insights by quantifying the impact of the distribution of values and the position in the intermediation chain on an agent’s profit. In particular, when the value distribution is uniform (shifted Pareto) the advertiser (the exchange) suffers more than the exchange (the advertiser) from intermediation. On the other hand, the exchange and the advertiser are equivalently affected when the value distribution is exponential. These observations are aligned with the optimal intermediation position in multi-stage intermediation settings identified in our analysis.

21.4 Concluding Remarks

We study two models which shed light on the problems related to intermediation in online advertising. In the first part, we characterize the optimal mechanism for an intermediary which offers a contract to an advertiser with a private budget and targeting criterion to acquire impression on her behalf. We show that the presence of the intermediary does not harm the advertiser surplus, even further leads to a simpler bidding policy. For this model, Balseiro and Candogan (2017) also show that the profits of the intermediary are maximized for markets where budgets of advertisers are neither exceedingly small nor exceedingly large by using a combination of theoretical results and numerical experiments. In the second part, we provide a game theoretic model to understand the strategic interaction between intermediaries organized in a chain network, and derive economic insights via numerical analyses. As an immediate extension of this model, more general network structures can be considered in addition to the chain network. To this end, Balseiro et al. (2017) study symmetric tree networks, and formally establish that the results shown in the second part of this chapter, such as the impact of the value distribution on the most profitable position in a network, are generalized to those networks. Moreover, they also analyze the incentives of intermediaries to merge horizontally (within the same tier) and vertically (across different tiers). Although the exchange considered herein is assumed to have no reserve price, Balseiro et al. (2017) also model the strategic behavior of the exchange, and show that the intermediation network structure plays a key role in the profit of the exchange.

References

- Balseiro SR, Candogan O (2017) Optimal contracts for intermediaries in online advertising. *Oper Res* 65(4):878–896
- Balseiro SR, Besbes O, Weintraub GY (2015) Repeated auctions with budgets in ad exchanges: approximations and design. *Manag Sci* 61(4):864–884
- Balseiro SR, Candogan O, Gurkan H (2017) Multi-stage intermediation in display advertising. Working paper
- Belavina E, Girotra K (2012) The relational advantages of intermediation. *Manag Sci* 58(9):1614–1631
- Bhattacharya S, Goel G, Gollapudi S, Munagala K (2010) Budget constrained auctions with heterogeneous items. In: Proceedings of the 42nd ACM symposium on theory of computing, STOC'10. ACM, New York, pp 379–388
- Borgs C, Chayes J, Immorlica N, Mahdian M, Saberi A (2005) Multi-unit auctions with budget-constrained bidders. In: Proceedings of the 6th ACM conference on electronic commerce, EC'05. ACM, New York, pp 44–51
- Bresnahan TF, Reiss PC (1985) Dealer and manufacturer margins. *RAND J Econ* 16:253–268
- Brusco S, Lopomo G (2008) Budget constraints and demand reduction in simultaneous ascending-bid auctions. *J Ind Econ* 56(1):113–142
- Brusco S, Lopomo G (2009) Simultaneous ascending auctions with complementarities and known budget constraints. *Econ Theory* 38(1):105–124
- Chawla S, Malec DL, Malekian A (2011) Bayesian mechanism design for budget-constrained agents. In: Proceedings of the 12th ACM conference on electronic commerce, EC'11. ACM, New York, pp 253–262

- Che YK, Gale I (1998) Standard auctions with financially constrained bidders. *Rev Econ Stud* 65(1):1–21
- Che YK, Gale I (2000) The optimal mechanism for selling to a budget-constrained buyer. *J Econ Theory* 92(2):198–233
- Condorelli D, Galeotti A (2016) Strategic models of intermediation networks. In: Bramoulle Y, Galeotti A, Rogers BW (eds) *The Oxford handbook of the economics of networks*. Oxford University Press, New York
- Feldman J, Mirrokni V, Muthukrishnan S, Pai MM (2010) Auctions with intermediaries: extended abstract. In: *Proceedings of the 11th ACM conference on electronic commerce, EC'10*. ACM, New York, pp 23–32
- Ghosh A, Rubinstein BI, Vassilvitskii S, Zinkevich M (2009) Adaptive bidding for display advertising. In: *Proceedings of the 18th international conference on World Wide Web, WWW'09*. ACM, New York, pp 251–260
- Gummadi R, Key PB, Proutiere A (2011) Optimal bidding strategies in dynamic auctions with budget constraints. In: *49th annual Allerton conference on communication, control, and computing (Allerton)*. IEEE, Piscataway, p 588
- Internet Advertising Bureau (2016) Internet advertising revenue report, 2016 full year results. PricewaterhouseCoopers. Technical report. Available at https://www.iab.com/wp-content/uploads/2016/04/IAB_Internet_Advertising_Revenue_Report_FY_2016.pdf. Accessed 6 Mar 2018
- Iyer K, Johari R, Sundararajan M (2014) Mean field equilibria of dynamic auctions with learning. *Manag Sci* 60(12):2949–2970
- Jiang C, Beck CL, Srikant R (2014) Bidding with limited statistical knowledge in online auctions. *SIGMETRICS Perform Eval Rev* 41(4):38–41
- Laffont JJ, Robert J (1996) Optimal auction with financially constrained buyers. *Econ Lett* 52(2):181–186
- Lariviere MA, Porteus EL (2001) Selling to the newsvendor: an analysis of price-only contracts. *Manuf Serv Oper Manag* 3(4):293–305
- Loertscher S, Niedermayer A (2007) When is seller price setting with linear fees optimal for intermediaries? Technical report, Department of Economics, Universität Bern, Discussion papers
- Loertscher S, Niedermayer AF (2012) Fee-setting mechanisms: on optimal pricing by intermediaries and indirect taxation. Available at SSRN 2172386
- Manea M (2018) Intermediation and resale in networks. *J Polit Econ* 126(3):1250–1301
- Mansour Y, Muthukrishnan S, Nisan N (2012) Doubleclick ad exchange auction. <https://arxiv.org/pdf/1204.0535v1.pdf>
- Maskin ES (2000) Auctions, development, and privatization: efficient auctions with liquidity-constrained buyers. *Eur Econ Rev* 44(4–6):667–681
- Myerson RB, Satterthwaite MA (1983) Efficient mechanisms for bilateral trading. *J Econ Theory* 29(2):265–281
- Nguyen T, Subramanian V, Berry R (2016) Delay in trade networks. *Oper Res* 64(3):646–661
- Niazadeh R, Yuan Y, Kleinberg R (2014) Simple and near-optimal mechanisms for market intermediation. In: *Web and Internet economics*. Springer, Cham, pp 386–399
- Pai MM, Vohra R (2014) Optimal auctions with financially constrained bidders. *J Econ Theory* 150(0):383–425
- Perakis G, Roels G (2007) The price of anarchy in supply chains: quantifying the efficiency of price-only contracts. *Manag Sci* 53(8):1249–1268
- Stavrogiannis LC, Gerding EH, Polukarov M (2013) Competing intermediary auctions. In: *Proceedings of the 2013 international conference on autonomous agents and multi-agent systems*, St. Paul, pp 667–674
- Tirole J (1988) *The theory of industrial organization*, 1st edn. The MIT Press, Cambridge/London
- Wu SD (2004) Supply chain intermediation: a bargaining theoretic framework. In: Simchi-Levi D, Wu SD, Shen Z-J (eds) *Handbook of quantitative supply chain analysis*. Springer, New York, pp 67–115