Suresh Chandra Satapathy
Vikrant Bhateja
Swagatam Das   *Editors*

# Smart Computing and Informatics

Proceedings of the First International Conference on SCI 2016, Volume 1

KES
International

Springer

# Smart Innovation, Systems and Technologies

Volume 77

**Series editors**

Robert James Howlett, Bournemouth University and KES International,
Shoreham-by-sea, UK
e-mail: rjhowlett@kesinternational.org

Lakhmi C. Jain, University of Canberra, Canberra, Australia;
Bournemouth University, UK;
KES International, UK
e-mails: jainlc2002@yahoo.co.uk; Lakhmi.Jain@canberra.edu.au

The Smart Innovation, Systems and Technologies book series encompasses the topics of knowledge, intelligence, innovation and sustainability. The aim of the series is to make available a platform for the publication of books on all aspects of single and multi-disciplinary research on these themes in order to make the latest results available in a readily-accessible form. Volumes on interdisciplinary research combining two or more of these areas is particularly sought.

The series covers systems and paradigms that employ knowledge and intelligence in a broad sense. Its scope is systems having embedded knowledge and intelligence, which may be applied to the solution of world problems in industry, the environment and the community. It also focusses on the knowledge-transfer methodologies and innovation strategies employed to make this happen effectively. The combination of intelligent systems tools and a broad range of applications introduces a need for a synergy of disciplines from science, technology, business and the humanities. The series will include conference proceedings, edited collections, monographs, handbooks, reference books, and other relevant types of book in areas of science and technology where smart systems and technologies can offer innovative solutions.

High quality content is an essential feature for all book proposals accepted for the series. It is expected that editors of all accepted volumes will ensure that contributions are subjected to an appropriate level of reviewing process and adhere to KES quality principles.

Suresh Chandra Satapathy
Vikrant Bhateja · Swagatam Das
Editors

# Smart Computing and Informatics

Proceedings of the First International
Conference on SCI 2016, Volume 1

*Editors*
Suresh Chandra Satapathy
Department of Computer Science and
    Engineering
PVP Siddhartha Institute of Technology
Vijayawada, Andhra Pradesh
India

Swagatam Das
Electronics and Communication Sciences
    Unit
Indian Statistical Institute
Kolkata, West Bengal
India

Vikrant Bhateja
Department of Electronics and
    Communication Engineering
Shri Ramswaroop Memorial Group of
    Professional Colleges
Lucknow, Uttar Pradesh
India

# Preface

The 1st International Conference on Smart Computing and Informatics (SCI) was organized successfully with the excellent support of Department of CSE, ANITS, Visakhapatnam on 3–4 March 2017. The aim of this International Conference was to present a unified platform for advanced and multi-disciplinary research towards design of smart computing and information systems. The theme was on a broader front focused on various innovation paradigms in system knowledge, intelligence and sustainability that is applied to provide realistic solution to varied problems in society, environment and industries. The scope was also extended towards deployment of emerging computational and knowledge transfer approaches, optimizing solutions in varied disciplines of science, technology and health care. The conference received huge quality submissions in direct track and special session tracks. After stringent quality check and review process, only good papers were accepted with an acceptance ratio of 0.38. Several eminent researchers and academicians delivered talks addressing the participants in their respective field of proficiency. Prof. Ganapati Panda, IIT Bhubaneswar; Dr. R. Logeswaran, Malaysia; Dr. C. Krishna Mohan, IIT Hyderabad; Dr. P.S. Grover, KIIT, Group of Colleges, Gurgaon; Dr. A.K. Nayak, Hon. Secretary, Computer Society of India and Director, Indian Institute of Business Management, Patna; Dr. Arunkumar Thangavelu, VIT Vellore; Dr. Ramchandra V. Pujeri, Director, MIT College of Engineering Pune; Dr. Nilanjan Dey, TICT Kolkota; and Dr. Prashant Kumar Pattnaik, KIIT Bhubaneswar were the eminent speakers and guests on the occasion.

We would like to express our appreciation to the members of the Programme Committee for their support and cooperation in this publication. We are also thankful to Team from Springer for providing a meticulous service for the timely production of this volume. Our heartfelt thanks to Chairman, ANITS for the support provided. Special thanks to all guests who have honoured us with their presence in the inaugural day of the conference. Our thanks are due to all special session chairs, track managers and reviewers for their excellent support. Profound thanks to Organizing Chair Prof. Pritee Parweker, ANITS, Visakhapatnam, for marvellous support. Sincere thanks to Honorary Chair, Dr. Lakhmi Jain, Australia, for his valuable inputs and support during the conference. Last, but certainly not the least,

our special thanks go to all the authors who submitted papers and all the attendees for their contributions and fruitful discussions that made this conference a great success.

Vijayawada, India                                                    Suresh Chandra Satapathy
Lucknow, India                                                              Vikrant Bhateja
Kolkata, India                                                                Swagatam Das
March 2017

# Organizing Committee

## Special Session Chairs

Dr. M. Bhanu Sridhar, GVP College of Engineering for Women, Visakhapatnam, AP, India
Dr. D.N.D. Harini, GVP College of Engineering, Visakhapatnam, AP, India
Dr. Tusar Kanti Mishra, ANITS, Visakhapatnam, AP, India
Prof. (Dr.) R. Sireesha, Professor, GITAM University, Visakhapatnam, AP, India
Prof. (Dr.) R. Sivaranjani, ANITS, Visakhapatnam, AP, India
Dr. Hari Mohan Pandey, Amity University, Delhi, India
Ankit Chaudhary, Truman State University, USA
Yudong Zhang, Nanjing Normal University, China, Research Scientist, MRI Unit, Columbia University, USA
Tanupriya Choudhury, Amity University, Uttar Pradesh, India
Praveen Kumar, Amity University, Uttar Pradesh, India
Dr. Sai Sabitha, Amity University, Uttar Pradesh, India
Dr. Suma V., Dean, Research and Industry Incubation Centre, Dayananda Sagar College of Engineering, Bangalore, India

## International Advisory Committee/Programme Committee

S.K. Udgata, UoH, Hyderabad, India
C.A. Murthy, ISI Calcutta, Kolkata, India
M.K. Tiwari, IIT Kharagpur, India
C. Chandra Sekhar, IIT Madras, Chennai, India
Suresh Sundaram, NTU, Singapore
Lipo Wang, NTU, Singapore
Amit Mitra, IIT Kanpur, India
Aruna Tiwari, IIT Indore, India

D. Nagesh Kumar, IISc, Bangalore, India
V. Sushila Devi, IISc, Bangalore, India
C. Hota, BITS Pilani, Hyderabad, India
Chilukuri Mohan, Syracuse University, Syracuse, USA
Debjani Chakraborty, IIT Kharagpur, India
P.K. Kalra, IIT Kanpur, India
Vasant Pandian, University Putra Malaysia, Malaysia
Oscar Castillo, Tijuana Institute of Technology, Chula Vista, CA, USA
Indranil Bose, IIM Calcutta, Kolkata, India
S. Bapi Raju, IIIT Hyderabad, India
Brijesh Verma, CQ University, Brisbane, Australia
C.R. Rao, UOHYD, Hyderabad, India
B.L. Deekshatulu, IDRBT, Hyderabad, India
Arun Agarwal, UOHYD, Hyderabad, India
Arnab Laha, IIM Ahmedabad, India
Biplav Srivastava, IBM Research, New Delhi, India
B.K. Mohanty, IIM Lucknow, India
M. Janga Reddy, IIT Bombay, Mumbai, India
M.C. Deo, IIT Bombay, Mumbai, India
Pankaj Dutta, IIT Bombay, Mumbai, India
Usha Anantha Kumar, IIT Bombay, Mumbai, India
Faiz Hamid, IIT Kanpur, India
S. Chakraverty, NIT Rourkela, Rourkela
H. Fujita, Iwate Prefectural University, Iwate, Japan
Dries Benoit, Ghent University, Ghent, Belgium
S.A. Arul, Philips Electronics Singapore, Singapore
Pawan Lingars, Saint Mary's University, Halifax, Canada
Amuelson Hong, Oriental Institute of Technology, Taiwan
Zhihua Cui, Taiyuan University of Science and Technology, Taiyuan, China
Balasubramaniam Jayaram, IIT Hyderabad, India
K. Saman Halgamuge, The University of Melbourne, Melbourne, Australia
Nischal Verma, IIT Kanpur, India
Laxmidhar Behera, IIT Kanpur, India
Prof. YaoChuJin, University of Surrey, Guildford, England
Vineeth Balasubramian, IIT Hyderabad, India
Atul Negi, Professor, University of Hyderabad, India
M. Naresh Kumar, NRSC, Hyderabad, India
Maurice Clerc, Franch Roderich Gross, England
Dr. Syed Basha, India
Kalyanmoy Deb, IIT Kanpur, India
Saman Halgamuge, Australia
Jeng-Shyang Pan, Talwan Peng Shi, UK Javier Del Ser, Spain
Leandro Dos Santos Coelho, Brazil
S Pattanaik, India
Gerardo Beni, USA

K. Parsopoulos, Greece
Lingfeng Wang, China
Athanasios V. Vasilakos, Sweden
Athens Pei-Chann Chang, Taiwan
Chilukuri K. Mohan, USA
SaeidNahavandi, Australia
Abbas Khosravi, Australia
Almoataz Youssef Abdelaziz, Egypt
K.T. Chaturvedi, India
M.K. Tiwari, India
Yuhui Shi, China
Dipankar Dasgupta, USA
Lakhmi Jain, Australia
X.Z. Gao, Finland
Juan Luis Fernandez Martinez, Spain
Oscar Castillo, Mexico
Heitor Silverio Lopes, Brazil
S.K. Udgata, India
Namrata Khemka, USA
G.K. Venayagamoorty, USA
Zong Woo Geem, USA
Ying Tan, China
S.G. Ponnambalam, Malaysia
Halina Kwasnicka, Poland
M.A. Abido, Saudi Arabia
Richa Singh, India
Manjaree Pandit, India
Hai Bin Duan, China
Delin Luo, China
V. Ravi, India
S. Basker, India
M. Rammohan, South Korea
Munesh Chandra Trivedi, ABES Engineering College, Ghaziabad, India
Alok Aggarwal, Professor & Director, JP Institute of Engineering and Technology, Meerut, India
Dilip Kumar Sharma, Institute of Engineering and Technology, GLA University, Mathura, India
K. Srujan Raju, CMR Technical Campus, Hyderabad, India
B.N. Biswal, BEC, Bhubaneswar, India
Sanjay Sengupta, CSIR, New Delhi, India
NaeemHanoon, Malaysia
Cirag Arora, India
Steven Fernades, India
Kailash C. Patidar, South Africa
K. Srujan Raju, CMR Group, Hyderabad

# Contents

# About the Editors

**Suresh Chandra Satapathy** is currently working as Professor and Head, Department of Computer Science and Engineering at PVP Siddhartha Institute of Technology, Andhra Pradesh, India. He obtained his Ph.D. in Computer Science and Engineering from JNTU Hyderabad and M.Tech. in CSE from NIT, Rourkela, Odisha, India. He has 26 years of teaching experience. His research interests are data mining, machine intelligence and swarm intelligence. He has acted as programme chair of many international conferences and edited six volumes of proceedings from Springer LNCS and AISC series. He is currently guiding eight scholars for Ph.D. Dr. Satapathy is also a Senior Member of IEEE.

**Vikrant Bhateja** is a Professor, Department of Electronics & Communication Engineering, Shri Ramswaroop Memorial Group of Professional Colleges (SRMGPC), Lucknow, and also the Head (Academics & Quality Control) in the same college. His area of research includes digital image and video processing, computer vision, medical imaging, machine learning, pattern analysis and recognition, neural networks, soft computing and bio-inspired computing techniques. He has more than 90 quality publications in various international journals and conference proceedings. Prof. Bhateja has been on TPC and chaired various sessions from the above domain in international conferences of IEEE and Springer. He has been the track chair and served in the core-technical/editorial teams for international conferences: FICTA 2014, CSI 2014 and INDIA 2015 under Springer-ASIC Series and INDIACom-2015, ICACCI-2015 under IEEE. He is associate editor in International Journal of Convergence Computing (IJConvC) and also serving in the editorial board of International Journal of Image Mining (IJIM) under Inderscience Publishers. At present, he is guest editors for two special issues floated in International Journal of Rough Sets and Data Analysis (IJRSDA) and International Journal of System Dynamics Applications (IJSDA) under IGI Global publications.

**Swagatam Das** received the B.E. Tel. E., M.E. Tel. E. (Control Engineering specialization) and Ph.D. degrees, all from Jadavpur University, India, in 2003,

2005 and 2009, respectively. Currently, he is serving as an Assistant Professor at the Electronics and Communication Sciences Unit of Indian Statistical Institute, Kolkata. His research interests include evolutionary computing, pattern recognition, multi-agent systems and wireless communication. Dr. Das has published one research monograph, one edited volume and more than 150 research articles in peer-reviewed journals and international conferences. He is the founding co-editor-in-chief of "Swarm and Evolutionary Computation", an international journal from Elsevier. He serves as associate editors of the IEEE Trans. on Systems, Man, and Cybernetics: Systems and Information Sciences (Elsevier). He is an editorial board member of Progress in Artificial Intelligence (Springer), Mathematical Problems in Engineering, International Journal of Artificial Intelligence and Soft Computing and International Journal of Adaptive and Autonomous Communication Systems. He is the recipient of the 2012 Young Engineer Award from the Indian National Academy of Engineering (INAE).

# Analysis of Existing Text Hiding Algorithms for Image Steganography Using TLNUS and AES

**Jagadish Gurrala and P. Sanyasi Naidu**

**Abstract** In the last few years, many researchers putting many efforts for getting good data hiding algorithms which was complex in design and undergo rigorous investigation on starting from secret text sizes ranges from 4 kB to 1 MB file has been embedded in it so far for the sake of more secure communication among mobile nodes and local area networks. As of now, several steganographic concepts were conceived on data hiding approaches deployed in insecure channel. In this paper, authors designed new data hiding algorithm approach proposed based on three protection layers has been used to maintain secrecy of the embedded message in a true color image. Here, the data is embedded randomly instead of sequentially by an image segmentation algorithm that uses two-level non-uniform segmentation. Advanced encryption standard algorithm has been used to encrypt the secret text. Different performance measures from the experimental results have shown the reasonable prototype of the proposed steganography algorithm. The result after comparing the proposed algorithm and the wide spectrum of steganographic schemes confirm that the stego image with medium perception ratio has been reached even if the stego image holds a large amount of data with good visual quality and working under jpeg and gray-scale images and also resistant to statistical and visual problems.

**Keywords** Steganography · AES · Data hiding process

J. Gurrala (✉)
Department of CSE, Anil Neerukonda Institute of Technology and Sciences,
Visakhapatnam, India
e-mail: gjagadish.cse@anits.edu.in

P. Sanyasi Naidu
Department of CSE, GITAM Institute of Technology GITAM University,
Visakhapatnam, India

# 1 Introduction

Steganography is the branch of data hiding process inside carrier image that enables hiding encrypted secret text in cover image to produce stego image. The secret text cannot be observed in the stego image over computer networks. Because, the entropy difference between cover image and stegao image is equal to zero. Then only intruder cannot anticipate the secret details in and around stego image. The profusion of our digital media (image, video, and audio) in our modern life has led to a rapid technological development of steganography with digital media files being the carrier contents (camouflage). In the few years, enormous papers are published on steganalysis of new data hiding algorithms, and their contribution is for embedding secret data in grayscale [1].

## 1.1 Importance of Steganography in Transmission

Steganography is derived from the Greek origin and means "to hide in plain sight." Steganography is an art and science of hiding the existence of information in an image file. The main goal of steganography is to hide information well enough such that the unintended recipients do not suspect the steganographic medium of containing hidden data. Simple steganographic techniques have been in use for hundreds of years. Most Steganography jobs have been carried out on different cover images like text, image, audio, or video.

Steganography has a vital role to put various secret messages (M) inside cover image (CI) using function $F_K$ to calculate the places where secret message kept inside CI to transmit stego image across networks. In the best case, no one can see that the secret message by both the parties are able to communicate message to stay hidden until it received by destination. Image Steganography is a process that involves hiding a message in an appropriate carrier file such as an image file. Unlike cryptography, which simply conceals the content or meaning of a message, steganography conceals the very existence of a message. The carrier can then be sent to a receiver without anyone else knowing that it contains a hidden message. Hence image steganography provides better security than cryptography (Fig. 1).

## 1.2 Importance of Cryptography over Steganographic Based Communication

Cryptography is having an important role in improving security standard across networks to make them non-readable form of given cover image against third-party vendors. It takes place between the sender and receiver for hiding text. The sender sends the message to the receiver through the communication channel. Encryption

Function ($F_k$=CI **XOR** M)

Cover Image (CI)

Secret text(M)

**Embedding Process**

Stego- image transit

**Extraction Process**

Stego key

Finding Stego image

**Fig. 1** Model for the data hiding algorithm on cover image for image steganography

is a process where the ordinary information also called as plain text is coded into some unrecognizable form, which is usually called cipher text, which is in the unrecognizable form of ordinary information that is plain text. The entire process is done for the sole purpose of protecting the message from being used or manipulated by the intruder who is the third person.

Crytography only allows the authenticated users to prevent secret data from unauthorized access by supplying a private/public key to the end user to read the proper information. Cryptography does enciphering/deciphering mechanism that protects our valuable information such as your documents, pictures, or online transactions from unwanted people accessing or changing it. Encryption works by using a mathematical formula called a cipher and a key to convert readable data (plain text) into a form that others cannot understand (cipher text). The cipher is the general recipe for encryption and the key makes the encrypted data images.

Cryptography alone could not solve the problem about access control to end user. Thus, authors are thinking to merge the different concepts of steganography and cryptography methods together to form and claiming that the secret data which is received to the end user is solely. An authorized user can only decrypt data because decryption requires a secret key or password. So cryptography alone cannot supplement security aspects. However, steganography alone cannot solve the issue regarding security. Therefore, we are thinking to solve the problem raised to correlate these ideas of steganography and cryptography to produce better security results what the secure image transmission demands.

In this paper, the remaining related work was discussed in second section, analysis of apply TLNUS and AES algorithms on steganography and their mathematical proofs were discussed in third section, the results and comparisons of given inputs were discussed in fourth section, and summary and conclusion and future direction in fifth section.

## 2 Analysis on Existing Algorithms

The two-level non uniform segmentation algorithm is to hide secret information in an image to transfer from source mobile to destination mobile and to carry out the task more effectively against statistical attacks while producing a high-quality image. The aim of present paper is to delivery the cover image in terms of stego image i.e., hide the data over an image using steganographic algorithms and ensure the quality of concealing data to the destination. However, the paper could try to apply method for embedding and encrypting the text in an image for normal transmission of stego image. The current process [2] provided successful delivery of the stego image to allow the services only to the authorized destination users without any copy right violation. The proposed method will help to secure the content with in the image and to make the personal information much secure because even though if the unauthorized person succeeds in being able to intercept the stego image, the intruder will not able to read the message as well as acquire the information during transit.

Whenever the data is encrypted using steganographic algorithms with in image, neither data nor the image is embedded in it should lose its originality. The main aim is to embed sufficient data in a gray image to make it invisible to end user.

### 2.1 Refined Model of Data Hiding Algorithm from Existing Approach

In this paper, the authors have given refined model of the algorithm process in paper [3, 4],which have given detailed information about 5 layers. From this five-layer model, authors reviewed the algorithm to shorten it into simple 3 layers (Fig. 2).

## 3 Steganography Implementation

In the steganography process discussed so far is in theoretical aspect. Hence, now it has been implemented in step process, an implementation through C++ environment.

### 3.1 TLNUS Algorithm

In this algorithm, image segmentation process is defined in this work; it is based on different sizes of segments of a CI.

The following steps are applied to display the proposed image segmentation:

**Fig. 2** Stages involved in the reviewed model

Step 1: Let $S$ be the key size cipher key of AES (cipher key generated from random generator) in Eq. 1.

$$S = |ck| \tag{1}$$

Step 2: In the Eq. 2, first level of segmentation process is used that horizontal segment length (SHS) and vertical segment length(SVS) is formulated;

$$\text{SHS} = W_{\text{CI}}/S \quad \text{and SVS} = H_{\text{CI}}/S \tag{2}$$

//where cover image breath is represented as $W_{\text{CI}}$ and cover image height is represented as $H_{\text{CI}}$, respectively.

Step 3: For the second level of segmentation defined in Eq. 3.
Given the length of the different lengths of segments for both the vertical and horizontal directions (Ver, Hor) according to the values of the switching index $(k)$;

$$K = ((i+j)\bmod S) + 1 \tag{3}$$

where $i,j$ are two points in the different lengths of next-level segmentation of whole image.

Step 4: If partition the whole block into small block contains different sizes of chunks when
$k$ is odd parity, then compute Ver, Hor in Eq. 4;

$$\text{Ver}_{h,v}^{i,j} = \left[\frac{\text{val}(ck_i) * \text{SVS}}{\sum_{m=1}^{s} \text{val}(ck_m)}\right] \tag{4}$$

$$\forall i = 1, \ldots, S; \ \forall j = 1, \ldots, s; \qquad \forall h = 1, \ldots, s; \quad \forall v = 1, \ldots, s$$

Else
Compute $\text{Ver}_{h,v}^{i,j}$, $\text{Hor}_{h,v}^{i,j}$ in Eq. 5;

$$\text{Ver}_{h,v}^{i,j} = \left[\frac{\text{val}(ck_k) * \text{SVS}}{\sum_{m=1}^{s} \text{val}(ck_m)}\right] ; \quad \text{Hor}_{h,v}^{i,j} = \left[\frac{\text{val}(ck_j) * \text{SHS}}{\sum_{m=1}^{s} \text{val}(ck_m)}\right] \tag{5}$$

$$\forall i = 1, \ldots, S; \ \forall j = 1, \ldots, s; \qquad \forall h = 1, \ldots, s; \quad \forall v = 1, \ldots, s$$

When block segmentation is completed.
//where $\text{Val}(ck_k)$ represents the ASCII value of the every $s$th character in the cipher key such that $s\{k, j, m\}$. The $(h)$ and $(v)$ are the height and the breath of the cover image, respectively.

Step 5: Start reading of each pixel [5] on each segment at the first and second levels of segmentation in Eq. 5;

Step 6: End of the main block.

## 3.2 Encryption and Decryption in AES

AES Encryption Process:

Here, we were using AES algorithm for better encryption algorithm in mobile platforms.

The AES cipher is almost identical to symmetric enciphering techniques where low-consuming mobile devices works. It does encryption on multiple bits of input plain text with 128 bits. It contains 10 rounds to mangled the input text into cipher text with random key size (Fig. 3).

### 3.3 Image Segmentation Algorithm Using TLNUS

In this paper, authors initiated selection of locations [2] over patches promptly for embedding secret messages randomly to achieve a promising solution in the proposed steganography algorithm; this approach is based on two-level non-uniform segmentation (TLNUS) [3] (Fig. 4).

## 4 Simulation and Testing

For cover image BMP file has been taken for data hiding process. This paper has chosen as 115 × 145 dimension matrix pixel image in Planet Gray, scale image for hiding simple text "hai" in cover image where we apply TLNUS and AES and apply these algorithms in Dev C++ for getting better results in this paper.

**Test 1**: **Gray scale image before segmentation:**

This is the input sky image contains all values ranges from 0 to 255, most of the pixels are black in nature given to AES process when cipher key given (Fig. 5).

**Test 2**: **Gray scale image after segmentation:**

This is the output sky image contains all values ranges from 0 to 7 only (Fig. 6).



**Fig. 3** Encryption process

**Fig. 4** Image segmentation
process using two-level
adaptive non-uniform image
segmentation



**Fig. 5** BMP raw image
before TLNUS



**Fig. 6** Screenshot image
after segmentation



## 4.1 Results for TLNUS

After simulation, the images get matrix representation, which contains most of the
values are occupied with 0's and 255 only (Fig. 7).

## 4.2 Output for Plain Text to Binary File

After encryption, the key size shared the cipher key with secret information in
binary form (Fig. 8).

**Fig. 7** Screen shot of pixel information from gray scale using TLNUS



**Fig. 8** Screen shot pixel information from BMP

## 4.3  AES Keys for Encryption

After stego image, all secret information is about to insert into plain text for transmission (Fig. 9).

## 4.4  Stego Image Matrix

After embedding secret text binary information into Plain text using process of TLNUS (Fig. 10).

**Fig. 9** Screen shot pixel information after adding AES keys



**Fig. 10** Screen shot pixel information of Stego image

## 5 Conclusion

This paper contributed on the way of hiding the secret data in the image using steganography and encrypting the steganographic part of the image through analysis. After surveying the concept, it involves in sending the image along with the data embedded in it, authors having ability to transfer it in reduced bandwidth, since there is no need to use other bandwidth to send the secret message we want to share. This paper has simulated exiting algorithm, which divide the image into sub blocks by using TLNUS; instead of selecting the whole image, a part of the image has been selected and the secret message is in embedding process. Encryption process using

AES is simple and fast in hiding process of secret data. As of now, image steganography still requires more inputs for organizing secret text in image to be done effectively provided user authentication in given location-based authentication. Hence, steganography achieve the secure transmission of image through obscurity and unique location parameters along with the help of essence of AES algorithm borrowed from modern symmetric block enciphering scheme. Therefore, authors suggested the importance of algorithm that chosen TLNUS and AES together to solve the fundamental image transfer through steganography methods in terms of embedding and extracting process lies secret data.

# References

1. Johnson, N.F., Jajodia, S.: Exploring steganography: seeing the unseen. IEEE Comput. **31**(2), 26–34 (1998)
2. El-Emam, N.: Hiding a large amount of data with high security using steganography algorithm. J. Comput. Sci. **3**(4), 32–223 (2007). doi:10.3844/jcssp.2007.223.232
3. El-Emam, N.: Embedded a large amount of information using high secure neural based steganography algorithm. Int. J. Inf. Commun. Eng. **4**(2), 95–106 (2008)
4. Sanyasi, P., Naidu, J.G.: Investigation and Analysis of Location Based Authentication and Security Services of Wireless LAN's and Mobile Devices. published in IJCA, **146**(8) (ISBN: 973-93-80893-84-7, July (2016)
5. El-Emam, N.: New Data-Hiding Algorithm Based on Adaptive Neural Networks with Modified Particle Swarm Optimization. Elsevier (2013)

# Automated System for Detection of White Blood Cells in Human Blood Sample

**Siddhartha Banerjee, Bibek Ranjan Ghosh, Surajit Giri and Dipayan Ghosh**

**Abstract** Determination of the WBC count of the body necessitates the detection of white blood cells (leukocytes). During an annual physical checkup, generally doctors prescribe for a complete blood count report. WBC count is required to determine the existence of disease for symptom like body aches, chills, fever, headaches, and many more. The existence of autoimmune diseases, immune deficiencies, blood disorders, and hidden infections within human body can also be alerted by the report of WBC count. The usefulness of chemotherapy or radiation treatment, especially for cancer patients, is also monitored by this report. This paper introduces an automated system to detect the white blood cell from the microscopic image of human blood sample using several image processing techniques.

**Keywords** WBC · RBC · Thresholding · Region labeling · Erosion · Dilation

## 1 Introduction

An important component of the immune system is leukocytes. Attacking of the body by viruses, bacteria, germs, and many others is controlled by these leukocytes. After the production of white blood cells in the bone marrow, it circulates through the bloodstream. The test of WBC count reveals the number of white blood cells in human body. A complete blood count (CBC) includes this test. A percentage of

S. Banerjee (✉) · B. R. Ghosh · S. Giri · D. Ghosh
Department of Computer Science, Ramakrishna Mission Residential College (Autonomous),
Narendrapur, Kolkata 700103, India
e-mail: sidd_01_02@yahoo.com

B. R. Ghosh
e-mail: bibekghosh2003@yahoo.co.in

S. Giri
e-mail: girisurajit@gmail.com

D. Ghosh
e-mail: d6ghosh@gmail.com

**Fig. 1 a** Blood cell **b** stem cell type

each type of white blood cell is present in human blood. However, white blood cell count can increase or decrease from the healthy range.

Stem cells get mature and create some kind of new blood cells. Each and every blood type has its own function. Blood components are shown in Fig. 1(a). It consists four parts **red blood cells** (erythrocytes), **white blood cells** (leukocytes), **platelets,** and **plasma.**

After becoming old or damage, the cells die and are replaced by new cells. After getting matured, the stem cells change to several components of blood, as shown in Fig. 1b. They grew up either as myeloid stem cell or as lymphoid stem cell. After getting mature, the myeloid stem cells become myeloid blast. Platelet, red blood cell, and several types of white blood cell are formed during this blast. A mature lymphoid stem cell can form lymphoid blast, and this blast creates a type of white blood cells. The characteristics of white blood cells formed from these two blasts are different. There are five major types of white blood cells, namely neutrophil, lymphocyte, eosinophil, monocyte, and basophil. The study will focus on detection of white blood cells in a given highly magnified microscopic blood smear images.

The rest of this paper is organized as follows. After this small introduction, some related works are presented in Sect. 2. Section 3 describes the proposed methodology. Experimental results of the proposed method are discussed in Sect. 4. Finally, Sect. 5 draws the conclusions and future work.

## 2   Related Work

Sonali C. Sonar et al. [1] have detected WBCs. Original image was converted to grayscale image. Contrast enhancement techniques such as linear contrast stretching ($L$) and histogram equalization ($H$) are applied. Three images R1, R2, and R3 are obtained, such that $R1 = L + H, R2 = L − H,$ and $R3 = R1 + R2$. A 3 × 3 minimum filter is implemented three times on the image R3. Global threshold value is determined by Otsu's method. Morphological opening with disk structuring element is used to remove small pixel groups. The radius of disk is considered to be

nine pixels. Finally, the neighboring pixels are connected. The objects that are less than half of average RBC area are finally eliminated.

The white blood cells are identified using knowledge base learning by Rajwinder Kaur et al. [2]. In the first approach, the input image was followed by Hough transform, snake body detection algorithm was applied, and the cells were counted. Then, in the second approach, the image went through k-means clustering method followed by histogram equalization, the blood cells were extracted in image segmentation, and the cells were counted.

Nurhanis Izzati et al. [3] have segmented white blood cell nucleus using active contour. First, image segmentation based on the partial differential equation is mainly carried out by active contour model or snakes algorithm, and the parameter of WBC is calculated. The segmented images are converted into binary image. In order to calculate circularity of the object roundness to classify the shape of the WBC nucleus, ratio of the area of an object to the area of the circle is calculated.

Miss. Madhuri G. Bhamare et al. [4] have converted RGB image to gray scale by eliminating the hue and saturation information while retaining the luminance and enhanced the image by using median filter. Image segmentations such as Otsu adaptive thresholding method, watershed transform method, as well as segmentation by K-means clustering followed by EM-algorithm are compared and have followed with Hough transform. The two fundamental morphological operations such as erosion and dilation are used, thereby filling holes and noise spikes and ragged edges are eliminated. Shape, color, and texture features are analyzed by local binary pattern and sequential forward selection algorithm or by artificial neural network and support vector machine.

S. Pavithra et al. [5] applied different image processing techniques to extract white blood cells. After converting a grayscale image, edge detection using Sobel operator is carried out. Median filter is used for image smoothing. After that, unsharp masking and gradient magnitude watershed techniques are applied. Finally, morphological operation and circular Hough transform are used for final count.

Some authors also used different morphological operations for identification and classification of WBC from microscopic image [6–9].

## 3    Proposed Method

The following steps are applied to detect the white blood cell from the microscopic image of the human blood cell.

### 3.1    *Image Acquisition*

The first step in the process requires an image sensor to digitize the signal for acquiring a digital image, or images can be obtained in RGB color format from

**Fig. 2**  **a** Microscopic image of human blood cell **b** histogram of the image

online medical library or hospital blood samples images and are converted to grayscale level. Figure 2a shows a representative microscopic image of human blood cell with four WBCs and many RBCs.

## 3.2  Segmentation

After image acquisition, partition is made from an input image into its constituent objects or parts using following steps.

**Nucleus Identification:** It is clear from the microscopic image of blood cell shown in Fig. 2a that the nucleus in WBC has the lower intensity compared to any other part. Figure 2b shows the histogram of the image. The histogram consists of three peaks. Two peaks at the two ends describe the low-intensity nucleus and high-intensity background. The peak in the middle shows the blood cells excluding the nucleus. From these observations, simple thresholding technique is applied to detect the nucleus. The intensity of the leftmost peak is identified and used as threshold value. The image after nucleus detection is shown in Fig. 3a.

**Blood Cells Identification:** From the histogram, it is quite obvious that the rightmost peak represents the white background of the microscopic image. To obtain the blood components, i.e., RBC, WBC, and platelets, leaving the background, the thresholding technique is also applied. To choose the threshold value, the valley between the rightmost two peaks is considered. The threshold value is chosen at the intensity where minimum frequency occurs in the valley. Figure 3b shows the identified blood cells.

**WBC Detection:** WBC of the blood contains a nucleus in its center. Thus, after detecting the blood cells, the WBC is identified by inspecting both the identified blood cells and the identified nucleus. First, each identified blood cells are assigned unique label using region labeling [10] algorithm. The labeled components that contain nucleus are considered as WBC. The detected WBCs are shown in Fig. 4a. But it is clear from the figure that the identified WBC may be attached with some RBC(s).

Fig. 3 Image of Fig. 2a after **a** nucleus identification **b** blood cell identification **c** WBC detection

## 3.3 Refinement of WBC

As shown in Fig. 3c, the detected WBCs may be connected with neighboring RBC(s). Thus, some refinement is needed to extract the WBCs by separating them from connected RBC(s). To achieve this, following steps are applied.

**Elimination of White Patches:** Detected WBCs contain small white patches inside them. These small patches should be removed before separating the RBCs from WBC by using erosion [11] operation, because these white patches will be increased after erosion operation. To do this, each WBC is considered separately. The size of each white patch is determined. If the size of the white patch is less than '*S,*' the patches are converted to black. '*S*' is experimentally chosen as 40. Figure 4a shows the WBCs after removing the white patches.

**Separation of the WBCs:** A layer of pixels from both the inner and outer boundaries of regions is stripped out by the morphological erosion operation. The holes and gaps between different regions become larger, and small details are eliminated. In this experiment, a disk-shaped structuring element is used with radius 27 to separate the RBCs from the WBCs. Figure 4b shows the result.

**Fig. 4** **a** Elimination of white patches **b** separation of WBCs **c** formation of WBCs

**Formation of WBCs:** From Fig. 4b, it is clear that the WBCs are now separated from the RBCs. But the sizes of the WBCs are reduced. It is also the fact that some small components of RBCs may be present after dilation. These small black regions are eliminated by criteria on size. The average size '*A*' of all regions is calculated. The region with size less than '*A*' is eliminated from the image.

To increase the size of WBCs, the dilation [11] operation is applied. Dilation has the opposite effect to erosion—it adds a layer of pixels to both the inner and outer boundaries of regions. Same disk-shaped structuring element with radius 27 is applied to increase the size. The result is shown in Fig. 4c.

## 4   Result

The proposed method is applied on 200 blood samples. The results generated by this method are compared with human visualization. It is found that the results match with the manual observation in high percentage (97.57%). Figure 5 shows

six blood samples as representatives. Figure 6 shows the detected WBCs for these six samples. In all cases, it detects the WBCs successfully. To test the strength of this method, some blood samples without WBC are given. The method satisfactorily produces output without detecting any WBC. One such sample is given in Fig. 5f, and the corresponding output is given in Fig. 6f.



**Fig. 5** **a–f** Six blood samples



**Fig. 6** **a–f** Detected WBCs of Fig. 5a–f

## 5    Conclusion

This technique may facilitate work flow in biomedical science by replacing tedious and monotonous work with automation. This can be a developmental step to create an automated solution in a larger scale to detect a mere human killing disease as fast as possible. The methodology achieves an automated system for the detection of WBCs and can be used for better and accurate classification of WBCs. The time associated with the pathologist's views can be decreased. The speed and accuracy can be increased by applying this automated analysis prior to the pathologist spending any time on it.

## References

 1. Mohamed, M., Far, B., Guaily, A.: An efficient technique for white blood cells nuclei automatic segmentation. In: IEEE International Conference on Systems, Man, and Cybernetics (SMC), pp. 220–225 (2012)
 2. Kaur, R., Kaur, H.: Comparative analysis of white blood cell by different segmentation methods using knowledge based learning. Int. J. Adv. Res. Electr. Electron. Instrum. Eng. **3** (9), 11720–11728 (2014)
 3. Marzuki, N.I.C., Mahmood, N.H.H., Razak, M.A.A.: Segmentation of white blood cell nucleus using active contour. J. Teknologi **74**(6), 115–118 (2015)
 4. Bhamare, M.G., Patil, D.S.: Automatic blood cell analysis by using digital image processing: a preliminary study. Int. J. Eng. Res. Technol. **2**(9), 3137–3141 (2013)
 5. Pavithra, S., Bagyamani, J.: White blood cell analysis using watershed and circular hough transform technique. Int. J. Comput. Intell. Inform. **5**(2), 114–123 (2015)
 6. Prabakaran, M.K., Khan, F.I., Abrar, N.M., Abbas, F.M., Khrshid, S.S.: A smart sensing and quantification of platelets, red blood cells (RBC), white blood cells (WBC) and classification of WBC's using microscopic blood image. Int. J. Appl. Med. Sci. Res. **1**(1), 1–9 (2015)
 7. Ravikumar, S., Shanmugam, A.: WBC image segmentation and classification using RVM. Appl. Math. Sci. **8**(45), 2227–2237 (2014)
 8. Othman, M.Z., Ali, A.B.: Segmentation and feature extraction of lymphocyte WBC using microscopic images. Int. J. Eng. Res. Technol. **3**(12), 696–701 (2014)
 9. Shivhare, S., Shrivastava, R.: Automatic bone marrow white blood cell classification using morphological granulometric feature of nucleus. Int. J. Sci. Technol. Res. **1**(4), 125–131 (2012)
10. Chanda, B., Majumder, D.D.: Digital Image Processing and Analysis. Prentice Hall, India (2009)
11. Gonzalez, R.C., Woods, R.E.: Digital Image Processing, 3rd edn. Prentice Hall, India (2009)

# A Multiobjective Ideal Design of Rolling Element Bearing Using Metaheuristics

**S. N. Panda, S. Panda and P. Mishra**

**Abstract** Longest fatigue life is one of the most decisive criteria for design of rolling element bearing. However, the lifetime of bearing will depend on more than one numbers of explanations like fatigue, lubrication, thermal traits. Within the present work, two main goals specifically the dynamic load capability and elasto-hydrodynamic lubrication minimal film thickness have been optimized simultaneously utilizing a multiobjective optimization algorithm centered upon particle swarm optimization. The algorithm accommodates the generalized approach to control combined integer design variables and penalty function strategy of constraint dealing with. The outcomes obtain are encouraging in view of objective function value and computational time. A convergence learns has been applied to make certain the most desirable factor in the design. The most suitable design outcome shows the effectiveness and efficiency of algorithm without constraint violations.

**Keywords** Multiobjective · Particle swarm optimization · Constraint violations

## 1 Introduction

The design disorders of rolling element bearing impacts on its performance, fatigue life, wear life, and durability. Thus, the effective design of bearing can affect the high-quality operation further affecting financial system of machines. The accountability of choosing a most beneficial design from all feasible replacement designs is a tedious job for a bearing designer. The foremost objective of

S. N. Panda (✉) · S. Panda · P. Mishra
Veer Surendra Sai University of Technology, Burla, Odisha 768018, India
e-mail: suryanarayan.uce@gmail.com

S. Panda
e-mail: sumanta.panda@gmail.com

P. Mishra
e-mail: priya.punya@gmail.com

rolling-aspect bearing design is to ensure better efficiency, minimum frictional loss, and further trustworthy operation under variable operating conditions. It is essential that the bearing designer and manufacturer ought to cooperate so as to efficiently achieve the design objective with low life cycle cost for customer [1].

Very less numbers of research have been done on optimization for design of rolling-aspect bearings. Changsen [2] used gradient-centered numerical optimization procedure for curler detail bearing design and formulated a nonlinear constrained multi-objective optimization problem making use of five objective functions such as maximum fatigue life, greatest static load capacity, minimum wear, low frictional minute, and minimal spring to roll proportion. Aramaki [3] developed a PC program to resemble the kinematics attainment of rolling-aspect bearing with few constraints and validated the model with empirical results.

Choi and Yoon [4] applied evolutionary approach for design improvement of automobile wheel bearing assembly design and escalation of its life. Chakraverty et al. [5] evaluated a constraint nonlinear optimization statement for longest fatigue life of the bearing for five design variables by utilizing evolutionary algorithm as GA. Tiwari and Rao [6] have proposed a constraint nonlinear optimization process utilizing evolutionary computing as GA being the design of roller bearing. Waghole et al. [7] suggested novel hybrid procedure for highest quality design of needle roller bearing for maximization of dynamic load capability.

Extensive research has been proposed for improvement of dynamic load capacity using evolutionary methods like GAs. However, in the present work, multiobjective optimization approach with a penalty function method for constraint handling is used to achieve higher dynamic load capacity and optimum minimum film thickness, so that performance of optimally designed bearing can be enhanced without constraint violations. With the intention to attain effective results, ten design parables and ten practical constraints used by the bearing fabricator are utilized in this design optimization problem.

## 2    Problem Formulation

The detail geometry (Fig. 1) of a deep score ball bearing is characterized by the different geometric variables like, bore diameter ($d$), outer diameter ($D$), bearing width ($w$), ball diameter ($D_b$), pitch diameter ($D_m$), inner and outer raceway curvature coefficients ($f_i$ and $f_o$), and number of rolling component ($Z$).

In this article design for entire inner geometry (i.e., $D_b$, $D_m$, $f_i$, $f_o$ and $Z$) of a bearing, at the same time optimizing its efficiency attributes and global fatigue life is addressed. The presence of multiple design requirements make the aforesaid problem come under the area of multiobjective optimization. Any constrained optimization problem composed of three add-ons, particularly design parameters, goal features, and constraints (viable outline parameter space). These factors for the present design problem mentioned below.

**Fig. 1** Radial deep groove ball bearing internal geometries

## 2.1 Design Variables

The design variables are fundamentally geometric specifications and different factors, called primary specifications. Stated specifications were to be resolved for bearing outline. Detailed information specifications as

$$X = [D_b, Z, D_m, f_o, f_i, K_{D\min}, K_{D\max}, \varepsilon, e, \zeta] \tag{1}$$

where $f_o = r_o/D_b$ and $f_i = r_i/D_b$

The parameters for bearing interior geometries are $D_m$, $D_b$, $Z$, $f_i$, and $f_o$ whereas $K_{D\min}$, $K_{D\max}$, $\varepsilon$, $\zeta$, and $e$ are parts of imperatives [6] and are typically saved steady even as designing bearings [6]. For reward work, theses are also handled as variables. Assembly angle ($\phi_0$) of a bearing is a principal constraint on the performance of rolling elements. Established on the analytical formulation in [6], the assembly angle represented as,

$$\phi_0 = 2\Pi - 2\cos^{-1} \frac{\left[\{(D-d)/2 - 3(T-4)\}^2 + \{(D/2 - (T/4) - D_b\}^2 \{(d/2 + (T/4)\}^2\right]}{2\{(D-d)/2 - 3(T/4)\}\{D/2 - (T/4) - D_b\}} \tag{2}$$

where $T = D - d - 2D_b$

## 2.2 Objective Function

There are two performance measures of rolling element bearing (deep groove ball bearing), namely dynamic load capacity ($C_d$) and elastohydrodynamic (EHD) minimum film thickness ($H_{min}$), to be synchronously optimized for best performance of bearing.

The dynamic load capacity, which directly influence the fatigue life of bearing. It is expressed as,

Fatigue life in millions of revolutions,

$$L = \left(\frac{C_d}{F}\right)^a \tag{3}$$

where $F$ is applied load and $a = 3$ for ball bearings

$$C_d = \begin{cases} \max\left[-f_c Z^{2/3} D_b^{1.8}\right] & D_b \leq 25.4 \text{ mm} \\ \max\left[-3.647 f_c Z^{2/3} D_b^{1.4}\right] & D_b > 25.4 \text{ mm} \end{cases} \tag{4}$$

$$f_c = 37.91 \left\{ 1 + \left[1.04\left(\frac{1-\gamma}{1+\gamma}\right)^{1.72}\left(\frac{f_i(2f_o - 1)}{f_o(2f_i - 1)}\right)^{0.41}\right]^{10/3} \right\}^{-0.3} \left[\frac{\gamma^{0.3}(1-\gamma)^{1.39}}{(1+\gamma)^{1/3}}\right]\left[\frac{2f_i}{2f_i - 1}\right]^{0.41} \tag{5}$$

where $\gamma = D_b \cos\alpha / D_m$ is not an impartial parameter; thus, it does not show up in the aim of plan parameters. Observed $\alpha$ is the free contact angle (in present case zero) that relies on bearing geometry. The dynamic load capacity is being determined on the basis of extreme octahedral stress developed between rolling element and races.

Another performance measure of bearing is the longest wear which is predominantly affected by minimum film thickness ($H_{min}$). As per the theory of elastohydrodynamic lubrication, minimum film thickness is formulated [6] considering the minimum value for outer and inner rings calculated separately.

$$H_{min, ring} = 3.63 a_1^{0.49} R_{x, ring}^{0.466} E_o^{-0.117} Q^{-0.073} \left\{\frac{\pi n_i D_m \eta_o (1-\gamma^2)}{120}\right\}^{0.68} \left[1 - \exp\left\{-0.703\left(\frac{R_{(y, ring)}}{R_{(x, ring)}}\right)^{0.636}\right\}\right] \tag{6}$$

$$H_{min} = \min\left(H_{min, inner}, H_{min, outer}\right) \tag{7}$$

where $i$ represents number of rows and it is equal to one for unit row deep score rolling bearing. Some subdefinitions for these aimed functions are

$$Q = F_r/iZ \cos \alpha \tag{8}$$

$$R_{(x, \text{ inner})} = D_b/2(1 - \gamma), \quad R_{(y, \text{ inner})} = f_i D_b/(2f_i - 1) \tag{9}$$

$$R_{(x, \text{ outer})} = D_b/2(1 + \gamma), \quad R_{(y, \text{ outer})} = f_o D_b/(2f_o - 1) \tag{10}$$

Thus, the multiobjective function simultaneously can be formulated using weighted parameters.

$$\max [f(X)] = W_1 \times (C_g/C_{g_0}) + W_2 \times (H_{\min}/H_{\min_0}) \tag{11}$$

$$W_1 + W_2 = 1 \tag{12}$$

where $C_{g_0}$, $H_{\min_0}$ are maximum dynamic load capacity and maximum minimum elastohydrodynamic film thickness. $C_{g0}$, $H_{\min 0}$ are the dynamic load capacity and minimum elastohydrodynamic film thickness, when these were optimized as single objective function.

## 2.3 Constraints

As scope of design of rolling-aspect bearing, various practical design requirements are given by scientists so to cut back the parameter space for ease of design optimization. For the benefit of the bearing assembly, nos. and diameter of rolling element should accompanying prerequisite as below.

$$2 (Z - 1) \sin^{-1}(D_b/D_m) \leq \phi_0 \tag{13}$$

$$g_1(X) = \frac{\phi_0}{2 \sin^{-1}(D_b/D_m)} - Z + 1 \geq 0 \tag{14}$$

Rolling element diameter must be chosen from certain point of confinement, that is,

$$K_{Dmin} \frac{D - d}{2} \leq D_b \leq K_{Dmax} \frac{D - d}{2} \tag{15}$$

where $K_{Dmin}$ and $K_{Dmax}$ are constants as minimum and maximum ball diameter limiters. The corresponding constraint conditions are derived from Eq. (16) and can be given as

$$g_2(X) = 2D_b - K_{D\min}(D - d) \geq 0 \tag{16}$$

$$g_3(X) = K_{D\,max}(D - d) - 2D_b \geq 0 \tag{17}$$

The following two requirements are to be persuaded, so as to ensure the strolling movability of bearings as the distinction between the pitch breadth and the normal diameter in a bearing should be not over shoot a prescribed standard.

$$g_4(X) = D_m - (0.5 - e)(D + d) \geq 0 \tag{18}$$

$$g_5(X) = (0.5 + e)(D + d) - D_m \geq 0 \tag{19}$$

where $e$ is constants named as parameter for mobility condition. In practice, the stress induced in the internal ring is continuously greater than stress induced in the outer ring; this requirement loads to imposition of constraint on the inner ring thickness that it must be greater than or equivalent as outer ring thickness, given as

$$g_6(X) = \frac{d_i - d}{2} - \frac{D - d_o}{2} \geq 0 \tag{20}$$

Terms $d_i$ and $d_o$ represent the inner and outer raceway diameters at the grooves. As reported in [6], thickness of bearing ring at outer raceway bottom should not be less than $\varepsilon D_b$ where $\varepsilon$ is a constant named as parameter for outer ring strength consideration. so the restraint case is

$$g_7(X) = 0.5(D - D_m - D_b) - \varepsilon D_b \geq 0 \tag{21}$$

The constraint on the diameter of the ball given by width of bearing $w$ is

$$g_8(X) = \beta w - D_b \geq 0 \tag{22}$$

For standard bearing specification, dynamic load rating deteriorates if the groove curvature radii of inner and outer raceways in a bearing are not exactly $0.515D_b$ so is continuously greater than $0.515D_b$. Accordingly other restraint cases are as follows:

$$g_9(X) = f_i \geq 0.515D_b \tag{23}$$

$$g_{10}(X) = f_o \geq 0.515D_b \tag{24}$$

## 3   Method for Constraints Handling

The most commonly used approach of constraint handling for constrained optimization problems is the penalty function strategy [8]. A penalty function $P(x)$ is added to original objective function to constrain the constraint violation and convert the constrained optimization problem into an unconstrained optimization statement.

This new output function thus obtained is termed as the pseudo-objective function. Further target objective function is being summed up by the gross constraint violation after normalization. Keeping in mind the end goal to have a superior productivity of the optimization algorithm, a greater value of penalty factor $r_p$ ought to be utilized. Representation of pseudo-objective function is,

$$\phi(X) = F_c(X) + r_p * P(X) \tag{25}$$

$$P(x) = \left[ \sum_{i=1}^{m} \left( \max[0, g_i(X)]^2 \right) + \sum_{i=1}^{p} [h_i(X)]^2 \right] \tag{26}$$

The initial objective function represented as $F_c(x)$ imposed penalty function $P(x)$ as by Eq. (26); $g_i$ and $h_i$ are the inequality and equality constraints, respectively. The multiplier $r_p$ is chosen by taking a few trial runs as in case an optimum result is obtained. Traditionally, it is rather tough to decide on $r_p$ for this reason, the essential trouble of handling constraints is decreased to the decision of the penalty time period so as to strike a steadiness among the knowledge upkeep and the strain of decision of $r_p$.

## 4 Optimization Algorithm (PSO)

A population centered evolutionary algorithm proposed & developed by Kennedy and Eberhart [9] named as Particle swarm optimization (PSO). Here bound to the search space, each and every particle maintains track of its positions, which is correlated to the most effective solution (fitness) it has observed up to now, pBest. Another best esteem followed by the global best version of the particle swarm optimizer is the global best value, gBest and its position, obtained thus far by any particle in the population. The strategy for imposing the PSO is, initialize a population of particles, evaluation the fitness worth of each and every particle, evaluation of each and every particle's evaluated fitness and the fitness evaluation with the population's total prior best, update the velocity and position of the particle as Eqs. (27) and (28), once more evaluation of fitness value of each and every particle until the stopping criterion is met commonly as the highest number of iterations.

$$\begin{aligned} v[] = v[] + c1 * rand() * (pBest[] - present[]) \\ + c2 * rand() * (gBest[] - present[]) \end{aligned} \tag{27}$$

$$present\,[] = present[] + v\,[] \tag{28}$$

**Table 1** Optimization outcomes

| D | d | w | $D_b$ | $D_m$ | Z | $f_i$ | $f_o$ | $\phi_0$ | $K_{Dmin}$ | $K_{Dmax}$ | $\varepsilon$ | $e$ | $\beta$ | $H_{min}$ (μm) | $C_g$ |
|---|---|---|-------|-------|---|-------|-------|------|--------|--------|--------|--------|--------|--------|--------|
| 30 | 10 | 9 | 6.97 | 20.68 | 7 | 0.515 | 0.519 | 3.12 | 0.4447 | 0.6543 | 0.3008 | 0.0445 | 0.7658 | 0.171 | 6912.3 |
| 35 | 15 | 11 | 6.96 | 26.35 | 7 | 0.523 | 0.515 | 3.44 | 0.4564 | 0.6238 | 0.0302 | 0.0658 | 0.7356 | 0.197 | 7015.1 |
| 47 | 20 | 14 | 8.56 | 34.60 | 8 | 0.515 | 0.515 | 3.68 | 0.4128 | 0.6782 | 0.300 | 0.0426 | 0.7240 | 0.234 | 10,546.0 |
| 62 | 30 | 16 | 10.50 | 45.46 | 8 | 0.515 | 0.515 | 3.69 | 0.4469 | 0.6246 | 0.300 | 0.0479 | 0.7891 | 0.316 | 17,698.0 |
| 80 | 40 | 18 | 12.44 | 61.50 | 9 | 0.515 | 0.515 | 3.12 | 0.4378 | 0.6456 | 0.300 | 0.0354 | 0.7354 | 0.435 | 26,854.0 |
| 90 | 50 | 20 | 12.68 | 70.55 | 11 | 0.515 | 0.515 | 3.73 | 0.4129 | 0.6364 | 0.300 | 0.0691 | 0.7769 | 0.512 | 28,752.0 |
| 110 | 60 | 22 | 17.46 | 86.56 | 11 | 0.515 | 0.515 | 3.95 | 0.4398 | 0.6458 | 0.332 | 0.0801 | 0.8231 | 0.627 | 49,957.0 |
| 125 | 70 | 24 | 17.92 | 98.243 | 11 | 0.515 | 0.515 | 3.36 | 0.4436 | 0.6433 | 0.300 | 0.0578 | 0.8971 | 0.711 | 50,567.0 |
| 140 | 80 | 26 | 18.87 | 118.37 | 11 | 0.515 | 0.515 | 3.58 | 0.4129 | 0.6758 | 0.324 | 0.0982 | 0.8862 | 0.812 | 63,568.0 |
| 160 | 90 | 30 | 21.66 | 125.25 | 11 | 0.515 | 0.515 | 3.12 | 0.4589 | 0.6468 | 0.301 | 0.0458 | 0.8497 | 0.924 | 79,436.0 |
| 170 | 95 | 32 | 23.08 | 134.02 | 11 | 0.515 | 0.515 | 3.46 | 0.4897 | 0.6456 | 0.300 | 0.0652 | 0.8346 | 1.121 | 89,565.0 |

## 5  Outcomes and Analysis

Table 1 demonstrates an improvement of dynamic load capacity and EHL minimum film thickness as compared to validated available research results reported in [10]. Table 2 gives an imperative perception of the advancement of dynamic capacity of bearing composed utilizing PSO when contrasted with GA and catalog, likewise the rate of lessening of least film thickness when contrasted with GA. Further, it is observed that there is 36.98% of the average reduction in minimum film thickness value as compared to GA. Time consumed by the algorithm for computation of outcome on a dual-core 1.8 GHz (2 GB RAM) Window 7 platform was 0.198 min. Figure 2 depicts the convergence traits of the algorithm. Figure 3 depicts the convergence characteristics of all design variables. It may be insured from Fig. 3 that the variables $f_i$, $f_o$, $Z$, and $e$ need better control (as these variables

**Table 2** Comparisons of outcomes

| $(C_g)$ (PSO) | (Tiwari et al.GA) $(C_d)$ | (Catalog) $(C_s)$ | $\lambda_1 = Cg/Cs$ | $\lambda_2 = C_d/C_s$ | $H_{min}$ (µm) | (Gupta et al. GA) $H_{min}$ (µm) | % Reduction for $H_{min}$ (µm) |
|---|---|---|---|---|---|---|---|
| 6912.3 | 5942.36 | 3580 | 1.93 | 1.66 | 0.171 | 0.2096 | 18.42 |
| 7015.1 | 6955.35 | 5870 | 1.2 | 1.18 | 0.197 | 0.2581 | 23.67 |
| 10,546 | 10,890.9 | 9430 | 1.12 | 1.15 | 0.234 | 0.362 | 35.36 |
| 17,698 | 16,387.4 | 14,900 | 1.19 | 1.1 | 0.316 | 0.5037 | 37.26 |
| 26,854 | 26,678.4 | 22,500 | 1.19 | 1.19 | 0.435 | 0.6748 | 35.54 |
| 28,752 | 28,789.3 | 26,900 | 1.07 | 1.07 | 0.512 | 0.777 | 34.11 |
| 49,957 | 42,695.3 | 40,300 | 1.24 | 1.06 | 0.627 | 0.9763 | 35.78 |
| 50,567 | 51,117.4 | 47,600 | 1.06 | 1.07 | 0.711 | 1.1323 | 37.21 |
| 63,568 | 59,042.9 | 55,600 | 1.14 | 1.06 | 0.812 | 1.2894 | 37.02 |
| 79,436 | 75,466.8 | 73,900 | 1.07 | 1.02 | 0.924 | 1.5026 | 38.51 |
| 89,565 | 89,244.7 | 83,700 | 1.07 | 1.07 | 1.121 | – | – |



**Fig. 2** Convergence traits of PSO

**Fig. 3** Convergence traits of design parameters

undergo larger variation in the optimum design). All other variables are converging rapidly to the optimal solution point.

## 6   Conclusions

The PSO algorithm is practiced on stated constraint multiobjective optimization problem involving dynamic load capacity and EHL minimum film thickness for deep groove ball bearing. The algorithm successfully handles mixed integer variables, and a penalty function approach has been used for efficient constraint handling. The reported result indicates the preeminence of proposed PSO algorithm over GA. The convergence characteristics show that $f_i$, $f_o$, $Z$, and $e$ converge to a very narrow region after undergoing large variation. This shows that these are the key design variable and need better control in optimum design. The very narrow region thus can be designed as key design variables. The EHD minimum film thickness reported is less as compared to GA, hence the manufacturer to take in to account the surface finish of the balls, outer raceway, and inner raceway (a tight tolerance is to be maintained during manufacturing of bearing).

# References

1. Asimow, M.: Introduction to Engineering Design. McGraw-Hill, New York (1966)
2. Changsen, W.: Analysis of Rolling Element Bearings. Mechanical Engineering Publications Ltd., London (1991)
3. Aramaki, H.: Basic technology research and development center, motion and control no. 3. NSK, rolling bearing analysis program package BRAIN (1997)
4. Choi, D.H., Yoon, K.C.: A design method of an automotive wheel bearing unit with discrete design variables using genetic algorithms. Trans. ASME J. Tribol. **123**(1), 181–187 (2001)
5. Chakraborthy, I., Vinay, K., Nair, S.B., Tiwari, R.: Rolling element bearing design through genetic algorithms. Eng. Optim. **35**(6), 649–659 (2003)
6. Tiwari, R., Rao, B.R.: Optimum design of rolling element bearings using genetic algorithms. Mech. Mach. Theory **42**(2), 233–250 (2007)
7. Waghole, V., Tiwari, R.: Optimization of needle roller bearing design using novel hybrid method. Mech. Mach. Theory **72**(2), 71–85 (2014)
8. Deb, K.: Multi-objective Optimisation Using Evolutionary Algorithms. Wiley-Interscience Series in Systems and Optimisation, New York (2001)
9. Kennedy, J, Eberhart, R.: Particle swarm optimization. In: Proceedings of the IEEE International Conference on Neural Networks (ICNN), pp. 1942–1948 (1995)
10. Gupta, S., Tiwari, R., Nair, S.B.: Multi objective design optimization of rolling bearing using genetic algorithm. Mech. Mach. Theory **42**, 1418–1443 (2007)

# Predicting Binding Affinity Based on Docking Measures for Spinocerebellar Ataxia: A Study

**P. R. Asha and M. S. Vijaya**

**Abstract** An obsessive stipulation impairs the regular function or structure of an organ in humans. Spinocerebellar ataxia disorder is a hereditary genetic disorder which is originated by the massive number of sequence variants found in large sets of genes. The mutation in the genes causes many of these disorders. There are certainly no effective drugs to treat those disorders. There are many types of spinocerebellar ataxia, and a better knowledge is required to forecast binding affinity. Binding affinity is crucial to screen the drugs for spinocerebellar ataxia disorder. Accurate identification of binding affinities is a profoundly demanding task. To overcome this issue, a new approach is to be designed in identifying the binding affinity effectively. Due to rapid growth of biological data, there is an increase in the processing time and cost efficiency. This paves the way for challenges in computing. The purpose of machine learning is to excavate beneficial knowledge in distinct to corpus of information and data by constructing effective feasible designs. In this paper, a preface to spinocerebellar ataxia, conventional and innovative strategies involved in predicting binding affinity are discussed.

**Keywords** Proteins · Protein structure · Homology modeling
Docking and affinity

## 1 Introduction

Spinocerebellar ataxia is a hereditary anarchy portrayed by deviations in gray matter handling its tasks. The disorder is due to mutations in the genes results in brain and spinal cord degeneration. Individuals with SCA encounter progressive

P. R. Asha (✉) · M. S. Vijaya
Department of Computer Science (PG), PSGR Krishnammal College for Women,
Coimbatore, India
e-mail: ashamscsoft@gmail.com

M. S. Vijaya
e-mail: msvijaya@psgkrc.com

atrophy, or muscle wasting. Negative consequences of the illness might include the person's hands and wrists, lower limbs, and also speech patterns. Till now, 36 types of spinocerebellar ataxia are discovered. Each type of SCA features its own symptoms [1]. Spinocerebellar ataxia is considered as an inherited disorder flow in a family, even when one parent is influenced. Autosomal recessive, dominant and X-linked are the three patterns of inheritance that causes spinocerebellar ataxia. The recessive pattern of a disease involves two copies of passed down faulty genes, one from each parent in which both will be carriers of the disease but usually not subject to the disease. The dominant pattern requires one copy of faulty gene to trigger the disorder. In X-linked pattern, the disease is passed only from mother to their offspring. In female, two pairs of X chromosomes exist thereby the daughters turn up into carriers, and usually not influenced by the disease, on the other hand, the male comprises of only one X-chromosome. It arises mainly as a result of autosomal dominant (AD) inheritance pattern [2]. SCA1, SCA2, SCA3, SCA6, SCA7, and SCA8 are some of the common forms of spinocerebellar ataxia [3]. Some of the common types of spinocerebellar ataxia and its clinical features, mutations, age onset, number of expansions occurred in the gene due to mutation are presented in Table 1.

## 2 Homology Modeling

Homology modeling, otherwise called as comparative modeling of protein, which tells about creating a model of the "target" protein from its sequence, and 3D structure of a protein from its suitable homologs known as template. It is an essential computational approach, to figure out the 3D structure. This method uses high-end resolution protein structures to create a model of unknown structure [4].

In Fig. 1 the simplified elucidation of modeling procedure is given. Initial process to retrieve the homologous sequences is performed in order to make an alignment. Based on this alignment the model is built. Then, adjusting the alignment using external data such as secondary structure information should be done by known motifs and conserved features. With this information, model is built using software. After constructing the model, it has to be refined and inspected.

Template identification is the crucial initial step. It lays the basis by determining ideal homolog(s) of recognized structure, referred to as templates, which are adequately parallel to the objective sequence to be modeled. The next phase engages generating an alignment with the template structure. The objective, as well as the template sequence, is generally connected with a protein domain family retrieved from Pfam [5], or a custom alignment can be generated by means of BLAST. In segment matching methods, the target is divided into short segments, for example, will split the sequence into 30 amino acids, and alignment is conducted over fragments instead of over the whole protein [6]. The technique is executed making use of the preferred program Modeler [7], which includes the CHARMM [8] energy conditions that guarantee valid stereochemistry is combined with spatial restraints.

Table 1 Details of common types of spinocerebellar ataxia

| Types | SCA1 | SCA2 | SCA3 | SCA6 | SCA7 | SCA8 |
|---|---|---|---|---|---|---|
| Age onset | Onset over age 60 years | Fourth decade | Second to fifth decade | Age onset is 43–52 yrs | Second to fourth decade | Age onset is from one to 73 years |
| Clinical features | Speech and swallowing difficulties, spasticity and weakness in the muscles which control the movement of eye | Difficulty faced in speech and swallowing, rigidity, tremors | Dystonia, spasticity, bulging eyes | Nystagmus, double vision | Muscle weakness wasting, hypotonia, poor feeding, failure to thrive | Muscle spasticity, drawn-out slowness of speech and reduced vibration sense |
| Gene mutation and chromosome no. | ATXN1 gene located on the chromosome 6 [29] | ATXN2 gene found on the chromosome 12 [30] | ATXN3 gene situated on the chromosome 14 [31] | CACNA1A gene being on the chromosome 19 [32] | ATXN7 gene found on the chromosome 3 [33] | ATXN8 and ATXN8 OS (reverse strand) genes found on the chromo Some 13 [34] |
| No. of repeats mutation occurred in the gene | Victims possess alleles with 39 or more CAG trinucleotide repeats [35] | Victims possess alleles with 33 or even more CAG trinucleotide repeats [36] | Victims carry alleles from 552 to 86 CAG trinucleotide repeats [37]. | Victims contain 20–33 CAG repeats [38] | Victims will often have higher than 36 CAG repeats [39] | Individuals carry from 80–250 CTG repeats [40] |

Target Sequence

Template identification and selection

Alignment

Structural alignment

Multiple sequence alignment

Iterate

Build model

Refine model

Evaluate model

Good?

No          Yes

Use model

**Fig. 1** Simplified illustration of the modeling process

A universal principle is that any specific insertion/loop prolonged than about five residues must not be considered. Once the model is built, it needs to be validated. Modeling involves a high-resolution experimental structure as a template, the precision which sprightly impacts the quality of the model. Progressively, the quality will depend on the degree of sequence uniqueness between the template and the protein to be modeled [9]. Alignment mistakes rise swiftly, whenever the sequence uniqueness is low that is below 30%. Moderate models possess between 30 and 50% sequence uniqueness to the template. Greater accuracy models are usually achieved if there is more than 50% sequence identity. There will be lot of information regarding biological function that may be produced from a three-dimensional structure [10]. Structures can be utilized to describe the impact of mutations in drug confrontation and genetic disorders [11]. Some of the tools which are used for homology modeling are listed in Table 2.

**Table 2** Tools used for homology modeling

| Tools used | Description | Pioneered papers |
|---|---|---|
| BLAST | Basic search tool in which the alignment of sequences is performed either for DNA or protein | Altschul et al. [41] |
| PSI–BLAST | Similar to BLAST, but find distant homologs and it can be performed only for protein sequences | Altschul et al. [42] |
| Scanps | Multiple sequence alignment server | Barton [43] |
| Clustal | Multiple sequence alignment program | Higgins and Sharp [44] |
| Muscle | Gives better sequence alignments, especially for larger alignments | Edgar [45] |
| T-coffee | More accurate alignment than other methods. It is equipped with many tools for evaluation and alignments | Notredame et al. [46] |

## 3 Docking

Molecular docking is always to forecast the prominent binding mode(s). When compared with standard experimental high-throughput screening (HTS), virtual screening is definitely a straighter and a reasonable drug discovery strategy [12] and also possesses the benefit of inexpensive and valuable screening. VS could be categorized into ligand-based and structure-based methods. Regarding structure-based drug design, molecular docking is considered the most frequent approach that has been commonly used from the time the early 1980s [13].

Recognizing the hot spot before docking procedures considerably improves the docking performance. Without the presence of information about the binding sites, cavity identification applications are available, otherwise online servers, e.g., GRID [14], POCKET [15], SurfNet [16], PASS [17] and MMC [18] can be utilized to recognize putative effective sites within proteins. Docking without knowing the idea regarding hot spot is known as blind docking. An earlier illumination for the ligand–receptor-binding procedure is lock-and-key principle [19], wherein the ligand sits into the macromolecule just like lock-and-key. After that induced-fit concept [20], it carries lock-and-key theory a phase more, proclaiming that the energetic site of the protein is constantly reshaped by interactions with the ligands since the ligands communicate with the macromolecule.

Docking can be accomplished via two interlinked procedures: one is by sampling conformations of ligand in the active site; after that rating these types of conformations by means of a scoring function [21]. Scoring functions might be categorized in force field-based, empirical, and knowledge-based scoring functions [22]. Traditional force field-based scoring functions [23] evaluate the binding energy by computing the amount of the nonbonded communication. There are two methodologies available in docking. They are rigid-ligand and rigid-receptor docking and flexible-ligand and rigid-receptor docking [21]. Some of the docking programs which are classified as freeware are listed in Table 3 and which are

classified as commercial is listed in Table 4. Table 5 shows the docking programs which are academic

## 4 Binding Affinity

Receptor is a protein molecule that acquires a signal by binding to a chemical called ligand. Affinity is described as measure of the strength of attraction between a receptor and its ligand. The ability of ligand is to form coordination bond with a receptor known as binding affinity [24].

The affinity of a ligand is determined by the force of attraction between macromolecule binding sites with the ligand. When there is high-affinity binding, it illustrates a lengthy dwelling time at the binding site of the macromolecule when compared with low-affinity binding. Table 6 portrays the review of works done for binding affinity prediction using general approaches and computational approaches.

The main drawback of binding affinity prediction using general approach is scoring function. Scoring functions are not able to forecast binding affinity or binding free energy for a couple of reasons: They determine mainly enthalpic terms, and disregard entropy, especially of the protein. Entropy is needed, certainly, to calculate binding free energy. Scoring functions merely are familiar with the bound state of the protein–ligand system, not the unbound states of the macromolecule and the ligand. Binding free energy can only be anticipated through knowledge of the bound state and the unbound states of the binding partners.

Machine learning techniques can automatically learn the model by taking intelligent hints from the data and predict the output more accurately. Several models were built using machine learning techniques to predict binding affinity. The authors Krammer et al. [25] proposed affinity prediction by using 118 protein–ligand complexes. Surface descriptors features and two empirical scoring functions methods were used for prediction. In another literature, Darnell et al. [26] anticipated protein–protein affinity prediction using decision tree by the method knowledge-based models. Shape specificity and biochemical contact features were used to predict affinity. Li et al. [27] projected two-layer support regression model for protein–protein affinity prediction. The suggested features were heavy atoms

**Table 3** List of docking programs (freeware)

| Program | Description |
|---------|-------------|
| I-click docking | Identifies binding orientation and affinity of a ligand |
| AADS | Automated active site detection, docking, and scoring (AADS) protocol |
| AutoDock | Automated docking of ligand to macromolecule |
| AutoDock Vina | New generation of AutoDock |
| BetaDock | Based on Voronoi diagram |
| HADDOCK | Mainly for protein–protein docking |
| Score | It is to calculate docking scores of ligand–receptor complex |

**Table 4** Docking programs (commercial)

| ADAM | Prediction of binding mode |
|---|---|
| Docking server | Integrates a number of computational chemistry software |
| FlexX | Incremental build-based docking program |
| Glide | Exhaustive search-based docking program |
| GOLD | Genetic algorithm-based docking program |
| ICM-Dock | Pseudo-Brownian sampling and local minimization docking program |

**Table 5** Docking programs (academic)

| Hammerhead | Fully automated docking of flexible ligands to protein |
|---|---|
| PSI-DOCK | Pose-sensitive inclined (PSI)-DOCK |
| PSO@AUTODock | Particle swarm optimization (PSO) algorithms are used for rapid docking of flexible ligands |

interacted with protein. Li et al. [28] proffered affinity prediction using random forest by RF score. Totally 36 features were considered to predict affinity.

The further work can be suggested by building machine learning models to predict affinity, rather than extracting from available opensource tools.

## 5 Observations and Discussions

From the study, it is observed that few machine learning techniques are being adopted to predict binding affinity based on docking measures. Machines learning technique reduces the time of processing large data and provides high accuracy by building regression models to predict binding affinity. Binding affinity of ligand to receptor is very essential since certain level of binding energy is commonly employed in the macromolecule to carry out a conformation in the transformation. Binding affinity is utilized for drug investigation, through which the receptor is going to block the progress of ailments otherwise trigger the enzyme to execute certain functions on target protein.

Ligands similar to drugs possess some specificity for the binding site on the macromolecule similar to objective protein. When the specificity is significantly fewer, this implies binding affinity is low; an increased dosage of medicine might be preferred, which can be connected with certain side effects. Therefore, the efficiency of a treatment depends upon the binding affinity, for which the binding site along with their binding capacity to produce the desired effects.

**Table 6** Review of work binding affinity prediction

| Objectives | Author | Data set | Method | Features |
|---|---|---|---|---|
| *General approaches* | | | | |
| To estimate the binding constant using empirical scoring function | Böhm [47] | 45 protein–ligand complex | Free energy function | – |
| Identifying binding modes and hot spots using knowledge-based scoring functions | Gohlke et al. [48] | PDB | DrugScore | – |
| Predicting interactions of a complex | Huang and Zou [49] | 100 protein–ligand complex | IT score | – |
| Identification of a protein-ligand Interactions | Plewczynski et al. [50] | PDBbind | Consensus approach | – |
| Prediction of a ligand binding poses | Tanchuk et al. [51] | PDBbind | Hybrid scoring function | – |
| *Machine learning approaches* | | | | |
| Affinity prediction | Krammer et al. [25] | 118 protein–ligand complexes | Two empirical scoring functions | Surface descriptors |
| Protein–protein affinity prediction using decision tree | Darnell et al. [26] | PDB | Two knowledge-based models | Shape specificity and biochemical contact features |
| Protein–protein affinity prediction | Li et al. [27] | PDBbind-CN | Two-layer support regression model | Heavy atoms were considered as the feature vectors |
| Predicting affinity using random forest | Li et al. [28] | PDBbind database | RF score was used | RF Score features with 6 AutoDock vina Features |

# 6   Conclusion

In our recent survey research, predicting binding affinity based on docking mea-
sures for spinocerebellar ataxia disease identification through computational intel-
ligence is reviewed. Proposed work will be focused on finding affinity by machine

learning models, rather than extracting the values from opensource tools. This paper elucidates the introduction of protein mutations, the vulnerability of spinocerebellar ataxia and various techniques involved in predicting the affinity briefly. From the observations, it is concluded that efficiency of predicting the binding affinity is very tough to process for huge complex. To deal with a large number of complexes, new disease identification model should be designed and developed based on the advanced learning techniques like deep learning. A distributed environment should be created with the big data technologies like Hadoop and its components that support in predicting the affinity effectively using a large number of mutated protein sequences.

# References

1. Weiss, T.C.: Ataxia Spinocerebellar: SCA Facts and Information (2010)
2. Whaley, N.R., Fujioka, S., Wszolek, Z.K.: Autosomal dominant cerebellar ataxia type I: a review of the phenotypic and genotypic characteristics (2011). doi:10.1186/1750-1172-6-33
3. Bird, T.D.: Hereditary Ataxia Overview, March 3 (2016)
4. Bishop, A., de Beer, T.A., Joubert, F.: Protein homology modelling and its use, Feb 2008, South Africa
5. Sonhammer, E.L., Eddy, S.R., Durbin, R.: Pfam: a comprehensive database of protein domain families based on seed alignments. Proteins **28**, 405–420 (1997)
6. Levitt, M.: Accurate modelling of protein conformation by automatic segment matching. J. Mol. Biol. **226**, 507–533 (1992)
7. Sali, A., Blundell, T.L.: Comparative protein modelling by satisfaction of spatial restraints. J. Mol. Biol. **234**, 779–815 (1993)
8. Brooks, B.R., Bruccoleri, R.E., Olafson, B.D., States, D.J., Swaminathan, S., Karplus, M.: CHARMM: a program for macromolecular energy, minimization, and dynamics calculations. J. Comput. Chem. **4**, 187–217 (1983)
9. Baker, D., Sali, A.: Protein structure prediction and structural genomics. Science **294**, 93–96 (2001)
10. Thornton, J.M., Todd, A.E., Milburn, D., Borkakoti, N., Orengo, C.A.: From structure to function: approaches and limitations. Nat Struct. Biol. Suppl. 991–994 (2000)
11. Marbotti, A., Facchiano, A.M.: Homology modeling studies on human galactose-l-phosphate uridylytransferase and on its galactosemia-related mutant Q188R provide an explanation of molecular effects of the mutation on homo- and heterodimers. J. Med. Chem. **48**, 773–779 (2005)
12. Walters, W.P., Stahl, M.T., Murcko, M.A.: Virtual screening—an overview. Drug Discov. Today. 160–178 (1998)
13. Kuntz, I.D., Blaney, J.M., Oatley, S.J., Langridge, R., Ferrin, T.E.: A geometric approach to macromolecule-ligand interactions. J. Mol. Biol. **161**(2), 269–288 (1982)
14. Goodford, P.J.: A computational procedure for determining energetically favorable binding sites on biologically important macromolecules. J. Med. Chem. **28**(7), 849–857 (1985)
15. Levitt, D.G., Banaszak, L.J.: POCKET: a computer graphics method for identifying and displaying protein cavities and their surrounding amino acids. J. Mol. Graph. **10**(4), 229–234 (1992)
16. Laskowski, R.A.: SURFNET: a program for visualizing molecular surfaces, cavities and intermolecular interactions. J. Mol. Graph. **13**(5), 323–330 (1995)
17. Brady Jr., G.P., Brady Jr., G.P., Stouten, P.F.: Fast prediction and visualization of protein binding pockets with PASS. J. Comput. Aided Mol. Des. **14**(4), 383–401 (2000)

18. Mezei, M.: A new method for mapping macromolecular topography. J. Mol. Graph. Model. **21**(5), 463–472 (2003)
19. Fischer, E.: Einfluss der configuration auf die wirkung derenzyme. Ber. Dt. Chem. Ges. **27**, 2985–2993 (1894)
20. Koshland Jr., D.E.: Correlation of structure and function in enzyme action. Science **142**, 1533–1541 (1963)
21. Meng, X.-Y., Zhang, H.-X., Mezei, M., Cui, M.: Molecular Docking: A Powerful Approach for Structure-Based Drug Discovery (2012)
22. Kitchen, D.B., Decornez, H., Furr, J.R., Bajorath, J.: Docking and scoring in virtual screening for drug discovery: methods and applications. Nat. Rev. Drug Discov. **3**(11), 935–949 (2004)
23. Aqvist, J., Luzhkov, V.B., Brandsdal, B.O.: Ligand binding affinities from MD simulations. Acc. Chem. Res. **35**(6), 358–365 (2002)
24. http://chemistry.tutorvista.com/inorganic-chemistry/binding-affinity.html
25. Krammer, A., Kirchhoff, P.D., Jiang, X., Venkatachalam, C.M., Waldman, M.: LigScore: a novel scoring function for predicting binding affinities, Nov (2004)
26. Darnell, S.J., Page, D., Mitchell, J.C.: An automated decision-tree approach to predicting protein interaction hot spots (2007)
27. Li. X., Zhu, M., Li, X., Wang, H.-O., Wang, S.: Protein–protein binding affinity prediction based on an SVR ensemble. Intelligent Computing
28. Li, H., Leung, K.-S., Wong, M.-H., Ballester, P.J.: The use of random forest to predict binding affinity in docking. Bioinform. Biomed. Eng. Ser, **9044**, 238–247 (2015)
29. Spinocerebellar Ataxia 1; SCA1, Online Mendelian Inheritance in Man (OMIM)
30. Spinocerebellar Ataxia 2, SCA2; Online Mendelian Inheritance in Man (OMIM)
31. Spinocerebellar Ataxia, Type 3, SCA3, Machado-Joseph Disease; Online Mendelian Inheritance in Man (OMIM)
32. Spinocerebellar Ataxia 6, SCA6; Online Mendelian Inheritance in Man (OMIM)
33. Spinocerebellar Ataxia 7, SCA7; Online Mendelian Inheritance in Man (OMIM)
34. Spinocerebellar Ataxia 8, SCA8; Online Mendelian Inheritance in Man (OMIM)
35. Subramony, S.H., Ashizawa, T.: Spinocerebellar Ataxia Type 1, July 3 (2014)
36. Pulst, S.M.: Spinocerebellar Ataxia Type 2, November 12 (2015)
37. Paulson, H.: Spinocerebellar Ataxia type 3, Sept. 24 (2015)
38. Gomez, C.M.: Spinocerebellar Ataxia type 6, July 18 (2013)
39. Garden, G.: Spinocerebellar Ataxia type 7, December 20 (2012)
40. Ayhan, F., Ikeda, Y. et al.: Spinocerebellar Ataxia type 8, April 3 (2014)
41. Altschul, S.F., Gish, W., Miller, W., Myers, E.W., Lipman, D.J.: Basic local alignment search tool. J. Mol. Biol. **215**, 403–410 (1990)
42. Altschul, S.F., Madden, T.L., Schaffer, A.A., Zhang, J., Zhang, Z., Miller, W., Lipman, D.J.: Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. Nucl. Acids Res. **25**, 3389–3402 (1997)
43. Barton, G.J.: Computer speed and sequence comparison. Science **257**, 1609–1610 (1992)
44. Higgins, D.G., Sharp, P.M.: CLUSTAL: a package for performing multiple sequence alignment on a microcomputer. Gene **73**, 237–244 (1988)
45. Edgar, R.C.: MUSCLE: a multiple sequence alignment method with reduced time and space complexity. BMC Bioinform. **5**, 113 (2004)
46. Notredame, C., Higgins, D., Heringa, J.: T-Coffee: a novel method for multiple sequence alignments. J. Mol. Biol. **302**, 205–217 (2000)
47. Böhm, H.-J.: The development of a simple empirical scoring function to estimate the binding constant for a protein-ligand complex of known three-dimensional structure, 15 Dec (1993)
48. Gohlke, H., Hendlich, M., Klebe, G.: Predicting binding modes, binding affinities and 'hot spots' for protein-ligand complexes using a knowledge-based scoring function, pp. 115–144, vol. 20, Dec (2000)
49. Huang, S.-Y., Zou, X.: An iterative knowledge-based scoring function to predict protein–ligand interactions: II. Validation of the scoring function (2006)

50. Plewczynski, D., Łażniewski, M., Grotthuss, M.V., Rychlewski, L., Ginalski, K. VoteDock: Consensus docking method for prediction of protei–ligand interactions. Sept (2010)
51. Tanchuk, V.Y., Tanin, V.O., Vovk, A.I., Poda, G.: A new, improved hybrid scoring function for molecular docking and scoring based on autodock and autodock vina. Chem. Biol. Drug Des. Dec (2015)

# A Novel Image Hiding Technique Based on Stratified Steganography

M. Radhika Mani and G. Suryakala Eswari

**Abstract** In current Web-based era, the security is becoming crucial factor for exchange of digital data. Other than cryptography, steganography is widely used due to its applications to multimedia information. Among various multimedia content, images are considered for steganography due to its large amount of redundant space availability. In steganography, bit stream-based information can be hidden in any carrier. The present paper proposes a novel algorithm for hiding gray image in the input cover image called as "Logical Stratified Steganography Technique." The proposed algorithm uses OR operation to combine the cover and hidden images. To test the integrity and robustness, the proposed method is applied on various images like human faces, textures, and medical images. Various performance measures are used to compare the quality of the proposed method. The strength of the proposed method is clearly exhibited with the evaluated results.

**Keywords** Quality · Shift · Pixel and quality measures

## 1 Introduction

Nowadays, information security is widely used in various areas of real-time applications. Information security deals with the exchange of secret information hidden in cover media. Various forms of information security are cryptography, steganography, and digital watermarking. Among them, steganography is found to be efficient for hidden communication. In cryptography, the intruder can identify the presence of communication but they fail to understand the original message.

M. Radhika Mani (✉)
Department of CSE, Pragati Engineering College, Surampalem, India
e-mail: radhika_madireddy@yahoo.com

G. Suryakala Eswari
Department of IT, Pragati Engineering College, Surampalem, India
e-mail: padmakala1@gmail.com

This problem is conquered by the steganography. In the process of steganography, presence of the hidden message is not identified by the intruders. The cryptography is applicable to only plain text whereas the steganography is applicable to any form of the multimedia, viz., image, audio, and video. Currently, major research is going in the image-based steganography methods. Asymmetric and symmetric methods represent the categorization [1]. The aim of these algorithms should be the security.

The original meaning is to hide the data in a plain sight. It is usually used for hiding information within information. The primary goal is to not to be suspected by any intruder for carrying information. It can use various forms of carrier such as video, audio, image, and text. For any type of carrier, the encryption and decryption techniques should ensure the quality. Though there are numerous techniques available, they are failed in maintaining the confidentiality. Thus, a novel approach is required for higher confidentiality measure.

Virus-evolutionary genetic algorithm [2] is using ensemble classifier and genetic optimization with an additional virus population and is found to be efficient than other universal classifiers for steganography. A message with double bits [3] is encrypted in the bit planes of second and fourth, and the embedding process is performed at one bit per plane. This method yields acceptable capacity of performance. X-Box mapping [4] is used in image steganography to enhance the security for Web-based data transmission. The injection steganography [5] technique combines the watermarking and encryption techniques. An optimal approach [6] uses cover image's hue color represented with 24 bits. This will embed the information [7] in higher LSB layers which results an increase in imperceptibility. The ALP-PC (Adaptive Linear Programming of Polar Codes) [8] is used to minimize distortion function for steganography [9, 10]. Color table expansion [11] technique is used for indexed images. This will yield high efficiency for the security and hiding capacity. Edge processing-[12] based method consists of various stages, viz., marginalization, reconstruction, block markers, and embedding. The hybrid crypto-steganography [13] is developed based on hierarchical access to the security.

The current algorithms are further required to enhance the confidentiality measure for hiding image. Thus, the present work proposes an innovative approach for efficient image steganography. The present work is structured into four sections. The proposed methodology is arranged in Sect. 2, the results and discussions in Sect. 3, and the conclusions in Sect. 4.

## 2  Methodology

In this technique, an image is hidden in the cover image using 4-bit LSB steganography. In the process of hiding the input hidden image, 4 LSBs of each pixel are replaced by the 4 MSBs of the corresponding pixel which is used to form the stego image as shown in Fig. 1. As 4 MSBs of the cover image pixels remain

same, the stego image does not differ greatly from the cover image. Human perception cannot recognize the change. Thus formed stego image is transmitted to receiver side. The flow chart of Fig. 3 details the process of hiding the image. During the retrieving process, all stego image pixels are scanned.

The formation of stego pixel from the cover and hidden pixels is depicted in Fig. 2. The 4 LSBs are retrieved and are shifted left 4 times, so that they form the most significant bits. Then, the corresponding intensity is displayed to get the hidden image. As the 4 MSBs of the hidden image remain same, the received hidden image doesn't vary greatly from the original hidden image. Fig. 3 depicts the retrieval procedure. The present paper evaluated various performance measures [14], viz., RSNR (Root Signal-to-Noise Ratio), MAE (Mean Absolute Error), SNR (Signal-to-Noise Ratio), PSNR (Peak Signal-to-Noise Ratio), RMSE (Root Mean Square Error), and MSE (Mean Square Error).



**Fig. 1** Flow chart for logical stratified steganography technique

Cover image pixel

| 1 | 0 | 1 | 1 | 0 | 0 | 1 | 1 |

Hidden image pixel

| 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 |

| 1 | 0 | 1 | 1 | 0 | 1 | 1 | 0 |

Stego image pixel

**Fig. 2** Formation of stego pixel in the proposed LSST

Stego image pixel

| 1 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | → Left shift 4 times → | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |

Hidden image pixel

**Fig. 3** Hidden image retrieval in the proposed LSST

## 3  Results and Discussions

The Logical Stratified Steganography Technique (LSST) uses 4 MSBs of the corresponding hidden pixel. The proposed method is experimented on various images shown in Fig. 4a–i. The hidden image of Fig. 5 is embedded in the input images, and the stego images are shown in the 6a–i. The stego images of LSST represent the robustness and integrity of the proposed method.

The performance measure for the proposed LSST is evaluated and given in Table 1. From Table 1, it is apparent that the LSST has higher values of MSE, MAE, and PSNR. The SNR, RSNR, and PSNR values are comparatively very low indicating lesser perceptibility of the message.

## 4  Conclusions

Presence of redundant data in images makes use of them as efficient cover media for steganography. An innovative method for hiding information within cover for reliable communication is focused. The proposed LSST technique is found to maintain the perceptual quality.

**Fig. 4** Original images **a** Human 1 (H1) **b** Human 2 (H2) **c** Human 3 (H3) **d** Texture 1 (T1) **e** Texture 2 (T2) **f** Texture 3 (T3) **g** Medical 1 (M1) **h** Medical 2 (M2) **i** Medical 3 (M3)

**Fig. 5** Hidden image for LLST



**Fig. 6** LSST stego images **a** Human 1 (H1) **b** Human 2 (H2) **c** Human 3 (H3) **d** Texture 1 (T1)
**e** Texture 2 (T2) **f** Texture 3 (T3) **g** Medical 1 (M1) **h** Medical 2 (M2) **i** Medical 3 (M3)

**Table 1** Estimated measures of proposed LSST

| Image name | MAE | MSE | RMSE | RSNR | SNR | PSNR |
|---|---|---|---|---|---|---|
| H1 | 5.9008 | 43.2550 | 6.5769 | 0.4564 | 0.2083 | 15.4101 |
| H2 | 6.5105 | 50.7387 | 7.1231 | 0.4214 | 0.1776 | 14.0240 |
| H3 | 5.9641 | 53.3017 | 7.3008 | 0.4112 | 0.1691 | 13.5960 |
| T1 | 5.6370 | 43.6508 | 6.6069 | 0.4544 | 0.2064 | 15.3310 |
| T2 | 5.4448 | 35.2065 | 5.9335 | 0.5059 | 0.2560 | 17.1983 |
| T3 | 7.9559 | 70.1532 | 8.3758 | 0.3584 | 0.1285 | 11.2099 |
| M1 | 6.6286 | 51.2655 | 7.1600 | 0.4193 | 0.1758 | 13.9343 |
| M2 | 3.8485 | 40.5868 | 6.3708 | 0.4712 | 0.2220 | 15.9631 |
| M3 | 7.5072 | 79.9823 | 8.9428 | 0.3357 | 0.1127 | 10.0709 |

# References

1. Kaushik, A., Kumar, A., Barnela, M.: Block encryption standard for transfer of data. In: International Conference on Networking and Information Technology, pp. 381–385 (2010)
2. Fuqiang, D., Minqing, Z., Jia, L.: Virus-evolutionary genetic algorithm based selective ensemble for steganalysis. In: Ninth International Conference on Parallel, Grid, Cloud and Internet Computing (3PGCIC), pp. 553–558 (2014)
3. Daneshkhah, A., Aghaeinia, H., Seyedi, S.H.: A more secure steganography method in spatial domain. In: Second International Conference on Intelligent Systems, Modelling and Simulation, pp. 189–194 (2011)
4. Dagar, E., Dagar, S.: LSB based image steganography using X-Box mapping. In: International Conference on Advances in Computing, Communications and Informatics, pp. 351–355 (2014)
5. Al-Hadidi, B.: A novel data hiding algorithm for all types of file based on Injection process. In: International Conference on Control, Decision and Information Technologies (CoDIT), pp. 872–875 (2013)
6. Bobate, R.V., Khobragade, A.S.: Optimal implementation of digital steganography in an true color images for the secrete communication. In: 3rd International Conference on Advances in Recent Technologies in Communication and Computing (ARTCom 2011), pp. 91–95 (2011)
7. Singla, D., Juneja, M.: An analysis of edge based image steganography techniques in spatial domain. In: Recent Advances in Engineering and Computational Sciences (RAECS), pp. 1–5 (2014)
8. Diouf, B., Diop, I., Keita, K.W., Farssi, S.M., Khouma, O., Diouf, M., Tall, K.: Adaptive linear programming of polar codes to minimize additive distortion in steganography. In: SAI Computing Conference (SAI), pp. 1086–1092 (2016)
9. Guan, Q., Dong, J., Tan, T.: Evaluation of feature sets for steganalysis of JPEG image. In: IEEE Youth Conference on Information Computing and Telecommunications (YC-ICT), pp. 359–362 (2010)
10. Lee, C.-F., Li, K.-T.: VQ-based image embedding scheme using adaptive codeword grouping strategy. In: 5th International Conference on Future Information Technology, pp. 1–6 (2010)
11. Wang, H., Fei, R.: Hiding data in indexed images. In: IEEE International Conference on Software Engineering and Service Science (ICSESS), pp. 833–836 (2013)
12. Chen, G., Cao, M., Fu, D., Ma, Q.: Research on an steganographic algorithm based on image edge. In: International Conference on Internet Technology and Applications (iTAP), pp. 1–4 (2011)

13. Vegh, L., Miclea, L.: Securing communication in cyber-physical systems using steganography and cryptography. In: 10th International Conference on Communications (COMM), pp. 1–4 (2014)
14. Kumar, B.V., Mani, M.R., NesaKumari, G.R., Kumar, D.V.V.: A new marginal color image watermarking method based on logical operators. Int. J. Secur. Appl. **3**(4), 1–8 (2009)

# Efficient System for Color Logo Recognition Based on Self-Organizing Map and Relevance Feedback Technique

**Latika Pinjarkar, Manisha Sharma and Smita Selot**

**Abstract** This paper proposes an efficient system for color logo recognition. The proposed logo retrieval approach uses self-organizing map (SOM) and relevance feedback based on three types of query improvement policies: new query point, query rewriting, and query development. Feature extraction techniques are implemented for representing the visual features of the images from the data set as well as the query image. For color feature extraction, color histogram, color moments, and color correlogram techniques are used. For texture feature extraction, Gabor wavelet and Haar wavelet are implemented. And for shape feature extraction, Fourier descriptor and circularity features are used. The topological mapping property of the SOM is used to map these database features. This map contains the comparable images arranged closer to each other. The data set consists of about 3000 color logo images. Euclidian distance is used for similarity computation. The proposed approach is able to recognize the logo images with good retrieval efficiency. The results obtained are remarkable after the integration of SOM with relevance feedback technique.

**Keywords** Content-based image retrieval (CBIR) · Feature vector
Logo · Relevance feedback · Self-organizing map (SOM)

## 1 Introduction

There is very fast boost in the number of registered logo images around the globe because logos are playing very important role in any business or organization nowadays. Therefore, it is very important to design an efficient logo retrieval system for its recognition to avoid duplicate logos with the already registered logos stored

L. Pinjarkar (✉) · S. Selot
Shri Shankaracharya Technical Campus, Bhilai, CG 490020, India
e-mail: latikabhorkar@gmail.com

M. Sharma
Bhilai Institute of Technology, Bhilai, CG 490020, India

in the database. Logo image retrieval is an important application of content-based image retrieval (CBIR) framework. A logo can be differentiated into three categories as follows: (1) logo having just text, (2) logo having just symbol, and (3) logo having combination of both text and symbol. A number of logo image recognition frameworks have been proposed to deal with all categories of logo images. However, the recognition results of these frameworks are poor when judged against the frameworks that are mainly intended to deal only one category of logo [1]. The proposed system is designed based on self-organizing map and relevance feedback technique to handle all the three types of logos.

The self-organizing map is an unsupervised learning method. Many data analysis tasks make use of this technique. The nonlinear mapping of a high-dimensional input space into two-dimensional grid of artificial neural units is the important property of self-organizing map. During the training phase of SOM, a topological map is formed which involves the mapping of input vectors near to each other in close by map units. In this way, the samples situated on SOM which are similarly comparable in regard to the given feature extraction plan are found close to each other [2].

Relevance feedback (RF) is the feedback taken from the user regarding the relevance of the recognized images by the CBIR system as an output [3]. The data associated with this feedback is used for improvement of the query. The query improvement process may help to enhance the results of the image retrieval framework in terms of image recognition. A number of RF-based CBIR systems are proposed to obtain better results. But still these systems have some problems like redundant browsing and exploration convergence. Redundant browsing means many of the RF schemes concentrate upon getting the user's contentment in one query scheme. The query is developed repeatedly by exploring the information associated with the accurate retrieved images marked by the user as relevant.

To solve the above problem, the navigation pattern mining method and the query improvement approach are applied to the conventional logo retrieval system (based on visual feature extraction). The navigation samples mined from the user query log are used as the smallest route to attain the attention of the user. It is used as a best searching route to meet the searching room nearer to the user's attention efficiently. The query improvement process results in more accurate relevant images as output. In this way, the difficulty of redundant browsing can be resolved using this mining process. Three query improvement approaches, new query point (NQP), query rewriting (QRW), and query development (QDE), are implemented to handle the difficulty of exploration convergence [4].

In this paper, a system for color logo retrieval based on shape, color, and texture features by using self-organizing map and relevance feedback technique is proposed. The organization of the paper is as follows. Related work in this area is discussed in Sect. 2. The proposed approach and details about data set used are discussed in Sect. 3. Section 4 deals with the experimentation and results obtained. Finally, conclusions are drawn in Sect. 5.

## 2 Related Work

Laaksonen et al. [2] have proposed self-organizing maps (SOMs) as a relevant feedback mechanism for CBIR. They have proposed the PicSOM CBIR scheme which illustrates the applicability of SOMs as a relevance feedback method. An exclusive tree-structured SOM for every individual feature kind was present in the PicSOM. Laaksonen et al. [5] have suggested a scheme intended for CBIR in big data sets. The basis of the framework was tree-structured self-organizing maps (TS-SOMs) called as PicSOM. A group of sample images were given to this system as a reference. The PicSOM was able to retrieve the group of images similar to these given samples. The image feature description such as color, texture, or shape was used to form TS-SOM. The relevance feedback technique was implemented for taking feedback from the user upon the preference of the retrieved images. The query images were given by means of the World Wide Web and were redefined continuously. Suganthan [6] has suggested a shape indexing method using self-organizing map. Pairwise relational attribute vectors were used to extract the structural information contained in the geometric shape. These vectors were quantized using a SOM. Two SOMs were used in this proposed work named as SOM1 and SOM2. Global histogram of relational attribute vectors was contained in SOM1. These histograms were given as input vectors to SOM2. SOM1 was used to confine the shape properties of the objects, and topology conserving mapping for the structural shapes was generated by SOM2. The database used was from the UK TradeMark Registry office containing the trademark image database of over 10,000 images. Testing was done upon a subset of this database, having 990 trademark images. Pakkanen and Iivarinen [7] have examined the use of the self-organizing map (SOM) as a means for CBIR. The experimentation was done on the data set of 1300 images and the MPEG-7 features. The proposed system performs with good retrieval results. Vesanto et al. [8] have used the SOM as a vector quantization technique which puts the sample vectors on a regular low-dimensional grid in a sorted style referred as SOM Toolbox. The article covers the SOM Toolbox and its usage. The performance of the system was evaluated with the means of computational load and was evaluated against a similar C program. Hussain and Eakins [9] have presented an approach for visual clustering of multi-component images of trademarks, based on the topological characteristics of the SOM. The proposed approach was implemented in two stages. In the first step, the features extracted from image components were used to construct a 2-D map. In the second step component, similarity vector from a query image was derived. A 2-D map of retrieved images was derived using this. The experimentation was done using a database of 10,000 trademark images, and the performance measure was precision and recall. The results prove that the component-based matching procedure achieves improved matching as compared with using whole-image clustering. The proposed system was not sensitive to alterations in input parameters like network size. Alnihoud [10] has proposed a CBIR system based on SOM. Fuzzy color histogram and subtractive fuzzy color clustering algorithms were used for marking

the class of the query image. The results proved that the proposed approach has reduced the computational complexity of the CBIR system. Huneiti and Daoud [11] have suggested a CBIR method by extracting both color and texture feature vectors using the discrete wavelet transform (DWT) and the SOM. Euclidean distance was utilized for computing the resemblance of the query image with the data set images using the Wang database. The work proposed by the above researchers is not based upon multiple features extraction. Also, none of them have used the SOM technique integrated with RF for logo image retrieval. The proposed approach is based upon multiple features extraction and SOM integrated with RF for logo image retrieval, making the framework very efficient.

## 3   Proposed Approach

The experimentation was done using an image database of 3000 color logo images. Some sample logo images from the database are as represented in Fig. 1. These images are saved in JPG format, and each resized to the size 256 × 256. The visual description of the image is depicted as far as its low-level features, i.e., color, texture, and shape [12]. The features of all the logo images in the data set are extracted by implementing color feature extraction, texture feature extraction, and shape feature extraction. Color histogram, color moments, and color correlogram techniques are used for color feature extraction; Gabor wavelet and Haar wavelet are implemented for texture feature extraction and for shape feature extraction; Fourier Descriptor and circularity features are implemented. These feature vectors are utilized as input training data to prepare the SOM. The feature vectors are fed to the SOM as sample data, and the best matching unit (BMU) is determined, and finally the neighbors are scaled using this BMU.

When the system is given a query image, it searches for the similar images in the map generated and retrieves the similar logo images present in the data set. Then the feedback is taken from the user about relevancy of the retrieved images (relevance feedback). The information from the user feedback is used to improve the query. Performance measures used for the assessment of the system are precision and recall.



**Fig. 1**  Typical logo images in the database

The proposed approach is implemented as per the following steps:

Step 1: Preprocessing of the data set images.

  (a) Extract the color feature by:
      (i) color histogram, (ii) color correlogram, (iii) color moments.
  (b) Extract the texture feature by:
      (i) Gabor wavelet, (ii) Haar wavelet.
  (c) Extract the shape feature by:
      (i) Fourier descriptor, (ii) circularity feature.

Step 2: Training the self-organizing map using these feature vectors.
Step 3: Querying an image to the system.
Step 4: Feature extraction of the query image (color feature, texture feature, and shape feature extraction).
Step 5: Retrieval of images from the data set relevant to the query image.
Step 6: Taking feedback from the user whether the retrieved images are relevant or not relevant? (relevance feedback).
Step 7: Categorize into positive ($P$) and negative images ($N$).
Step 8: Determine modified query point (new$_{qp}$)
Step 9: Retrieval of the images from the database as final output images using this new query point.

The following section describes the implementation details of SOM and RF techniques.

## 3.1 Self-Organizing Map (SOM)

The SOM defines the mapping from the input data set $D^n$ onto a standard two-dimensional array of nodes (map network). Every node $j$ in the map is associated with a parametric reference vector $n_j \in D^n$. In the proposed approach, the array of nodes is projected onto a rectangular lattice. For mapping the input, every input vector $y$ is evaluated against the ratio $n_j : s$, and the top equivalent is derived. Every input vector $y \in D^n$ is compared with the $n_j : s$ using the Euclidean distance. The winning node $d$ is calculated using the following equation:

$$\|y - n_d\| = \min\{\|y - n_j\|\} \tag{1}$$

where $y$ is mapped onto $d$ relative to the parameter values $n_j$. Nodes which are topographically near to another in the array learn from the same input. The update formula is given as:

$$n_j(t+1) = n_j(t) + h_{d,j}(t)[y(t) - n_j(t)] \qquad (2)$$

where $t$ is the discrete-time coordinate and $h_{d,j}$ is the region defining function. The early values of $n_j : s$ are random [13].

## 3.2  Relevance Feedback Implementation

The various data structures needed for the implementation of the proposed relevance feedback approach are described as following:

1. Unique record table—contains query image name, iteration, query point, and relevant image name
2. Record table—contains query image name, cluster number
3. Query position table—contains query image name and query point
4. Navigation operation table—contains query image name, iteration number, and cluster number (e.g., if query q is present in cluster 1 of first iteration, then the entry will be c11)

   This table is constructed as following:

(a) The retrieval results of every iteration are clustered,
(b) The query points are clustered,
(c) The cluster Ids called as item set are present in this table, and
(d) These clusters are traversed for constructing a pattern for every query. For example given a query image, if the retrieved images are available in cluster 2 of iteration 1, cluster 1 of iteration 2, and cluster 2 of iteration 3, then the pattern is c21, c12, and c23.

   These sequential patterns constructed are used for the mining process, and the frequently occurred patterns are mined.

5. Record partition table—contains query point and relevant image name

   The proposed logo retrieval system performs the iterative search procedure as per the following steps:

   Step 1: The new query point is determined by taking the average of features of positive images.
   Step 2: The closest query seeds (root) are determined to get the matching sequential patterns.
   Step 3: These matching navigation pattern trees are used to determine the closest leaf nodes.
   Step 4: Then the top '$r$' relevant query points from the group of the closest leaf nodes are determined.
   Step 5: The top '$t$' relevant images are retrieved and given as output to the user.

In the proposed methodology for NQP generation, the step 1 is implemented. Steps 2–5 are followed for QDE. The new feature weights are determined using the features of positive images given by user at every feedback for executing QRW process [4].

The details of the NQP generation, QRW, and QDE are described as follows:

### New Query Point Generation (NQP)

Suppose in the preceding feedback the images retrieved by the query point denoted by $old_{qp}$. A new query point $new_{qp}$ is determined by taking the average of the features of the positive images. $P$ is given by the user as a feedback. Suppose the positive images are given as $P = \{p_1, p_2, \ldots, p_k\}$ and $m$ dimensions of the $j$th feature $R_j = \{r_1^Y, r_2^Y \ldots r_m^Y\}$ extracted from the $Y$th positive image. Then the new query point $new_{qp}$ indicated by $P$ can be given as [4]:

$$new_{qp} = \{\overline{R_1}, \overline{R_2}, \ldots, \overline{R_c}\} \quad \text{where } 1 \leq j \leq c \tag{3}$$

$$\overline{R_j} = \{\overline{r_1}, \overline{r_2}, \ldots, \overline{r_m}\} \text{ and } \overline{r_t} = \frac{\sum 1 \leq y \leq s, r_t^y \in R_j r_t^y}{s}$$

where $new_{qp}$ and the positive images are saved into the unique record table.

### Query Rewriting (QRW)

Suppose a set of positive images is given as $P = \{p_1, p_2, \ldots, p_k\}$ determined by the old query point $old_{qp}$ in the earlier feedback. The new weight of the $j$th feature $R_j$ is given as [4]:

$$T_j = \frac{\dfrac{\sum_{x=1}^{a} \alpha_x}{\alpha_j}}{\sum_{y=1}^{a} \dfrac{\sum_{x=1}^{a} \alpha_x}{\alpha_y}} \tag{4}$$

where

$$\alpha = \sum_{z=1}^{m} \frac{\sqrt{\sum_{i=1}^{b} \left(r_i^z - r_i^{old_{qp}}\right)^2}}{b} \tag{5}$$

and $1 \leq j \leq a$.

**Query Development (QDE)**

The weighted KNN search is done by QDE method: first the closest query seed to every of *P* is determined, called positive query seed, and the closest query seed to each of *N*, called negative query seed [4].

There may be some query seeds present in both positive query seed group and the negative query seed group. To handle this situation, a token pn.chk is assigned to every seed. If the seed is having greater number of negative examples than the positive examples, then pn.chk = 0 otherwise pn.chk = 1.

After this, the relevant query seeds are determined. The navigation pattern tree is traversed to find a set of matching leaf nodes. The new feature weights as calculated in Eq. (4) are used in the search procedure to find the required images.

The search procedure is classified into two steps.

In the first step, the relevant visual query points are generated, and in second step, the relevant images are determined.

The query position table is used to determine the related query points. The top related query points are determined using the KNN search method. Top '*r*' similar query points are determined in step 1. The top '*t*' images which are closer to $new_{qp}$ are retrieved as a result.

# 4 Experiments and Results

The experimentation is done using an image database of 3000 color logo images. Feature vector's database is formed by extracting color, texture, and shape features of all 3000 images. The SOM is trained using these feature vectors. The experiment includes giving a query image, extracting features of the image, retrieving relevant logo images from the data set as output, taking feedback from the user about the relevancy of these output images, and finally retrieving the relevant images based on the feedback given by the user and new query point.

The execution of the proposed framework is assessed as far as standard assessment parameters precision and recall which can be defined as follows:

$$\text{Precision} = \frac{\text{Number of retrieved relevant images}}{\text{Total number of retrieved images}}$$

$$\text{Recall} = \frac{\text{Number of retrieved relevant images}}{\text{Total number of relevant images in the database}}$$

Results in terms of precision and recall for sample 10 query images are depicted in Table 1. Table 2 provides the comparison of the proposed approach with Su et al. [4].

The performance of the proposed approach has been improved after integration of SOM with the RF technique. The precision value is improved considerably as

**Table 1**  Results in terms of precision and recall for sample 10 query images

| S no. | Query image | Precision | Recall | S no. | Query image | Precision | Recall |
|---|---|---|---|---|---|---|---|
| 1 | | 0.9567 | 0.7121 | 6 | ŚKODA | 0.9481 | 0.696 |
| 2 | Sun microsystems | 0.9528 | 0.7712 | 7 | Microsoft | 0.9579 | 0.7987 |
| 3 | Infosys® | 0.9586 | 0.7267 | 8 | Disney | 0.9641 | 0.7321 |
| 4 | hike | 0.9745 | 0.8101 | 9 | Signature | 0.9698 | 0.7648 |
| 5 | iGATE | 0.9623 | 0.7534 | 10 | adidas | 0.9598 | 0.7251 |
| Average | | Precision: 0.96037 | | | Recall: 0.74902 | | |

**Table 2**  Comparison for Su et al. [4] with the proposed approach

| S no. | Method | Parameter | | | |
|---|---|---|---|---|---|
| | | Data set used | Precision | Execution time | |
| 1 | Su et al. (using only RF) | Seven classes of different categories of 200 images each (1400 images) | 0.910 | 1.184667 s (for each class of 200 images) | |
| 2 | Proposed approach (using SOM + RF) | 3000 images (of logos) | 0.9603 | 1.2 min (for total 3000 images) | |

depicted in Table 2, though the improvement is needed further for reducing the execution time required by the framework.

## 5  Conclusion

In the proposed approach, an efficient system for color logo retrieval is proposed based on self-organizing map and relevance feedback technique. The color, texture, and shape features of the logo images are extracted and are used as sample data to train the SOM. The relevance feedback mechanism helps to improve the query. The performance of the system is measured using precision and recall. The results obtained are remarkable and proved that the performance of the system has been improved after integrating SOM with relevance feedback technique. In the future

scope of the work feature, optimization technique can be combined with proposed system to identify the best features from the database feature set which can further reduce the execution time of the retrieval process.

# References

1. Wei, C.-H., Li, Y., Chau, W.-Y., Li, C.-T.: Trademark image retrieval using synthetic features for describing global shape and interior structure. Pattern Recogn. **42**, 386–394 (2009)
2. Laaksonen, J., Koskela, M., Laakso, S., Oja, E.: Self-Organising Maps as a Relevance Feedback Technique in Content-Based Image Retrieval, No. 2–3, pp. 140–152. Springer, Berlin (2001)
3. Pinjarkar, L., Sharma, M., Mehta, K.: Comparative evaluation of image retrieval algorithms using relevance feedback and it's applications. Int. J. Comput. Appl. (IJCA) **48**(18), 12–16 (2012)
4. Ja-Hwung, S., Huang, W.-J., Yu, P.S., Tseng, V.S.: Efficient relevance feedback for content-based image retrieval by mining user navigation patterns. IEEE Trans. Knowl. Data Eng. **23**(3), 360–372 (2011)
5. Laaksonen, J., Koskela, M., Laakso, S., Oja, E.: PicSOM -content-based image retrieval with self-organizing maps. Pattern Recogn. Lett. **21**, 1199–1207 (2000)
6. Suganthan, P.N.: Shape indexing using self-organizing maps. IEEE Trans. Neural Netw. **13** (4), 835–840 (2002)
7. Pakkanen, J., Iivarinen, J.: Evaluating SOM as an Index in Content-Based Image Retrieval. Helsinki University of Technology Laboratory of Computer and Information Science, Helsinki (2003)
8. Vesanto, J., Himberg, J., Alhoniemi, E., Parhankangas, J.: Self-organizing map in MATLAB: the SOM Toolbox. In: Proceedings of the MATLAB DSP Conference Nov 16–17, Espoo, Finland, pp. 35–40 (1999)
9. Hussain, M., Eakins, J.P.: Component-based visual clustering using the self-organizing. Neural Netw. **20**, 260–273 (2007)
10. Alnihaud, J.: Content based image retrieval system based on self organizing map, fuzzy color histogram and subtractive fuzzy clustering. Int. Arab J. Inf. Technol. **9**(5), 452–458 (2012)
11. Huneiti, A., Daoud, M.: Content-based image retrieval using SOM and DWT. J. Softw. Eng. Appl. **8**, 51–61 (2015)
12. Long, F., Zhang, H., Feng, D.D.: Multimedia information retrieval and management. In: Chapter-1 Fundamentals of Content-Based Image Retrieval, pp. 1–26. Springer, Berlin (2003)
13. Kohonen, T.: Self Organinzing Maps, 2nd edn. Springer, New York (1997)

# Congestion-Aware Opportunistic Routing Protocol in Wireless Sensor Networks

**Maya Shelke, Akshay Malhotra and Parikshit N. Mahalle**

**Abstract** It is expected that 50 billion devices in the world will be connected on IOT by 2025. The importance of wireless sensor networks cannot be overstated in this scenario. Network becomes more beneficial to an application when it can be used to its full potential, which is difficult to achieve because of limitations of resources (processor, memory, and energy). There are many existing routing mechanisms which deal with these issues by reducing number of transmissions between sensor nodes by choosing appropriate path toward base station. In this paper, we propose a routing protocol to select the optimized route by using opportunistic theory and by incorporating appropriate sleep scheduling mechanisms into it. This protocol focuses on reduction of congestion in the network and thus increases an individual node's life, the entire network lifetime, and reduces partitioning in the network.

**Keywords** Wireless sensor networks · Opportunistic routing protocol
Congestion control · Sleep scheduling mechanisms

M. Shelke (✉) · A. Malhotra
Symbiosis Institute of Technology (SIT) Affiliated to Symbiosis International University
(SIU), Pune, India
e-mail: mayabembde07@gmail.com

A. Malhotra
e-mail: dydirector@sitpune.edu.in

P. N. Mahalle
Smt. Kashibai Navale College of Engineering Affiliated to Savitribai Phule Pune University,
Pune, India
e-mail: AALBORG.PNM@gmail.com

# 1   Introduction

Wireless sensor networks (WSNs) are an essential technology which has the potential to change our lives. Hence, WSN is considered as one of the very important elements of Internet of Things. WSN holds huge count of applications where human approach cannot reach easily. Recent advances in technology have made manufacturing of tiny and low-cost sensors technically and economically practical. WSN comprises hundreds or thousands of sensor nodes. These tiny sensor nodes are capable of sensing, processing, and communicating environment or application parameters to the base station. These nodes are heavily deployed either inside the device or in its perimeter.

Sensor nodes can be deployed using random or deterministic approach [1]. In random deployment approach, nodes are deployed in inaccessible terrains or disaster relief operation with no fixed topology. In deterministic approach, nodes are deployed with a fixed topology. WSNs are used in wide variety of application areas such as military, ambient monitoring, weather monitoring, security, inventory control, disaster management, health, forest fire detection. In all these applications, the sensor nodes sense the required parameters and transmit them to the sink node. Sensor nodes can communicate with each other and directly or indirectly with the base station. Communication of data from sensor node to base station requires discovery of efficient path known as routing. Routing is a very challenging task in WSN because of the following reasons [2, 3]:

1. Nodes are deployed in unattended areas and in ad hoc manner in many applications.
2. Large number of sensor nodes make use of global addressing scheme infeasible.
3. Sensor nodes are constricted in terms of processing power, memory, and battery.
4. Existing IP-based protocols may not be easily extended and used in WSNs.

Hence there is need of an efficient routing protocol.

Opportunistic routing (OR) has recently drawn attention of researchers as it has proved to enhance the performance of wireless ad hoc and sensor networks over traditional routing protocols that perform best path routing [4, 5]. This protocol selects one or more fixed routes in advance from source to base station. This strategy does not work in dynamic WSNs because of unreliable links. Unreliable links cause frequent transmission failures which require retransmission of data causing inefficient use of network resources or even a system breakdown [6]. Opportunistic routing takes advantage of broadcast nature of WSN by involving all nodes overhearing communication into the routing process. It helps to build routes dynamically and makes optimal use of energy.

WSN is an indivisible part of IOT applications [2, 7]. Hence, huge and continuous amount of data generation poses a big challenge in the form of congestion which adversely affects application performance. Congestion is the situation in the network where incoming traffic is more than what the network can handle. Congestion degrades performance of a network drastically. Various factors that are

responsible for congestion include insufficient memory, bursty traffic, slow processor, fast packet arrival rate, and congestion feeds itself. It is necessary to take preventive steps to avoid congestion because once it occurs it reinforces itself to worsen the problem.

To further reduce energy consumption in WSN, sleep scheduling mechanisms can be used [8, 9]. Sensor nodes consume more energy in idle stage than transmitting or receiving stage. Hence, sleep scheduling mechanism intelligently makes some nodes to sleep, while others awake to avoid performance from degrading. Sleep mode switches off radio for some time period to reduce idle listing period.

To satisfy the purpose of IOT, huge number of sensors are required to be deployed for a particular application to provide accurate information. A system is useful only if it provides correct data in time. It is possible to achieve this goal by incorporating context awareness into routing process. It is not always possible to provide explicit information to nodes in WSN; hence, it is preferred if nodes can acquire the required information by themselves. System is called context-aware if it can understand the situation and accordingly take automatic decisions dynamically. Context can be defined as any information that can be used to illustrate the situation of entities, for example, a person, place, or object that are considered relevant to the interaction between a user and an application, including the user and the application themselves [10, 11].

## 2 Overview of Opportunistic Routing Protocols

Energy efficiency is still the most crucial issue in WSN [12, 13]. A number of routing techniques have been proposed to deal with this constraint. Most of the energy-aware protocols do not consider lossy nature of wireless links [14]. Multipath routing protocols have been put forward to counteract the unreliability of wireless channels. All these protocols follow a traditional design principle in network layer of wired networks. In wireless networks, selecting good paths in advance does not work because of characteristics such as interference, fading, and multipath effects. These factors cause heavy packet loss and make it difficult for these protocols to achieve stable performance.

Liu et al. [15] provide insights into OR protocols. This paper compares OR with conventional routing protocols in regard to energy consumption and throughput. The key to success of OR protocols is selection of appropriate metric and candidate selection algorithm. Stepwise working of OR protocol is shown in Fig. 1.

OR protocols can be classified into different types based on routing efficiency; use of network state information, compatibility with existing MAC protocols, protocol overhead, location information, and use of coding function are shown in Fig. 2.

Opportunistic protocols such as ExOR, GeRaF, and EQGOR address the above-mentioned challenge by selecting routes on the fly. OR protocols take advantage of broadcast nature of WSN by including all nodes that overhear

**Fig. 1** Stepwise opportunistic routing process



**Fig. 2** Classification of OR protocols

transmission into routing (forwarding) process. This new concept has an issue that multiple nodes hear a packet and unnecessarily forward the same packet. We compare OR protocols in Table 1.

**Table 1** Comparison of OR protocols

| Protocol name | Parameters to calculate route | Candidate coordination method | Evaluation parameters | Limitations |
|---|---|---|---|---|
| ExOR [16] | ETX 1 Number of transmissions 2 Number of hops 3 Internode packet loss, S/N ratio | Ack-based | 1 Distance per transmission 2 Batch size 3 Throughput | 1 For wireless mesh networks not for WSN 2 Concentrates on transmission of large files than energy consumption |
| GeRaF [17] | 1 Average energy consumed 2 Geographic distance of each node to destination | RTS/CTS-based | 1 Energy consumed 2 Latency 3 Number of nodes in the network | 1 Assumed constant traffic model and high density 2 More communication overhead |
| EEOR [6] | 1 Energy of nodes 2 Link error probability | Ack-based | 1 Energy consumption 2 Packet loss ratio 3 End-to-end delay 4 Packet duplication ratio | 1 Agreement cost can be calculated in more realistic way instead of number of acknowledgements sent 2 Sleep Scheduling modes are not considered |
| ENS_OR [18] | 1 Optimum distance of sensor node to sink 2 Residual energy of each node | Ack-based | 1 Receiving packets ratio 2 Network lifetime 3 Avg. of residual energy 4 First dead node | 1 Interference of signals is not considered 2 Sleep scheduling mechanism is not considered |

## 3 Proposed Congestion-Aware Opportunistic Routing Protocol

Hereby, we propose congestion-aware opportunistic routing mechanism which aims to increase network lifetime and reduce energy consumption by controlling congestion in the network and incorporating sleep scheduling mechanism into it. We attempt to integrate congestion control and routing mechanism together which will automatically cut down processing power and energy consumption. Stepwise working of proposed mechanism is illustrated in Fig. 3.

**Fig. 3** Congestion-aware opportunistic routing protocol

This algorithm unlike the works presented in Table 1 considers opportunistic routing mechanism, congestion control, and sleep scheduling modes together.

Residual energy, distance, density, traffic, and link quality are considered as metrics for calculating cost of relay node.

According to [19], distance can be calculated using

- Time of flight [TOF]:

$$\Delta t_{\text{RTT}} - \Delta t_{\text{latency}}$$

- Received signal strength:

$$r = \lambda/4\pi . \sqrt{P_t/Pr . G_t . G_r}.$$

$r$ is equivalent distance of sender and receiver
$Pr$ is maximum possible power in receiving node
$P_t$ is maximum transmitting power

These parameters can be measured in terms of each other and can be correlated using regression correlation techniques.

## 3.1 Congestion Control Techniques

Each node in the network has two types of traffic. One is generated by the node itself, and second is transit traffic which is passing through this node. So traffic of an individual node can be calculated as [20, 21]

$$T(n_i) = t(n_{ai}) + t(n_{i-1})$$

where $T(n_i)$ is total traffic at node $n_i$, $t(n_{ai})$ is the traffic generated by the application layer of node $n_i$, $t(n_{i-1})$ is the traffic coming from previous node.

Congestion threshold is calculated by using packet inter-arrival time and packet service time. Packet inter-arrival time is the time interval between arrivals of two packets at MAC layer. Whereas packet service time is the time taken by an algorithm to send packet out of node which includes time period for which packet resides in buffer, time for collision resolution, and time at which last bit is successfully sent out.

Congestion threshold can be calculated as

$$C(i) = \text{Pkt\_Ser\_time}/\text{Pkt\_IntArr\_time}$$

$$\text{Pkt\_Ser\_time} = \text{Pkt\_wait\_time} + \text{Collision\_Resolution\_time} + \text{Pkt\_Trans\_time}$$

Rate adjustment algorithm is called when traffic at a particular node crosses congestion threshold. This algorithm regulates rate of packets from network layer to MAC layer (Fig. 4).

## 3.2 Sleep Scheduling Mechanism

Sensor node consumes more energy in idle phase than transmitting or receiving phase. Considerable amount of energy can be saved if large number of nodes are made to sleep [8, 18].

Assume node remains in idle state for $T$ time period and spends $e_{id}$ energy per second. Therefore, total energy spent during idle state is

$$E_{id} = e_{id}T.$$

**Fig. 4** Congestion control
mechanism [21]



If nodes are made to sleep intelligently, then $E_{id}$ amount of energy can be saved per node without affecting network performance (QOS parameters). Hence, efficient and intelligent sleep scheduling mechanisms are desirable in WSN applications.

## 4  Conclusion

Hereby, we propose congestion-aware opportunistic routing protocol in WSN. Increase in IOT applications is giving rise to strong need for congestion control mechanisms to reduce traffic in the network to achieve stable performance. These mechanisms are required in WSN as well as at the interface of IOT. It is impractical to implement heavy mechanisms in WSN; hence, this is an area of considerable research.

Here, we stated current state of the art of opportunistic protocols, congestion control mechanisms, and sleep scheduling modes along with advantages and limitations of few of them. Further incorporation of sleep scheduling mechanisms along with proposed idea will minimize energy consumption and delay while maximizing network lifetime and packet delivery ratio. We plan to validate results of proposed idea on simulator and implement it on real-time network.

# References

1. Akyildiz, I.F., Weilian, S., Sankarasubramaniam, Y., Cayirci, E.: Wireless sensor networks: a survey. Comput. Netw. **38**(4), 393–422 (2002)
2. Gubbi, J., Buyya, R., Marusic, S., Palaniswami, M.: Internet of Things (IoT): a vision, architectural elements, and future directions. Future Gener. Comput. Syst. **29**(7), 1645–1660 (2013)
3. Al-Karaki, J.N., Kamal, A.E.: Routing techniques in wireless sensor networks: a survey. IEEE Wirel. Commun. **11**(6), 6–28 (2004)
4. Liu, D., et al.: Duplicate detectable opportunistic forwarding in duty-cycled wireless sensor networks. IEEE/ACM Trans. Netw. **24**(2), 662–673 (2016)
5. Kumar, N., Singh, Y.: An energy efficient opportunistic routing metric for wireless sensor networks. Ind. J. Sci. Technol. **9**(32), 1–5 (2016)
6. Mao, X., Tang, S., Xu, X., Li, X.-Y., Ma, H.: Energy-efficient opportunistic routing in wireless sensor networks. IEEE Trans. Parallel Distrib. Syst. **22**(11), 1934–1942 (2011)
7. Singh, D., Tripathi, G., Jara, A.J.: A survey of Internet-of-Things: future vision, architecture, challenges and services. In: 2014 IEEE World Forum on Internet of Things (WF-IoT), pp. 287–292. IEEE, Washington (2014)
8. Baba, S.B., Mohan Rao, K.R.R.: Improving the network life time of a wireless sensor network using the integration of progressive sleep scheduling algorithm with opportunistic routing protocol. Indian J. Sci. Technol. **9**(17), 1–6 (2016)
9. Zhang, Z., et al.: Energy-efficient and low-delay scheduling scheme for low power wireless sensor network with real-time data flows. Int. J. Ad Hoc Ubiquitous Comput. **22**(3), 174–187 (2016)
10. Dey, A.K., Abowd, G.D., Salber, D.: A conceptual framework and a toolkit for supporting the rapid prototyping of context-aware applications. Hum. Comput. Interact. **16**(2), 97–166 (2001)
11. Perera, C., Zaslavsky, A., Christen, P., Georgakopoulos, D.: Context aware computing for the internet of things: a survey. Commun. Surv. Tutor. IEEE **16**(1), 414–454 (2014)
12. Shelke, M., et al.: Fuzzy-based fault-tolerant low-energy adaptive clustering hierarchy routing protocol for wireless sensor network. Int. J. Wirel. Mob. Comput. **11**(2), 117–123 (2016)
13. Yao, Y., Cao, Q., Vasilakos, A.V.: EDAL: an energy-efficient, delay-aware, and lifetime-balancing data collection protocol for heterogeneous wireless sensor networks. IEEE/ACM Trans. Netw. **23**(3), 810–823 (2015)
14. Chachulski, S., et al.: trading structure for randomness in wireless opportunistic routing vol. 37(4). ACM, New York (2007)
15. Liu, H., Zhang, B., Mouftah, H.T., Shen, X., Ma, J.: Opportunistic routing for wireless ad hoc and sensor networks: present and future directions. IEEE Commun. Mag. **47**(12), 103–109 (2009)
16. Biswas, S., Morris, R.: ExOR: opportunistic multi-hop routing for wireless networks. In: ACM SIGCOMM Computer Communication Review vol. 35(4), pp. 133–144. ACM, New York (2005)
17. Zorzi, M., Rao, R.R.: Geographic random forwarding (GeRaF) for ad hoc and sensor networks: energy and latency performance. IEEE Trans. Mob. Comput. **2**(4), 349–365 (2003)
18. Luo, J., Hu, J., Wu, D., Li, R.: Opportunistic routing algorithm for relay node selection in wireless sensor networks. IEEE Trans. Ind. Inform. **11**(1), 112–121 (2015)
19. Jörger, T., Höflinger, F., Gamm, G.U., Reindl, L.M.: Wireless distance estimation with low-power standard components in wireless sensor nodes. arXiv preprint arXiv:1601.07444 (2016)
20. Kaur, J., Grewal, R., Singh Saini, K.: A survey on recent congestion control schemes in wireless sensor network. In: Advance Computing Conference (IACC), 2015 IEEE International, pp. 387–392. IEEE, Washington (2015)

21. Wang, C., Li, B., Sohraby, K., Daneshmand, M., Hu, Y.: Upstream congestion control in wireless sensor networks through cross-layer optimization. IEEE J. Sel. Areas Commun. **25** (4), 786–795 (2007)

# A Multi-level Secured Approach Using LBP and Spiral Scan Path

**N. Subramanyan, S. Kiran, R. Pradeep Kumar Reddy and P. Manju Yadav**

**Abstract** The world is becoming more interconnected with the advent of Internet and new networking technologies. Security is becoming an important factor and a vital component while transmitting data in the communication channel. As the data is highly confidential, the most popular approach is cryptography which deals with the techniques of secret writing. The goal is to allow the intended recipients of a message to receive the message securely while preventing eavesdroppers from understanding the message. The techniques of cryptography and network security are fully growing, and these lead to the process of development of practical and readily available robust and dynamic techniques. The proposed method encrypts the data in three levels—using dynamically generated key, local binary pattern, gray code, and spiral scan path. Local binary pattern is a nonparametric descriptor which efficiently summarizes the local bits of data. Gray code along with spiral scan path provides high security.

**Keywords** Cryptography · Eavesdropper · Gray code · Internet
Local binary pattern · Security · Spiral scan path

N. Subramanyan (✉) · S. Kiran · R. Pradeep Kumar Reddy
Department of Computer Science and Engineering, Y.S.R. Engineering College of Y.V.U, Proddatur, India
e-mail: subramanyam.neelam@gmail.com

S. Kiran
e-mail: rkirans125@gmail.com

R. Pradeep Kumar Reddy
e-mail: pradeepmadhavi@gmail.com

P. Manju Yadav
Accenture Technologies Pvt. Ltd, Hyderabad, India
e-mail: manjup125@gmail.com

# 1 Introduction

Nowadays, security is becoming an important factor and a vital component while transmitting data in the communication channel. Cryptography [1] is important for secure transmission of information through Internet. Many algorithms or methods are available for the integrity and authenticity of information transferring through the network. Goals of cryptography are confidentiality or privacy, data integrity, authentication, and non-repudiation. Cryptographic systems are classified along three independent properties, based on the type of methods used for converting the plaintext into ciphertext, techniques used for processing the plaintext, and the number of keys used for encryption and decryption. Two types of cryptographic schemes, symmetric cryptography and asymmetric cryptography, differ based on the number of keys.

Fig. 1 shows the conventional encryption model [1]. The plaintext $X$ is converted into ciphertext $Y$ using the encryption algorithm and the key $K$. Here, key $K$ can be any value, which is independent of the plaintext. After getting the ciphertext, it may be transmitted to the receiver.

After receiving the ciphertext, the plaintext is reproduced by the decryption algorithm and same key that was used for encryption. Several factors decide the security of conventional encryption model. The encryption algorithm and secrecy of the key are two important factors. For message $X$ and for encryption key $K$, the encryption algorithm produces the ciphertext $Y$:



**Fig. 1** Conventional encryption model

$$Y = E_k(X). \tag{1}$$

The intended receiver after possession of the same key used for encryption is able to invert the transformation to get the original plaintext:

$$X = D_k(Y). \tag{2}$$

An intruder, observing $Y$ but not having access rights the $K$ or $X$, intruder may try to recover plaintext $X$ or key $K$ or both of them. Assume that the intruder knows both the encryption algorithm ($E$) and decryption algorithm ($D$). If the intruder is showing interest in knowing only for this particular message, then intruder puts effort to get back plaintext $X$ by generating a plaintext estimate $X^{\wedge}$. However, if the intruder is interested in being able to read all future messages as well, he tries to recover key $K$ by generating an estimate $K^{\wedge}$. The remaining part of this paper is organized as follows. Section 2 presents the existing system, i.e., "Automatic Key Generation of Caesar Cipher," the background work. Section 3 presents the proposed method. Section 4 discusses the performance comparisons. Finally, Sect. 5 concludes this paper.

## 2 Background Work

Automatic Key Generation of Caesar Cipher [2] is a simple symmetric encryption technique proposed by B. Bazith Mohammed. It uses the combination of substitution and transposition techniques. The algorithm can accept the plaintext which contains alphabets, numbers, and special characters. The encryption is performed in two phases: In phase one, Caesar cipher is used, and phase two applies the Rail fence technique.

### 2.1 Encryption Algorithm

Phase 1: Using Caesar cipher substitution technique

Step 1: Find the ASCII value for each plaintext character.
Step 2: Generate the key

$$\text{Key} = (\text{sum of all ASCII values}) \bmod 256. \tag{3}$$

Step 3: Apply Caesar cipher

$$E = ((P + K) \mod 256). \tag{4}$$

where P is the plaintext and K is the key.

Phase 2: Using Rail fence technique

Step 4: Consider key as 3 and apply Rail fence technique for the ciphertext obtained in step 3.

## 2.2 Limitations in the Existing System

1. Generating same ciphertext for different plaintexts. For example, the ciphertext for "welcome" and "WELCOME" is same as "c[QOYXQ".
2. Ciphertext is almost same as plaintext for some cases. For example, for plaintext welcomehello, after the Caesar encryption the intermediate ciphertext is "welcomehello" and after applying the Rail fence technique the final ciphertext is woeecmhlolel.
3. It is a stream cipher model.

## 3 Proposed Method

The proposed method provides three levels of security. At first, key is dynamically generated. In the first level, local binary pattern [3] code for the key is generated, and then, logical XOR operation is performed between the LBP code [4] and each plaintext character. The second level takes care of binary-to-gray code conversion. Spiral scan path is applied in the third level.

## 3.1 Key Generation Process

The proposed method uses a dynamic key for encryption and decryption. The bits of plaintext are involved in key generation process. The final key Kn is generated as follows:

$$K_n = ((K_{n-1} * i)\%255) + i. \tag{5}$$

where $i = 1,2,\ldots$ $K_n$ is the final key used for encryption and decryption. $K_0$ is the initial seed value found by using the bits of the plaintext.

### 3.2 Local Binary Pattern

Local binary pattern (LBP) [3] is a nonparametric method; LBP utilizes the bit positions of a decimal number and generates the local binary pattern codes or LBP codes, which encode the decimal value (Fig. 2). Rule sequence for generating the LBP code [4] is shown in Fig. 3.

The local binary pattern code from rule sequence I and II is "g a b c f e d h".

### 3.3 Gray Code

Gray code [5] uses a binary number system invented by Frank Gray. There is only one bit difference between two successive numbers. For example, the gray code for binary number 0011 (3 in decimal) is 0010, whereas for binary number 0100 (4 in decimal) it is 0110. Gray code is primarily used as an encoding code. The conversion from binary to gray and gray to binary is simple process.

### 3.4 Spiral Scan Path

In the spiral scan path [6, 7], the characters or text is scanned from the outside to the inside, tracing out a spiral curve starting from the left corner of the matrix to the center cell. The characters or text is placed in a $5 \times 5$ matrix with a predefined order before the scanning proceeds. The filling of text or characters in the spiral scan path proceeds as follows. Filling of characters follows the shapes one by one: First 12 characters are filled in cross shape, then next 8 characters are filled in diamond shape, then next 4 characters are filled in four corners, and finally, the last character is filled in central cell (Figure 4).

### 3.5 Key Generation

Step 1: Finding the seed value $K_0$.

Take the least significant bits of every plaintext character (exclude the last plaintext character if a number of plaintext characters are 25). Arrange the bits sequentially into three groups having eight bits each.

**Fig. 2** LBP pattern

| $2^7$ | $2^6$ | $2^5$ |
|---|---|---|
| $2^0$ | Key | $2^4$ |
| $2^1$ | $2^2$ | $2^3$ |

**(a)**

| Rule/Base Pairs | Sequence Bits |
|---|---|
| i<j(01,02,12) | a,b,c |
| i>j(10,20,21) | d,e,f |
| i=j(00,22) | g,h |

**(b)**

| Bit Number | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|
| Bits Obtained from Rule sequence I | g | a | b | c | f | e | d | h |

**Fig. 3** **a** Rule sequence I, **b** rule sequence II

**(a)**          **(b)**          **(c)**          **(d)**

**Fig. 4** **a** Positions of first 12 characters of *cross shape*, **b** the positions of next 8 characters of *diamond shape*, **c** the positions of next four characters in the *corners,* and **d** the position of last character (i.e., center)

Step 2: Find the decimal values for each group of bits.

$K_0$ is the sum of these group values.

Step 3: Generating final key $K_n$

$$K_n = ( (K_{n-1} * i)\%255) + i, \qquad (6)$$

where $i = 1, 2, ....$ and $K_n$ is the final key used for encryption and decryption. $K_0$ is obtained in the previous step.

Step 4: $K_n$ is the final key used for encryption and in decryption methods.

## 3.6  Encryption Process

Step 1: Find the ASCII values of each plaintext character.
Step 2: Generate the key using key generation algorithm.
Step 3: Find the LBP code for key.
Step 4: Perform logical XOR between each plaintext character and the key LBP code (obtained in step 3).

Step 5:  Perform significant bit swapping in the result obtained in step 4.
Step 6:  Convert the binary result obtained in step 5 to gray code.
Step 7:  Convert each binary number to ASCII code.
Step 8:  Apply spiral scan path for getting the ciphertext.

## 4  Performance Analysis

The parameters such as computation time and avalanche effect help us in analyzing the performance of cryptographic algorithms.

### 4.1  Avalanche Effect

Avalanche effect [8] is the significant characteristic for encryption algorithm. When a bit of change occurs in either plaintext or key, there should be a significant change in the generated ciphertext:

$$\text{Avalanche Effect} = \frac{\text{No. of flipped bits in the ciphertext}}{\text{No. of bits in the ciphertext}} \times 100\% \qquad (7)$$

Table 1 shows the avalanche effect comparisons for the existing and proposed methods. The results clearly show that proposed method is giving high avalanche effect.

### 4.2  Computation Time

Computation time [9] is another criterion for measuring the performance of cryptography algorithm. The processor speed, algorithm complexity, etc., influence the computational speed. The smallest computation time [9] algorithm is preferred. Encryption computation time is the time required for the encryption algorithm for generating the ciphertext from plaintext. Decryption computation time is the time required by the decryption algorithm for generating the plaintext from ciphertext. Table 2 shows the comparison of encryption and decryption times for different input size files.

**Table 1** Avalanche effect for existing and proposed methods

| Plaintext | | Ciphertext | Ciphertext in binary | No. of flipped bits |
|---|---|---|---|---|
| welcome | Existing Method | c[QOYXQ | 0110001101011010101000101001111010110010101100001010001 | 17 |
| welcomf (e is changed with f) | Existing Method | d\RPZYS | 0110010001011100010100100101000000101101001011001010010011 | |
| welcome | Proposed Method | σ■±2√■≥ | 111001011111101111000100110010111101111111011110010 | 28 |
| welcomf (e is changed with f) | Proposed Method | ■σΩ)α&Θ | 1111111011100101111010100010010011110000000110011011101001 | |

**Table 2** Comparison of encryption and decryption execution times

| Input file size (KB) | Execution times (m s) | | | |
|---|---|---|---|---|
| | Encryption execution time | | Decryption execution time | |
| | Existing method | Proposed method | Existing method | Proposed method |
| 20 | 22,815 | 18,507 | 22,815 | 18,507 |
| 25 | 25,311 | 18,995 | 25,311 | 18,995 |
| 30 | 28,269 | 21,322 | 28,269 | 21,322 |
| 35 | 32,099 | 24,032 | 32,099 | 24,032 |

## *4.3 Encryption and Decryption Throughputs*

$$\text{Encryption throughput (KB/s)} = \frac{\sum \text{input files}}{\sum \text{EET}} \qquad (8)$$

$\sum$input files (existing) = 110 KB, $\sum$EET (existing) = 108,167 ms.
Encryption throughput (existing) = 1.0169 KB/s.
$\sum$EET (proposed) = 83,778 ms.
Encryption throughput (proposed) = 1.3129 KB/s.

$$\text{Decryption throughput (KB/s)} = \frac{\sum \text{input files}}{\sum \text{DET}} \qquad (9)$$

$\sum$DET (existing) = 108,494 ms.
Decryption throughput (existing) = 1.0138 KB/s.
$\sum$DET (proposed) = 82,856 ms.
Decryption throughput (proposed) = 1.3276 KB/s.

The proposed method is giving high throughputs when compared to existing method. For analyzing the above-said performance measures, the existing and proposed algorithms are implemented in Java (Java 7) programing language and executed on the system having AMD six-core processor with 3.50 GHz speed and installed with Windows 8.1 Pro 64-bit operating system.

## 5   Conclusion

Security plays an important role in the communication channel via an untrusted media such as Internet. In order to protect the data from hackers, intruders, and unauthorized persons, one needs better security standards. Proposed method is providing better security when compared with existing Caesar cipher method with three levels of security using LBP codes, gray code, and spiral scan path along with dynamic key generation, which is difficult to break.

Proposed method is showing best performance when compared with existing Caesar cipher schemes in terms of avalanche effect, encryption and decryption times, and throughput. Present encryption can further extended to the system having the increased matrix size for spiral scan path with UNICODE system support and is able to generate ciphertext characters containing only alpha numerals not having special symbols.

## References

1. Stallings, W.: Cryptography and Network Security: Principles and Practices, 4th edn, pp. 30–39. Prentice Hall, Upper Saddle River (2006)
2. Mohammed, B.B.: Automatic key generation of Caesar Cipher. Int. J. Eng. Trends Technol. **6**(6), 2231–5381 (2013)
3. Huang, D., Shan, C., Ardabilian, M.: Local binary patterns and its application to facial image analysis: a survey. IEEE Trans. Syst. Man Cybern. Part C Appl. Rev. **41**(6), 765–781 (2011)
4. Pradeep Kumar Reddy, R., Nagaraju, C.: Text encryption through level based privacy using DNA steganography. Int. J. Eng. Trends Technol. Comput. Sci. **3**(3), 2278–6856 (2014)
5. Wakerly, J.F.: Digital Design: Principles and Practices, 4th edn, pp. 48–50. Pearson Education India, New Delhi (2009)
6. Singh, A., Goswami, L., Ur Rahman Khan, A.: A novel image encryption algorithm for enhanced security. Int. J. Eng. Res. Technol. (IJERT) **2**(9), 2942–2949 (2013)
7. Paul, M., Kumar Mandal, J.: A novel symmetric key cryptographic technique at bit level based on spiral matrix concept. In: International Conference on Information Technology, Electronics and Communications (ICITEC—2013), Bangalore, India, 30–31 March 2013, pp. 6–11 (2013)
8. Ramanujam, S., Karuppiah, M.: Designing an algorithm with high avalanche effect. Int. J. Comput. Sci. Netw. Secur. (IJCSNS) **11**(1), 106–111 (2011)
9. Singh, S., Maakar, S.K., Kumar, S.: Performance analysis of DES and RSA cryptography. Int. J. Emerg. Trends Technol. Comput. Sci. (IJETTCS) **2**(3), 418–422 (2013)

# A Diametrical Association to SkipCloud for Consistent Matching Services in Publishing/Subscribing Framework

**Yerininti Venkata Narayana and Dalli Varun Prasad**

**Abstract** The arrival pace of data is increasing day-by-day to the applications which pose a great confront to disseminate large scale of live content in a scalable and consistent manner to the most interested users. For data distribution, the publish/subscribe model is used, which has capacity of mounting the system to enormous size. The chuck of intricate computing and consistent communication, cloud computing affords an immense chance. In this paper, a mechanism called SkipCloud Management System (SCMS) is proposed, which affords a direct interaction for publishing/subscribing an event in a scalable and consistent manner in cloud environment. To achieve very less routing inactivity and consistent links among different datacenters, a SkipCloud is used to organize the datacenters of SCMS. A hybrid space partitioning method called HPartition is used to map a large-scale pertinent subscription into multiple subspaces which ensures a high relevant matching and provides a number of datacenters for a piece of event.

**Keywords** Pace · Confront · Enormous · Mounting · Intricate · Chuck
Immense · Pertinent · Supple · Secluded · Substance · Doled · Feeble
Inactivity · Intrigued · Fizzled · Acquaintance

## 1 Introduction

Generally, the SkipCloud is additionally a cloud utilized for sending the memberships by directing to its comparing datacenters. The user has to interact first with the datacenters nearby, as they are relevant for subscribing and publishing the content. If they are not relevant for any of those, they forward the event to the

Y. V. Narayana (✉) · D. Varun Prasad
Department of Computer Science & Engineering, DVR & Dr. HS MIC College of
Technology, Kanchikacherla, Krishna District 521180, Andhra Pradesh, India
e-mail: naarayanaa808@gmail.com

D. Varun Prasad
e-mail: dallivarun@gmail.com

SkipCloud and it has to check the relevant datacenter for subscribing and publishing the content. This is a crucial part of the SkipCloud which take care of the routing to a particular datacenter on demand. For this, the user has to connect the nearest datacenter for publishing or subscribing the content. If the data center is not pertinent then the request will be forwarded to the SkipCloud. The availability of the actual data center is verified and content will be subscribed or published successfully to the actual user. So, this is a time-consuming process for the user to wait until the need is fulfilled. However, they are other factors which are influencing the SkipCloud in turns of performance. But, still the throughput can be increased by simple modifications to the system.

## 1.1 Distributed Computing

Distributed computing [1] is put into practice by employing a scheme of secluded servers assisted on Internet to accumulate, regulate, handle instruction, instead of a locale server. The distributed computing system defined in Fig. 1 is storing some content, if needed or subscribing some content when required is based on the user's tasks.

Distributed computing is rapidly evolving that has been widely used in the following ways. Consistency, scalability and sustainability, Effortless and supple deployment, $24 \times 7$ carry, device and location independent, Frees up internal resources, secure storage management, lower capital expenditure. Currently, there are three types of cloud computing are available, namely public, private, and hybrid clouds. Depending upon the type of usage these are used. A cloud is used in services namely: software-as-a-service, platform-as-a-service, and infrastructure-as-a-service.

## 1.2 Existing Process

A number of pub/sub services offered by each of the datacenter are a quick punishment to send the pub/sub information to the SkipCloud. This chiefly originates



**Fig. 1** Distributed computing system

from the accompanying truths: The majority of them are wrong to coordinate the live substance with high information magnitude because of the insufficiency of membership space partition strategies that brings any low coordinating throughput or soaring memory overhead. Subscriptions might be allocated to erroneous coordinating servers, which acquire accessibility issue, substance of coordinating servers furthermore, the subscriptions might be doled out to the wrong coordinating servers.

Firstly, Publishers/Subscribers connecting to the nearest datacenter for event matching. If it is not the pertinent datacenter, it checks for the availability of matching servers in the SkipCloud. Secondly, the time diminution for the association to the SkipCloud ie., from pub/sub to the datacenter and from datacenter to the SkipCloud.

## 2 Literature Survey

Applying a series of steps for increasing the performance of the SkipCloud and also reducing the time for publishing or subscribing the content is needed for the user. This report is well ordered into different parts. Part 1 describes how the connectivity to the datacenter is done. Part 2 describes how the event assignment and forwarding of relevant content are been published or subscribed from the particular datacenter to the actual user. Part 3 describes how the neighbor list is maintained in the SkipCloud. Part 4 describes how routing is done from the SkipCloud to the pertinent datacenter.

### 2.1 Connectivity to Datacenter

The publisher/subscriber [2–6] who is in need of publishing or subscribing some live content from the datacenter will first connect to the nearest available datacenter [2, 7] irrespective of the content may or may not contain in that particular datacenter.

### 2.2 Event Assignment and Forwarding

The publisher/subscriber is feeble with the nearest datacenter; the subscriptions are sent to the SkipCloud, i.e., after accepting a subscription, the broker advances this subscription to its relating subspaces [4, 8, 9].

## 2.3  Neighbor List Maintenance

Clusters of non-peak stages, uses a incompetent associate testing convention in view of Cyclon [10], which gives vigorous network of every cluster. In each and every(m) stage of SkipCloud, every cluster darts Cyclon to remain consistent network along with all the brokers [9]. Since every broker drops into a cluster at every stage, as it keeps up m-neighbor lists. At cluster of stage s, the neighbor list tests parallel number of acquaintance from their relating brood clusters at stage s + 1. This guarantees the routing from the underneath stage can simply locate a broker indicating a superior stage. Since, brokers keep up stages of acquaintance list, also upgrade the neighbor rundown of one stage and successive ones to trim down the obstruction charge. This topology of the multistage acquaintance lists like Tapestry [4]. Contrasted and Tapestry, SkipCloud utilizes various brokers as a part of peak clusters as focuses to guarantee consistent directing.

## 2.4  Prefix Directing from SkipCloud to Datacenter

Prefix directing [3, 7, 9] is chiefly inured to proficiently direct memberships and events to peak clusters. Note that the cluster identifiers at stage s + 1 are produced by joining one k-array to consequent clusters at stage s. The connection of identifiers between clusters is the establishment of directing to target clusters. Quickly, while accepting a directing solicitation to particular group, a broker looks the acquaintance lists of every stage and picks the acquaintance that imparts highest basic prefix to aim clusterid as following hop. Directing action rehashes an acquaintance whose identifier is nearer than itself, until a broker can't discover.

## 3  Proposed Method

### 3.1  Framework Process

From Fig. 2, we come to know that all the brokers in SCMS are exposed to Internet, so that whichever subscriber or publisher can attach to them directly. To guarantee consistent connection and less directing inactivity, the brokers are associated with a distributed overlay, called SkipCloud. This gives an easy maintenance of individual datacenters for data retrieving. The intact substance space is segregated into disjoint subspaces, each managed by a number of brokers. If a subscriber needs some event to dispatch to the datacenter, the event is placed in the particular datacenter through SkipCloud. The subscription and events drops into equivalent subspace are coordinated on the equivalent broker. Subsequent to coordinating procedure is finished, and the events are telecasted to the comparing intrigued subscriber.

**Fig. 2** SCMS framework

From Fig. 2, the publisher P interacts through SkipCloud directly and it finds the exact availability of the event in the corresponding broker identity given and forwards the event to broker B7 which interacts with the datacenter 4 to publish some data. The datacenter, which is relevant for publishing the subscription, publishes the content and sends the successful response to the user.

S1 is another subscriber ie., the interested user starts subscribing the data and interacts with the SkipCloud for relevant datacenter availability. The availability of the event is present in broker B4 and thus the event forwarded to that particular broker who is in datacenter 2 will be transmitted to subscriber S1.

S2 is another subscriber ie., the interested user starts subscribing the data and interacts with the SkipCloud for relevant datacenter availability. The availability of the event is present in broker B6 and thus the event forwarded to that particular broker who is in datacenter 3. And again the event is will be transmitted to subscriber S2.

## 3.2 Component Description

After watchful investigation, the framework has been recognized to have the accompanying components:

1. SkipCloud Organization.
2. SkipCloud Performance.

3. Hybrid multidimensional partition Technique.
4. Publisher/Subscriber Module.

## SkipCloud Organization

Every one of the brokers as the forepart is presented to Internet and in the least subscriber or publisher can interface with them straightforwardly for interaction. Keeping in mind the end goal to accomplish consistent availability and low routing inactivity, every one of these brokers is associated through a distributed overlie called SkipCloud. Memberships and events are posted to the subspaces covering and events drops into the equivalent subspace are coordinated on to equivalent broker. After coordinating procedure, events are shown to comparing intrigued subscribers. Below algorithm1 gives the organization of SkipCloud, and algorithm 2 checks whether the datacenter is available at that particular instant or not.

### Algorithm 1: SkipCloud Management

Step 1: String    urlToRedirect="";
Step 2: String    typeOfData=request.getParameter("typeOfData");
Step 3:  Check if (typeOfData equals to the typeofDataCenter1)
               Array of Urls=Get the relevant datacenters available;
               for(String url: Urls ) then
                  Check if(checkUrl(url))
                     urlToRedirect=url;
                     break;
Step 4: Else check if(typeOfData equals to the typeofDataCenter2)
               Array of Urls=Get the relevant datacenters available;
               for(String url: Urls ) then
                  Check if(checkUrl(url))
                     urlToRedirect=url;
Step 5: redirectUrl(urlToRedirect);

### Algorithm 2: CheckUrl

Step 1: Initialize a boolean active;
Step 2: Assign a String variable for obtaining the Urls
Step 3: Establishing a connection to the url
Step 4: Get the response message from the datacenter
Step 5: Check if(message.equals ("ok")) then
             Datacenter is alive or active;
          else
             Datacenter is not alive or active;

**SkipCloud Performance**

The SkipCloud composes every one of the brokers into various stages of clusters. At the peak stage, brokers are sorted out into different clusters [8], where topography is complete graphs. At this stage, every cluster is a peak cluster [9]. It contains a superior broker which produces a novel k-array identifier length utilizing a jumble function cluster are responsible for the equivalent content subspaces, which gives numerous coordinating candidates to every event. Since, brokers in equivalent peak produce successive correspondence among themselves, for example, overhauling memberships and posting events are sorted out to each additional in single hop to achieve a complete graph [3]. After sorting out, the clusters at the rest stages can be created stage by stage which is depicted in Fig. 3. This identifier is called clusterid.

Comparison Graph of SkipCloud in Terms of Its Performance

The time taken for existence methodology is about 4 s and the time taken for proposed methodology is about 2 s which is represented in Fig. 4. So, from the proposed methodology we reduce the time taken for an event to happen.

Scalability Test for Existence System

We assess the scalability [5] of all approaches through measuring the changing of matching rate with different values of number of brokers, $N_b$ as described in Table 1. The change in number of brokers $N_b$, the matching pace of each approach increases linearly with the growth of $N_b$, from 1 to 4 as shown in Fig. 5.



**Fig. 3** Cluster stages

**Fig. 4** Experimental graph



**Table 1** Scalability test for existence system

| No. of brokers | Matching rate for event doll pics | Matching rate for event building pics | Matching rate for event project pics |
|---|---|---|---|
| 1 | 0.3 | 1.7 | 3.4 |
| 2 | 0.4 | 2.1 | 3.9 |
| 3 | 0.5 | 2.5 | 4.2 |
| 4 | 0.6 | 2.9 | 4.8 |



**Fig. 5** Scalability graph for existence

Scalability Test for Proposed System

We evaluate the scalability of all approaches through measuring the changing of matching rate with different values of number of brokers, $N_b$ as described in Table 2. The change in number of brokers $N_b$, the matching pace of each approach increases linearly with the growth of $N_b$, from 1 to 4 is shown in Fig. 6

**Table 2** Scalability test for proposed system

| No. of brokers | Matching rate for event of doll pics | Matching rate for event of building pics | Matching rate for event project pics |
|---|---|---|---|
| 1 | 0.2 | 1.5 | 3.1 |
| 2 | 0.3 | 1.9 | 3.4 |
| 3 | 0.4 | 2.3 | 3.6 |
| 4 | 0.5 | 2.7 | 4.4 |



**Fig. 6** Scalability graph for proposed system

Reliability

Evaluation of reliability is by testing its capability to recover from server failures. But in our proposed system, even if the server is failed due to some reasons, the events are retrieved from the next datacenter.

**Hybrid Multidimensional Partition Technique**

To accomplish tensile and consistent event coordinating among numerous servers, we proposed a Hybrid multidimensional space partition strategy, call HPartition [9]. It permits comparative memberships to be isolated to equivalent server, gives different candidate coordinating servers to every event. In addition, it attentively assuages hotspots that put workload equalization among the entire servers. HPartition isolates the whole substance space into dislodge subspaces. Memberships and events with covering subspaces are posted and coordinated on the equivalent peak group. To keep workload equalization among servers, it isolates hotspots [9, 11, 12] to numerous coldspots in a versatile way.

**Publisher/Subscriber**

Every subscriber sets up similarity with home broker, and occasionally throws its membership as a pulse note, keeps up a clock for its each cradled membership. On the off chance, if the broker hasn't got pulse note from subscriber over $T$ out time, subscriber should be disconnected. Next, the home broker expels memberships and informs the brokers including the fizzled membership to evacuate it.

## 4   Conclusion

SCMS interfaces the brokers during a distributed overlie SkipCloud, guarantees consistent connectivity among brokers during its multistage groups, brings a less routing inactivity through routing steps. In the course of a Hybrid multidimensional space partition technique, SCMS achieves high scalable skewed memberships and every event is permitted to coordinate on every applicant servers. Despite the fact that event matching service can productively shift through unessential clients from enormous information volume, there are still various issues we have to explain.

Firstly, we don't give flexible reserve provisioning systems to acquire decent concert value proportion. We sketch to actualize the flexible procedures of altering the size of servers in light of the agitate workloads. Furthermore, it doesn't promise that the brokers disseminate substantial live content with different information sizes to the relating subscribers in a continuous way. For the dispersal of mass content, the transfer limit turns into the primary bottleneck. Event matching service is used as a cloud-helped system to understand a common and scalable information broadcasting service over live substance of different information sizes. Apart from these, we didn't examine with respect to the security of the SkipCloud anyplace in the paper. So, as a future work we can include some security related activities to the SkipCloud.

## References

1. Ma, X., Wang, Y., Qiu, Q., Sun, W., Pei, X.: Scalable and elastic event matching for attribute-based publish/subscribe systems. Future Gener. Comput. Syst. **36**, 102–119 (2013)
2. Cao, F., Singh, J.P.: Efficient event routing in content-based publish/subscribe service network. In: INFOCOM, 2004, pp. 929–940 (2004)
3. Gupta, A., Sahin, O.D., Agrawal, D., El Abbadi, A.: Meghdoot: content-based publish/subscribe over p2p networks. In: Middleware, 2004, pp. 254–273 (2004)
4. Kermarrec, A.-M., Massoulié, L., Ganesh, A.J.: Probabilistic reliable dissemination in large-scale systems. IEEE Trans. Parallel Distrib. Syst. **14**(3), 248–258 (2003)
5. Hakiri, A., Berthou, P., Gokhale, A.: Publish/subscribe-enabled software defined networking for efficient and scalable. IEEE Trans. IoT Commun. **53**(9), 48–54 (2015)
6. Daubert, J., Fischer, M., Grube, T., Schiffner, S., Kikiras, P., Mühlhäuser, M.: AnonPubSub: anonymous publish–subscribe overlays. Comput. Commun. **76**, 42–53 (2016)

7. Voulgaris, S., Riviere, E., Kermarrec, A., Van Steen, M., et al.: Sub-2-sub: self-organizing content-based publish and subscribe for dynamic and large scale collaborative networks. In: Research Report RR5772. INRIA, Rennes, France (2005)
8. Rao, W., Chen, L., Hui, P., Tarkoma, S.: Move: a large scale keyword-based content filtering and dissemination system. In: ICDCS, 2012, pp. 445–454 (2012)
9. Ma, X., Wang, Y., Pei, X.: A scalable and reliable matching service for content-based publish/subscribe systems. IEEE Trans. Cloud. Comput. (99), 1–13 (2014)
10. Voulgaris, S., Gavidia, D., van Steen, M.: Cyclon: inexpensive membership management for unstructured p2p overlays. J. Netw. Syst. Manag. **13**(2), 197–217 (2005)
11. Kazemzadeh, R.S., Jacobsen, H.-A.: Reliable and highly available distributed publish/subscribe service. In: 28th IEEE International Symposium on Reliable Distributed Systems, 2009. SRDS'09. IEEE (2009)
12. Zhao, Y., Wu, J.: Building a reliable and high-performance content-based publish/subscribe system. J. Parallel Distrib. Comput. **73**(4), 371–382 (2013)

# Multiple Images Defending Process by Leaning Techniques

**Yerininti Venkata Narayana, Atmakuri Prashant, Dalli Varun Prasad and Chitturi Satya Pavan Kumar**

**Abstract** Providing security to data in the form of text is cryptography, but visual cryptography provides defense for visual images. In this, we intended a mechanism to encipher an image into different levels, i.e., intensity variation, pixel swapping, steganography, and randomization, which seizes network as a means to attain the actual receiver who does the retrieval procedure in reverse order. Different methods are developed, namely intensity variation on each and every pixel of an image and color pixel value swapping within an image, calculating the pixel intensity values and embedding them using steganography, and randomization process. Instead of using external images for manipulation, it provides a better encryption process. Multiple images can be used, which is a good parameter of the proposed procedure. Hence, multiple images can be transformed at a time and more security can be achieved by this proposed method.

**Keywords** Enormous · Retain · Treatment · Imbricated · Decree
Arbitrarily · Gloom

Y. V. Narayana (✉) · A. Prashant · D. V. Prasad
Department of Computer Science & Engineering, DVR & Dr. HS MIC College
of Technology, Kanchikacherla, Krishna District 521180, Andhra Pradesh, India
e-mail: naarayanaa808@gmail.com

A. Prashant
e-mail: prashant.atmakuri@gmail.com

D. V. Prasad
e-mail: dallivarun@gmail.com

C. S. P. Kumar
SCOPE, VIT University, Vellore 632014, Tamilnadu, India
e-mail: pavan540.mic@gmail.com

# 1 Introduction

Visual cryptography, an encryption process which enables visual information such as pictures, text to be encrypted such that decryption can be done via human naked eye [1]. This technique was developed by Moni et al. [1]. They demonstrated a visual secret sharing scheme, where an image was divided into $n$ shares so that only someone with all $n$ shares could decrypt the image. It is not possible to reveal the secret information from any one of the images. Every share reproduced in an individual clarity, and deciphering is done by stacking those shares together. Grouping all shares together, the primary image is appeared.

Different types of notions in the fundamental strategy include m-out-of-n in the visual cryptography [2, 3]. In Fig. 1, you can observe that the two layers grouped each other are correctly aligned and the data hidden will appear [4]. Two shares are obtained from the concrete image. Every share image has a couple of pixels for each pixel in the concrete image. The couple of pixels may be black or may be white depending on the decree: The couple of pixels in the share images are white, and then, the concrete image pixels should be equivalent, and if the couples are imbricated, they look like gleam gray. On the next, if the concrete image pixel was black, the couples in the share images are reciprocal: arbitrarily gloom one with ■□, and other with □■. When both couples are imbricated, they look like gloomy gray.

Subsequently, if both shares are placed over each other, actual has appeared. A single share image will not reveal any data regarding the original and is equivalent from an arbitrary pattern of both couples. Additionally, if anyone of the share is with you, using some strategic rules you can construct an assumed share that merges it with that to reproduce by any means.

# 2 General Cryptography Techniques for Key Generation

Providing security to the text data by using some mathematical equations is called cryptography. This was categorized into two types.

1. Public key cryptography and
2. Secret key cryptography.

**Fig. 1** Demonstration of visual cryptography



Share 1

Share 2

Share 1+Share 2

Public key cryptography is also known as asymmetric cryptography. In this, every user consists of two keys generally called as public key and private key: public key, available to everyone in the group, and private key, which is known user itself. The sender can use any of the keys and send the data to the receiver, with security. If the sender uses public key for encryption, the receiver should use private key for decryption and vice versa. In this, if the intruder is aware of the secret key, the data will be revealed easily to the intruder which is an unauthorized access of the data. Secret key cryptography, also termed as symmetric cryptography, uses only one common key mutually by using some set of procedures for encryption procedure and decryption procedure. In this, both the sender and receiver should be aware of the secret key which is also in confidential. Key distribution is the difficult task in this approach.

In order to overcome the problems in the above approaches, visual cryptography lays a gradient for securing the data. This approach was introduced by Naor and Shamir in 1994 [5]. Visual cryptography is also an encryption method which is used to hide the data in the form of images, such that it does not need any mathematical computations for decrypting the actual data. Indeed, it is decrypted by the human eye by stacking the shared images.

Visual cryptography means encrypting a secret image into n nonsensical shares which cannot reveal a few information with any n combination of shares. Though share is overwhelmed, some procedures are designed in special, wherever pixel values that appear in a concrete hold the actual secret and additionally the secret is being shared among the three shares. If single share is lost, the secret will not be acquired not possible to decipher the original data. If the sent key is numeric tracing the secret is easy for the powlers.

## 3  Basic Methodology of (2, N) Shares Generation

Secret sharing, with unlimited generous people N, as a minimum of 2 needed for decrypting the secret, is one among visual secret sharing techniques introduced in 1994 by Naor and Shamir [5]. Here, we had a secret which encrypted to form N shares that become visible as casual and no legible facts about hidden secret, although any of 2 shares are grouped together the secret happen to legible to human naked eye.

Proposed method considers a bit map to encrypt by satisfying some procedures. A bit map takes 3-byte and 4-byte color form, i.e., bits for each pixel. In the 3-byte color form, we have red, green, and blue colors, and in the 4-byte form, additional factor is transparency. Each color is prepared of 1-byte building up a sum of 3/4 bytes. The value of each color ranges between 0 and 255, which also facilitate the succor of randomly generated numbers on each run of the appliance and consequently they operate as secret inputs [6, 7].

# 4    Proposed System

In this, different encryption techniques are used for an input image which underwent, namely variation in pixel intensity, swapping of pixel values, and randomization of shares.

## 4.1    Variation in Pixel Intensity

First category of encryption, in which 3 arbitrary values are made for every pixel such as RB, RG, and RR identical to every pixel color [8]. The 3 arbitrary values are maintained in temp_B, temp_G, and temp_R; the concrete color values of a pixel are B, G, and R. In these, either the modified or actual values are stored in b, g, r which are used as input color values for the next type of encryption, i.e., pixel swapping. The image shown in Fig. 2 is undergone intensity variation encryption, and output is shown in Fig. 3.

## 4.2    Swapping of Pixel Values

Swapping of pixel values is the second category of encryption, which is done after the output of pixel intensity variation process [4]. The concrete input is separated

**Fig. 2**  Concrete input



**Fig. 3**  Intensity-varied image

**Fig. 4** Possibilities of pixel swapping process



into 8 chunks C1, C2, C3 … C8 as shown in Fig. 4. While implementing this technique, the exact position of the pixel in the last swap should be saved, because pixels are processed in an order. The C1 chunk pixels are changed totally; then, the task is changed to C2 chunk, where interchanging of pixels is persist from the location of last interchanged with C1 chunk. The same process is persisted for C3 chunk; however, it is not replicated for C4 chunk because it includes pixels from C1, C2, and C3 chunks.

Pixel swapping means single swap for each pixel, i.e., each pixel underwent swapping only single time. An exception condition is existed, i.e., in what way pixels of chunk C1 are swapped with pixels of C3 block, if C2 chunk wants to swap the pixels of C2 chunk with the C3 chunk, which it is not possible because C3 chunk already gets swapped with C1 chunk. So, for this instant pixel in C2 chunk is considered to have a zero value and the process continued further.

## 4.3 Generation of Shares

From second form of encryption, swapping of pixel values is performed depending on the arbitrary values of each and every pixel. Pixel color values, namely blue, green, and red, are interchanged, but not the arbitrary values of a color. The arbitrary values of a color of each and every pixel created during the initial encryption type continue undamaged in the pixel position itself and are not changed during the pixel value swapping process.

The arbitrary values and flag values, should be protected, are considered as secret inputs for both encryption procedure and decryption procedure. These values are numeric which makes tracing the values difficult for the intruders. In order to avoid this threat and to improve better security, the values are transformed into an image, i.e., numerical values are converted into an image.

Therefore, the arbitrary and the flag values of each and every pixel are constructed into three different bit map images called as "shares" [3] as shown in

**Fig. 5** Generated shares for concrete image—Share1, Share2, and Share3

Fig. 5. In encrypted image, each share has the same number of pixels, since the length and breadth of obtained shares are equal as the concrete and three factors are obtained from random numbers of the encrypted images. The three shares generated are of 4-byte size, i.e., 32 bits including the encoded image constructed during encryption procedure are set free to the decryption procedure at the receiver.

## 4.4 Average Pixel Intensity

For each generated share, we calculate the average pixel intensity value and also we calculate the average of the three average values of pixel intensities and steganized in the shares of a corresponding image. The same is repeated for every image which we want to send it to the receiver. Here, another protection technique is involved called steganography [9, 10] which is the practice of hiding messages or information within other text.

The below formula is used to find out the average pixel intensity of an concrete image of three shares, which comprises of three values red, green, blue ranging from 0 to 255. $N1$, $N2$, and $N3$ are the total number of pixel values in the each obtained shares of a concrete image, and $i$, $j$, and $k$ are the three planes of an RGB-colored image. If the chosen is a gray-colored image, then there will be only two planes. Here, (1) is used to calculate the RGB-colored image in which there are 3-D plane and (2) is used to calculate the gray-colored image which is a 2-D plane.

$$\frac{1}{N1}\left[\frac{1}{N2}\left[\frac{1}{N3}\sum_{i,j=1,k}^{255} f(i,j,k)\right]\right] \quad \text{where } k = 1, 2, 3 \tag{1}$$

$$\frac{1}{N1}\left[\frac{1}{N2}\left[\sum_{i,j=1}^{255} f(i,j)\right]\right]. \tag{2}$$

## 4.5   Randomization of Shares

Randomization refers to the practice of using possibility methods to assign shares to treatments. In this way, the available data lurking shares are distributed at different chance levels across the treatment conditions. After this, the shares are transmitted to receiver.

## 4.6   Sender (Transmitter) Side Process

Step 1:  Start.
Step 2:  Select an image from the list.
Step 3:  Send selected image to the intensity variation process refer paper (generation of random numbers and flag values) and after go to the next step.
Step 4:  Pixel swapping process and factor calculation and go to the next step.
Step 5:  Generation of 3 shares and go to the next step.
Step 6:  Now find the average pixel intensity value of individual shares and also the average of three average individual pixel intensity values and go to the next step.
Step 7:  By using steganography, we add these values, i.e., generated in step 6 to the 3 individual shares of an image and go to the next step.
Step 8:  Repeat step 2 until the required number of images should be sent. If finished, go to step 9.
Step 9:  Randomizing the shares of each individual shares.
Step 10:  Send all these randomized shares to the receiver using cryptography techniques and go to the next step.
Step 11:  Stop.

## 4.7   Receiver Side Process

Step 1:  Start.
Step 2:  Select any share from the randomized group.
Step 3:  Decrypt the share image using steganography, so that we get the average pixel intensity value of individual share and also the average of three average individual pixel intensity values and go to the next step.
Step 4:  Repeat step 2 for two or more times and then find the average of any 3 shares matching with the total average pixel intensity value of 3 shares and go to the next step.

Step 5: If it is equal, now we can stack the three shares and form an image. Repeat steps 2 for finding all the other images until all the shares has been finished and then go to the next step.

Step 6: All the images one by one are continued for pixel swapping process and then go to the next step.

Step 7: Next, the images are followed with the intensity variation process. Here, each image is recovered to original image from random values of RGB and then go to the next step.

Step 8: Finally, original encrypted images have been obtained.

Step 9: Stop.

## 5 Conclusion

The defense methods employed in this paper are enhanced in affording a good confidentiality for sharing of files than the other secret key cryptographic techniques. Encryption of an image is done at different levels, i.e., intensity variation, pixel value swapping, steganizing, and randomization. Distinct kinds of mechanisms are employed for data to transform it into multiple image securing process. Secret is encrypted into an image; a distinct kind of concept of converting data to image is kept in advance, which is also a fictional method in other available cryptographic procedures.

The generated arbitrary numbers are retained by flag values for the respective random generation. An attacker is not able to break the secret key which are the arbitrary values used as input, because we had pixel values ranging between 0 and 255, an enormous task to attacker to find the accurate value of sender used. About 255 values for three arbitrary values of each pixel in an image and n numbers of pixels values to be originated out for n number of images cast-off in sender process. Another type of method called steganography, used for embedding the calculated pixel intensity values of each share and the total image, is a special procedure for providing more security to data, and also, it is not easy to extract the pixel intensity values from the shares which are unique. Randomization of shares is to confuse the attacker to find the data embedded in it.

In the other visual cryptography methods, the output at the last step is, however, acquired, but it is not as clear enough which is the foremost problems of other methods. But, despite of generous number of treatments made on each pixel value, we acquire the concrete again, which was fed as an input to algorithms with same pixel color value. This is possible by using plain mathematics, which is used for treatments.

# References

1. Gnanaguruparan, M., Kak, S.: Recursive hiding of secrets in visual cryptography. Cryptologia **26**, 68–76 (2002)
2. Verheul, E.R., van Tilborg, H.C.A.: Constructions and properties of k out of n visual secret sharing schemes. Des. Codes Cryptogr. **11**(2), 179–196 (1997)
3. Ateniese, G., Blundo, C., De Santis, A., Stinson, D.R.: Extended capabilities for visual cryptography. Theor. Comput. Sci. **250**, 143–161 (2001)
4. Anitha, R., Sasirekha, N.: Image securing mechanism by gradient techniques. Int. J. Comput. Eng. Appl. **VIII**(I), 55–64 (2014)
5. Naor, M., Shamir, A.: Visual cryptography. In: Advances in Cryptology EUROCRYPT "94. Lecture Notes in Computer Science
6. Koppu, S., Madhu Viswanatham, V.: A survey on security issues: digital images. Int. J. Pharm. Technol. **8**(2),13420–13427 (2016)
7. Addin Ahmed Salem Al-Maweri, N.: State-of-the-art in techniques of text digital watermarking: challenges and limitations. Int. J. Comput. Sci. **12**(2) (2016)
8. Liu, F.: Improving the visual quality of size invariant visual cryptography scheme. J. Vis. Commun. Image Represent. **23**(2), 331–342 (2012)
9. Yu, J., Li, F., Cheng, H., Zhang, X.: Spatial steganalysis using contrast of residuals. IEEE Signal Process. Lett. **23**(7), 989–992 (2016)
10. Sah, H.R., Gunasekaran, G.: Privacy preserving data mining using image slicing and visual cryptography. In: IEEE Transactions, pp. 1–7 (2016)

# Implementing Decision Tree in Air Pollution Reduction Framework

**Anindita Desarkar and Ajanta Das**

**Abstract**  Air pollution, which is one of the biggest threats to the civilization, refers to the contaminated air. It happens due to occurrence of harmful gases, dust, and smoke into the atmosphere which is vulnerable to almost every living creature. It poses serious threats to environmental and social well-being. This paper proposes a layered air pollution reduction framework through implementing the decision tree approach. The proposed framework recommends suggestive measures for reducing air pollution level with the help of an innovative rule base along with mining proper data from the massive dataset. It also discusses the experimental results based on the decision tree approach which shows the implementation of the rule base depending on the pollution level by analyzing the impact factors like holiday, festival, political gathering, etc.

**Keywords**  Predictive analysis · Air pollution · Knowledge discovery
Decision tree · Machine learning

## 1 Introduction

Environmental pollution is nothing but the introduction of harmful substances into the environment which has the effect of adversely affecting every life form on the Earth. Also there are various forms of environment pollution like—air pollution, water pollution, etc. It possesses significant threats to our lives like causing respiratory disorders, global warming and greenhouse effect, various types of cancer, destruction of habitats, climate change, etc.

A. Desarkar (✉) · A. Das
Department of Computer Science and Engineering, Birla Institute of Technology Mesra,
Deemed University, Kolkata Campus, Ranchi, India
e-mail: aninditadesarkar@gmail.com

A. Das
e-mail: ajantadas@bitmesra.ac.in

There are numerous causes of air pollution like fossil fuel burning, agricultural activities, industrialization which includes increased number of vehicles [1]. So, one of the solutions can be to reduce carbon emission level so that it remains within a safer limit. A study from "Nature" can be reported that the outdoor air pollution kills 3.3 million people every year across the world, and the number will rise in the next 35 years. This study also found that China and India have the highest rate of deaths from air pollution. So, air pollution control measures should be taken up in emergency basis.

The objective of this paper is to propose a layered air pollution reduction framework where appropriate rule set is suggested to the citizen by predicting the forthcoming pollution level through implementing the decision tree approach. Organization of the paper is as follows. Section 2 describes the related work across the globe. Section 3 gives the brief description of the decision tree approach which is a well-known supervised machine learning algorithm. Section 4 explains smart air pollution reduction approach along with the decision tree. This framework can be used for detecting and taking immediate steps for reducing air pollution with minimum human intervention. Sample experimental results are listed in Sect. 5, and Sect. 6 concludes the paper.

## 2 Related Work

Several measures are taken worldwide to reduce the air pollution. As of 2014, the US government took help of an organization called Environment Protection Agency which partnered with private businesses to decrease industry emissions. Another joint initiative was made by USA and China government where resolution was taken to control industrial air pollution. A specific limit was set by China government on the use of coal and vehicles with high emissions on the road. Several incentives were declared to the cities that can reduce pollution level. As the next step, Beijing shut down hundreds of factories responsible for emitting high pollutants and performed clean audits of hundreds more as of 2014. The central pollution board of India has declared an initiative names as NAQMP—National Air Quality Monitoring Program. It measures the air quality of 341 stations across 126 cities/towns along with 4 union territories in India including Kolkata [2, 3]. Various air pollutants like $SO_2$, $NO_2$, SPM, and RSPM/PM10 are regularly monitored as a part of the initiative [4]. Delhi government has proposed the odd—even rule to reduce congestion as well as to reduce pollution resulting from vehicular emissions. The rule says that cars with odd number plate will be on the road on odd dates only and so on the even date. The government has decided to add 1000 more buses gradually along with 9000 CNG contract carriages as an alternative transport in Delhi. The Supreme Court mandated the use of masks for traffic police as they are exposed to pollution for long hours. A whopping 100% hike has been asked by the Supreme Court for the green cess on commercial vehicles entering Delhi, and new boards are also formed by the Delhi government to notify the new cess in 125 toll

booths across the state. Registration of luxury SUVs and diesel cars above 2000 cc as national capital has been prohibited by the order of Supreme Court as diesel cars are the major source of pollution. Adequate initiative has been taken to convert the existing taxis plying in the city to CNG. The National Green Tribunal (NGT) has requested central and state government not to use diesel vehicles for its employees. An order has been issued by NGT which bans the burning of crop residues in few states like Delhi, Punjab, Rajasthan, Haryana, and Uttar Pradesh as it was observed in the past that there is a sharp rise of air pollution in the NCR zone due to this phenomenon [5].

## 3 Decision Tree: A Supervised Machine Learning Approach

Machine learning is a branch of artificial intelligence which can process petabytes of data, analyses with minimal human intervention and learns from the past experience to achieve the reliable and effective end result. Supervised learning is a machine learning task which infers a function from labelled training dataset. Here the training dataset is fed as input to the supervised learning algorithm for data analysis and deriving the inferred function. This inferred function is used to map new datasets. Supervised learning is broadly classified into two major types of problems: classification and regression. And decision tree is a very well-known method used for classification and regression problems across the globe. The primary objective of this method is to create a model which is able to predict the value of a target variable by learning simple decision rules inferred from the data features [6, 7, 8]. Decision trees have the benefit for knowledge-based systems—they are easily comprehensible by human experts and can be transformed directly into production rules. Moreover, when a specific case is assigned to the decision tree, it provides the reason behind the choice along with the solution. Another advantage of decision tree is its simplicity and efficiency of their construction compared to other classifiers such as neural networks [9, 10].

## 4 Decision Tree Based Air Pollution Reducing Approach

The public health of metropolitan cities is significantly affected due to the increased pollution rate. The city people largely depend on its public transport, buses and non-motorized vehicles which cause the pollution. The air pollution increased through the pollutants generated from these various vehicles, particularly ground-level ozone and particulate matter can worsen respiratory diseases and trigger asthma attacks [11].

## 4.1  Proposed Air Pollution Reducing Framework

This section proposes a data analytics based layered framework towards reducing air pollution in various metropolitan cities. We are collecting pollution related data from various parts of the city, applying various transformation rules and finally store in the pollution database which is basically a warehouse. This warehouse will contain huge volume of data as its additive in nature based on the time period. Meaningful insight can be derived with the help of proper data mining techniques and predictive analysis on this huge pollution related data. This proposed framework consists of four layers, [12] and the layered wise functionalities are described in the following.

## 4.2  Source of Data

Different air pollution measurement parameters like $SO_2$, $NO_2$, RSPM, and SPM will be calculated from the samples collected through various sensors, cameras, and other devices placed across the city. These collected data work as the input parameters of this framework and will be passed to the next ETL layer.

## 4.3  Extraction, Transformation and Loading (ETL)

Further analysis will be carried out based on the aggregated air samples, and the required transformation will be applied on this aggregated set to match with the monitoring requirement. The transformed dataset will be finally stored into the pollution master database. Few such attributes involved in this analysis have been captured in the Table 1.

## 4.4  Analytics and Knowledge Discovery by Creating Decision Tree

This is the main process where a decision tree is formed to forecast the pollution of forthcoming days. *Pollution master database* will act as the primary source on which detailed analysis will be performed to form the decision tree which is

**Table 1**  Sample statistics analysis report

| Location | Pollution parameters | Weekly average | Standard deviation | Valid monitoring days | Percent violation |
|----------|----------------------|----------------|--------------------|-----------------------|-------------------|

basically a graphical representation of possible outcomes to a decision depending on specific conditions. The decision tree is able to forecast the pollution level of the forthcoming days based on the new input parameters.

*Historical database* also provides important inputs to generate the predicted pollution level. It contains the pollution history and their corresponding incidents and phenomenon, presented in following Table 2, like two patients expired as the ambulance got stuck due to heavy traffic and could not reach hospital timely. Finally, *Rule base* provides the rule sets which are suggested as the solution to reduce the pollution. The rule set is created by analyzing various phenomenons from pollution master database. New rules can be added or modified in the table based on the current scenario. It comes as the leaf nodes in the decision tree. *Knowledge database* keeps information regarding the future events of the city which are responsible to affect the pollution level. The decision tree can feed new datasets along with corresponding inputs from the knowledge database to guide us in implementing appropriate rules in critical situations where pollution threshold is exceeded.

Adequate laws and regulations should be in place for successful implementation of the above rules across the city to make it an efficient and effective initiate. The following section describes few such rules.

- **Rule 1** No private vehicle should be allowed to pass through the region except for the emergency situation—adequate measures should be taken while passing in case of emergency situation like using the pollution mask—academic institutions should remain closed—situation like board examinations, etc., can be excused—offices should allow work from home for their employees if possible.
- **Rule 2** Private vehicles carrying four people should only be allowed with adequate precautions and proper justifications, like it will not be allowed in case of casual shopping or any entertainment purposes.
- **Rule 3** Only senior citizens would be allowed to pass through with adequate precautions. Emergency situation will be treated as exceptions.
- **Rule 4** Entry should be given only to the odd numbered private vehicles with adequate precautions. Emergency situation will be treated as exceptions.
- **Rule 5** Entry should be given only to the even numbered private vehicles with adequate precautions. Emergency situation will be treated as exceptions.
- **Suggestive Measure 1** Extra tax benefit would be provided to the shared service cab providers.
- **Suggestive Measure 2** Special discount on road tax would be provided to the solar energy enabled vehicles.

**Table 2** Sample structure of historical dataset

| Location | Pollution parameter | Date | Pollution percent violation | Overall pollution percent violation | Incident name | Incident desc. | Imp. phenomenon |
|---|---|---|---|---|---|---|---|

Here the following Fig. 1 represents the decision tree where OPPV denotes the overall pollution percent violation. Based on OPPV level, it can be categorized into three levels. The first level constitutes where the OPPV level is zero which means the environment is currently safe. The second level is created where OPPV is less than forty which leads to moderate pollution. Based on other conditions, corresponding rules are applied. The third level is the extreme case where the OPPV is greater than forty. Here more strict rules will be applied based on various conditions.

The next part is *Knowledge Discovery* which would be performed based on the prediction generated in the decision tree with the help of Rule Base and knowledge database. Knowledge Database (as presented in Table 3) will store information regarding the forthcoming events which may affect the pollution level of the city, like information about the general holiday list—assuming that city will encounter less traffic in the holidays which leads to lesser pollution level.

The decision tree will feed new set of input data and corresponding information from the knowledge database. Based on that input, the decision tree will guide us to apply the appropriate rule set considering the pollution level.

## 4.5 *Visualization and Interpretation*

This is the last but very important layer as it will showcase the decisions to the citizens as well as proper support is required from the Government authorities for adherence of the decided rules.



**Fig. 1** Decision tree for predicting pollution level

**Table 3** Sample structure of knowledge database

| Location | Date | Major incident | Holiday |
| --- | --- | --- | --- |

# 5 Experimental Results

The following two cases depict two different scenarios which can arise while predicting the forthcoming pollution for taking appropriate steps against it.

*Case 1*

- Suppose we want to predict the pollution level for tomorrow for place A and the current OPPV (overall pollution percent violation) of the place is 50 which is received from the Historical dataset.
- The next step is to check in the knowledge database about important phenomenon of the date. Suppose it comes as a festive holiday.
- As a final step, the decision tree needs to be checked to decide the implementation of the appropriate rule set based on the current conditions. Here the required inputs are current OPPV = 50, Holiday = Yes and Festival = Yes.
- According to the decision tree, Rule 1 which is "No private vehicle should pass touching the region. Emergency situation would be excused.", should be applied on that locality to handle the present pollution level and try to lower it down.

*Case 2*

- Suppose tomorrow's pollution level needs to be predicted for place B where the current OPPV is 18.
- After checking in the knowledge database, suppose we receive information that it's a working day along with a planned political gathering.
- The required inputs for decision tree are current OPPV = 18, Holiday = No and Political Gathering = Yes.
- Based on the decision tree, Rule 3 which is "Only senior citizens would be allowed to pass through. Emergency situation would be excused", should be applied to face the challenges.

The following Table 4 summarizes the above two cases.

Core java is used for implementing the above decision tree. The following screenshots describe the predicted rule set along with the given inputs of the above two cases where Fig. 2 represents Case 1 and Fig. 3 represents Case 2.

**Table 4** Evaluation of sample inputs based on decision tree

| Location | Pollution prediction date | Current OPPV level | Important phenomenon | Rule set—derived from decision tree |
|---|---|---|---|---|
| Place A | 12/9/2016 | 50 | Festive holiday | If current OPPV > 40 and holiday = Yes and festival = Yes, then apply Rule 1 |
| Place B | 12/9/2016 | 18 | Working day with political gathering | If current OPPV($\geq 0$ and $\leq 40$) and holiday = No and political gathering = Yes, then apply Rule 3 |

ffort

8. Kotsiantis, S.B., Zaharakis, I.D., Pintelas, P.E.: Machine learning: a review of classification and combining techniques. Artif. Intell. Rev. **26**(3), 159–190 (2006)
9. Pach, F.P., Abonyi, J.: Association rule and decision tree based methods for fuzzy rule base generation. Int. J. Comput. Electr. Autom. Control Inf. Eng. **2**(1) (2008)
10. Julián, C.I.F., Ferri, C.: Airvlc.: an application for real-time forecasting urban air pollution. In: Proceedings of the 2nd International Workshop on Mining Urban Data, Lille, France (2015)
11. Air Quality Assessment. Retrieved 17 June 2016 from http://www.cpcb.nic.in/15-44.pdf
12. Desarkar, A., Das, A.: A smart air pollution analytics framework. Presented in the International Conference on ICT for Sustainable Development (ICT4SD 2016), held in Bangkok, December 12–13, 2016 (2016)

# Identification of Subgroups in a Directed Social Network Using Edge Betweenness and Random Walks

K. Sathiyakumari and M. S. Vijaya

**Abstract**  Social networks have obtained masses hobby recently, largely because of the success of online social networking Web sites and media sharing sites. In such networks, rigorous and complex interactions occur among several unique entities, leading to huge information networks with first rate commercial enterprise ability. Network detection is an unmanaged getting to know challenge that determines the community groups based on common place hobbies, career, modules, and their hierarchical agency, the usage of the records encoded in the graph topology. Locating groups from social network is a tough mission because of its topology and overlapping of various communities. In this research, edge betweenness modularity and random walks is used for detecting groups in networks with node attributes. The twitter data of the famous cricket player is used here and network of friends and followers is analyzed using two algorithms based on edge betweenness and random walks. Also the strength of extracted communities is evaluated using on modularity score and the experiment results confirmed that the cricket player's network is dense.

**Keywords**  Edge betweenness · Random walks · Modularity
Community detection · Social network

## 1   Introduction

The developing use of the Internet has brought about the development of networked interaction environments consisting of social networks. Social networks are graph structures whose nodes represent people, corporations, or other entities, and whose

K. Sathiyakumari (✉) · M. S. Vijaya
PSGR Krishnammal College for Women, Coimbatore 641004, Tamilnadu, India
e-mail: sathiyakumari@psgrkc.ac.in

M. S. Vijaya
e-mail: msvijaya@psgrkc.ac.in

edges represent courting, interaction, collaboration, or have an effect on among entities. The edges in the network connecting the entities may have a direction indicating the flow from one entity to the other, and a strength denoting how much, how often, or how important the relationship is. Researchers are increasingly interested in addressing a wide range of challenges exist in these social network systems.

In recent years, social network studies has been completed the use of massive amount of statistics gathered from online interactions and from explicit courting links in online social community systems such as Facebook, Twitter, Linkedin, Flickr, Instant Messenger. Twitter is highly rated as a new shape of media and utilized in various fields, such as corporate marketing, education, broadcasting. Structural characteristics of such social networks can be explored the usage of sociometrics to recognize the shape of the network, the properties of links, the roles of entities, information flows, evolution of networks, clusters/communities in a network, nodes in a cluster, center node of the cluster/network, and nodes on the periphery, etc. To find out functionally related items from communities, [1, 2] allow us to observe interplay modules, lacking characteristic values and expect unobserved connections among nodes [3]. The nodes have many relationships among themselves in groups to proportion commonplace residences or an attributes. Figuring out network community is a trouble of clustering nodes into small corporations and a node can be belonging to a couple of communities immediately in a network structure.

Unique resources of statistics are used to perform the clustering challenge, first is about nodes and its attributes and the second is ready the relationship among nodes. The attributes of nodes in community structure are known properties of users like network profile, author publication, publication histories which helps to determines similar nodes and community module to which the node belongs. The connection between the nodes provides information about friendships, authors collaborate, followers, and topic interactions.

A few clustering algorithms [4, 5] employ node attributes, however, ignores the relationships among nodes. But the network detection algorithms make use of businesses of nodes which can be densely related [6, 7] but ignore the node attributes. By way of using these sources of records, sure set of rules fails to explain vital structure in a community. For instance, attributes may additionally tell about which community node with few links belonging to and it is far hard to decide from network structure on my own. On the opposite, the community offers detail approximately two nodes belong to identical community even someone of the node has no attribute values. Node attributes can stabilize the community structure which ends up in more correct detection of communities. Thus, community detection becomes difficult undertaking when thinking of both node attributes and network topology.

The proposed method overcomes the above problem by identifying communities based on node and its attributes by implementing Girvan–Newman edge betweenness and random walks algorithm.

## 2  Related Work

A network is a densely related subset of nodes; this is carefully related to the last network. Social networks are a mixture of essential heterogeneities in complicated networks, together with collaboration networks and interplay networks. Online social networking packages are used to represent and version the social ties among people. Finding communities within an arbitrary community may be a computationally difficult challenge. Numerous research dealings in recent past years have been carried out in the subject matter of network detection, and a number of the crucial research works are point out beneath.

Nicola Barbieri et al. [8] provided network-cascade community (CCN) model, which produced overlapping communities based totally on its interest and level of authority. Major drawback of this model changed into sluggish in getting to know section and also slows in estimate impact energy.

Xie et al. [9] determined numerous lessons of overlapping communities the usage of special algorithms like clique percolation, label propagation [10, 11], agent-based and debris-based totally fashions. This method is used to link partitioning and stochastic generative models.

Evans and Lambiotte [12] used everyday node partitioning to line graph for obtaining hyperlink portioning in authentic community. Ahn et al. [13] used Jaccard coefficient of the neighborhood node to discover similarity among hyperlinks. Kim and Jeong [14] used Infomap approach to encode random stroll path in line network.

Hughes and Palen [15] analyzed how twitter customers react and unfold statistics on social and political issues and located out that those massive activities entice new customers to twitter. Diakopoulos and Shamma [16] accumulated tweets concerning the US presidential election candidates of 2008 and analyzed public emotional response, visible expression, and so on. Kwak et al. [17] analyzed a first rate quantity of twitter dialogues and consumer relationships.

Fortunato [6] affords a entire evaluate in the vicinity of network detection for undirected networks from a statistical physics attitude, while Schaeffer [18] particularly specializes in the graph clustering problem as an unmanaged studying mission. Both surveys in short discuss the case of directed networks; but their consciousness is on the undirected case of the problem.

In this study, the Girvan–Newman algorithm based totally on edge betweenness and random stroll set of rules is applied for coming across communities in networks with node attributes. The twitter facts of the well-known cricket player are taken for take a look at and network of friends and fans is analyzed based on modularity score.

## 3  Girvan–Newman Algorithm

### 3.1  *Community Detection Framework*

The Girvan and Newman is a trendy community finding algorithm. It performs herbal divisions some of the vertices without requiring the researcher to specify the

numbers of groups are present, or placing limitations on their sizes, and without displaying the pathologies obvious within the hierarchical clustering techniques. Girvan and Newman [19] have proposed an set of rules which has three definitive functions: (1) Edges are steadily eliminated from a community, (2) the edges to be removed are selected with the aid of computing betweenness rankings, and (3) the betweenness scores are recomputed for removal of each aspect.

As a degree of traffic float, Girvan and Newman use place betweenness, a generalization to edges of the renowned vertex betweenness of freeman [20, 21]. The betweenness of a part  is described because the wide variety of shortest paths between vertex pairs. This amount can be calculated for all edges in the time complexity of $o$ $(mn)$ on a graph with $m$ edges and $n$ vertices [22, 23].

Newman and Girvan [24] define a measure known as modularity, that is a numerical index that shows ideal separation between nodes. For a separation with $g$ corporations, define as $g \times g$ matrix $e$ whose difficulty $e_{ij}$ is the fraction of edges within the authentic network that be part of vertices in group $i$ to those in organization $j$. Then, the modularity is defined as

$$Q = \sum_i e_{ii} - \sum_{ijk} e_{ij}e_{ki} = \text{Tr}\, e - \left\| e^2 \right\|,$$

which suggests the sum of all factors of $x$, $q$ is the fraction of all edges that lie inside communities minus the predictable fee of the identical quantity in a graph in which the vertices have the identical degrees; however, edges are positioned at random without look upon the companies. The $q$ = zero indicates that network form is not any stronger than could be anticipated through randomness and values aside from 0 represent deviations from randomness. Limited peaks in the modularity at some stage in the development of the network shape set of rules endorse suitable divisions of the community.

## 3.2  Girvan–Newman Partitioning Algorithm

**Successively Deleting Edges of High Betweenness**

Step 1: Find the edge or multiple edges with maximum betweenness; if there can be tie in betweenness, then put off those edges from graph. This system may spilt the graph into numerous additives; it make first degree partition of graph.

Step 2: Recalculate all betweenness values and then remove the edges/edge with high betweenness value. Again split the first-level region into several components such that there are nested within larger regions of graph.

Step 3: Repeat steps (1) and (2) till edges remain in graph.

**Computing Betweenness Values**

For each node A:

Step 1: Do breadth first search starting at node *A*.

Step 2: Count the number of the range of shortest paths from *A* to every different node.

Step 3: Decide the quantity of waft from *A* to all other nodes.

## 3.3  *Random Walks Algorithm*

Random walks is a mathematical idea formalizing a method consisting of a chain of random steps. In case of graphs, given a node that corresponds to a starting point, a random walks is defined due to the fact the collection of nodes fashioned with the resource of a repeating technique beginning from the initial node and randomly transferring to network nodes. At every step, the random walker is positioned on a node of the graph and jumps to a present-day node selected randomly and uniformly among its friends.

Mathematically, allow $GU = (V, E)$ be an undirected graph and $v_0$ be the starting node of the random walk. At the $t$th step, the random stroll is located at node $i$. At $t + 1$ step, the random walks is transferring from node $i$ to node $j$ (neighbor of $i$) with transition chance $1/k_i$. This defines the transition matrix $p$ of the random walks as

$$\begin{cases} \frac{A_{ij}}{k_i}, & if\,(i,j) \in E, \\ 0, & otherwise \end{cases} \tag{1}$$

This matrix may be written as $P = D - 1A$, in which $D - 1$ is the inverse of the diagonal degree matrix $d$. This matrix can also be considered as a diploma normalized model of the adjacency matrix. Random walks are considered to be Markov chains 1, wherein the set of feasible states corresponds to the vertex set of the graph.

Any distribution on a graph $G$ may be represented by the useful resource of a row vector $\pi = [\pi_{1,} \ldots \pi_n]^T$, wherein the $i$th entry that captures the amount of the distribution is residing at node $i$. In case of random walks, the chance distribution over the graph $g$ for each node $i \in v$ at any time step offers the opportunity of the random stroll of being at node $i$. As a result, if $\pi$ is the preliminary distribution, then $\pi_1 = \pi_P$ is the distribution after one step and $\pi_t = \pi_{Pt}$ is the distribution after $t$ steps. Therefore, it is able to outline a stationary distribution $\pi_s$, because the distribution where $\pi_s = \pi_{sPt}$, $\forall t$. The stationary distribution corresponds to a distribution that does not alternate through the years and describes the opportunity that the stroll is being at a selected node after a sufficiently long time. The combination time is the time wanted via the random stroll to reach its stationary distribution. The spectrum of the transition matrix $p$ can be used to sure the combination time of a random walks on a graph and in particular the second largest eigenvalue [25].

# 4   Experiments and Results

The proposed framework includes four phases: Twitter data, directed network, community detection algorithm, and modularity score. Every phase is described in the following sections, and the architecture of the proposed system is shown in Fig. 1.

   A directed network is created using Twitter friends/followers listing as the graph. In this community detection, two algorithms are used for the stage of evaluation of network Girvan–Newman, and random walks algorithm is used to detect communities and subgroups. The size of subgroups is found using Girvan–Newman algorithm and random walks of this network. The algorithm also detects modularity score of community. The real-time data is collected using the twitter application programming interface 1.1 for this research work. Nine thousand records of friends and followers list of the famous cricket player have been crawled from his twitter account. The data is collected at run time from twitter network using R3.2.4, a statistical tool. The cricket player's initial community network as shown in Figs. 2



**Fig. 1**  Community detection framework

**Fig. 2**  Cricket player's initial network

and 3 depicts the relationship types such as friends, followers, and friends and followers. This network has 7095 edges and 6831 vertices.

A community is a densely related institution of vertices, with only sparser connections to different groups. Girvan–Newman algorithm and random walks algorithm are employed here to detect communities from cricket player's twitter network, since it is a directed network. The modularity score for this network is obtained as 0.91. Thirty-nine different communities are extracted for this network based on edge betweenness modularity measure and demonstrated in different colors as shown in Fig. 4. These 39 communities are clustered based on followers, friends, and both followers and friends in the network. The distribution of nodes in various communities is shown in Fig. 6. The membership of size of community 1 is 69, community 2 has the highest size with 166 memberships. Communities 3 and 4 have the membership sizes 42 and 39, respectively. Communities 5 and 7 have the same membership size 37 and so on. Eight different communities are extracted from the same network based on random walks modularity measure and demonstrated in different colors as shown in Fig. 5. These eight communities are clustered based on followers, friends, and both followers and friends in the network. The distribution of nodes in various communities is shown in Fig. 7. The membership of size of community 1 is 207, community 2 has 166, and community 3 has the highest size with 237 memberships. Communities 3 and 4 have the membership sizes 193 and 209, respectively. Communities 1, 3, and 5 have the high membership size of other



**Fig. 3** Friends and followers network

**Fig. 4** Communities identified based on edge betweenness algorithm



**Fig. 5** Communities identified based on random walks algorithm



**Fig. 6** Membership distribution of communities (edge betweenness algorithm)

**Fig. 7** Membership distribution of communities (random walks community algorithm)



**Fig. 8** Modularity scores of community detection algorithms

communities. The modularity score of the network is shown in Fig. 8. The comparative study of community detection algorithms is made in terms of number of communities, membership distribution, and modularity score.

## 5   Discussion and Findings

The aim of network detection in graphs is to discover the subgroups by using the use of the statistics encoded inside the graph topology. In this research work, the modularity score obtained through edge betweenness algorithm is 0.91, which proves that the cricket player's friends and followers network is dense.

The Girvan–Newman algorithm has detected 39 different communities from the cricket player's network and found five communities dense out of total communities. The modularity score found through random walks algorithm is 0.7, which also confirms that the cricket player's friends and followers network is dense. The random walks algorithm has found eight communities from the cricket player's network which are all highly dense. Girvan–Newman algorithm has discovered more number of sparse communities than random walks algorithm and has eliminated them during clustering. Random walks algorithm finds less number of communities with high communication between the nodes.

## 6 Conclusion and Future Work

This work elucidates the application of Girvan–Newman algorithm and random walks for detecting communities from networks with node attributes. The real time twitter directed network of a cricket player is used to carry out network analysis. Modularity score is evaluated and subgroups are detected using two community detection algorithms of directed network. The membership distributions of the subgroups generated by two algorithms were discussed. The experimental results indicate that the network is dense and communication between the nodes is high. As scope for further work, analysis of nested communities can be carried out with cliques and subgroups.

## References

1. Coscia, M., Rossetti, G., Giannotti, F., Pedreschi, D.: Demon: a local-first discovery method for overlapping communities. In: KDD '12 (2012)
2. Girvan, M., Newman, M.: Community structure in social and biological networks. In: PNAS (2002)
3. Yang, J., Leskovec, J.: Overlapping community detection at scale: a non-negative factorization approach. In: WSDM '13 (2013)
4. Blei, D.M., Ng, A., Jordan, M.: Latent Dirichlet allocation. MLR **3**, 993–1022 (2003)
5. Johnson, S.: Hierarchical clustering schemes. Psychometrika **32**(3), 241–254 (1967)
6. Fortunato, S.: Community detection in graphs. Phys. Rep. **486**(3–5), 75–174 (2010)
7. Xie, J., Kelley, S., Szymanski, B.K.: Overlapping community detection in networks: the state of the art and comparative study. ACM Comput. Surv. (2013). doi:10.1145/2501654.2501657
8. Barbieri, N., Bonchi, F., Manco, G.: Cascade-based community detection. In: WSDM'13, February 4–8, Rome, Italy (2012)
9. Xie, J., Kelley, S., Szymanski, B.K.: Overlapping community detection in networks: the state of the art and comparative study. ACM Comput. Surv (2013). doi:10.1145/2501654.2501657
10. Gregory, S.: Finding overlapping communities in networks by label propagation. New J. Phys. **12**(10), 103018 (2010)
11. Padrol-Sureda, A., Perarnau-Llobet, G., Pfeie, J., Munes-Mulero, V.: Overlapping community search for social networks. In: ICDE (2010)

12. Evans, T.S., Lambiotte, R.: Line graphs, link partitions, and overlapping communities. Phys. Rev. E **80**, 016105 (2009)
13. Ahn, Y.-Y., Bagrow, J.P., Lehmann, S.: Link communities reveal multiscale complexity in networks. Nature **466**(7307), 761–764 (2010)
14. Kim, Y., Jeong, H.: The map equation for link communities. Phys. Rev. E **84**, 026110 (2011)
15. Hughes, A.L., Palen, L.: Twitter adoption and use in mass convergence and emergency events. Int. J. Emerg. Manag. **6**(3–4), 248–260 (2009)
16. Diakopoulos, N.A., Shamma, D.A.: Characterizing debate performance via aggregated twitter sentiment. In: Proceedings of the 28th International Conference on Human Factors in Computing Systems (CHI '10), pp. 1195–1198 (2010)
17. Kwak, H., Lee, C., Park, H., Moon, S.: What is twitter, a social network or news media? In: Proceedings of the 19th International World Wide Web Conference (WWW '10), pp. 591–600 (2010)
18. Schaeffer, S.E.: Graph clustering. Comput. Sci. Rev. **1**(1), 27–64 (2007)
19. Girvan, M., Newman, M.E.J.: Community structure in social and biological networks. Proc. Natl. Acad. Sci. USA **99**, 7821–7826 (2002)
20. Freeman, L.C.: A set of measures of centrality based upon betweenness. Sociometry **40**, 35–41 (1977)
21. Anthonisse, J.M.: The rush in a directed graph. Technical Report BN9/71, Stichting Mathematicsh Centrum, Amsterdam (1971)
22. Newman, M.E.J.: Scientific collaboration networks: II. Shortest paths, weighted networks, and centrality. Phys. Rev. E **64**, 016132 (2001)
23. Brandes, U.: A faster algorithm for betweenness centrality. J. Math. Sociol. **25**, 163–177 (2001)
24. Newman, M.E.J., Girvan, M.: Finding and evaluating community structure in networks. Preprint cond-mat/0308217 (2003)
25. Newman, M.E.J., Girvan, M.: Finding and evaluating community structure in networks. Phys. Rev. E **69**, 026113 (2004)

# Handover Techniques in New Generation Wireless Networks

**Vinodini Gupta and Padma Bonde**

**Abstract** In the present age of telecommunication, the growth of wireless technologies and its integration with telecommunication standards have made wireless internetworking highly ubiquitous. With increasing diversity and dynamicity, performance of traditional wireless technologies degrades. To enhance performance, existing technologies need to be integrated. This has led to the evolution of new generation wireless networks (NGWNs). Mobility, seamless migration, end-to-end connectivity and QoS requirements demanded the need of handover amongst NGWNs. This paper presented a survey over basics of handover, its mechanism and behavioural pattern, classification and analysis of existing approaches, existing issues and factors affecting handoff performance.

**Keywords** NGWN · Handover · Seamless migration · Mobility
QoS · Service continuity · Network security

## 1 Introduction

With the evolution of digital era, various technological revolutions and the growing need for faster communication have led to the evolution of many generations in ubiquitous networks from 1G to 4G networks, and further to 5G networks in the forthcoming advancements. With integration of wireless technology and cellular communication, no wireless technology can fulfil the QoS requirements all alone. This has led to the proliferation of NGWNs. These technologies support proactive approach for improvising user experience and ensure better performance, higher

V. Gupta (✉) · P. Bonde
Computer Science and Engineering Department, Shri
Shankaracharya Technical Campus, Bhilai, Chhattisgarh 490020, India
e-mail: gupta.vinu9@gmail.com

P. Bonde
e-mail: bondepadma@gmail.com

data rates, diverse coverage area and wide range of services. However, cost issues still need to be considered.

As the popularity of NGWNs is increasing, attention is now being paid over newly emerging 5G networks. Here, focus will be basically on enhancing coverage areas and user experience. Although 5G networks cope up traffic growth efficiently and provide required QoS standards, researchers are still looking forward to incorporate popular technologies like software-defined networking (SDN), machine-to-machine communication (M2M), etc., with newly emerging technologies like big data, leading to the evolution of 5G networks with desirable standards.

Handoff is a technique of achieving uninterrupted network services during channel transition from one network to another having either different or same air interfaces. Several handover techniques have been implemented till date to enhance the network performance. However, security, mobility, seamlessness, handoff and communication latency, quality of experience (QoE) and QoS standards are the issues of prime concern during handoff in NGWNs.

Paper presented a detailed literature store over the basics of handover. Section 2 focussed over major handover issues and challenges. The parameters affecting the handover mechanism were highlighted in Sect. 3. Section 4 discussed various handovers mechanisms in wireless networking. Section 5 presented discussion and analysis of the currently existing handoff mechanisms. Finally, conclusions of the survey were drawn in Sect. 6.

## 2  Challenges and Issues in Handover

Various challenges and issues during handoff in ubiquitous networking are briefly discussed below.

**Mobility:** Achieving desired mobility standards especially during handoff degrades network performance and security. Therefore, maintaining mobility along with performance and security is challenging.

**Reliability:** The network involved during handover must provide reliable services during handover. It must be robust in nature and ensure effective data delivery without any data loss

**Seamless Connectivity:** As the mobility increases, service continuity degrades. Therefore, to ensure always best connected (ABC) concept, handoff must be seamless such that the transition remains transparent to the user.

**Load Balancing:** Increased load over network degrades the network efficiency. Overloaded network results in handover failure, incurs communication latency and handoff delays and increase interference in microcellular networks. Therefore, minimal network load is desirable during handover.

**Security:** To prevent the terminals from malicious attacks, encryption and authentication mechanisms are needed. However, incorporating security credential in power-constrained scenario degrades network performance. Hence, maintaining security and network performance simultaneously is quite challenging.

**Handover Delays:** Minimizing handoff delay is quite challenging especially in delay sensitive applications. It degrades the throughput and QoS standards. Various latency reduction handoff techniques like Pre-active scanning, Selective channel scan, zone-based interleaved scanning [1] have been proposed to reduce handoff delays.

**Communication Latency:** Communicational latency occurs due to different operational standards of the networks. It degrades the performance. So it should be minimized

**Enhancing Interoperability:** Under highly dynamic and pervasive environment, the entities operate under different protocol standards. Therefore, interoperability must be enhanced to increase the overall performance

**Enhancing QoS:** QoS is the most important factor which affects the network performance and throughput. For enhancing the QoS factors, user mobility and network conditions must be handled properly.

**Enhancing QoE:** QoE can be defined as the factor for measuring the level of user satisfaction. Achieving end user satisfaction is yet another primary goal of any communication technology. Hence, it must be equally considered during handover.

**Ping-Pong Effect:** Ping-pong effect occurs in microcellular networks. It reduces the throughput and user satisfaction. So reducing the ping-pong effect is desirable.

**Corner Effect:** Corner effect occurs due to loss of signal strength when mobility of the terminal changes frequently. It is hard to detect. Hence, it is desirable to reduce the corner effect.

**Reduction in Interference:** Interference increases due to unfavourable environmental conditions like rain, wind, etc., and also with the mobility. It degrades QoS especially in multimedia and voice applications. So, interference must be reduced for achieving better throughput.

**Scalability:** In ubiquitous networking, the terminal nodes change their position frequently. Due to this, topology of network also changes frequently. Therefore, handover algorithms must be scalable to meet the real-time requirements.

**Better Resource Utilization:** Unmanaged resources cause handover failures and reduce network performance. However, resource utilization in highly mobile and power-constrained atmosphere is quite challenging.

**Maximizing Battery Lifetime:** The wireless networks are highly power constrained. So maximizing battery lifetime is highly desirable to provide seamless services during handoff.

**Maximizing Network Performance:** Network Performance can be evaluated over various factors like throughput, QoS, user satisfaction, fault tolerance, etc. However, in ubiquitous environment, maintaining performance standards is quite challenging. So handoff mechanisms should aim at maximizing the performance.

**Optimizing Handoff Techniques:** As the diversity and complexity of the communication network increase, the performance of the network degrades. Therefore, it is desirable to optimize the existing handoff techniques such that the system remains unaffected.

## 3  Parameters Affecting the Performance During Handover

Handover metrics are the parameters which decide the need of handover. Based on cause of handover and the frequency of occurrence, these metrics can be either dynamic or non-dynamic [2]. Moreover, these metrics can also be classified as network related, terminal related, user related or service related [3].

### 3.1  Dynamic Metrics

The value of these metrics changes frequently and therefore greatly affects the handover-related decisions. Some of the prominent dynamic factors are discussed below.

**Capability of Network:** Different networks possess different capacity in terms of bandwidth support, protocol support, interoperability standards, etc. Network capacity affects the performance; therefore, handoff decisions are highly dependable over capabilities of underlying networks.

**Network Conditions:** Handover performance is greatly affected by the dynamic changes occurring in the vicinity. So, network topology and underlying conditions play vital role during handoff decision-making.

**Network Security:** Security is the major issue that needs to be considered during handoff process to prevent malicious attacks. Security policies regarding integrity, authorization, authentication, confidentiality and resource modification must be infused properly during handover decision phase.

**Network Throughput:** Network throughput is the indicator of successful data delivery. It is therefore one of the major factors affecting the handover process and needs fair consideration.

**Traffic Balancing:** Frequent variations in network loads reduce the traffic-carrying capacity of the cells and degrade the QoS standards. So network load needs to be considered during handover.

**Bandwidth:** Higher bandwidth results in lower call dropping and reduced call blocking. Hence, higher bandwidth ensures better throughput, seamless connectivity and better QoS during handover.

**Received Signal Strength (RSS):** Received signal strength (RSS) plays significant role in reducing ping-pong effect. Lower values of RSS increase network load whereas higher value results in higher call dropping. Hence, optimal RSS value is desirable during handover initiation phase.

**Signal-to-Noise Ratio (SNR):** This ratio is affected by environmental factors like wind temperature, rain, etc., and is desirable when signal power is greater than noise power.

**Velocity:** In microcellular networks, higher velocity leads to frequent handovers increasing the overall handover counts. Hence, velocity of terminal needs to be considered for effective handover.

**Quality of Service (QoS):** QoS signifies the standard of service provided by the network. QoS levels certify the network performance. So higher value of QoS is desirable during handover.

**Handover Latency:** Handover latency affects the QoS factors and degrades the network throughput and performance. Hence, minimal value of handover latency is desirable.

**Handover Failure:** Lack of resources at target station and mobility are the main causes for handoff failure. Handover failure degrades QoS parameters and QoE standards. Therefore, minimization of handoff failure is desirable. To prevent failures, handover prioritization schemes are proposed [4] which reduce the failure probability to enhance performance and throughput.

**Handover Counts:** Minimal value of handover counts is desirable as frequent handovers increase ping-pong effect and cause resource wastage thereby degrading network performance.

**Occurrence of Unnecessary Handovers:** Unnecessary handovers result in ping-pong effect which further incurs both communication and handoff latency overheads. Therefore, unnecessary handovers must be minimized.

**Interference Prevention:** Interference during handover is highly unfavourable as it degrades the QoS standards which consequently reduce user satisfaction levels. So, interference prevention is highly desirable during handoff.

## 3.2   Non-Dynamic Metrics

Contrary to dynamic metrics, these metrics change very seldom and therefore having lesser impact over handover mechanism. Few of the non-dynamic factors are listed below.

**User Preferences:** Users may have various options as per their preferences and application requirements. Since user preferences influence QoE standards, it should be considered during handover.

**Power Consumption:** Handover process incurs battery consumption during decision phase due to interface activations. Hence, power consumption must be considered to improve energy efficiency and network performance.

**Network Cost:** It refers to the overall cost of accessing the network during handover. It is calculated on the basis of call arrival rates using a cost function. Least network cost ensures better QoS parameters. So, it is an important parameter to be considered during handover

## 4   Overview of Existing Handover Mechanisms

Handoff performance is one of the leading matters of concern prevailing in the realm of ubiquitous networking. Handoff classification is a well-heeled hierarchal structure founded over various perspectives. Considering network domain, handoff can be either horizontal handoff or vertical handoff [5]. Former type occurs between two base stations operating under same wireless technologies whereas latter occurs if the base stations operate under different wireless technologies. Moreover, horizontal handoff can be further classified as link-layer handoff and intra-system handoff [3].

During handoff, the terminal node breaks connection with the existing base station and gets attached to the target base station. Regarding connections during handoff [5], if the terminal node connects to only one base station at a time, it is termed as hard handoff. However, if the terminal maintains connection with more than one access point simultaneously, it is known as soft handoff. Former type is also called break before make handoff (BBM) as the old connection breaks new connection is set whereas latter type is called make before break handoff (MBB) as new connection is set before breaking the old connection. Moreover, as a special case of soft handoff, if connections can be switched only over the links belonging to same access points, it is termed as softer handoff [6]. Regarding frequency involved, if handoff occurs across access points operating on same frequency, it is called intra-frequency handoff. However, it is termed as inter-frequency handoff when handoff occurs between access points operating on different frequencies.

With viewpoint of administrative domain, handoff can be intra-administrative, if the participating networks are managed by the same administrative domain, or

inter-administrative if participating networks are governed by different administrative domains. Similarly, for subnetting involved, handovers can be either inter-subnet or intra-subnet [7].

From decision-making perspective, handover can be network-controlled handover (NCHO) if network controls decision-making, mobile-controlled handover (MCHO) [3], if terminal takes handoff-related decisions, mobile-assisted handover (MAHO), if terminal node gives support to the network in decision-making and network-assisted handover (NAHO) and if handoff decisions are taken by terminal under network assistance.

Based on necessity, handoff can be obligatory when handover is unavoidable or voluntary, if occurrence of handover is optional. Furthermore, as per user control allowance, handover can be proactive [5] if user has the authority of decision-making or passive if user has no control over occurrence of handover.

References [2, 3] discussed the mechanism involved in vertical handover. The various phases involved in VHO mechanism are briefly explained below.

**Network Discovery Phase:** In this phase, handoff process initiates and the terminal gathers information about the adjoining networks to search for networks that can be used for handoff.

**Network Analysis Phase:** The partaking networks advertise their data rates and QoS parameters; on the basis of which terminal decides which network should be selected.

**Network Selection Phase:** It is decision-making phase where current network decides when and where the handoff should be triggered.

**Handover Execution Phase:** In this phase, network-binding process starts, and connections are re-routed to new network ensuring seamless connectivity. It also deals with security issues and handles user's context (Fig. 1).

The categorization of various existing vertical handoff decision (VHD) algorithms is presented in [3, 5, 6, 8]. Based on this, the various VHD algorithms are discussed below.

**RSS-Based Algorithms:** These algorithms consider received signal strength (RSS) for deciding the need of handoff. These algorithms achieve minimal handoff failures, reduce unnecessary handoff and connection breakdowns, but suffer an increase in handoff delay.

**Bandwidth Based Algorithms:** Bandwidth serves as the basic criteria for handoff in these algorithms. These methods increase the throughput, reduce handover latency and provide effective load balancing. However, these algorithms suffer ping-pong effect.

**Cost Function-Based Algorithms:** These algorithms follow a priority-based approach where a cost function is evaluated over handoff metrics to select the best suited network. These algorithms increase user satisfaction and overall network throughput, reduce handoff blocking probability and handoff delay.

**Fig. 1** Handover mechanism



**Context Awareness-Based Algorithms:** For effective decision-making, these algorithms gather user-specific as well as application-specific context information. These algorithms improve system flexibility and efficiency. However, these algorithms are time-consuming and cause computational delays.

**Multiple Attribute Decision-Making (MADM) based Algorithms:** These algorithms evaluate various handover parameters for effective decision-making, improvising network throughput and maximizing resource utilization.

**Combination Algorithms:** These algorithms are based on fuzzy logic or artificial neural networks (ANN) to ensure effective handover decision-making. These algorithms increase the network performance by reducing unnecessary handovers and avoid ping-pong effect.

**User-Centric Algorithms:** These algorithms focus on user preferences for maximizing QoE factors during handoff. For attaining highest level of user satisfaction, these algorithms consider application-specific details along with user preferences.

**QoS-based Algorithms:** These algorithms aim at achieving better quality standards. These algorithms enhance network performance and throughput by reducing handoff latency, packet loss, signalling overhead and other computational costs.

**Authentication based Algorithms:** Security is the main criteria for the algorithms of this category. These algorithms infuse encryption algorithms with authentication mechanism to achieve highly secured handovers thereby increasing the robustness

and reliability of network. However, due to security overheads, handoff latency increases resulting in handoff failures and ping-pong effect.

## 5  Discussion and Analysis

During handover, decision-making complexity increases with the mobility of nodes thereby affecting the network performance and throughput. The simulation in [9] enhanced the performance of handover by merging the context information during decision-making. It did not consider the trajectory. However, for real-time scenario, the trajectory and context information can be incorporated in the algorithm through machine learning approach.

In [10], the algorithm resulted in higher bandwidth, lower cost and lower power consumption. However, mobility and security issues can be considered to further enhance the performance. The three-level optimization scheme in [11] reduced the computational complexity thereby increasing the handover accuracy. It achieved better network performance as compared to traditional techniques. However, the algorithm can be further enhanced by incorporating lookup tables for hysteresis values based on velocity and condition of network.

MPA [7] proved to be highly secured, optimized and loosely coupled to network layer and application layer mobility management protocols. It achieved significant reduction in handover delays and increased the quality of application performance. However, performance can be further enhanced by formulating a global interface for faster network selection.

BIGAP architecture in [12] achieved network performance along with mobility and efficient load balancing. However, the overall throughput degraded to some extent as the architecture was prone to collisions. Hence, a collision-free channel assignment is needed to increase the throughput of the network. The handover management model presented in [13] provided unified framework handover optimization and network analysis. It achieved favourable network performance under various mobility models and therefore proved favourable for real-time scenario.

The stochastic model presented in [14] implemented two stage optimal solutions for enhancing handover performance in multi-mobility scenario. It resulted in reducing the handover counts when implemented under different traffic models. The decision algorithm proposed in [15] modelled vehicular mobility for achieving better connectivity. It incorporated user preferences as well as application requirements to enhance handover decision-making. However, the performance can be further enhanced by adopting VANET technologies and related protocols.

The simulations done at [16] proved the effectiveness of soft handovers over hard techniques. It provided better connectivity and QoS standards in both MCHO and NCHO configurations. In [17], the proposed handover mechanism dynamically balanced the load over microcellular networks thereby reducing the interference during handover. User mobility and handover overheads were well incorporated to enhance the network throughput. The simulation results in [18] showed that the

proposed algorithm enhanced the QoS as well as QoE standards at minimal communication cost. It thereby ensured smooth handover through effective load balancing. The scanning algorithm in [19] proved to be effective in reducing energy consumption during handoff. It proved highly energy efficient in unpredictable mobility scenario, thereby making it suitable for microcellular communication. However, self-optimization techniques can be further incorporated to enhance the adaptive nature of the algorithm.

The algorithm in [20] enhances network performance in terms of bandwidth, handover delay, throughput and reliability. It reduced the signalling overload and provides QoS support for mobility model, thereby enhancing the QoS and QoE parameters. The algorithm presented in [21] effectively reduced the battery drainage. Moreover, battery enhancer software can be developed by different vendors to support battery life at the terminal end.

The handover algorithm in [22] improved the QoS standards by reducing handoff latency and packet loss during handover thereby achieving higher throughput. The MAC protocol in [23] considerably reduced handoff latency and packet loss. Higher throughput and live monitoring were achieved by implementing burst transmission along with seamless handoff. The wireless architecture presented in [24] effectively reduced the handoff latency as well as uplink delay. The architecture was further analysed for queuing delays in X2 packets. It was found backward compatible. The handoff-aware proactive approach for congestion control presented in [25] reduced the handover delay occurring at transport layer. Due to cross-layer assistance, the algorithm performed well in multiple handover scenarios thereby achieving seamless mobility.

Although the mobility framework proposed in [26] reduced communication as well as handoff latency, the performance can still be enhanced by exploiting the storage and operational capabilities of mobile nodes. The tunnel establishment approach for mobility management in [27] reduced the handoff latency and efficiently mitigated the signalling burden to enhance network throughput and performance. The simulation study in [28] revealed that the proposed algorithm reduced the processing delay and handoff delay during handover. It minimized the packet loss thereby increasing the TCP performance during handover in satellite networking.

Existing handoff techniques cause considerable amount of delay thus affecting the throughput. Authentication mechanisms cause delay and degrade the network performance. Hence, maintaining performance along with security is quite challenging. Pair-hand authentication protocol proposed in [29] achieved network efficiency along with security. It further incorporated batch signature verification scheme to increase the reliability of network, making it more feasible for real-time applications.

In [24], an open access network architecture was proposed to achieve interoperability in a secured fashion. The analytical discussion over authentication protocol in [30] revealed that EAP-ERP provided integrity but lacks in trust relationship. Contrary to this, EAP-CRA provided better integrity and trust relationship.

Context awareness approach is capturing attention day by day. Tamijetchelvy et al. [22] focussed on multiple criteria-based context aware handover algorithms to improve network performance. However, as the handover delay increased, there is a need of developing unified VHD algorithms which can effectively handle the complexity of mobile networking. The algorithm in [10] was computationally cheaper and free from ranking abnormalities. Furthermore, it proposed that to enhance the QoE levels, mobile communication networks need to be interconnected.

## 6    Conclusion

The main objective of this survey is to explore the potential handover challenges for providing framework for further research in the field of new generation wireless networks (NGWNs). The paper presented a brief overview of the basics of handover, their classification and behavioural pattern. It discussed the various existing handoff algorithms, their outcomes and limitations to provide path for future research. Various existing issues in handoff are also discussed to find out the potential still existing in this field. The survey presented in this paper will provide basic guidelines to researchers for devising novel handover algorithms for enhancing the performance of NGWNs.

## References

1. Tetarwal, M.L., Kuntal, A., Karmakar, P.: A review on handoff latency reducing techniques in IEEE 802.11 WLAN. In: National Seminar on Recent Advances in Wireless Networks and Communications, NWNC-2014. Published in International Journal of Computer Applications (IJCA)
2. Fachtali, I.E., Saadane, R., Koutbi, M.: A survey of handovers decision algorithms for next generation wireless networks. Int. J. Adv. Res. Comput. Commun. Eng. 4(1), 159–165 (2015)
3. Seth, A.: Vertical handoff decision algorithms for next generation wireless networks: some issues. In: International Journal of Advanced Research in IT and Engineering (IJARIE), Vol. 2, No. 8. ISSN: 2278-6244 (2013)
4. Sgora, A., Dimitrios, D., Vergados, D.D.: Handoff prioritization and decision schemes in wireless cellular networks: a survey. IEEE Commun. Surv. Tutor. 11(4), 57–77 (2009)
5. Chandavarkar, B.R., Reddy, G.R.M.: Survey paper: mobility management in heterogeneous wireless networks. In: International Conference on Communication Technology and System Design (2011)
6. Ravichandra, M., Gowda, H.N., Kumar, C.A.U.: A survey on handovers literature for next generation wireless networks. Int. J. Adv. Res. Comput. Commun. Eng. 2(12), 4671–4677 (2013)
7. Dutta, A., Famolari, D., Das, S.: Media-independent pre-authentication supporting secure interdomain handover optimization. IEEE Wirel. Commun. 15, 55–64 (2008)
8. Rajule, N., Ambudkar, B., Dhande, A.P.: Survey of vertical handover decision algorithms. Int. J. Innov. Eng. Technol. 2(1), 362–368 (2013)

9. Guidolin, F., Pappalardo, I., Zanella, A., Zorzi, M.: Context-aware handover policies in HetNets. IEEE Trans. Wirel. Commun. **15**(3), 1895–1906 (2016)
10. TalebiFard, P., Leung, V.C.M.: A dynamic context-aware access network selection for handover in heterogeneous network environments. In: IEEE (2011)
11. Fischione, C., Athanasiou, G., Santucci, F.: Dynamic optimization of generalized least squares handover algorithms. IEEE Trans. Wirel. Commun. **13**(3), 1235–1249 (2014)
12. Zubow, A., Zehl, S., Wolisz, A.: BIGAP—seamless handover in high performance enterprise IEEE 802.11 networks. In: 15th IEEE/IFIP Network Operations and Management Symposium (IEEE NOMS, 2016)
13. Nguyen, V.M., Chen, C.S., Thomas, L.: A unified stochastic model of handover measurement in mobile networks. IEEE/ACM Trans. Netw. **22**(5), 1559–1576 (2014)
14. Fang, B., Zhou, W.: Handover reduction via joint bandwidth allocation and CAC in randomly distributed HCNs. IEEE Commun. Lett. **19**(7), 1209–1212 (2015)
15. Marquez-Barja, J.M., Ahmadi, H., Tornell, S.M., Calafate, C.T., Cano, J.C., Manzoni, P., DaSilva, L.A.: Breaking the vehicular wireless communications barriers: vertical handover techniques for heterogeneous networks. IEEE Trans. Veh. Technol. **64**(12), 5878–5890 (2015)
16. Vegni, A.M., Natalizio, E.: A hybrid (N/M)CHO soft/hard vertical handover technique for heterogeneous wireless networks. Elsevier (2013)
17. Wang, Y., Haas, H.: Dynamic load balancing with handover in hybrid Li-Fi and Wi-Fi networks. J. Lightw. Technol. **33**(22), 4671–4682 (2015)
18. Sarma, A., Chakraborty, S., Nandi, S.: Deciding handover points based on context-aware load balancing in a WiFi-WiMAX heterogeneous network environment. IEEE Trans. Veh. Technol. **65**(1), 348–357 (2016)
19. Vondra, M., Becvar, Z.: Distance-based neighborhood scanning for handover purposes in network with small cells. IEEE Trans. Veh. Technol. **65**(2), 883–895 (2016)
20. Kadusic, E., Zivic, N., Kos, A.: QoS-aware dynamic MAP selection in HMIPv6 architectures. In: IEEE ACCESS (2016)
21. Habeib, M.Z., Elsayed, H.A.E., Elramly, S.H., Ibrahim, M.M.: Battery Based Vertical Handover Between WiMAX and WLAN Technologies. In: IEEE (2011)
22. Tamijetchelvy, R., Sivaradje, G., Sankaranarayanan, P.: Dynamic MAPT approach for vertical handover optimization in heterogeneous network for CBR and VBR QoS guarantees. In: International Conference on Information and Communication Technologies (ICICT) (2014)
23. Dargie, W., Wen, J.: A seamless handover for WSN using LMS filter. In: 39th Annual IEEE Conference on Local Computer Networks (2014)
24. Mukhopadhyay, A., Das, G.: A ring-based wireless optical network to reduce the handover latency. J. Lightw. Technol. **33**(17), 3687–3697 (2015)
25. Sinky, H., Hamdaoui, B., Guizani, M.: Proactive multipath TCP for seamless handoff in heterogeneous wireless access networks. IEEE Trans. Wirel. Commun. **15**(7), 4754–4764 (2016)
26. Wang, X.: A mobility frame for 6LoWPAN WSN. IEEE Sensors J. **16**(8), 2755–2762 (2016)
27. Kim, M.S., Lee, S.: Enhanced network mobility management for vehicular networks. IEEE Trans. Intell. Transp. Syst. **17**(5), 1329–1340 (2016)
28. Hu, H., Yuan, D., Liao, M., Liu, Y.: Packet cache-forward method based on improved bayesian outlier detection for mobile handover in satellite networks. China Commun. **13**(6), 167–177 (2016)
29. He, D., Chen, C., Chan, S., Bu, J.: Secure and efficient handover authentication based on bilinear pairing functions. IEEE Trans. Wirel. Commun. **11**(1), 48–53 (2012)
30. Ramezani, K., Sithirasenan, E., Su, K.: Formal security analysis of EAP-ERP using casper. IEEE Access **4**, 383–396 (2016)

# Defected Ground Structure Switchable Notch Band Antenna for UWB Applications

**Vamseekrishna Allam and B. T. P. Madhav**

**Abstract** A defected ground structure notch band antenna is proposed in this work. The switchable characteristics for the designed notch band antenna are achieved through open-end slots on/off positions. The proposed DGS notch band antenna is capable of notching the frequency bands 3–4, 5.5–6.5 GHz, respectively. A high notch band rejection with VSWR greater than 2 and the return loss greater than −10 dB is achieved at the notching frequencies. The defected ground structure is providing balance in the impedance bandwidth to the designed models. By sorting the slots on the radiating structure, the tunability in the notching frequencies is attained in this paper. The antenna radiation characteristics and the surface current distributions at operating bands as well as at notch bands are presented in this work. The proposed notch band antenna is providing high rejection of gain in the notch band and average gain of 2.8 dB in the operating band.

**Keywords** Defected ground structure · Monopole antenna · Notch band Switchability · Ultra-wideband

## 1 Introduction

The accelerated generation of wireless communication systems has constituted a demand for reconfigurable or tunable filters and antennas [1–3]. These type of devices reduces the compulsion to yield extravagant charges correlated with the refitting of wireless infrastructures, considering that an adjustment in the frequency, bandwidth, or other conditions of the hardware can be attained over electronic/ mechanical reconfiguration. Reconfigurable equipment further produces the appropriate hardware for a highly capable management and adoption of a spectrum

V. Allam (✉) · B. T. P. Madhav
Department of ECE, K L University, Guntur, Andhra Pradesh, India
e-mail: vamseekrishnaa3@kluniversity.in

B. T. P. Madhav
e-mail: btpmadhav@kluniversity.in

through the theory of dynamic spectrum approach and cognitive radio [4–8]. Utilization of reconfigurable devices again allows the time allocation of hardware, which in turn advantages to mass and size contraction of the communication system. This is a great condition in compact devices and of significant influence in satellite communication systems [9–12].

UWB is a radio technology that can use a very low energy level for short-range, high-bandwidth communications over a large portion of the radio spectrum. Ultra-wideband is a technology for transmitting information spread over a large bandwidth ranging from 3.1 to 10.6 GHz. Ultra-wideband characteristics are well suited to short-distance applications, such as peripherals of PC, personal area network (PAN), medical and radar imaging. Due to short duration of UWB pulses, it is easier to provide high data rates [13]. Ultra-wideband characteristics are well suited to short-distance applications, due to its low emission levels permitted by regular agencies. The main challenge while designing UWB antennas is to achieve good band notch characteristics with reduced size and cost [14]. Band notch filtering functionally can be achieved using various methods such as using radiating patch with different shapes like rectangular, circular, hexagonal and with different slots, such as V slot, U slot, W slot, inverted U slot, and elliptical slot.

In this paper, notch band antenna model is designed, and switches are incorporated in the slots to fine-tune the notch band frequencies. Monopole structure with defected ground model is proposed in this work. Frequency tunability with the on and off conditions of the switches at slots is examined and presented in this work. The antenna modeling is carried with CST microwave studio electromagnetic tool, and the parametric results are analyzed for the optimized dimensions in this work. Simulation-based results are presented, and analysis is well organized for understanding the performance characteristics of the designed antenna.

## 2   Antenna Design and Geometry

A notch band monopole antenna with inverted u-shaped radiating element is presented in this work. A U slot is also placed on the feed line to achieve notch band characteristics. The backside of the antenna consists of partial ground and an open-ended stub for impedance matching. The radiating element consists of another slot on the lower edge nearer to feed line. Depending on the slot on and off conditions, the notch band characteristics are analyzed. The geometric aspects of the designed antenna models are presented in Table 1. The proposed dual notch band antenna is prototyped on FR4 substrate with dielectric constant 4.4 and loss tangent 0.02.

**Table 1** Antenna dimensional characteristics in mm

| Antenna parameter | Dimensions (mm) | Antenna parameter | Dimensions (mm) |
|---|---|---|---|
| $L_s$ | 34 | $L_2$ | 4.4 |
| $W_s$ | 34 | $L_3$ | 1.6 |
| $L_f$ | 11.7 | $L_4$ | 9.87 |
| $W_f$ | 2.84 | $L_5$ | 60.02 |
| H | 1 | $L_6$ | 9.4 |
| $W_1$ | 13.8 | $L_7$ | 7.3 |
| $W_2$ | 2.2 | $L_8$ | 3.2 |
| $W_3$ | 0.4 | $L_9$ | 3.25 |
| $W_4$ | 0.5 | $L_{10}$ | 11.9 |
| $W_5$ | 0.6 | $E_s$ | 0.05 |
| $L_1$ | 4 | AR | 2 |
| G | 0.8 | $r_2$ | 6 |
| $G_x$ | 0.2 | $r_x$ | 2 |
| $r_1$ | 5 | $r_{x1}$ | R1/AR |
| $r_{x2}$ | $r_2$/AR | | |

## 3  Results and Discussion

The designed models are analyzed using commercial electromagnetic tool CST, and the results are presented in this section. The reflection coefficients of the designed models are presented in this work to know the operating bands and notch bands for the antenna models. Four models are designed in this work, and by placing switches on the slots, we examined the change in resonant frequency and the notch bands. When switch s1 and switch s2 are in off condition, then a single notch band can be observed at 5.8–6.2 GHz from Fig. 1. When s1 is off and s2 is on, then the notch band is widened, and the antenna is notching the band from 4.5 to 6.5 GHz. The same is the case with s1 on and s2 off. When both the switches are in on condition, the proposed model is notching dual band (3.2–4 GHz) and (5.5–6.5 GHz).

The band rejection is also very high at these two notch bands. Figure 2 is providing similar information like VSWR with respect to notch bands and operating bands. The on and off conditions of the switches lead to shift in resonant frequency and notch bands.

Parametric results of VSWR by changing radius r1 from 5 to 5.5 mm are presented in Fig. 3. It is being observed that there is no change in the notch band from 5.5 to 6.5 GHz, but at fundamental notch band, a shift of 0.5 GHz is obtained with change in r1.

The far field radiation of the dual notch band antenna with two switches in on condition is presented in Fig. 4. The three-dimensional and polar coordinates-based radiation pattern of the antenna is also shown in Fig. 4. Antenna is producing gain more than 5 dB in the operating band and negative values in the notch band.

(a) Antenna Model 1                    (b) Antenna Model 2

(c) Antenna Model 3                    (d) Antenna Model 4

(e) Ground plane

**Fig. 1** Notch band antenna model geometry

The fabricated prototype of the proposed antenna is shown in Fig. 5. The measured return loss plot is shown in Fig. 5c. The plot shows the notch band from 4.6 to 6.46 GHz and another notch band obtained at 6.9–7.6 GHz band. The

**Fig. 2** VSWR versus frequency



**Fig. 3** Parametric analysis with change in r1



**Fig. 4** Radiation pattern of dual notch band antenna at operating frequency 4.5 GHz

**(a)** Front view          **(b)** back view



**(c)** Measured result of return loss

operating band of the antenna extends up to 12.43 GHz which shows an ultra-wide bandwidth with notch band characteristics. The shift in notch band frequencies response is obtained by measured pattern which can be due to connection losses.

## 4  Conclusion

A notch band antenna has been designed with cutting slots and analyzed with monopole configuration in this paper. Antenna with slots is demonstrated with on and off conditions to fine-tune the notch band in this work. All the antenna parameters with respect to switching are analyzed and presented. The radiation characteristics and their corresponding directivity and gain parameters are simulated and presented. Antenna with closed switches is providing gain more than 5 dB in the operating band and antenna with open switches providing gain more than 5.2 dB. Time domain analysis shows the pulse fidelity value more than 0.5 and VSWR < 2 in the operating band. The designed model is simple in structure and easy to fabricate and satisfying the UWB communication requirements.

# References

1. Dastranj, A., Abiri, H.: Bandwidth enhancement of printed E-shaped slot antennas fed by CPW and micro strip line. IEEE Trans. Antenna Propag. **58**, 1402–1407 (2010)
2. Lottici, V., D'Andrea, A., Mengali, U.: Channel estimation for ultra-wideband communications. IEEE J. Sel. Areas Commun. **20**(9), 1638–1645 (2002)
3. Schantz, H.: The Art and Science of Ultra-Wideband Antennas. Artech House, Norwood, MA (2005)
4. Ojaroudi, N.: Bandwidth improvement of monopole antenna using p-shaped slot and conductor-backed plane. Int. J. Wirel. Commun. Netw. Mobile Comput. **1**, 14–19 (2014)
5. Ojaroudi, N.: Compact UWB monopole antenna with enhanced bandwidth using rotated L-shaped slots and parasitic structures. Microw. Opt. Technol. Lett. **56**, 175–178 (2014)
6. Chung, K., Yun, T., Choi, J.: Wideband CPW fed monopole antenna with parasitic elements and slots. Electron. Lett. **40**(17), 1038–1040 (2004)
7. Zeng, X., He, J., Wang, M., Abdulla, M.: New closed-form formula for series inductance and shunt capacitance based on measured TDR impedance profile. IEEE Microw. Wirel. Comput. Lett. **17**, 781–783 (2007)
8. Raman, Y.S.V., Madhav, B.T.P., Mounika, G., Sai Teja, K., Sai Kumar, S.B.V.N., Sri Harsha, K.: Analysis of circularly polarized notch band antenna with DGS. ARPN J. Eng. Appl. Sci. **11**(17), 10140–10150 (2016)
9. Madhav, B.T.P., Kumar, K.V.V.: Analysis of CPW fed step serrated ultra-wide band antenna on Rogers RT/Duroid substrates. Int. J. Appl. Eng. Res. **9**(1), 53–58 (2014)
10. Mehdipour, A., Mohammadpour-Aghdam, K., Faraji-Dana, R.: Complete dispersion analysis of vivaldi antenna for ultra-wideband applications. Progr. Electromagn. Res. PIER **77**, 85–96 (2007)
11. Jang, Y.W.: Experimental study of large bandwidth three-offset micro strip line-fed slot antenna. IEEE Microw. Wirel. Comput. Lett. **11**, 425426 (2001)
12. Ojaroudi, M., Ojaroudi, N., Ghadimi, N.: Dual band-notched small monopole antenna with novel coupled inverted U-ring strip and novel fork-shaped slit for UWB applications. IEEE Antennas Wirel. Propag. Lett. **12**, 182–185 (2013)
13. Almalkawi, M.J., Devabhaktuni, V.K.: Quad band-notched UWB antenna compatible with WiMAX/INSAT/lower-upper WLAN applications. Electron. Lett. **47**(19), 1062–1063 (2011)
14. Zhang, Y., Hong, W., Yu, C., Kuai, Z.-Q., Don, Y.-D., Zhou, J.-Y.: Planar ultra wideband antennas with multiple notched bands based on etched slots on the patch and/or split ring resonators on the feed line. IEEE Trans. Antennas Propag. **56**(9), 3063–3068 (2008)

# IoT-Based Multimodal Biometric Identification for Automation Railway Engine Pilot Security System

K. Sujatha, R. S. Ponmagal, K. Senthil Kumar, R. Shoba Rani and Golda Dilip

**Abstract** Railways are the most convenient mode of transport, but safety precaution is lagging. Train accidents, due to an unknown person operating the engine, will lead to the end of many lives and also loss of railway property. The optimal solution to meet this problem here proposes the effective system of "Automation of Railway Engine Pilot Security System using Multimodal Biometrics Identification" (AREPSS using MBI). Iris and Fingerprint inputs are given by engine pilot from cabin to control room using Internet of things (IoT). In control room, identifications take place by fusing the inputs and then pass the decision signal to automatically start the engine. The common unimodal biometric system can be seen in most of the places due to its popularity. Its reliability has decreased because it requires larger memory footprint, higher operational cost, and it has slower processing speed. So, we are introducing multimodal biometric identification system which uses iris and fingerprint for security reason. The major advantage of this several modality method is that as both modalities utilized the same matcher component, the reminiscence footprint of the system is reduced. High performance is achieved by integrating multiple modalities in user verification and identification causing high dependability and elevated precision. So this procedure improves the safety in engine and thus helps in saving lives and property.

**Keywords** Iris · Fingerprint · Unimodal · Multimodal · Wavelet transform Neural network · Internet of things and biometric

K. Sujatha (✉) · R. S. Ponmagal · K. Senthil Kumar · R. Shoba Rani · G. Dilip
Department of EEE/CSE/ECE, Center for Electronics Automation and Industrial Research (CEAIR), Dr. MGR Educational and Research Institute, Chennai, India
e-mail: drksujatha23@gmail.com

# 1   Introduction

Our Indian railway is the third largest, and to control the network of whole railway is more complicated. So there is a need to increase the security in railways. Automation is one of the best advanced solutions to increase the security in railways. Here introducing one of the automation approaches for pilot security in railway engine using effective multimodal biometric recognition [1, 2].

The basic idea is that the iris and fingerprint biometric inputs are given by the driver from the engine cabin to the control room. In the control room, the verification process of fusing both the inputs takes place and verifies the detail about the driver and then passes on the signal for the engine to operate. Misidentification rate of iris is very less. So here fingerprint is fused with iris so as to accelerate the performance [3, 4]. Here password-based authentication is replaced by multimodal biometric authentication.

# 2   Existing System

Presently, pilot security in railway engine is not automated. According to the ongoing trend, the driver inserts the key which in turn produces the required voltage necessary to run the train [5]. The driver and the guard are the decision makers while the train is running, though they get the required signal from the control room. So when an unknown person operates the engine, it leads to the end of many lives and also the loss of railway property [6]. So there is a need to increase the security system in railways.

## 2.1   Problems in Existing System

The problems in the existing system include manual control, high risk for train passengers, lose of railway property and less security because pilot is the decision maker to operate the engine [7]. These problems can handled if IoT based railway automation scheme is introduced.

## 2.2   Existing Methodologies

Iris and fingerprint biometric authentication technique is more effective to improve the pilot security. The existing biometric identification techniques are unimodal biometric identification and multimodal biometric identification with two complete unimodal systems [8]. Unimodal biometric system consists of its own unique feature extractor and classifier. Its reliability is decreased because it requires

memory footprint, less accuracy, and it has slower processing speed. Figures 1 and 2 show the schematic representation for the unimodal and multimodal biometric systems, respectively.

## 2.3 Drawbacks in Existing Methodologies

It requires larger memory footprint, high operational cost, and multiple algorithms. Each unimodal system contains separate matcher (classifier), moderate fusion score level, and slower processing speed.

## 3 Proposed System Using IoT

The main idea behind the projected structure is to engage the fully illicit and automated pilot security in railway engine using fusion-based multi-factor identification system devoid of the accessibility of two single mode structures by incorporating the technique in Internet of things (IoT). Here iris and fingerprint is used as a part of biometric database. Fusing these two input images decrease the probability of hacking. The basis behind the proposed real-time system is schematically represented in Fig. 3.

First, the input is given by the driver from the engine cabin to the control room and in the control room, the verification process of fusing both the inputs, i.e., iris and fingerprint, takes place. For better safety, both the biometric inputs are send to control room. Thus, the priority is given to the control room operator, who checks in first and then verifies the detail about the driver and then passes on the signal for the engine to start. This process is repeated for each driver in break journey for long-distance covering trains, and whereas for short distance traveling like local trains, there are shifts for each drivers or they may be used only for one way journey.



**Fig. 1** Architecture for unimodal biometric

**Fig. 2** Architecture for existing multimodal biometric



**Fig. 3** Proposed automated pilot security system using IoT

## 3.1 Advantages of Proposed System

Automated checkout and loss prevention with safe journey for passengers eliminate the loss of railway property and high security and maintenance free.

## 3.2 Proposed Methodology

This proposed plan is used to defeat the problems occurred in existing methodologies and progress the pilot security in existing system. The block diagram for the proposed methodology is shown in Fig. 4. Euclidean distance or Hamming distance is used for matching the patterns with minimum distance of separation. Minimum distance indicates that the patterns are more similar to each other.

Primarily, fingerprint biometric input is obtained and approved by the fingerprint feature extractor. The processed indication is judged against the reference pattern in the database via the provided matcher. In the superseding time, the following input is acquired and sent to the iris feature extractor. In the time that the matcher completes the processing of the first biometric and generates the matching output, the second biometric input is processed and is ready for matching. The Euclidean distance matcher is now used to evaluate the iris biometric reference with the templates and generate the output. The fusion takes place once both matching scores are available. Proposed fusion-based multimodal biometric approach (Fig. 4) is for both



Fig. 4 Proposed multimodal biometric methodology

modalities utilized same classifier is that both output scores will be same format. In general, the proposed structure will advance the security among pilot in railways.

### 3.3  Advantages of Proposed Methodology

Eliminating additional normalization functions, improves the processing speed, reduces the memory footprint, simple design process, and both modalities utilize a single classifier.

## 4  Recognition Process

Implementation of multimodal biometric system consists of the following processes (1) Iris recognition process. (2) Fingerprint recognition process, and (3) Fusion process. The realization approach involves the process of acquiring a representation of the region containing the transcript, preprocessing that image, extracting (segmenting) the individual characters, describing the characters in a form suitable for computer processing, and recognition.

## 5  Result and Discussion

Identification of engine pilot is based on the classifier output and fusion score level. In this section, we are going to discuss simulation output for classifier using MATLAB. The details regarding the identification database and false accept error rates are provided in MATLAB coding for all the three processes. In this identification process, one fingerprint/iris image is selected as the verification image input from database. Then input image is matched against the entire database images. If the given input matched with the database, then the user is identified as an authorized person to start the engine, and the maximum matching score value will be displayed. If the given input is not matched with the database then the user is identified as unauthorized person, and the minimum matching score value is displayed with error percentage. The simulation results of fingerprint and iris are shown in Figs. 5 and 6.

When the fusion process score value is compared with the individual score values of both input images, it will improve the efficiency of the average score value. The simulation result of fusion process is shown in Fig. 7.

The simulation result of the proposed fusion-based multimodal biometrics identification clearly shows the improved recognition rate compared with existing identification method.

**Fig. 5** Simulation matching result for fingerprint



**Fig. 6** Simulation matching result for Iris

## 6 Conclusion

The proposed work in this paper has the concept of combining the features of both iris and fingerprint, and we can attain very high efficiency, and the performance is also improved. The major advantage of this multimodal biometric identification is that both modalities utilizes the same matcher, low cost with a small memory easier for hardware implementation. This recognition can be implemented in high security areas like airports, war fields, and ATM centers instead of using passwords.

**Fig. 7** Simulation result for fusion process



The simulation results clearly show that the proposed multimodal biometrics identification is more secure and efficient for automation of pilot security in railway engine. So this technique enhances high security in railway engine and thus saves lives and property.

# References

1. Sujatha, K., Pappa, N.: Combustion monitoring of a water tube boiler using a discriminant radial basis network. ISA Trans. **50**, 101–110 (2011)
2. Jain, A.K., Hong, L., Kulkarni, Y.: A multimodal biometric system using finger prints, face and speech. In: 2nd International Conference Audio and Video-based Biometric Person Authentication, pp.182–187. Washington, March 22–24 (1999)
3. Jain, A.k., Prabhakar, S., Chen, S.: Combing multiple matchers for a high security fingerprint verification system. **20**(11–13), 1371–1379 (1999)
4. Roli, F., Kittler, J., Fumera, G., Muntoni, D.: An experimental comparison of classifier fusion rules for multimodal personal identity verification system, pp. 76–82 (2002)
5. Lumini, A., Nanni, L.: When fingerprints are combined with iris—a case study: FVC2004 and CASIA. Inter. J. Net. Sec. 4, 27–34 (Jan 2007)
6. Nandakumar, K.: MultiBiometric systems: fusion strategies and template security. Ph.D. Thesis, Michigan State University (2008)
7. Masek, L., Kovesi, P.: MATLAB source code for a biometric identification system based on iris patterns. The School of Computer Science and Software Engineering, the University of Western Australia (2003)
8. Baig, A., Bouridane, A., Kurugollu, F., Qu, G.: Fingerprint iris fusion based identification system using single hamming distance matcher. Inter. J. Biosci. BioTech **1**, 47–57 (Dec 2009)

# Architectural Outline of GIS-Based Decision Support System for Crop Selection

**Preetam Tamsekar, Nilesh Deshmukh, Parag Bhalchandra, Govind Kulkarni, Vijendra Kamble, Kailas Hambarde and Vijay Bahuguna**

**Abstract** For Indian farmers, the crop selection decision is a very crucial task as number of factors need to be taken into consideration. To drive out from this situation, a solution is proposed which apply analytic hierarchy process (AHP) and GIS in terms of a crop selection decision support system in Indian scenario. The framework was outlined, created and executed over a selected farm. This paper portrays the implemented systems architectural framework.

P. Tamsekar · N. Deshmukh · P. Bhalchandra · G. Kulkarni · V. Kamble · K. Hambarde
S.C.S., Swami Ramanand Teerth Marathwada University, Nanded, Maharashtra 431606,
India
e-mail: pritamtamsekar@gmail.com

N. Deshmukh
e-mail: nileshkd@yahoo.com

P. Bhalchandra
e-mail: srtmunparag@gmail.com

G. Kulkarni
e-mail: govindcoolkarni@gmail.com

V. Kamble
e-mail: vijendrakamble5@gmail.com

K. Hambarde
e-mail: kailassrt@gmail.com

V. Bahuguna (✉)
Department of Geography, DBS PG College, Dehradun, Uttarakhand 248001, India
e-mail: vijaybahugunadbs@gmail.com

# 1   Introduction

Since long, the GIS (Geographical Information System) is based on computer system containing a spatial database mapped with geographic location. The key feature of GIS is the interpretation of geometric and thematic attributes of spatial data. GIS is capable of handling and organising huge scale of data including manipulations on it. In the current scenario, GIS has tremendous potential in applied use across agricultural, industrial, business logistics, etc. domains. Our underlined work cites use of GIS in agriculture domain. The countries on the top edge of agriculture development list make use of GIS, which helps to bridge the difference that the developing countries lack. The farmers in the developing nations heavily rely on the information passed on by the predecessors for selecting the crop for cultivation; on the other hand, the lack of resources or information, the facts based on scientific ground which could result in higher yield is overlooked. The hurdle to overcome and challenge for Indian agriculture sector is precision farming. With a specific end goal to compose the scientific process for selection of crop, based on expert's advice, the integration of GIS can come up with better yield efficiency. GIS which is considered as a boon for decision support system in various sectors can also play a vital role in agricultural sector for crop selection process; however, not a wide study has been done in support of crop selection's decision support system using GIS. It thus opens a wide opportunity for research in this sector. In the present age, the huge challenge around which the globe is revolving has the geographic dimensions such as overpopulation, natural disaster, deforestation; whereas if we see the local challenges, it also has the geographic components that can be visualised using the GIS technology [1]. Changing worldwide environment and increase in the demand of food requirement have also raised a concern over ensuring the food security. Moreover, the environment challenges are demotivating the man force from farming and which is resulting in monetary feasibility of agricultural practises. This all raised a concern for research in precision agricultural practise [2] (Fig. 1).

Frequently, we arrive at a point when we need to take a decision. Learning about the situation is the key to make decision. It is easy to take a decision on small issue, but when the issue or situation is complex, it is a tough call to conclude to a decision and we end-up to a decision with partial data or information. But on a contrary note, GIS makes the task much easier to take a decision call on complex scenario by using the spatial database and graphical or statistical products [3]. Subsequently, the interaction among different parameters is straightforward using GIS. In concern to the current theme to model a decision support system for crop selection, we require soil texture, biological details of soil, irrigation, climate, etc. The collected information all alone is of no use to make decision, as well as the geo-referenced locations are nothing if not plotted on a map to understand, with the use of map, it is easy to visualise the scenario. In the mid-nineteenth century, in London, to control the outbreak of cholera GIS was effectively used to take the decision [4].

**Fig. 1** GIS elements



## 2 Agriculture Pattern in Maharashtra

The primary and basic occupation of the Maharashtra people is agriculture. The area covered by cash crop and food crop is large in state.

## 3 System Outline

It is speculated that improvement in the productivity can be accomplished by utilising different parameters like soil texture, meteorology, irrigation in GIS-based decision support system. The research work comprises of the enlisted objectives like to draft a thematic cartography of soil, meteorology, irrigation, well inventory and road network followed by the development of a model framework for decision support system for precision farming and at the end to validate the approach through an applied case study.

During the investigation, we came across a contemporary work [5] that has explained the role of GIS and RS technologies in crop estimation. This supports for crop forecasting at local levels. The purpose of this system was to design a supporting system for crop estimation. It stated that the main objective of this study was to design the GIS system which estimates the crop area from ground location based on AFS and prototype system for the area estimation. For the accuracy of the system, there is cross comparison with the results with agricultural statistics information division (ASID) of the Ministry of Agriculture of I.R. of Iran from the same data. The result of this study has shown that spatial analysis has provided the

natural solution for the spatial problem, which support for decision-making management.

In another cited work [6], a GIS-based DSS framework was used for the optimisation and identification of location to implement rain water harvesting (RWH) management strategies effectively and efficiently. During this study, sensitivity analysis on percentage influence was conducted which has provided an insight into the influence of each RWH indicator in weights assigned to each indicator. To determine the potential RWH, site methodology was developed by using GIS and remote sensing technology. The model used is only valid for this type of relation. However, it does not deny the possibility of other kind of relations within our study system. In this regard, the modelling approach should be considered as a guide in identifying essential interactions on which empirical efforts should be targeted to confront the results obtained with purely observational approaches such as that employed in this study.

Another study [7] is based on the Beijing as a case study which introduced the way to set up a capable IT service system for processing, gaining, analysis of data, transmission, storage and application of the GFPB site selection with the help GPS and RS which are as tools to obtain data, by using network-based information-transferring and sharing channel. Author also thought that the selection of site must depend upon soil quality assessment air quality, water quality, soil organic matter, and the surrounding traffic and pollution situation and many other factors. Through this study, the researcher presented a MSHAF-based analysis program which was combined with agricultural multi-service collaborative technology which has a aim of multi-sectoral, multi-platform spatial data, attribute data, structured data, unstructured data and other information resources to enable them to be grouped hierarchically as per their types and logical relations, which proved appropriate specifications and feasible methods in order to achieve effective integration and fusion of cross-platform, cross-sector and multi-resource heterogeneous data sources and build a unified decision management support system.

The applied architecture of DSS for crop selection is depicted in the Fig. 2. Soil's biological information, seasons, irrigation details and crop-related facts are taken into concern in order to keep the model efficient as well as compact. Every layer is weighted in accordance with the Saaty's AHP method [8]. Each and every class is assigned an intensity value with expert's advice. And finally, the results are obtained in tabular as well as graphical form.

## 4   Proposed Architecture

To develop the crop selections decision support system, the following method has been adopted (Fig. 2).

1. **Collection of Data**: The first and foremost thing is the collection of data after the recognition of the essential fields, i.e. biological and chemical details of soil,

**Fig. 2** Proposed DSS for crop selection

which were acquired after the test in lab, next was the irrigation details, which were gathered from tube well inventory and past 5 years rainfall record. These all collected details were placed into the attribute tables of GIS database. The

following table depicts the soil pH of the collected samples, which were tested in lab by the experts (Table 1).

2. **Thematic Map**: The basic requirement of the suitability analysis is the thematic map of the study area and then the primary step is the cartography of the topographic guide map of the survey region.

3. **Geo-referencing**: Geo-referencing is the basic process before digitisation. In the geo-referencing process, the coordinates are assigned to the scanned map in order to relate the attributes on the map to the actual position on the earth surface. The coordinate system adopted here is latitude and longitudinal degree. In the current scenario, the GCP (ground control points) has been collected using hand-held GPS machine and then synchronised into the GIS software, and then by using these GCP, the map is geo-referenced.

4. **Data Quality Check**: The quality of data plays an important part in the quality of graphical and statistical product. The data quality element needs to be considered are completeness, logical consistency, potential accuracy, temporal accuracy and thematic accuracy.

5. **Plot GCP of Soil Sample**: GCP which gives the exact location of a physical feature on map. In this study, we used the GCP to show the location of the collected soil samples. To get the GCP of soil sample, location hand-held device is used. In the present study, a theme is generated to show the GCP of soil sample and also used to for Kriging.

6. **Data Layers**: In GIS, the raw data is of no use unless and until it is being processed and hence the process of digitization is an essential one. In this phase, the collected attributes are digitised on layers using the features like point, line, polyline and polygon. In the present study, point is used to depict the tube well; polygon is used to show the border of the area taken for study and the irrigated zone is also shown using the polygon. The final outcome after AHP process is also shown using point on a map layer.

7. **Data of Land and Land Owner**: The detail data related to the land and land holders are acquired from the legal government agencies. This information covered the land holding area of respective farmer, details of irrigation facility in farm, last cultivated crop pattern and the personal info of farmers. This collected info has been used to generate the respective layers to process in AHP for DSS.

8. **Expert's Advice**: To address the agriculture concerned knowledge, the expert from respective field where called on. The selection of crops as options and respective parameters such as soil composition details, season and irrigation was taken in, with concern of the expert. And the vital thing of AHP process is to assign the values of intensity of importance to criteria, and alternatives are also completed under the guidance and help of expert.

**Table 1** pH of collected soil sample

| Soil sample | S1 | S2 | S3 | S4 |
|---|---|---|---|---|
| pH | 8.51 | 8.13 | 8.14 | 7.82 |

9. **Building Database**: In this step, the database required for the decision support system is build and organised. In the current study, the land holders details and related fields are recorded in GIS, and data collected from expert are also added to the database.

10. **Integration of Spatial and Non-spatial Data**: The vital step after data collection and digitization of layer in spatial decision support system is the spatial and non-spatial data integration. In this process, the spatial and non-spatial data's database is attached to the particular map.

11. **Import Layers in AHP**: In order to process the layer in analytic hierarchy process tool to get the result map and table, it is important to import the required raster layers in the AHP process tool. The digitised data layers of the components are imported; firstly, the layers of criteria and then the options or alternatives are imported.

12. **Process AHP**: It is one of the important step in GIS-based DSS for crop selection, here the layers of various components imported in AHP tool are assigned with the respective values as called for in by the expert. Then, the layers with the values are computed using the available feature in AHP tool, after the successful computation, the result map is visualised in graphical form in canvas area and the tabular data are saves in (.txt) note pad extension file.

13. **Results**: After processing the layers in AHP, the results are displayed on map and also generated in a notepad text file.

## 5 Conclusion

This paper introduces a novelty use of integrating GIS and AHP technologies for Indian agriculture scenario. This work has outlined a DSS framework to support multi-criteria selection of Indian crops. Our experimental outcomes addressed the enhancement in ease of crop analysis and decision-making process. The area opted for this study was a village Khadkut under the jurisdiction of Nanded District in Maharashtra state. The present study is vital to improve and bridge the knowledge vacuity of farmers in decision-making procedure regarding selection of respective crops.

## References

1. Website source at http://www.peer.eu/fileadmin/user_upload/opportunities/metier/course2/c2_updating_sessions_GIS.pdf as on date 10 May 2015
2. Peeters, A., Ben-Gal, A., et al.: Developing a GIS-based spatial decision support system for automated tree crop management to optimize irrigation inputs. In: iEMSs (2012)

3. Seppelt, R., Voinov, A.A., Lange, S., Bankamp, D. (Eds.): International Congress on Environmental Modelling and Software Managing Resources of a Limited Planet, Sixth Biennial Meeting, Leipzig, Germany
4. Website source at http://www.iemss.org/society/index.php/iemss-2012-proceedings
5. Pradhan, S.: Crop area estimation using GIS, remote sensing and area frame sampling. **3**(1) (2001)
6. Nketiaa, A.K., Forkuob, E.K., Asamoaha, E., Senayaa, J.K.: A GIS based model as a decision support framework for identifying suitable rain water harvesting sites. Int. J. Adv. Technol. Eng. Res. **3**(4) (2013)
7. Yan, X., Wang, W., Liang J.: How to design and apply a DSS based on multi-source heterogeneous data fusion (MSHDF) technology for the site selection of green food production base (GFPB). In: 2010 Second IITA International Conference on Geoscience and Remote Sensing (IITA-GRS), vol. 2, pp. 627, 633, 28–31 Aug 2010. doi:10.1109/IITA-GRS.2010.5602324
8. Saaty, T.L.: Decision making with the analytic hierarchy process. Int. J. Serv. Sci. **1**(1) (2008)

# PROMETHEE-Based Analysis of HCWM Challenges in Healthcare Sector of Odisha

**Sasanka Sekhar Mishra, Kamalakanta Muduli, Manoranjan Dash and Devendra K. Yadav**

**Abstract** Waste produced from hospitals is throwing threat to environment as well as generating precarious disease. Improper waste management not only is a major cause of environment pollution but also hazardous to the human resources associated with this industry, i.e. allied healthcare professionals, patients, paramedical, non-paramedical staff, and to the community health. Previous research found factors that slow down usefulness of HCWM. Factors affecting HCWM and their significance if known in advance will be supportive for the practioners and analyst to effectively handle these barriers in enhancing the effectiveness of healthcare waste management. PROMETHEE methodology for multiattribute decision-making (MADM) was used to prioritize barriers, according to the degree of their adverse effect. Most significant barriers being perceived and ranked 1 are reuse of healthcare waste not authorized by law followed by improper segregation practices, hospital administrators' lack of accountability and poor training and awareness programs which are alleged as the least significant barriers.

**Keywords** Barriers of waste management · Health care waste (HCW)
MADM · Healthcare unit · PROMETHEE

S. S. Mishra · M. Dash
Siksha O Anusandhan University, Bhubaneswar, India

K. Muduli (✉)
Papua New Guinea University of Technology,
LAE, Morobe Province 411, Papua New Guinea
e-mail: kamalakantam@gmail.com

D. K. Yadav
NIT Calicut, Kozhikode, India

163

# 1   Introduction

Wastes generated from hospitals and allied units are termed as healthcare waste which includes materials, syringes, dressings, body tissue and samples. Poor management of healthcare waste generated from hospitals has an adverse effect on patients, health workers and the poetical risk of making environment polluted [1, 2]. A major portion around 75–90% of the medical waste is considered equivalent to household waste as it consists of food remains, paper and packaging material. The remaining portion consisting of human tissues, syringes, bandages, used culture media containing micro-organisms [3] is considered as hazardous waste as it contains harmful, carcinogenic, toxic and infectious materials [4]. As per statistics revealed by WHO in the year 2000, millions of hepatitis B, C and HIV infections were caused due to the infected syringes [1]. Improper waste management practices such as directing untreated waste into water bodies or dumping in low-lying areas or roadside municipal bins have been observed to be a common practice in many countries. Besides these, poor HCWM practices end up with growth and multi-plication of insects, rodents and worms as well as spread of diseases such as cholera, HIV, typhoid [3, 5]. Over the last few years, due to rapid mushrooming of hospitals in urban areas and migration of population from rural areas to urban areas, these hospitals generate a large amount of healthcare waste and disposable health products [6, 7]. This raises concerns for proper healthcare waste management which if not done in advance will have an impact on the health community, and the impact has to be minimized [8]. Realizing the potential negative impact on community, many developed as well developing countries have taken steps in addressing theses issues and had adopted many commendable steps [5]. However, in India, the intuitive and awareness regarding HCWM is in infancy stages and the results are not satisfactory [9, 10]. This raises the question on knowledge and awareness level about various factors affecting the HCWM in India. Information about these factors particularly the barriers will help the healthcare practiioners in formulating their proactive strategies and to minimize the impact on the environment.

In this context, this research is carried out with the following objectives

1. to explore the parameters that obstruct HCWM practices in different HCUs in India.
2. to use PROMETHEE to rank the parameter identified as per their degree of inhibiting strength.

# 2   Literature Review

Worldwide poor management of healthcare wastes is considered as a serious issue for the society as well as for the organizations. Hence, researchers from several countries such as Indonesia [8, 11], Greece [12, 13], USA [14, 15], Turkey [16],

**Table 1** In Indian context identified barriers of healthcare waste management (HCWM)

| Sl. No. | Barriers | References |
|---|---|---|
| 1 | Improper segregation practices | [14, 18, 19] |
| 2 | Implementation of poor HCWM operational practices | [20, 21, 19] |
| 3 | Insufficient help from govt. agencies | [20, 21, 17] |
| 4 | Less importance to green purchasing | [20, 22, 17] |
| 5 | Reuse of HCW not authorized by law | [20, 23, 21] |
| 6 | Hospital administrators lacking accountability | [24, 21, 19] |
| 7 | Inadequate HCWM infrastructure | [25, 2, 6] |
| 8 | Financial constraints | [20, 25, 17] |
| 9 | Low level of awareness and training programs | [12, 2, 18] |
| 10 | Resistance to change and adoption | [21] |
| 11 | Low level of acquaintance between patients, paramedical staff and general public | [24, 26] |
| 12 | Improper coordination among pollution control board (PCB), HCUs and municipality | [24, 3] |
| 13 | Less importance to waste management issues as per policy drafting of health care units | [24] |
| 14 | No stringent implementation of infection control measures | [24, 6] |

Italy [17], Cameroon [18], Croatia [4], England [19], Palestine [2] have studied healthcare waste management issues. In Indian context, few researchers Thakur and Anbanandam [7, 20], Sharma et al. [10], Madhukumar and Ramesh [21], Das and Biswas [22], Muduli and Barve [23] have also studied HCWM practices.

Several authors designed questionnaires, conducted field research and used interview techniques in surveying wastes generated by healthcare institutions [16] which is summarized in Table 1.

# 3 Methodology

Based on the literature, a structured questionnaire was designed consisting of 14 items relating to barriers causing hindrance in the implementation of HCWM. Survey was conducted among respondents working in different healthcare units, physicians and healthcare professionals. Five-point Likert scale was used (1 = totally disagree, 5 = completely agree), and values produced were used for prioritization of healthcare waste management barriers using PROMETHEE II.

## 3.1 PROMETHHE-Based Evaluation of Barriers

Brans [27] first developed the Preference Ranking Organization Method for Enrichment Evaluation (PROMETHEE) as one of the most prevalent multicriteria decision-making techniques [28], which has been employed in this study to analyse various HCWM challenges. PROMETHEE II has been chosen in this study for being interactive and its ability to rank the variables which are complex and difficult to compare [29]. For computation of a preference function for each variable, it involves a pairwise comparison of variables. The preference function is used as a basis to determine a preference index for the variable $i$ over $i'$. This preference index is the measure to support the hypothesis that variable i is given priority over $i'$. The steps in PROMETHEE II technique are listed as below [27–30].

Step-1:   Normalization of decision matrix using the following equation:

$$R_{ij} = [X_{ij} - \min(X_{ij})]/[\max(X_{ij}) - \min(X_{ij})] \quad (i = 1, 2, \ldots, n; j = 1, 2, \ldots, m) \tag{1}$$

where $X_{ij}$ is the performance measure of $i$th factor with respect to $j$th criterion.

Step-2:   Computation of evaluative differences of $i$th factor with respect to other factors. This step involved the calculation of differences in criteria values between different variables pairwise.

Step-3:   Calculation of preference function $P_j(i, i')$ using the following equation:

$$P_j(i, i') = 0 \quad \text{If } R_{ij} \leq R_{i'j} \tag{2}$$

$$P_j(i, i') = (R_{ij} - R_{i'j}) \quad \text{If } R_{ij} > R_{i'j} \tag{3}$$

Step-4:   Aggregated preference function was calculated taking into account the weight criterion. Aggregated preference function:

$$\Pi(i, i') = \left[ \sum_{j=1}^{m} W_j \times P_j(i, i') \right] \Big/ \sum_{j=1}^{m} W_j \tag{4}$$

where $W_j$ was the relative importance (weight) of $i$th criterion.

Step-5:   Determination of leaving and entering outranking flows. Leaving (or positive) flow for $i$th barrier:

$$\phi^+(i) = \frac{1}{n-1} \sum_{i'=1}^{n} \Pi(i, i') \quad (i \neq i') \tag{5}$$

Entering (or negative) flow for $i$th barrier:

$$\phi\text{-}(i) = \frac{1}{n-1} \sum_{i'=1}^{n} \Pi(i', i) \quad (i \neq i') \tag{6}$$

Here, $n$ was the number of barriers.

Step-6: Calculation of net outranking flow for each variable using Eq. 7.

$$\phi(i)=\{\phi^+(i)\}-\phi^-(i)\} \tag{7}$$

Step-7: Depending on the values of $\phi(i)$, rank the variables. The more important was the barrier to be looked upon depending upon the higher value of $\phi(i)$. Thus, the most important barrier to be looked upon was the one having the highest $\phi(i)$ value.

## 4 Result Analysis

Prioritization of the HCWM barriers needs expert opinion regarding the degree of influence of each barrier on HCWM practices. This needed a questionnaire-based survey using which included respondents from two hospitals operating in Odisha. The respondents provided their judgment in using numerals between 1 and 5. The survey responses were then normalized using Eq. 1 and shown in Table 2. Next, aggregated preference function was calculated from these normalized values using steps 2–4 and represented in Table 3. Finally, evaluations of net outranking flow for each barrier are done using steps 5 and 6 and are shown in Table 4. The barriers are then ranked depending upon their net outranking values as mentioned in step 7 and represented in Table 4.

**Table 2** Normalized table

| Barriers | Averaged values |
|----------|-----------------|
| 1 | 0 |
| 2 | 0.25 |
| 3 | 0.25 |
| 4 | 0.25 |
| 5 | 1 |
| 6 | 0 |
| 7 | 0.5 |
| 8 | 0.75 |
| 9 | 0 |
| 10 | 0.75 |
| 11 | 0.75 |
| 12 | 0.25 |
| 13 | 0.75 |
| 14 | 0.75 |

**Table 3** Aggregated preference function for HCWM

|   | A | B | C | D | E | F | G | H | I | J | K | L | M | N |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| B | 0.25 | 0 | 0 | 0 | 0 | 0.25 | 0 | 0 | 0.25 | 0 | 0 | 0 | 0 | 0 |
| C | 0.25 | 0 | O | 0 | 0 | 0.25 | 0 | 0 | 0.25 | 0 | 0 | 0 | 0 | 0 |
| D | 0.25 | 0 | 0 | 0 | 0 | 0.25 | 0 | 0 | 0.25 | 0 | 0 | 0 | 0 | 0 |
| E | 1 | 0.75 | 0.75 | 0.75 | 0 | 1 | 0.5 | 0.25 | 1 | 0.25 | 0.25 | 0.75 | 0.25 | 0.25 |
| F | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| G | 0.5 | 0.25 | 0.25 | 0.25 | 0 | 0.5 | 0 | 0 | 0.5 | 0 | 0 | 0.25 | 0 | 0 |
| H | 0.75 | 0.5 | 0.5 | 0.5 | 0 | 0.75 | 0.25 | 0 | 0.75 | 0 | 0 | 0.5 | 0 | 0 |
| I | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| J | 0.75 | 0.5 | 0.5 | 0.5 | 0 | 0.75 | 0.25 | 0 | 0.75 | 0 | 0 | 0.5 | 0 | 0 |
| K | 0.75 | 0.5 | 0.5 | 0.5 | 0 | 0.75 | 0.25 | 0 | 0.75 | 0 | 0 | 0.5 | 0 | 0 |
| L | 0.25 | 0 | 0 | 0 | 0 | 0.25 | 0 | 0 | 0.25 | 0 | 0 | 0 | 0 | 0 |
| M | 0.75 | 0.5 | 0.5 | 0.5 | 0 | 0.75 | 0.25 | 0 | 0.75 | 0 | 0 | 0.5 | 0 | 0 |
| N | 0.75 | 0.5 | 0.5 | 0.5 | 0 | 0.75 | 0.25 | 0 | 0.75 | 0 | 0 | 0.5 | 0 | 0 |

**Table 4** Net outranking flow for each barrier

| Sl. No. | Barriers | $\phi^+$ | $\phi^-$ | $\phi$ | Rank |
|---|---|---|---|---|---|
| 1 | Improper segregation practices | 0 | 0.481 | −0.481 | 4 |
| 2 | Implementation of poor HCWM operational practices | 0.03846 | 0.269 | −0.231 | 3 |
| 3 | Inadequate assistance from government agencies | 0.03846 | 0.269 | −0.231 | 3 |
| 4 | Low priority to green purchasing | 0.03846 | 0.269 | −0.231 | 3 |
| 5 | Reuse of healthcare waste illegitimately | 0.51923 | 0 | 0.5192 | 1 |
| 6 | Hospital administrators lacking accountability | 0 | 0.481 | −0.481 | 4 |
| 7 | Inadequate HCWM infrastructure | 0.15385 | 0.1346 | 0.0192 | 3 |
| 8 | Financial constraints | 0.28846 | 0.0192 | 0.2693 | 2 |
| 9 | Low level of awareness and training programs | 0 | 0.481 | −0.481 | 4 |
| 10 | Resistance to change and adoption | 0.28846 | 0.0192 | 0.2693 | 2 |
| 11 | Low level of acquaintance between patients, paramedical staff and general public | 0.28846 | 0.0192 | 0.2693 | 2 |
| 12 | Improper coordination among pollution control board (PCB), HCUs and municipality | 0.03846 | 0.2692 | −0.231 | 3 |
| 13 | Less importance to waste management issues as per policy drafting of healthcare units | 0.28846 | 0.0192 | 0.2693 | 2 |
| 14 | No stringent implementation of infection control measures | 0.28846 | 0.0192 | 0.2693 | 2 |

# 5    Conclusion

Healthcare units experience a number of challenges during implementation of waste management practices. Hence, it is not only required to identify these challenges but also required to quantify their degree of negative influence. In this respect, fourteen challenges of healthcare waste management have been identified in this research. Further, PROMETHEE II has been employed to evaluate the adverse impact of these challenges. The present research reveals that 'improper segregation practices', 'hospital administrators accountability' and 'low level of awareness and training programs' are the least significant factors, while 'reuse of healthcare waste illegitimately' is the most important challenge of HCWM experienced by the Indian HCUs. Prioritization of these HCWM challenges will give an insight to the decision-makers to improve the HCWM effectiveness in their organizations.

**Declaration**    Consents of the respondents have been taken to use the data in this research. If any issue arises hereafter, then the authors will be solely responsible. Neither the editors nor the publisher will be responsible for any misuse or misinterpretation of the data.

# References

1. Liu, H.C., Wu, J., Li, P.: Assessment of health-care waste disposal methods using a VIKOR-based fuzzy multi-criteria decision making method. Waste Manag. **33**(12), 2744–2751 (2013)
2. Al-Khatib, I.A., Sato, C.: Solid health care waste management status at health care centers in the West Bank-Palestinian Territory. Waste Manag. **29**(8), 2398–2403 (2009)
3. Dwivedi, A.K., Pandey, S., Shashi, : Fate of hospital waste in India. Biol. Med. **1**(3), 25–32 (2009)
4. Marinković, N., Vitale, K., Holcer, N.J., Džakula, A., Pavić, T.: Management of hazardous medical waste in Croatia. Waste Manag. **28**(6), 1049–1056 (2008)
5. Swain, S., Muduli, K., Biswal, J.N., Tripathy, S., Panda, T.K.: Evaluation of barriers of health care waste management in India—a gray relational analysis approach. In: Proceedings of the 5th International Conference on Frontiers in Intelligent Computing: Theory and Applications, vol. 515, no. 1, pp. 181–188. Springer, Singapore (2017)
6. Gupta, S., Boojh, R., Dikshit, A.K.: Environmental education for healthcare professionals with reference to biomedical waste management—a case study of a hospital in Lucknow, India. Int. Res. J. Environ. Sci. **1**(5), 69–75 (2012)
7. Thakur, V., Ramesh, A.: Healthcare waste management research: a structured analysis and review (2005–2014). Waste Manag. Res. 1–16 (2015). doi:10.1177/0734242X15594248
8. Chaerul, M., Tanaka, M., Shekdar, A.V.: A system dynamics approach for hospital waste management. Waste Manag. **28**(2), 442–449 (2008)
9. Biswal, M., Mewara, A., Appannanavar, S.B., Taneja, N.: Mandatory public reporting of healthcare-associated infections in developed countries: how can developing countries follow? J. Hosp. Infect. **90**(1), 12–14 (2015)
10. Sharma, A., Sharma, V., Sharma, S., Singh, P.: Awareness of biomedical waste management among health care personnel in Jaipur, India. Oral Health Dent. Manag. **12**(1), 32–40 (2013)

11. Zurbrügg, C., Gfrerer, M., Ashadi, H., Brenner, W., Küper, D.: Determinants of sustainability in solid waste management—the Gianyar waste recovery project in Indonesia. Waste Manag. **32**(11), 2126–2133 (2012)
12. Tsakona, M., Anagnostopoulou, E., Gidarakos, E.: Hospital waste management and toxicity evaluation: a case study. Waste Manag. **27**(7), 912–920 (2007)
13. Komilis, D., Fouki, A., Papadopoulos, D.: Hazardous medical waste generation rates of different categories of health-care facilities. Waste Manag. **32**(7), 1434–1441 (2012)
14. Conrardy, J., Hillanbrand, M., Myers, S., Nussbaum, G.F.: Reducing medical waste. AORN J. **91**(6), 711–721 (2010)
15. Berwick, D.M., Hackbarth, A.D.: Eliminating waste in US health care. JAMA **307**(14), 1513–1516 (2012)
16. Dursun, M., Karsak, E.E., Karadayi, M.A.: A fuzzy multi-criteria group decision making framework for evaluating health-care waste disposal alternatives. Expert Syst. Appl. **38**(9), 11453–11462 (2011)
17. Cagliano, A.C., Grimaldi, S., Rafele, C.: A systemic methodology for risk management in healthcare sector. Saf. Sci. **49**(5), 695–708 (2011). ISSN 0925-7535
18. Manga, V.E., Forton, O.T., Mofor, L.A., Woodard, R.: Health care waste management in Cameroon: a case study from the Southwestern Region. Resour. Conserv. Recycl. **57**, 108–116 (2011)
19. Loveday, H.P., Wilson, J., Pratt, R.J., Golsorkhi, M., Tingle, A., Bak, A., et al.: epic3: national evidence-based guidelines for preventing healthcare-associated infections in NHS hospitals in England. J. Hosp. Infect. **86**, S1–S70 (2014)
20. Thakur, V., Anbanandam, R.: Healthcare waste management: an interpretive structural modeling approach. Int. J. Health Care Qual. Assur. **29**(5), 559–581 (2016)
21. Madhukumar, S., Ramesh, G.: Study about awareness and practices about health care waste management among hospital staff in a medical college hospital, Bangalore. Iran. J. Basic Med. Sci. **3**, 7–11 (2012)
22. Das, S.K., Biswas, R.: Awareness and practice of biomedical waste management among healthcare providers in a Tertiary Care Hospital of West Bengal, India. Int. J. Med. Public Health **6**(1), 19 (2016)
23. Muduli, K., Barve, A.: Barriers to green practices in health care waste sector: an Indian perspective. Int. J. Environ. Sci. Dev. **3**(4), 393–399 (2012)
24. Gupta, S., Boojh, R.: Report: biomedical waste management practices at Balrampur Hospital, Lucknow, India. Waste Manag. Res. **24**, 584–591 (2006)
25. Abdulla, F., Qdais, H.A., Rabi, A.: Site investigation on medical waste management practices in norther Jordan. Waste Manag. **28**(2), 450–458 (2008)
26. Athavale, A.V., Dhumale, G.B.: A study of hospital waste management at a rural hospital in Maharastra. J. ISHWM **9**(1), 21–31 (2010)
27. Brans, J.P.: Lingenierie de la decision. Elaboration dinstrumentsdaide a la Decision. Method Promethee. In: Nadeau, R., Landry, M. (eds.) Laide a la decision: Nature, instruments et perspectives davenir, pp. 183–214. Presses de Universite Laval, Quebec (1982)
28. Kilic, H.S., Zaim, S., Delen, D.: Selecting 'the best' ERP system for SMEs using a combination of ANP and PROMETHEE methods. Expert Syst. Appl. **42**(5), 2343–2352 (2015)
29. Amaral, T.M., Costa, A.P.C.: Improving decision-making and management of hospital resources: an application of the PROMETHEE II method in an emergency department. Oper. Res. 3(1), 1–6 (2014)
30. Sen, D.K., Datta, S., Patel, S.K., Mahapatra, S.S.: Multi-criteria decision making towards selection of industrial robot. Benchmark Int. J. 22(3), 465–487 (2015)

# Performance Analysis of NSL_KDD Data Set Using Neural Networks with Logistic Sigmoid Activation Unit

**Vignendra Jannela, Sireesha Rodda,**
**Shyam Nandan Reddy Uppuluru, Sai Charan Koratala**
**and G. V. S. S. S. Chandra Mouli**

**Abstract** Network intrusion detection system (NIDS) is a software tool that scans network traffic and performs security analysis on it. NIDS performs match operations upon passing traffic with a pre-established library of attacks in order to identify attacks or abnormal behavior. One of the standard data sets used widely for network intrusion systems is the NSL_KDD data set. The current paper aims to analyze the NSL_KDD data set using artificial neural network with sigmoid activation unit in order to perform a metric analysis study that is aimed at discovering the best fitting parameter values for optimal performance of the given data. Evaluation measures such as accuracy, *F*-measure, detection rate, and false alarm rate will be used to evaluate the efficiency of the developed model.

**Keywords** Network intrusion detection system · NSL_KDD
Neural networks · Logistic sigmoid activation unit

## 1 Introduction

Network intrusion refers to the illegal access of data housed within a network by unauthorized entities. A network intrusion detection system (NIDS) is a security management system employed to detect intrusion attempts within the traffic circulating in a network. An NIDS performs detection through an analysis of all active traffic within a network, constantly matching the traffic with a pre-established library of known attacks. Whenever a match occurs, an alert is usually issued to the system administrator who is responsible to take the necessary next steps.

Techniques of Detection
NIDSs in traditional networks work primarily on two types of techniques:

V. Jannela (✉) · S. Rodda · S. N. R. Uppuluru · S. C. Koratala ·
G. V. S. S. S. Chandra Mouli
Department of Computer Science and Engineering, GITAM Institute of Technology, GITAM
University, Visakhapatnam, India
e-mail: vignendraj@gmail.com

1.  Signature-based intrusion detection techniques

Signature-based detection refers to a set of techniques that search for special patterns or byte combinations known as "signatures" that have already been established as malicious in nature. A signature is treated as a threat only if it matches a precise set of pre-established conditions—thus reducing false positive rates. However, these techniques cannot detect threats beyond those already identified in a pre-established library of threat signatures.

2.  Anomaly-based intrusion detection techniques

Anomaly-based intrusion detection refers to a set of techniques that work within a well-defined scope or "rules" that are used to classify instances of activity as either "normal" (harmless) or "anomalous" (potentially malicious). They are usually built in two phases: The first is the "training phase" where scope of operation is constructed; the second is "testing phase" where actual network traffic is subject to evaluation within scope conditions. Although there are higher chances of encountering false positives, these techniques are much more potent at detecting previously un-encountered malicious instances.

Due to their "Learn and classify" methodology, supervised machine learning techniques can be employed toward building anomaly-based network intrusion detection. In our paper, we aim to explore logistical regression, a supervised learning algorithm, to evaluate the NSK_KDD data set in order to evaluate various evaluation parameters.

In this paper, a performance analysis has been performed on the NSL_KDD data set using neural networks with logistic sigmoid activation units. The remainder of the paper is divided as follows: Section 2 houses a discussion on related works. Section 3 houses a general outlook of the workings of the performance analysis and the neural network in addition to an overview of the NSL_KDD data set. Section 4 houses the actual results from experiments performed on the neural network. Section 5 outlines the findings of the performance analysis and provides inputs for future expansion of this paper's ideas.

## 2  Related Work

There have been significant contributions related to the performance analysis of the NSL_KDD data set. An assortment of machine learning algorithms has been applied in the past to study the performance of the NSL_KDD data set. Most performance analyses concentrate on studying dedicated issues that affect the performance of NSL_KDD in detection network intrusion. These include issues such as the significance of data quality, effectiveness of specific algorithms, and methods of discretization.

Revathi and Malathi [1] proposed a new model that aimed at performing dimensionality reduction using CFS subset. Initially, the unmodified data was evaluated using various classification algorithms. Then, a modified data set based

on the proposed model was evaluated with the same set of classifiers. The results indicated a clear reduction in computational time while seeing a significant surge in accuracy.

Aggarwal and Sharma [2] discuss the importance of data for improving the accuracy and overall performance of network intrusion detection systems. As such, the authors have classified the data into three categories: basic, content traffic, and host. Analysis of data emphasizes on two primary test metrics: FAR and DR. These metrics have been used to perform a well-founded analysis by testing over permutations of categorized data. In conclusion, it has been established that "content" data holds the highest FAR and "host" data holds the lowest DR values.

Dreiseitl and Machado [3] aimed at performing a comparative analysis of Artificial Neural Networks and Logistic Regression. The focus has been triangulated on biomedicine data sets and the quality of the model developed by the above methods.

Aziz et al. [4] proposed a new algorithm for discretization of data. Applying an equal width interbinning model, the algorithm performs discretization of data while simultaneously ensuring that a degree of homogeneity exists between data values. The proposed algorithm boasts simplicity in implementation. However, the overall findings indicate that the algorithm does not perform efficiently when compared to the most commonly used existing models.

Hussain and Lalmuanawma [5] highlighted the challenge posed by the existence of noise in real-time environments for network intrusion detection models. The authors apply various machine learning algorithms to analyze how each performs under noisy conditions. Their findings show that certain classifiers such as J48 and JRip that give high-quality performance in test environments fail to recreate the same performance efficiency in noisy environments. In these cases, neural networks outperform these algorithms, thereby establishing the robustness of the neural network-based model.

Singh and Bansal [6] in this paper applied four algorithms on the NSL_KDD data set to evaluate their performances in NIDS models. The algorithms applied are as follows: BSF, multilayer perceptron, logistic regression, and voted perception. Through the tests, the multilayer perceptron algorithm was deemed to possess highest accuracy lowest error rate, thereby implying it as the most suited algorithm for NIDS models.

## 3 Methodology

The network intrusion detection system we study as part of this paper will be using the NSL_KDD data set. The previously used KDD99 data set had certain inherent flaws that were attempted to be rectified in the NSL_KDD data set. These flaws are discussed in [7]. The newer KDD data set does still possess some of the issues aired by McHugh [8], thereby questioning its authority as the absolute representative of networks in real-time environments. However, due to the scarcity of publicly

available data sets for NIDSs, it is our best belief that the NSL_KDD data set can still be utilized fruitfully as a benchmarking data set in order to perform analyses and comparative studies upon various methods of network intrusion. In addition to this, the number of records available in the NSL_KDD data set is sizable to a practical degree. This results in increased affordability in running experiments on the entire data sets and eliminates the need to perform random selection on snippets of the overall data. As a consequence, the results showcased by works of research from different sources will still hold overall consistency and will facilitate comparisons. A brief description of the attacks in NSL_KDD data set (Table 1).

Multilayer Perceptron

A multilayer perceptron (MLP) is an artificial neural network model working on a "feed forward" principle. The model works by mapping input data sets onto output data sets. The structure of the MLP is characterized by the presence of multiple layers of nodes in the form of a directed graph, where each layer is connected to the next layer in totality. Every node in the model is a processing element (known as neuron) that has a nonlinear activation function. The only exceptions to this are the input nodes. The multilayer perceptron model trains the network using back propagation technique. Each MLP contains one input layer, one or more hidden layers, and one output layer, consisting of nonlinearly activating nodes. The activation units in the hidden layers are modeled using logistic sigmoid function which can be used to model probability. Ease of usage and the ability to approximate any input/output map are this model's primary advantages.

Each perceptron of the neural network is governed by the following equations:

$$Z = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n \tag{1}$$

Hypothesis Function $h_\theta(x)$

$$h_\theta(x) = g(Z) \tag{2}$$

$$g(Z) = \frac{1}{1 + e^{-Z}} \text{ (sigmoid/Logistic function)} \tag{3}$$

$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}} \tag{4}$$

**Table 1** Attacks in NSL_KDD data sets [2]

| Attacks in data set | Attack type |
|---|---|
| DOS | Back, Land, Neptune, Teardrop, Mailbomb, Processable, Udpstorm, Apache2, Worm |
| Probe | Satan, IPsweep, Nmap, Mscan, Saint. |
| R2L | Guess_password, Ftp_write, Imap, Phf, Multihop, Warezmaster, Xlock, Xsnoop, Snmpguess, Snmpgetattack, Httptunnel, Sendmail, Named. |
| U2R | Buffer_overflow, Loadmodule, Rootkit, Perl, Sqlattack, Xterm, Ps |

Minimization Objective (Cost-$J(\theta)$):
Let

$L$ be the total number of layers in the network.
$S_l$ be the number of units (not counting bias unit) in layer "$l$".
$K$ be the number of output units/classes.
$\lambda$ be the regularization parameter to avoid over-fitting.

$$h_\emptyset(x) \in R^K \; (h_\emptyset(x))_i = i\text{th output}$$

$$
\begin{aligned}
J(\theta) = \frac{-1}{m} &\left[ \sum_{i=1}^{m} \sum_{K=1}^{K} y_K^{(i)} \log\left( h_\theta(x^{(i)})_K \right) + \left(1 - y_K^{(i)}\right) \log\left(1 - h_\theta(x^{(i)})_K\right) \right] \\
&+ \frac{\lambda}{2m} \sum_{L=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} \left( \theta_{ji}^{(l)} \right)^2
\end{aligned}
\tag{5}
$$

Feed forward propagation is used to compute the cost.
Backward propagation is used to update the weights of activation unit of each perceptron.

## Algorithm: Back Propagation

(1) Training set $\{(x^1, y^1), \ldots, (x^m, y^m)\}$
(2) Set $\Delta_{ij}^{(l)} = 0$ (for all $l, i, j$)
(3) for $i = 1$ to $m$

    (a) Set $a^1 = x^1$
    (b) Perform forward propagation to compute $a^{(l)}$ for $l = 2, 3, 4, \ldots, l$
    (c) Using $y^i$, compute $S^{(l)} = a^{(l)} - y^{(l)}$
    (d) Compute $\delta^{(l-1)}, \delta^{(l-2)}, \ldots, \delta^{(2)}$
        Using Eqs. (6) and (7)
    (e) $\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_{ij}^{(l)} \delta_i^{(l+1)}$

(4)

    (a) $D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)}$ if $(j \neq 0)$
    (b) $D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)}$ if $(j = 0$ correspondence to biased terms)

(5) $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = D_{ij}^{(l)}$

For the computation of $\delta^{(l)}$:

$$\delta^{(l)} = (\theta^{(l)})^T \delta^{(l+1)}. * g'(z^{(l)}) \tag{6}$$

Sigmoid gradient

$$g'(z^{(l)}) = a^{(l)}. * (1 - a^{(l)}) \tag{7}$$

The performance of the model built using training data set is analyzed using test metrics on the test data set of NSL_KDD. The test metrics used for evaluation are as follows: (a) accuracy, (b) precision, (c) recall, (d) $F$-measure, and (e) false alarm rate (FAR).

## 4   Results

A number of training models are built for different values of regularization parameter and varying number of hidden units. These models are tested on NSL_KDD test data set for different test metrics noted above. The model's performance for varying values of regularization parameters and a number of hidden layer units are tabulated in Tables 2 and 3, respectively. The results tabulated in Table 3 correspond to varying number of hidden units for a neural network with a single hidden layer architecture. Figures 1, 2, 3, 4 and 5 correspond to the graphs for values tabulated in Table 2. Figures 6, 7, 8, 9, and 10 correspond to the graphs for values tabulated in Table 3.

**Table 2** Lambda versus test metrics

| Lambda value | Accuracy | Precision | Recall | $F$-measure | False alarm rate |
|---|---|---|---|---|---|
| 0 | 99.4361 | 99.4361 | 99.4173 | 99.4252 | 0.3263 |
| 0.15 | 99.4561 | 99.4561 | 99.4374 | 99.4449 | 0.3272 |
| 0.2 | 99.4321 | 99.4321 | 99.4133 | 99.4211 | 0.3285 |
| 0.25 | 99.4521 | 99.4521 | 99.4331 | 99.441 | 0.3149 |
| 0.3 | 99.4401 | 99.4401 | 99.4208 | 99.4288 | 0.3142 |
| 1 | 99.3761 | 99.3761 | 99.3592 | 99.3656 | 0.3794 |
| 3 | 99.1802 | 99.1802 | 99.1644 | 99.1701 | 0.503 |
| 10 | 97.6645 | 97.6645 | 97.6531 | 97.6205 | 1.3276 |
| 30 | 95.8648 | 95.8648 | 95.1465 | 95.4891 | 2.0254 |
| 100 | 92.7335 | 92.7335 | 92.0715 | 92.0967 | 2.7552 |
| 300 | 85.5829 | 85.5829 | 76.2865 | 80.6027 | 3.6892 |

**Table 3** Hidden units versus test metrics

| Hidden units | Accuracy | Precision | Recall | F-measure | FAR |
|---|---|---|---|---|---|
| 1 | 87.8864 | 87.8864 | 78.5043 | 82.8591 | 2.1184 |
| 2 | 98.3563 | 98.3563 | 97.9276 | 98.1351 | 0.8477 |
| 3 | 98.3563 | 98.3563 | 97.9276 | 98.1351 | 0.8477 |
| 4 | 98.7682 | 98.7682 | 98.7619 | 98.7607 | 0.811 |
| 5 | 99.1842 | 99.1842 | 99.2042 | 99.188 | 0.524 |
| 10 | 99.4081 | 99.4081 | 99.3884 | 99.3968 | 0.344 |
| 20 | 99.4121 | 99.4121 | 99.3925 | 99.401 | 0.3307 |
| 30 | 99.4361 | 99.4361 | 99.4216 | 99.4271 | 0.3227 |
| 35 | 99.4401 | 99.4401 | 99.4211 | 99.4289 | 0.3185 |
| 40 | 99.4401 | 99.4401 | 99.421 | 99.4287 | 0.3084 |
| 42 | 99.3761 | 99.3761 | 99.3578 | 99.3652 | 0.3639 |
| 50 | 99.3361 | 99.3361 | 99.318 | 99.3252 | 0.3859 |
| 100 | 99.0402 | 99.0402 | 99.0193 | 99.0287 | 0.5506 |
| 200 | 97.1606 | 97.1606 | 97.1369 | 97.1447 | 1.7419 |

**Fig. 1** Lambda versus test accuracy



**Fig. 2** Lambda versus test precision

**Fig. 3** Lambda versus recall



**Fig. 4** Lambda versus *F*-measure



**Fig. 5** Lambda versus false alarm rate

**Fig. 6** Hidden layer versus test accuracy



**Fig. 7** Hidden layer versus precision



**Fig. 8** Hidden layer versus recall

**Fig. 9** Hidden layer versus *F*-measure



**Fig. 10** Hidden layer versus false alarm rate



# 5   Conclusion and Future Works

In this paper, various parameters relating to the ANN are constantly tried, tested, and readjusted in order to find the optimal values that lead to the best performance. The best values for test metrics are obtained when the regularization parameter (represented by "Lambda") has a value that is close to zero (Lambda = 0.25). This infers that the test data should be as close as possible to the training data in terms of the hypothesis space. In spirit of this very inference, U2R attacks have been deemed as the attacks with the best possibility of avoiding detection. This is owed to the fact that the proportion of data available for U2R type attacks is significantly lower when compared to other attacks. Finally, the optimal hidden layer size of the network architecture has been found to be 35 units in size in order to gain optimal performance with regard to the NSL_KDD training and test data sets. The current

model works with the full data set containing all 41 attributes. Future works can attempt the same experiments on modified data sets upon which dimensionality reduction has been performed. Reduced dimensions lead to decreased computing costs, and may potentially lead to different performance metrics. The current paper concentrates on building artificial neural networks—at present, promising research is being performed on support vector machines, which is a prospective approach to attempt for future works. It is a point of supreme importance that this paper has worked with a noiseless data set. The results related to noisy data sets may show significant variations. Attempting to study intrusion detection performance under noisy environments is thus another promising field to explore for future research.

# References

1. Revathi, S., Malathi, A.: A detailed analysis on NSL_KDD dataset using various machine learning techniques for intrusion detection. Int. J. Eng. Res. Technol. **2**(12), 1–6 (2013)
2. Aggarwal, P., Sharma, S.K.: Analysis of KDD dataset attributes—class wise for intrusion detection. In: 3rd International Conference on Recent Trends in Computing 2015 (ICRTC-2015)
3. Dreiseitl, S., Ohno-Machado, L.: Logistic regression and artificial neural network classification models: a methodology review. J. Biomed. Inform. **35**, 352–359 (2002)
4. Aziz, A.S.A., Azar, A.T., Hassanien, A.E, Hanafy, S.A.-O.: Continuous features discretization for anomaly intrusion detectors generation. In: Soft Computing in Industrial Applications, Volume 223 of the series Advances in Intelligent Systems and Computing, pp. 209–221
5. Hussain, J., Lalmuanawma, S.: Feature analysis, evaluation and comparisons of classification algorithms based on noisy intrusion dataset. In: 2nd International Conference on Intelligent, Computing, Communication & Convergence (ICCC-2016)
6. Singh, S., Bansal, M.: Improvement of intrusion detection system in data mining using neural network. Int. J. Adv. Res. Comput. Sci. Softw. Eng. **3**(9), 1124–1130
7. Tavallaee, M., Bagheri, E., Lu, W., Ghorbani, A.: A detailed analysis of the KDD CUP 99 data set. In: Second IEEE Symposium on Computational Intelligence for Security and Defense Applications (CISDA) (2009, submitted)
8. McHugh, J.: Testing intrusion detection systems: a critique of the 1998 and 1999 DARPA intrusion detection system evaluations as performed by lincoln laboratory. ACM Trans. Inf. Syst. Secur. **3**(4), 262–294 (2000)

# Distributed Task Scheduling in Cloud Platform: A Survey

**Debojyoti Hazra, Asmita Roy, Sadip Midya and Koushik Majumder**

**Abstract**  Cloud computing is a booming area in distributed computing and parallel processing. Cloud provides services to its customer on pay-per-use basis. It has gained a lot of attention due to its unique features—elasticity, scalability, and on-demand services. Cloud facilitates both computational and storage service to its customers. This reduces the cost of deployment and maintenance for any organization. As a result, demand for cloud computing has increased considerably. To provide the services, cloud service provider needs to utilize all resources in an optimal way. To utilize all resources efficiently, task schedule plays a significant role. It is responsible for scheduling users' tasks in the cloud environment. The task scheduler arranges tasks in a queue for the available connected resources. This arrangement benefits the cloud service provider to achieve maximum performance in a cost efficient manner. In this paper, an extensive study of some well-known task-scheduling algorithms in cloud environment is done while identifying the advantages and weaknesses of these existing algorithms. Future research areas and further improvement on the existing methodologies are also suggested.

**Keywords**  Cloud computing · Task scheduling · Task scheduler · Makespan

D. Hazra · A. Roy · S. Midya · K. Majumder (✉)
Department of Computer Science and Engineering, West Bengal University of Technology, Kolkata, India
e-mail: koushik@ieee.org

D. Hazra
e-mail: debojyoti.hazra22@gmail.com

A. Roy
e-mail: asmitaroy2002@gmail.com

S. Midya
e-mail: sadip20@gmail.com

# 1 Introduction

Cloud computing is a recent technology that provides services to its users on the basis of demand. These services are Software as a Service (SaaS), Platform as a Service (PaaS), and Infrastructure as a Service (IaaS). To enhance system performance and reliability of cloud service, different factors such as task scheduling, workflow scheduling, and job scheduling are responsible.

Task-scheduling policy has direct effect on the Quality of Experience (QoE) of cloud users. A pool of resources [1, 2] is available to cloud service providers (CSP) and those resources are allocated to the users on the basis of their demand. When user's need is over, then these resources are released. This mapping and releasing of resources are done by a task scheduler. As a result, CSP requires minimal effort for managing these resources. When no task scheduler is available in a cloud system, then incoming tasks are randomly mapped to cloud resources. As a result, some resources will be underutilized, whereas some will have immense load on them leading to a deadline situation. This will increase the makespan, which is the time period within which all the queued tasks will be executed, and the task queue will be emptied. So, it can be said that task scheduler is one of the most vital components for cloud computing paradigm.

# 2 Task Scheduling

User sends service requests to the cloud system, and these requests are termed as tasks. A task is a collection of jobs that include entering data to cloud system, processing the data, accessing some required software, and infrastructure or store data to some data center and then return the output to users. Scheduling a task to a particular resource considers various parameters like deadline of the task, its waiting time, arriving time, and cost to perform a task and energy consumption by a resource to execute the task. A task scheduler aims to allocate task to resources in a time and cost efficient manner that will result in optimal resource utilization. In Fig. 1a, basic task-scheduling architecture in cloud environment is shown.

Task scheduling [3–7] is divided into two broad categories—(a) centralized scheduling and (b) distributed scheduling. Figure 1b shows the entire taxonomy of task schedulers.

## 2.1 Centralized Task Scheduling

In centralized scheduling, one single scheduler is available for the entire system to schedule various tasks. Centralized scheduling is very easy to implement; however, the main drawback of centralized scheduling is that the whole system halts when

**Fig. 1** **a** Task-scheduling architecture, **b** task-scheduling classification

the scheduler stops working. Scalability and fault tolerance of a centralized task scheduling are low. This single point failure problem is overcome by distributed task scheduling.

## 2.2 Distributed Task Scheduling

In distributed scheduling, a collection of schedulers is connected with each other. When any task comes from user to the cloud system, the schedulers perform the scheduling collectively. That means the workload is distributed among multiple schedulers, thereby saving significant amount of process cycle. Advantage of distributed scheduling is that this scheduling type exchanges workload with neighbor nodes. As a result, the problem of single point break down is eliminated. However, distributed scheduling is highly complex as there are several interconnections between schedulers.

The remaining paper is divided into following sections. In Sect. 3, a literature review on different existing task-scheduling algorithms has been conducted. In Sect. 4, a comparison table and a comparative study among the algorithms are given. Finally, Sect. 5 draws the conclusion to this survey with some future scopes in this field.

# 3   Literature Review of Existing Task-Scheduling Algorithms

## 3.1   Activity-Based Costing Algorithm

In [8], an Activity-Based Costing (ABC) algorithm for task scheduling in cloud environment has been proposed. This algorithm calculates a priority list for each task that arrives in cloud system. This priority is calculated by considering number of resources required to complete the task, cost for accessing each resource, and the profit that cloud service provider can earn on successful completion of that task.

The scheduling of task is done in two phases. In Fig. 2, working strategy of this algorithm is explained. At first, this algorithm creates three priority activity lists namely LOW, MID, HIGH. In the first phase, the priority for all incoming tasks is calculated and they are put in an appropriate list based on their priority value. In cloud system CPU, storage and various other services are considered as resources. In the next phase of the algorithm, listed tasks are allocated to some resource. Tasks scheduled as HIGH list priorities are dispatched first followed by MID and finally from LOW list. This phase is in action until all the tasks are allocated to some resource. For different resource of cloud system, there is different cost for usage. The total cost for a task is calculated by summing up the cost of different individual resources needed. This algorithm sets a low-cost value for more frequent tasks. However, the total profit is high as the request for such task is high. On the other hand, high cost is charged for special tasks like providing customized firewall service to a user or giving specialized priority value to some special users. The service provider is able to earn a good profit though the frequency of such tasks. The ABC algorithm does not consider time as a parameter. As a result, if there is large number of tasks in the HIGH list, the tasks in the LOW list have to wait for a longer time period. Such tasks may miss deadline resulting in customer dissatisfaction and hampering the quality of service.



**Fig. 2** Task and resource mapping by activity-based costing algorithm

## 3.2  *Improved Activity-Based Costing Algorithm*

An improved version of Activity-Based Costing has been proposed by the authors in [9]. It is called Improved Cost-Based Algorithm.

This algorithm also considers cost of individual resources. Working strategy of this algorithm is shown in Fig. 3. Like ABC, even here service provider sets low cost for easier and frequently executed tasks and for special tasks the cost is set to high value. This algorithm executes in three phases. In first phase, all incoming tasks priority is calculated. Next, this algorithm arranges the tasks according to HIGH, LOW, and MID priority queue. In the last phase, all tasks are grouped. This grouping is done based on resource granularity size, i.e., total million number of instruction a resource can handle in a time period. At first, the tasks listed in HIGH queue are grouped. On completion of HIGH-priority tasks, MID and LOW priority listed tasks are grouped further. Total number of Million Instructions (MI) of every group must be less than or equal to the resource's execution capacity. This capacity is multiplication of resource's MIPS (Million Instruction per Second a resource can handle) and prefixed granularity size. When this defined criterion is matched, then the group is allocated to that resource. However, when MI of a group will become larger than a resource's execution capacity, then last task of that group will be discarded and resource allocation will be executed.

## 3.3  *Double Level Priority-Based Optimization Algorithm*

A new kind of task-scheduling strategy is proposed in [10] which works depending upon two levels of priorities. As per Fig. 4, at first level, all tasks are collected. Then these tasks are divided into two categories depending upon resource availability for tasks' execution in one data center. The task is put into the Fully Available category list if all the resources required by task are available in one data center. However, the task is put into the Partially Available category list if the required resources exist in more than one data center.



**Fig. 3** Resource allocation by improved activity-based costing algorithm

**Fig. 4** Resource allocation by double level priority-based optimization algorithm

For a task, if all the required resources are available in a single data center and one task is dependent on any previous activity, then that task is listed in Dependent Task list otherwise listed in Independent Task list. Thus, three types of lists, namely, Independent List, Dependent List, and various category lists are obtained. Upon all these lists, cost-based algorithm is applied to create LOW, MID, and HIGH queue. As a result, multiple low, mid, and high queues are generated. Comparing among the selected tasks, the one with highest priority is selected. In second level, the authors categorized available VMs depending on processing power and cost for computing a task. Resources are sorted based on turnaround time, and highest priority task is assigned a resource first. This process is repeated until all the task lists are emptied.

## 3.4 Cost–Deadline-Based Algorithm

An optimized scheduling algorithm is proposed in [11] which take into consideration the deadline and cost of a task in cloud system.

This algorithm has been proposed based on priority and admission control scheme. Figure 5 shows the scheduler assigns priority for each task arriving for execution and put them in a queue. Admission to this queue is determined by calculating tolerable delay and cost of service. This algorithm is also called Task Priority Delay (TPD) algorithm since three parameters are considered in this approach. They are deadline, cost, and task length. To generate the tasks' priority queue, deadline and task length is considered. The task queue is then sorted according to the calculated priority. In case there are two tasks with same priority, the task that yields maximum profit is selected first. Finally, the MIN–MIN heuristic approach [7] is followed to allocate the tasks to some virtual machine.

**Fig. 5** Resource allocation by cost–deadline algorithm



## 4 Comparative Analysis

A good task-scheduling algorithm plays a key role in providing better performance for user. From the above literature review, various advantages and disadvantages of existing task-scheduling algorithm are identified. Table 1 shows the comparative study of those algorithms. Activity-Based Costing [8] and Improved cost-based algorithm [9] consider only cost as a scheduling parameter. The waiting time of the task is not considered for this algorithm. So the number of tasks that misses its deadline is quite high compared to the scheduler proposed in [11]. Double priority-based optimization algorithm [10] considers only resource availability for a task in a VM, thus gives a cost-effective solution with an improved makespan value. Other QoS parameters, like task with critical deadline, power consumed by each scheduler, are not paid attention to. The scheduling policy in [11] does not give a good QoS for low-priority task and follows a non-preemptive scheduling resulting to a number of context switches. None of the above algorithms aims to minimize the power consumed by the task scheduler which is an important QoS parameter for providing a green environment.

## 5 Conclusion and Future Work

In cloud computing, resources are allocated virtually based on the need of user. In order to manage these resources, good scheduling algorithms are required. In this paper, a brief study is done on different task scheduling algorithms. Most of the existing algorithms consider only one or two parameters for scheduling policy. So these algorithms could be further developed to get better results.

To design a robust task scheduling algorithm, more number of parameters are needed to be considered. For cost optimization-based algorithms, time parameter and energy parameter can be included along with the cost parameter in the scheduling policy. A new kind of cost-based task-scheduling algorithm can be

**Table 1** Comparison table of task scheduling algorithms

| Scheduling name | Scheduling parameter | Advantage | Disadvantage |
| --- | --- | --- | --- |
| ABC [8] | Cost | 1 Both resource cost and overhead cost can be minimized | 1 When overhead costing is high this algorithm may not be efficient<br>2 There is no time parameter in ABC |
| Improved cost-based algorithm [9] | Cost | 1 Total profit is increased<br>2 Makespan is reduced | 1 Tasks with lower costs have to wait long<br>2 Quality Of Service attributes are not considered |
| Double level priority-based optimization [10] | Cost, dependency | 1 Cost is reduced<br>2 Completion time improved | 1 Deadline of task is not considered<br>2 Energy parameter for resources is not considered |
| Cost–deadline-based algorithm [11] | Deadline | 1 Deadline of task is considered<br>2 It predicts initial cost of task | 1 Cost optimization is not up to level<br>2 Low priority tasks may have to wait long due to availability of large tasks with high precedence<br>3 Non-preemptive type |

designed by considering waiting time of a task. For cost-based algorithm, tasks belonging to the LOW queue have to wait longer to gain access to CPU. If priority function is modified using the waiting time of the task, a new result can be obtained that will reduce completion time for low-priority tasks.

Deadline-based costing algorithm could be modified by considering resource utilization parameter while allocating task on different VMs. When resource allocation will be in an efficient way, energy consumed by a resource will automatically decrease. Since resources will be in working mode or alive mode only when a task appears. Otherwise, resources will go into an idle mode and from that to sleep mode resulting in minimum energy consumption.

A good task-scheduling algorithm always aims at both user and service provider satisfaction. For user level satisfaction, cost for a task needs to be low and response with desired result needs to be high. For service provider level satisfaction, profit should be high. So to fulfill these criterions, task-scheduling algorithm should try to pay attention on more and more number of parameters.

# References

1. Mathew, T., Sekaran, K.C., Jose, J.: Study and analysis of various task scheduling algorithms in the cloud computing environment. In: International Conference on Advances in Computing, Communications and Informatics (ICACCI), pp. 658–664 (2014)

2. Salot, P.: A survey of various scheduling algorithm in cloud computing environment. Int. J. Res. Eng. Technol. (2013). ISSN 2319-1163
3. Arya, L.K., Verma, A.: Workflow scheduling algorithms in cloud environment—a survey. In: Recent Advances in Engineering and Computational Sciences, pp. 1–4 (2014)
4. Dave, Y.P., Shelat, A.S., Patel, D.S., Jhaveri, R.H.: Various job scheduling algorithms in cloud computing: a survey. In: International Conference in Information Communication and Embedded Systems, pp. 1–5 (2014)
5. Fakhfakh, F., Kacem, H.H., Kacem, A.H.: Workflow scheduling in cloud computing: a survey. In: 18th IEEE International Enterprise on Distributed Object Computing Conference Workshops and Demonstrations, pp. 372–378 (2014)
6. Patil, S., Kulkarni, R.A., Patil, S.H., Balaji, N.: Performance improvement in cloud computing through dynamic task scheduling algorithm. In: 1st International Conference on Next Generation Computing Technologies, pp. 96–100 (2015)
7. Nagadevi, S., Satyapriya, K., Malathy, D.: A survey on economic cloud schedulers for optimized task scheduling. Int. J. Adv. Eng. Technol. **5**, 58–62 (2013)
8. Cao, Q., Wei, Z.B., Gong, W.M.: An optimized algorithm for task scheduling based on activity based costing in cloud computing. In: 3rd International Conference on Bioinformatics and Biomedical Engineering, pp. 1–3 (2009)
9. Selvarani, S., Sadhasivam, G.S.: Improved cost-based algorithm for task scheduling in cloud computing. In: IEEE International Conference on Computational Intelligence and Computing Research, pp. 1–5 (2010)
10. Parikh, S., Sinha, R.: Double level priority based optimization algorithm for task scheduling in cloud computing. Int. J. Comput. Appl. **62**(20) (2013)
11. Sidhu, H.S.: Cost–deadline based task scheduling in cloud computing. In: Second International Conference on Advances in Computing and Communication Engineering, pp. 273–279 (2015)

# Fuzzy Current Control of Grid Interactive Voltage Source Converter with Solar Energy

R. S. Ravi Sankar and S. V. Jaya Ram Kumar

**Abstract** This work deals with design and simulation of photovoltaic system connected to the grid. Maximum power is extracted from solar panel by means of boost convert which is operated by means of perturb and observe method. The output of the boost converter connected to the three-phase inverter is wired to the grid. Inverter is controlled by the fuzzy control technic. In this control technique, load current is taken as reference current and grid current is taken as actual current. Error in direct and quadrature axis currents is given to the fuzzy controller which generates the reference voltage. These are transferred into the 'abc' frame and compared with the carrier signal, and generated pulses are given to the inverter. The fuzzy current control method switches the power device in the inverter such that minimizes the error between the grid currents and the load current. Grid current tracks the reference current with less transient time, steady-state error, acceptable total harmonic distortion. This is implemented through the MATLAB simulation.

**Keywords** Fuzzy controller · Three-phase voltage source inverter
Photovoltaic (PV) · Total harmonic distortion (THD)
Maximum power point tracking (MPPT)

## 1 Introduction

Renewable energy resources are an interesting alternative to that of the non-renewable energy sources because they can be used on a large scale without polluting the environment. Photovoltaic (PV) power generation is a concept of converting solar energy into direct current electricity using semiconductor materi-

R. S. Ravi Sankar (✉)
Electrical and Electronic Engineering Department, Vignan's Institute of Information
Technology, Visakhapatnam 530046, India
e-mail: Satya_ravi2001@yahoo.com

S. V. Jaya Ram Kumar
Electrical and Electronic Engineering Department, GRIET, Hyderabad, India

als. A PV system consisting of the solar panel has a basic element of solar cells which are used to generate the required solar power. These renewable energies can be used efficiently with power electronic converters. When power converters are placed in the system, harmonics are generated due to the switching process, which may cause disturbances in the distribution grid, so these harmonics should be filtered, which requires the usage of filters like 'L, LCL, LC' to attain the grid current THD below 5% [1].

Control of one-phase inverter in synchronous (usually few kW) has been reported by [2–4]. PI controller gains are tuned by means of fuzzy for the adaptive control strategy and have given the fast transient response and better performance under various parametric conditions of one-phase inverter connected to the grid [5]. Comparison of the classical PI and fuzzy logic PI current controller of PWM rectifier implemented based on field orientation control has given better performance [6]. Fuzzy current controller used in the inner current control loop and fuzzy controller used in the outer voltage loop to control the power flow between the fuel cell and battery system connected to grid have given the better performance in terms of the THD value of the grid current under the different parametric conditions [7].

In this paper, Sect. 2 explains modeling of PV system and MPPT, and Sect. 3 describes about the fuzzy controller and overall block diagram of system. Section 4 projects the simulation results.

## 2 PV System and MPPT

### 2.1 PV Array Mathematical Modeling

PV array comprises the number of solar modules wired in cascaded and shunted to meet the required voltage and power rating. Figure 1 illustrates the basic element of PV array which is simple P–N junction that converts the solar radiation into electricity. $I_{pv}$ is the current generated by the light. D1 and D2 are two antiparallel diodes, $R_p$ is the shunt leakage resistance, and $R_s$ is the contact resistance, and $V$ is output voltage of single cell, and Now, current $I$ is given (1) [8, 9]

$$I = I_{pv} - I_{D11} - I_{D22} - \left(\frac{V + IR_s}{R_p}\right) \tag{1}$$

$$I_{D_{11}} = I_{o11}\left[\exp\left(\frac{q(V + IR_S)}{A_{11}KT}\right) - 1\right] \tag{2}$$

$$I_{D_{22}} = I_{o22}\left[\exp\left(\frac{q(V + IR_S)}{A_{22}KT}\right) - 1\right] \tag{3}$$

$$I_{\text{pv}} = [I_{scr} + K_i(T_k - T_{\text{ref}k})] \times \lambda/1000 \tag{4}$$

Here, $\lambda$ is the solar irradiation, $I_{scr}$ is cell shorted current, $T_k$ and $T_{\text{ref}k}$ are the actual and standard temperatures, $K$ is temperature coefficient of short-circuit current (A/K), $A_{11}$ and $A_{22}$ are the diode ideality factors, '$q$' is the charge of the electron, and $I_{o11}$ and $I_{o22}$ are the reverse saturation currents of $D_1$ and $D_2$, respectively. Equation (1) is modified for PV module as (5–6) [10]:

$$I = N_{pp}I_{\text{pv}} - N_{pp}I_{D1} - N_{pp}I_{D2} - \left(\frac{V + IR_s}{R_p}\right) \tag{5}$$

$$
\begin{aligned}
I = N_{pp} \times I_{\text{pv}} - N_{pp} \times I_{o11} &\left[\exp\left\{\frac{q \times (V + IR_s)}{N_s A_{11} kT}\right\} - 1\right] - N_{pp} \\
\times I_{o22} &\left[\exp\left\{\frac{q \times (V + IR_s)}{N_s A_{22} kT}\right\} - 1\right] - \left(\frac{V + IR_s}{R_p}\right)
\end{aligned}
\tag{6}
$$

Here, $N_{ss}$ and $N_{pp}$ are the number of cells that are cascaded and shunted in module. PV modules are connected in cascaded and shunted to meet the voltage and power rating of PV array. So the PV array current equation becomes (7)

$$
\begin{aligned}
I = N_{ppp} \times I_{\text{pv}} - N_{ppp} \times I_{o1} &\left[\exp\left\{\frac{q \times \left(V + IR_s\left(\frac{N_{ss}}{N_{pp}}\right)\right)}{N_{sss} A_{11} kT}\right\} - 1\right] - N_{ppp} \\
\times I_{o2} &\left[\exp\left\{\frac{q \times \left(V + IR_s\left(\frac{N_{ss}}{N_{pp}}\right)\right)}{N_{ss} A_{22} kT}\right\} - 1\right] - \left(\frac{V + IR_s\left(\frac{N_{ss}}{N_{pp}}\right)}{R_p\left(\frac{N_{ss}}{N_{pp}}\right)}\right)
\end{aligned}
\tag{7}
$$

where $N_{sss}$ and $N_{ppp}$ are the number of cascade- and shunt-connected modules, respectively [10]. In this paper, standard KC200GT data sheet parameters are used to do PV array simulation.



**Fig. 1**   PV cell equivalent circuit

## 2.2  Maximum Power Point Tracking

Temperature and irradiations values of Light is depending on climatic conditions, so to extract the maximum power from PV array. Perturb and observe algorithm is implemented in this paper, which is a simple method by means of which variable duty cycle is generated and given to the boost converter.

## 2.3  DC–DC Boost Converter

In this work, boost convert is used to extract the maximum power from the PV array. It consists of the inductor (L), one controlled switch (S), and one uncontrolled switch (D), as shown in Fig. 2. When $S$ is 'ON,' inductor stores the energy. When it is OFF, stored energy is released to the load. Capacitor is used at the out to reduce the replies in the voltage and gives the smooth voltage to the inverter.

To design the $L$ and C values, the PV array voltage is taken as input ($V_I$) 152 V. Switching frequency 10 kHz and output voltage ($V_o$) 300 V are considered. From the Eq. (8). The values of $C$ = 15 mF and $L$ = 0.1 mH.

$$\text{Duty cycle } D = 1 - \frac{V_I}{V_o}; \quad \text{Inductance } L = \frac{V_I D}{\Delta If}; \quad \text{Capacitance } C = \frac{I_0 D}{\Delta If} \quad (8)$$

## 3  Fuzzy Current Control Technique

The architecture of this controller is shown in Fig. 3. It has three units: fuzzifier unit, interface unit, and defuzzifier at the out terminals. It has input and output variables. Error and change in error currents are inputs, whereas voltage is an output variable. The current error means difference between the load current $I_L(k)$ and grid current $I_g(k)$, and process output is reference voltage $y(k)$. Change in error current means difference between the current values of two successive sample periods,



**Fig. 2**  Boost converter circuit

which are defined in (9) and (10), respectively. In this work, seven membership function levels are defined as high negative (HN3), moderate negative (MN2), small negative (SN1), zero (Z), small positive (SP1), moderate positive (MP2), and high positive (HP3) which are shown in Table 1. The input variables are processed by membership functions. Triangular or trapezoidal functions are simple and fast computation, which are used in this work.

$$ei\,(k) = I_l(k) - I_g(k) \qquad (9)$$

$$\Delta ei = ei(k) - ei(k-1) \qquad (10)$$

In this work, two fuzzy controllers are used as shown in Fig. 7, whose membership functions are shown as in Figs. 4, 5, and 6.

Membership Functions for $I_d$ and $I_q$

Rules represented given below:

Rule1: 'error($ei$)' is HN3, and error in change '($\Delta ei$)' is HN3; then, $\Delta u$ is HN3
Rule2: 'error($ei$)' is HN3, and error in change '($\Delta ei$)' is HN3; then, $\Delta u$ is HN3
- - - - - - - - - - -
Rule49: 'error($ei$)' is HP3, and error in change '($\Delta ei$)' is HP3; then, $\Delta u$ is HP3.
Overall System Block Diagram of the System



**Fig. 3** Architecture of the fuzzy controller

**Fig. 4** Error current membership function



**Fig. 5** Change in error current membership function

## 4  Simulation Results

In this paper, PV array is designed for a power of 12 kW with an output of 150 V. It has 5 modules in cascade and 10 strings in parallel. Voltage at PV array depends on the number of modules in cascade. It is shown in Fig. 8; initially, it has some transient behavior by means of the DC–DC converter, and the output voltage is increased to 300 V, as shown in Fig. 9; it acts as an input to the three-phase inverter circuit.

At Point of Common Coupleing (PCC), static load is connected which draws a peak current of 5 A as shown in first wave in Fig. 10, which is compared with grid current as shown in the second waveform in Fig. 10. The error in current is processed through fuzzy controller to generate the reference voltage is shown in Fig. 11. This voltage is compared with triangular carrier wave of frequency 10 kHz

**Fig. 6** Output voltage membership function

**Table 1** Rule matrix table

| Error Δerror | HN3 | MN2 | SN1 | Z | SP1 | MP2 | HP3 |
|---|---|---|---|---|---|---|---|
| HN3 | HN3 | HN3 | SN1 | SN1 | SN1 | SN1 | Z |
| MN2 | HN3 | MN2 | MN2 | SN1 | SN1 | Z | SP1 |
| SN1 | MN2 | MN2 | SN1 | SN1 | Z | SP1 | SP1 |
| Z | MN2 | SN1 | N1 | Z | SP1 | SP1 | MP2 |
| SP1 | SN1 | SN1 | Z | SP1 | SP1 | MP2 | MP2 |
| MP2 | SN1 | Z | SP1 | SP1 | MP2 | MP2 | HP3 |
| HP3 | Z | SP1 | SP1 | MP2 | MP2 | HP3 | HP3 |

## Overall System block diagram of the System



**Fig. 7** Overall system block diagram representation

to produce pulses to the PWM Inverter. These pulses are shown in Fig. 12. The inverter output voltage as shown in Fig. 13. The grid voltage shown in Fig. 14. This controller has less transit behavior where peak transient current is 9 A also within the limit.

**Fig. 8** PV voltage



**Fig. 9** PV voltage after MPPT



**Fig. 10** Load and grid current

Total harmonic distortion of the grid current as shown in Fig. 15 has very less value 0.99, which is in the limit.

**Fig. 11** Reference voltage generated by fuzzy controller



**Fig. 12** Pulses generated by PWM modulator



**Fig. 13** 3-phase inverter voltage



**Fig. 14** Grid voltage

**Fig. 15**   Analysis of grid current

## 5   Conclusion

This work presented the fuzzy current controller for a PV inverter connected to the grid. It takes the input as error in current, and change in error current gives the reference voltage to the PWM modulator. Pulse is produced which is applied to the grid-connected inverter. This controller that produced balanced grid current is same as the load current with less THD value and low transient behavior. It is simple current control and implementation without sacrificing quality and accuracy of the grid current.

## References

1. Ekanayake, B.J., Holland, M.P.: Exploiting PV inverters to support local voltage—small signal model. IEEE Trans. Energy Convers. **29**(2) (2014)
2. Kjaer, S.B., Pedersen, J.K., Blaabjerg, F.: A review of single-phase grid-connected inverters for PV modules. IEEE Trans. Ind. Appl. **41**(5), 1292–1306 (2005)
3. Rodrguez, J., Ponttz, J., Silva, A.C., Correa, P.: Predictive control of a voltage source inverter. IEEE Trans. Ind. Electron. **54**(1), 495–503 (2007)
4. Ravi Sankar, R.S., Jayaram Kumar, S.V., Deepika, K.K.: Model predictive current control of PV inverter connected to the grid. Indo Natl. J. Electr. Eng. Comput. Sci. **2**(2), 285–296 (2016)
5. Sefa, I., Altin, N., Ozdemir, S.: Fuzzy controller inverter for grid interactive renewable energy systems. IET Renew. Power Gener. **9**(7), 729–738 (2015)
6. Jansinski, M., Liserre, M., Blaabjerg, F., Cichowlas, M.: Fuzzy logic current controller for pulse width modulation rectifiers. In: IECON 02, 2002 IEEE Conference, vol. 2, pp. 1300–1305
7. Zambri, N.A., Ismali, N.M.: Performance comparison of PI and PI-fuzzy controller for grid interactive fuel cell inverter system. In: 2015 IEEE Conference, pp. 1–6
8. Kumar Das, S., Akil Raju, R.: Development of PV cell/module/array and non-uniform irradiance effect based on two diode model by using PSPICE. In: International Conference on 2015, Nascent Technologies in the Engineering Field (ICNTE), 9–10 Jan 2015, pp. 1–6 (2015)
9. Satarupa, B., Anup, A., Chitti Babu, B.: Comparative analysis of mathematical modeling of (PV) array. In: Annual IEEE India Conference (INDICON2012), 7–9 Dec 2012, pp. 269–274
10. Marcelo Gradella, V., Jonas Rafael, G.: Comprehensive approach to modeling and simulation of PV arrays. IEEE Trans. Power Electron. **24**(5) (2009)

# Improving the Map and Shuffle Phases in Hadoop MapReduce

J. V. N. Lakshmi

**Abstract** Massive amounts of data are needed to be processed as analysis is becoming a challenging issue for network-centric applications in data management. Advanced tools are required for processing such data sets for analyzing. As a proficient analogous computing programming representation, MapReduce and Hadoop are employed for extensive data analysis applications. However, MapReduce still suffers with performance problems and MapReduce uses a shuffle phase as a featured element for logical I/O strategy. The map phase requires an improvement in its performance as this phase's output acts as an input to the next phase. Its result reveals the efficiency, so map phase needs some intermediate checkpoints which regularly monitor all the splits generated by intermediate phases. MapReduce model is designed in a way that there is a need to wait until all maps accomplish their given task. This acts as a barrier for effective resource utilization. This paper implements shuffle as a service component to decrease the overall execution time of jobs, monitor map phase by skew handling, and increase resource utilization in a cluster.

**Keywords** MapReduce · Hadoop · Shuffle · Big data
Data analytics · HDFS

## 1 Introduction

The objective of data analytics is scrutinizing, cleansing, renovating, and molding of the data for extracting functional information, portentous termination and sustaining choice making [1]. Data analysis has various sides and looming methods beneath diverse identities in special business, science and social science fields [2].

Big data is a meticulous technique of data analysis that focuses on analyzing huge data sets which materialize from various fields of intensive informatics data

J. V. N. Lakshmi (✉)
AIMS Institutes of Higher Education, Peenya, Bengaluru, Karnataka, India
e-mail: jlakshmi.research@gmail.com

centers [3]. Big data typically comprises of data sets of massive volume beyond the skill of traditional software tools to analyze, handle, and process the data [4].

Procedures written in this practical way are mechanically parallelized and implemented on an immense cluster of commodity equipment [5, 6]. In program execution, runtime structures are concerned of the splits which are scheduled in handling many operations such as implementation across set of machines, managing failures, and handling inter-machine communications [7]. The crucial drawback is exhibited on Hadoop performance affecting the cluster.

The significant explanation of Hadoop is outlined as below:

(1) Distinct phases are leaped into a single task—the implementation of reduce function is CPU intensive and memory intensive as to segregate the map task data and produce the absolute outcome.
(2) Arbitrary requests from I/O effecting the shuffle phase—task tracker receives plenty of I/O reading requests. Each request will prompt plenty of I/O reading operations with different offset on the task tracker.

In this paper, an attempt is done to extricate shuffle phase from reduce task and instrument it as a standard resource provider. Integrate the shuffle service with sequential read policy and handling partitioning skew in reduce task to manage stragglers. Section 2 portrays the background and Hadoop MapReduce programming model, and Sect. 3 describes the problem statement. Section 4 discusses design process, and Sect. 5 analyzes on improvement in map phase. Section 6 involves the evaluation of algorithm, and Sect. 7 reviews the results. Finally, Sect. 8 concludes.

## 2 Background

The recent efforts from Hadoop MapReduce features are analyzed in improving performance are illustrated as follows:

(1) Map step, reduce step, the sort and merge step are included in Google MapReduce model implemented by Hungchih yang, Ali Dasdan et al.
(2) An architectural combination of MapReduce and database technologies resulted as HadoopDB is developed for analytical workloads.
(3) Hadoop MapReduce HDFS layer is replaced with concurrency optimized data storage layer which improves efficiency of data accessing concurrency, proposed by B Nicolae, G Antoniu et al.
(4) A pipeline architecture was proposed by N Conway, T Condie et al., which supports online streaming for many networking sites
(5) Resource manager and scheduler are alienated into separate components by YARN from Apache for solving the blockage of job tracker.

## 2.1 MapReduce

A data flow standard such as MapReduce is widely used for parallelizing the data on various applications [8]. This is a simple and open data flow programming model preferential when compared over usual high-level database approaches. This training model is used for processing large-scale datasets in computer clusters by exercising two function map ( ) and Reduce ( ). The functions Map ( ) and reduce ( ) are as follows:

$$\text{Map}(K1, \ V1) \rightarrow \text{list}(K2, \ V2) \quad \text{Reduce}(K2, \ \text{list}(V2)) \rightarrow \text{list}(V2)$$

The Map ( ) functions uses key/value pair as input generating the intermediate key/value pairs. The generated intermediate key/value pairs are the input given to reduce function to produce final output [9].

## 2.2 Hadoop

Hadoop executes shuffle as a component of reduce task because of which there is high utilization of bandwidth in the cluster, resulting low usage of processor and unproductive performance [10, 11].

Hadoop distributed file system (HDFS) provides high throughput access to application data, resource allocation task in cluster and high unsystematic disk I/O requests are suitable for application that has large data sets [12].

From the Fig. 1, data in a Hadoop cluster is busted down into minor portions and circulated all through the collection, where a job tracker keeps track of jobs in both parent and child segments. The map and reduce functions can be implemented on



**Fig. 1** MapReduce model

slighter subsets of your larger data sets, and this provide the scalability metrics that is needed for data processing [13].

## 3   Problem Statement

MapReduce Programming Model is very simple but as it processes, we come across many problems in map ( ) function. Map ( ) function is assigned with each split if one split cannot execute with any problem (or) if one split fails then we cannot compute the result of Map ( ) function. As combining the individual result of each map function is assigned as input to reduce function, Map function should perform in better way [14].

Two tasks associated with improved reduce phase are shuffle part and reduce part. The initial shuffle segment calls for transitional outcome from map phase. This necessitates more buffer area for various operations sorting and mapping to elicit output.

Numerous disk I/O requests from shuffle phase result in inefficient usage of resources. The above-specified reasons lead to cluster performance problems. The analysis shows an improvement in certain phases of Hadoop MapReduce specifically in terms of execution [15, 16].

## 4   Design Process

Shuffle and reduce as individual stages of tasks: Primarily remove copy and merge operations of shuffle from reduce as an entity splits.

### 4.1   Joining of Unusual Splits into Solitary Task

The shuffle phase fetches the transitional outcome from each and every map task where as the reduce function could not start its processing until shuffle phase releases the processed output data. This wastes the CPU resource time and decreases the network bandwidth.

### 4.2   Random I/O Request of Shuffle Task

Each map task needs to read facts from disk and transfer the response to defined reduce task instantly. This results in large amount of random disk I/O operations which in turn reduces the performance.

## 4.3 Design

Our features mainly involve the following stages to increase enhancement. Shuffle can process meager data improving the resource utilization efficiently within the same amount of time but the disk I/O request is progressively increased.

Fig. 2 describes the various stages of improved MapReduce architecture by implementing the technique of disjoint maps with skew in them. They are handled separately by the task tracker in slave node. The usage of generate function improves the shuffle phase and processes the data to reduce task.

Services from Shuffle: By implementing shuffle as service, resource utilization has been incremented because light weight common service relocates the on command for reduce task as a service [17].

Overcoming stragglers: To avoid blocking of slots, a skewed task is recognized and implemented. A skewed task is identified and accomplished them under similar task master. It detects partitioning skew before shuffling of data begins by monitoring data sizes produced by map and handles it by dynamically creating multiple reduce task per skewed partition [18].

Managing disk I/O requests in map phase: The I/O requests from different disk drives are processed within certain interval of time. These requests are sorted and grouped into a sequential list, forwarding them to respective output files of map tasks [19]. Task tracker reconstructs responses by reading data emerging from disk and sending their responses to reduce task in order.



**Fig. 2** Improved shuffle phase and handling straggler

## 5   Improving the Map Phase

From the problem statement discussed, there is a need to improve the map ( ). The function needs some checkpoints which monitor the function regularly and try to solve when a split intimates an interrupt.

---

*Algorithm*: Improving Shuffle service and avoiding stragglers in Map task.

---

1.  Map phase runs the map task
     //each map task does the job
     //assigned by the Job Tracker
2.  if(collision=true)
3.  Generate(); //Straggler detection and removal     and partition the task with skew
     //collision is an interrupt due to which map task cannot give an output
     // if there is a collision it calls generate method
4.  Map() /* After the map processing the task tracker initialize the shuffle task launching the services*/
5.  End if
6.  Regular Check_points()
     // to verify proper execution
7.  // Managing Disk I/O request by processing
         Read_ratio, read_time, file_size, read_total methods
8.  Map task output is transferred to Shuffle phase
9.  Shuffle sorts and merges the generated output
            $st= dr/bw$
       *st*: shuffle time for reduce task
       *dr*: data to be shuffle per reduce
       *bw*: band width between nodes.
10. This transfers the results to Reduce()

---

The generate function monitors the map phase at regular checkpoint and views the status of each map split. These checkpoints are arranged dynamically and access the needs of the splits. Distributed storage structure shares information among different tasks. The above algorithm specifies the design in Fig. 2, and various phases of handling the stragglers and handling the resources efficiently are elaborated.

The map jobs are scheduled in a queue, and reduce jobs use priority queue structure. In this way, interpretation of result from the map ( ) results intermediate key/value pairs [20]. These pairs are given as an input to the reduce function, and after interpretation, we generate the final output. The generate function is also used even in reduce phase. Dynamically arranged checkpoints monitor the reduce phase and split into smaller splits when an interrupt occurs. It finally combines all the split's output for obtaining the final result.

# 6 Evaluation

Presenting the performance and resource utilization of MapReduce jobs by implementing the shuffle service can be analyzed below:

## 6.1 Simulation Experiments

Intricate methods, routine calls, resource requirements, etiquette, and exchanges in the Hadoop cluster influence the ratio of disk read/write operations. Because of these random requests, there is a decrease in the regular reading ratio on disk and by which there is an increase in reduce task time.

## 6.2 MapReduce Job Experiment

Pi estimator utilizes more of CPU computations so it is CPU-intensive task where as word count and TeraSort are resource oriented. The resources are memory and band width.

## 6.3 Straggler Handling

The techniques are employed in distributing the reducers with even number of map outputs in parallel, ensuring there are no skews.

## 6.4 Settings

The configuration of our Hadoop with 0.24 version requires 12-node cluster among one is a master and remaining ten are slaves. Every node in a cluster uses core processor organizes 2 GHZs 4 GB of Ram and 500 GB disk drive.

*Read total*: Sum of read operations per each reduce task.

*File size*: Volume of information produced by each map task in core state.

*Read time*: Mean time between transfer requests and to obtain data for process.

*Read ratio*: Average ratio of read total to read time.

*Local read total*: Word count of reduce task with comparison of job time.

*Reduce_skew*: Stragglers count in reduce task.

# 7   Results

Reading performance is verified with an improved fetch phase with varied file sizes such as 128, 256, and 512 MB than the earlier fetch task. If sequential strategy on read operation is applied, then mean increase in read ratio is 94.17%, and if concurrent strategy is applied, then the ratio is 62.81%.

Figure 3 shows the data-read in local mode by varying in their speed of accessing. The data read is measured as 128, 256 and 512 MB per second avoiding stranglers. Drawing the result from word count showing the diminishing time utilization of reduce phase from 8.94% to 6.32%.

The reduce phase utilization of resources and the word count are low as the data from map phase has to release the entire output. Best word count is observed in Fig. 4 as it visualizes the read and write ratios on disk. After the necessary modifications done for the shuffle phase, the graph illustrates the improvement by showing 7% increase in resource utilization.

# 8   Conclusion

MapReduce programming model requires improvement in map phase as well as in shuffle phase. Though it is simple, but while implementation some complications are observed at map phase. If one map fails, it cannot compute the output as the result of map phase is an output for reduce phase.

The reduce phase adds a scheduler for every node. So, by using generate function which dynamically monitors the reduce phase will solve the basic problem in map phase. Cluster resources are well utilized efficiently when data is huge for processing transitional information then shuffle is determined as a service with minute amounts of time.

Hadoop MapReduce uses word count and TeraSort which acts as an added advantage for performance enhancement with different data structures. Resource deployment, absolute time usage are perfection features observed in skew handling technique.



Fig. 3   Read total (MB) local read

**Fig. 4** Word count
comparison of reduce task



# References

1. Arulmurugan, A., Srinivasan, R.: Enhanced task scheduling scheme for Hadoop MapReduce systems. In: IJETCSE, May 2015
2. Dimitris, F., Ioannis, M.: Scheduling MapReduce Jobs and Data Shuffle on Unrelated Process. MIT, Cambridge (2015)
3. Pavloet, A.: A comparison of approaches to large-scale data analysis. In: Proceedings of ACM SIGMOD, vol. 5, pp. 367–378 (2009)
4. Yandong, W., Yu, W., Que, X.: Virtual shuffling for efficient data movement in MapReduce. In: IEEE Transitions on Computers Conference, June 2015
5. Luiz, A.B., Jeffrey, D., Holzle, U.: Web search for a planet: the Google cluster architecture. IEEE Micro **23**(2), 22–28 (2003)
6. Huston, L., Wickremesinghe, R., SatyaNarayana, M.: Storage architecture for early discard in interactive search. In: FAST Conference Proceedings (2004)
7. Dean, J., Ghemawat, S.: MapReduce: simplified data processing on large clusters in Google, Inc OSDI (2004).
8. Lakshmi, J.V.N., Ananthi, S.: A theoretical model for big data analytics using machine learning algorithms. In: ICACCI Conference, Delhi, October 2015
9. Kwon, Y.C., Howe, B.: A study of skew in Map Reduce application. In: International Conference, USA (2014)
10. Alan, F.G., Olga, N., Shubham, C., Pradeep, K., Shravan, M.N.: Building a high level dataflow system on top of MapReduce: the pig experience. In: IEEE Conference (2009)
11. Yanfei, G., Jia, R., Xiaobo, Z: IShuffle—improving Hadoop performance with shuffle-on-write. In: USENIX ICAC, USA (2013)
12. Abouzeid, A., Bajda, P., Abadi, D.J., Rasin, A., et al.: HadoopDB: an architectural hybrid of MapReduce and DBMS technologies for analytical workloads. PVLDB **2**(1), 922–933 (2009)
13. Ananthi, S., Lakshmi, J.V.N.: A study on Hadoop architecture for big data analytics. In: Delhi Conference ICETSCET, September 2014
14. Herodotos, H., Lim, H., Luo, G.: StarFish—a self tuning system for Big Data Analytics, CIDR, USA (2011)
15. Ronnie, C., et al.: SCOPE: easy and efficient parallel processing of massive data sets. In: Proceedings of VLDB (2008)
16. Ashish, T., Joy deep Sen, S.: HIVE—a warehousing solution over a MapReduce framework. In: VLDB (2009)
17. Li, J., Ye, Y.: Improving the shuffle of Hadoop MapReduce. In: Proceedings of IEEE ICCCTS (2013)
18. Li, J., Yue, Y., Lin, X.: Improving the shuffle of Hadoop MapReduce. In: IEEE ICCCTS, Beijing, China (2013)

19. Prateek, D., Sriram, K., Janakiram, D.: Chisel: resource savvy approach for handling skew in MapReduce application. In: IEEE Conference on Cloud Computing, vol. 35, pp. 45–56 (2013)
20. Dean, J., Ghemawat, S.: MapReduce: a flexible data processing tool. ACM Commun. **53**, 72–77 (2010)

# Improved Decision Making Through IFSS

**T. R. Sooraj, R. K. Mohanty and B. K. Tripathy**

**Abstract**  Decision making has become a common feature in day-to-day activities. Uncertainty-based models are more efficient in handling such problems. In this chapter, we propose an algorithm in this direction by using a hybrid model formed by combining the two models of soft set (SS) and the intuitionistic fuzzy set (IFS). We follow the characteristic function approach used by Tripathy et al. for the purpose. The illustrative real-life example shows the efficiency of our algorithm over other such models.

**Keywords**  SS · Fuzzy set · IFS · Decision making

## 1  Introduction

In order to handle uncertainty-based data sets, several models [8, 24–26] have been found in literature among which FS [27] and its extension IFS [1] have been very fruitful in real life. SS [7] was introduced by Molodtsov in [7] as a parametric tool to handle uncertainty, which is its strong point over many other such models. Among the many applications of SS [2–6, 9–11], decision making has drawn the attentions of several researchers [13–23]. Hybrid algorithms along with SS have been used in many instances to develop such algorithms [13–23]. A new approach was introduced in defining the SS in [12], and also, the basic operations on them were rigorously defined. We follow this approach to redefine IFSS in this chapter, and also, the basic operations are revised. It is a common practice for all of us in real-life situations to take decisions. There are several algorithms in literature to

T. R. Sooraj (✉) · R. K. Mohanty · B. K. Tripathy
SCOPE, VIT University, Vellore, Tamilnadu, India
e-mail: soorajtr19@gmail.com

R. K. Mohanty
e-mail: rknmohanty@gmail.com

B. K. Tripathy
e-mail: tripathybk@vit.ac.in

support decision making by individually or otherwise [9–11, 13–23]. In this chapter, we propose an improved algorithm for data presented in the form of IFSS structure and establish its efficiency. The parameters used under the realm of SS help in setting the characteristics of domain of study under consideration. The working principle of the algorithm is explained and illustrated through a real-life application.

## 2 Definitions and Notions

Let $U$ & $K$ be the set of universal and parameters, respectively.

**Definition 2.1** A soft set, (Z, K) is defined as follows.

$$Z : K \rightarrow \text{PS}(U) \tag{2.1}$$

PS $(U)$ is the power set of $U$.

**Definition 2.2** A FSS, (Z, K) is defined as follows.

$$Z : K \rightarrow \text{FS}(U) \tag{2.2}$$

FS $(U)$ is the set of all fuzzy subsets of $U$.

**Definition 2.3** An intuitionistic FSS (IFSS), (Z, K) is defined as follows.

$$Z : K \rightarrow \text{IFS}(U) \tag{2.3}$$

IFS $(U)$ is the set of all intuitionistic fuzzy subsets of $U$.

Tripathy et al. [17] defined the membership ($\mu^e_{(Z,K)}$) and non-membership function ($v^e_{(Z,K)}$) for IFSS. By using $\mu^e_{(Z,K)}$ and $v^e_{(Z,K)}$, they redefined all the operations of IFSS.

## 3 Improved Decision Making Through IFSS

Tripathy et al. [16] discussed an application of IFSS. Here, we improve decision making algorithm with the help of a new score Eq. (3.1) and compare our algorithm with the algorithm proposed in [16] and show how this algorithm is better than the algorithm mentioned in [16]. As we know, to select a suitable house, the user has to consider the parameters which match to his needs. Let the parameters be $e_1$, $e_2$, $e_3$, $e_4$, and $e_5$. Here, $e_1$ = beautiful, $e_2$ = wooden, $e_3$ = green surrounded, $e_4$ = expensive, and $e_5$ = distance. As mentioned in [16], the parameters $e_4$ and $e_4$ are negative parameters. The algorithm is discussed in the following section.

## *3.1  Algorithm*

1. Input IFSS (*F*, *E*).
2. Input the priority given by customer for every parameter.
3. Construct priority table (PT) [17] for $\mu$, $v$ and $h$.
4. Construct comparison table (CT) [17] for $\mu$, $v$ and $h$.
5. Apply the Eq. 3.1. to find the score.

$$\text{Score} = \begin{cases} -\mu(1+h) & \text{if } m<0 \ \& \ h<-1 \\ \mu(1+h) & \text{Otherwise} \end{cases} \tag{3.1}$$

6. The object having the highest score is the best choice.

## 4  Illustration of the Algorithm

Here, we take the same example as we took in [16].

The Universal set, *U* = {house1, house2, house3, house4, house5, house6}

The parameter set *E* = {beautiful, wooden, green surrounded, expensive, distance}. Let (*F*, *E*) be the IVFSS which describe the "attractiveness of houses," shown in Table 1.

The customer gives priorities for parameters as 0.7, 0, 0.2, −0.5, and −0.2, respectively.

Next step is to find PT [16] for $\mu$, $v$ and $h$ values. It is shown in Tables 2, 3, and 4.

Next step is to find the comparison table for membership, non-membership, and hesitation values. It is shown in Table 5, 6 and 7.

By using the formula (3.1), the decision table can be formulated

**Table 1** Tabular representation of the IFSS (*F*, *E*)

| *U* | Beautiful | | Wooden | | Green surrounded | | Expensive | | Distance | |
|-----|------|------|------|------|------|------|------|------|------|------|
| $h_1$ | 0.10 | 0.70 | 0.00 | 0.90 | 0.20 | 0.70 | 0.10 | 0.80 | 0.80 | 0.20 |
| $h_2$ | 0.90 | 0.00 | 0.60 | 0.40 | 0.80 | 0.20 | 0.80 | 0.10 | 0.30 | 0.60 |
| $h_3$ | 0.30 | 0.50 | 0.10 | 0.80 | 0.20 | 0.70 | 0.30 | 0.40 | 0.40 | 0.60 |
| $h_4$ | 0.70 | 0.10 | 0.70 | 0.20 | 0.60 | 0.40 | 0.60 | 0.20 | 0.60 | 0.30 |
| $h_5$ | 0.30 | 0.40 | 0.40 | 0.50 | 0.40 | 0.50 | 0.50 | 0.20 | 0.10 | 0.90 |
| $h_6$ | 0.90 | 0.10 | 0.50 | 0.40 | 0.60 | 0.30 | 0.60 | 0.20 | 0.50 | 0.40 |

**Table 2** Priority table for membership values

| $U$ | Beautiful | Wooden | Green surrounded | Expensive | Distance | Row sum |
|-----|-----------|--------|------------------|-----------|----------|---------|
| $h_1$ | 0.07 | 0.0 | 0.04 | −0.05 | −0.16 | −0.10 |
| $h_2$ | 0.63 | 0.0 | 0.16 | −0.40 | −0.06 | 0.33 |
| $h_3$ | 0.21 | 0.0 | 0.04 | −0.15 | −0.08 | 0.02 |
| $h_4$ | 0.49 | 0.0 | 0.12 | −0.30 | −0.12 | 0.19 |
| $h_5$ | 0.21 | 0.0 | 0.08 | −0.25 | −0.02 | 0.02 |
| $h_6$ | 0.63 | 0.0 | 0.12 | −0.30 | −0.10 | 0.35 |

**Table 3** Priority table for non-membership values

| $U$ | Beautiful | Wooden | Green surrounded | Expensive | Distance | Row sum |
|-----|-----------|--------|------------------|-----------|----------|---------|
| hi | 0.21 | 0.9 | 0.56 | −0.4 | −0.16 | 1.11 |
| $h_2$ | 0 | 0.4 | 0.16 | −0.05 | −0.48 | 0.03 |
| $h_3$ | 0.15 | 0.8 | 0.56 | −0.2 | −0.48 | 0.83 |
| $h_4$ | 0.03 | 0.2 | 0.32 | −0.1 | −0.24 | 0.21 |
| $h_5$ | 0.12 | 0.5 | 0.4 | −0.1 | −0.72 | 0.2 |
| $h_6$ | 0.03 | 0.4 | 0.24 | −0.1 | −0.32 | 0.25 |

**Table 4** Priority table for hesitation values

| $U$ | Beautiful | Wooden | Green surrounded | Expensive | Distance | Row sum |
|-----|-----------|--------|------------------|-----------|----------|---------|
| $h_1$ | 0.72 | 0.1 | 0.4 | −0.55 | −0.68 | −0.01 |
| $h_2$ | 0.37 | 0.6 | 0.68 | −0.55 | −0.46 | 0.64 |
| $h_3$ | 0.64 | 0.2 | 0.4 | −0.65 | −0.44 | 0.15 |
| $h_4$ | 0.48 | 0.8 | 0.56 | −0.6 | −0.64 | 0.6 |
| $h_5$ | 0.67 | 0.5 | 0.52 | −0.65 | −0.26 | 0.78 |
| $h_6$ | 0.34 | 0.6 | 0.64 | −0.6 | −0.58 | 0.4 |

**Table 5** Comparison table for membership values

| $h_i$ | $h_j$ | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|---------|
| | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_6$ | Row sum |
| $h_1$ | 0 | −0.43 | −0.12 | −0.29 | −0.12 | −0.45 | −1.41 |
| $h_2$ | 0.43 | 0 | 0.31 | 0.14 | 0.31 | −0.02 | 1.17 |
| $h_3$ | 0.12 | −0.31 | 0 | −0.17 | 0.00 | −0.33 | −0.69 |
| $h_4$ | 0.29 | −0.14 | 0.17 | 0 | 0.17 | −0.16 | 0.33 |
| $h_5$ | 0.12 | −0.31 | 0.00 | −0.17 | 0 | −0.33 | −0.69 |
| $h_6$ | 0.45 | 0.02 | 0.33 | 0.16 | 0.33 | 0 | 1.29 |

**Table 6** Comparison table for non-membership values

| $h_i$ | $h_j$ | | | | | | |
|---|---|---|---|---|---|---|---|
| | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_6$ | Row sum |
| $h_1$ | 0 | 1.08 | 0.28 | 0.9 | 0.91 | 0.86 | 4.03 |
| $h_2$ | −1.08 | 0 | −0.8 | −0.18 | −0.17 | −0.22 | −2.45 |
| $h_3$ | −0.28 | 0.8 | 0 | 0.62 | 0.63 | 0.58 | 2.35 |
| $h_4$ | −0.9 | 0.18 | −0.62 | 0 | 0.01 | −0.04 | −1.37 |
| $h_5$ | −0.91 | 0.17 | −0.63 | −0.01 | 0 | −0.05 | −1.43 |
| $h_6$ | −0.86 | 0.22 | −0.58 | 0.04 | 0.05 | 0 | −1.13 |

**Table 7** Priority table for hesitation values

| $h_i$ | $h_j$ | | | | | | |
|---|---|---|---|---|---|---|---|
| | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_6$ | Row sum |
| $h_1$ | 0 | −0.65 | −0.16 | −0.61 | −0.79 | −0.41 | −2.62 |
| $h_2$ | 0.65 | 0 | 0.49 | 0.04 | −0.14 | 0.24 | 1.28 |
| $h_3$ | 0.16 | −0.49 | 0 | −0.45 | −0.63 | −0.25 | −1.66 |
| $h_4$ | 0.61 | −0.04 | 0.45 | 0 | −0.18 | 0.2 | 1.04 |
| $h_5$ | 0.79 | 0.14 | 0.63 | 0.18 | 0 | 0.38 | 2.12 |
| $h_6$ | 0.41 | −0.24 | 0.25 | −0.2 | −0.38 | 0 | −0.16 |

**Table 8** Decision table according to our algorithm

| Houses | Membership score | Non-membership score | Hesitation score | Score | Rank |
|---|---|---|---|---|---|
| $h_1$ | −1.41 | 4.03 | −2.62 | −2.2842 | 6 |
| $h_2$ | 1.17 | −2.45 | 1.28 | 2.6676 | 1 |
| $h_3$ | −0.69 | 2.35 | −1.66 | −0.4554 | 4 |
| $h_4$ | 0.33 | −1.37 | 1.04 | 0.6732 | 3 |
| $h_5$ | −0.69 | −1.43 | 2.12 | −2.1528 | 5 |
| $h_6$ | 1.29 | −1.13 | −0.16 | 1.0836 | 2 |

**Decision Making:** From the Table 8, we can conclude that the customer's best option is "$h_2$." The order of selection is $h_2$, $h_6$, $h_4$, $h_3$, $h_5$, and $h_1$.

The following Table 9 shows the results obtained using the algorithm in [6]. In both cases, we can see the first 3 ranks are same. But the ranks 4 and 5 are different in both cases. If we are comparing the values given in Table 1, we can see that $h_3$ is a better choice than $h_5$.

**Table 9** Decision table using the algorithm in [16]

| Houses | Membership score | Non-membership score | Hesitation score | Score | Final rank |
|---|---|---|---|---|---|
| $h_1$ | −1.41 | 0.87 | 0.54 | −0.37 | 6 |
| $h_2$ | 1.17 | −1.05 | −0.12 | 1.55 | 1 |
| $h_3$ | −0.69 | 0.75 | −0.06 | −0.25 | 5 |
| $h_4$ | 0.33 | −0.33 | 0.00 | 0.83 | 3 |
| $h_5$ | −0.69 | 0.33 | 0.36 | 0.17 | 4 |
| $h_6$ | 1.29 | −0.57 | −0.72 | 1.07 | 2 |

## *4.1 Comparison of Algorithms*

1. The formula in [16] is giving results if the value of membership score (m) is less than zero and the value of hesitation score is less than −1 which is corrected using the formula (3.1).
2. In [16], the priority for the non-membership values and hesitation values was same as the membership values. But in real-life situations, priority of non-membership value is normally one's complement of membership value. Here, we are using priority of non-membership value as one's complement of membership value, and it's giving better results than existing algorithms.
3. In [16], priority of hesitation values is same as priority of membership value. But in real-life situations when priority of membership is less than 1, the hesitation will increase.

## 5  Conclusion

In this chapter, we follow the membership function for IFSS which extends the notion of characteristic function for FSS introduced by Tripathy et al. in 2015. Also, we improved the existing algorithms with the help of score function. We compared the existing results with our improved algorithm, and we are getting better results.

## References

1. Atanassov, K.: Intuitionistic fuzzy sets. Fuzzy Set Syst. **20**, 87–96 (1986)
2. Maji, P.K., Biswas, R., Roy, A.R.: Fuzzy soft sets. J. Fuzzy Math. **9**(3), 589–602 (2001)
3. Maji, P.K., Biswas, R., Roy, A.R.: An application of soft sets in a decision making problem. Comput. Math. Appl. **44**, 1007–1083 (2002)
4. Maji, P.K., Biswas, R., Roy, A.R.: Soft set theory. Comput. Math. Appl. **45**, 555–562 (2003)

5. Mohanty, R.K., Sooraj, T.R., Tripathy, B.K.: An application of IVIFSS in medical diagnosis decision making. Int. J. Appl. Eng. Res. (IJAER) **10**(92), 85–93 (2016)
6. Mohanty, R.K., Sooraj, T.R., Tripathy, B.K.: IVIFS and decision-making. Adv. Intell. Syst. Comput. **468**, 319–330 (2017)
7. Molodtsov, D.: Soft set theory—first results. Comput. Math Appl. **37**, 19–31 (1999)
8. Narayanan, S.J., Bhatt, R.B., Paramasivam, I., Khalid, M., Tripathy, B.K.: Induction of fuzzy decision trees and its refinement using gradient projected-neuro-fuzzy decision tree. Int. J. Adv. Intell. Paradig. **6**, 346–369 (2014)
9. Sooraj, T.R., Mohanty, R.K., Tripathy, B.K.: Fuzzy soft set theory and its application in group decision making. Adv. Intell. Syst. Comput. **452**, 171–178 (2016)
10. Sooraj, T.R., Tripathy, B.K.: Interval valued hesitant fuzzy soft sets and its application in stock market analysis. Adv. Intell. Syst. Comput. **517**, 755–764 (2017)
11. Sooraj, T.R., Mohanty, R.K., Tripathy, B.K.: Hesitant fuzzy soft set theory and its application in decision making. Adv. Intell. Syst. Comput. **517**, 315–322 (2017)
12. Tripathy, B.K., Arun, K.R.: A new approach to soft sets, soft multisets and their properties. Int. J. Reason. Based Intell. Syst. **7**(3/4), 244–253 (2015)
13. Tripathy, B.K., Sooraj, T.R., Mohanty, R.K.: Advances decision making usisng hybrid soft set models. Int. J. Pharm. Technol. **8**(3), 17694–17721 (2016)
14. Tripathy, B.K., Sooraj, T.R., Mohanty, R.K., Parida, S.Ch.: Rough multisets and their properties. Int. J. Sci. Innov. Math. Res. **3**(2), 690–694 (2015)
15. Tripathy, B.K., Sooraj, T.R., Mohanty, R.K.: A new approach to fuzzy soft set theory and its application in decision making. Adv. Intell. Syst. Comput. **411**, 305–313 (2016)
16. Tripathy, B.K., Mohanty, R.K., Sooraj, T.R., Arun, K.R.: A new approach to intuitionistic fuzzy soft sets and its application in decision-making. Adv. Intell. Syst. Comput. **439**, 93–100 (2016)
17. Tripathy, B.K., Mohanty, R.K., Sooraj, T.R., Tripathy, A.: A modified representation of IFSS and its usage in GDM. Smart Innov. Syst. Technol. **50**, 365–375 (2016)
18. Tripathy, B.K., Mohanty, R.K., Sooraj, T.R.: On intuitionistic fuzzy soft sets and their application in decision-making. Lect. Notes Electr. Eng. **396**, 67–73 (2016)
19. Tripathy, B.K., Sooraj, T.R., Mohanty, R.K.: A new approach to interval-valued fuzzy soft sets and its application in decision-making. Adv. Intell. Syst. Comput. **509**, 3–10 (2017)
20. Tripathy,B.K., Sooraj, T.R, Mohanty, RK: A new approach to interval-valued fuzzy soft sets and its application in group decision making. In: Proceedings of International Conference on Computer systems, Data Communication and Security, CDCS-2015, Kochi, pp. 1–8
21. Tripathy, B.K., Mohanty, RK, Sooraj, T.R.: On intuitionistic fuzzy soft set and its application in group decision making. In: Proceedings of ICETETS-2016, pp. 1–5 (2016). doi:10.1109/ICETETS.2016.7603002
22. Tripathy, B.K., Mohanty, R.K., Sooraj, T.R.: Application of uncertainty models in bioinformatics. In: Dash, S., Subudhi, B. (eds.) Handbook of Research on Computational Intelligence Applications in Bioinformatics, Chapter-9, pp. 169–182. IGI Global, Hershey (2016)
23. Tripathy, B.K., Sooraj, T.R., Mohanty, R.K., Arun, K.R.: Parameter reduction in soft set models and application in decision making. In: Sangaiah, A.K., Gao, X.-Z., Abraham, A. (eds.) Handbook of Research on Fuzzy and Rough Set Theory in Organizational Decision Making, pp. 331–354. IGI Global, Hershey (2016). doi:10.4018/978-1-5225-1008-6
24. Tripathy, B.K.: Rough sets on intuitionistic fuzzy approximation spaces. In: IEEE Intelligent Systems, UK, pp. 776–779 (2006)
25. Tripathy, B.K.: Rough sets on fuzzy approximation spaces and intuitionistic fuzzy approximation spaces. Stud. Comput. Intell. **174**, 3–44 (2009)
26. Tripathy, B.K., Satapathy, M.K., Choudhury, P.K.: Intuitionistic fuzzy lattices and intuitionistic fuzzy boolean algebras. Int. J. Eng. Technol. **5**(3), 2352–2361 (2013)
27. Zadeh, L.A.: Fuzzy sets. Inf. Control **8**, 338–353 (1965)

# Design and Analysis of Compact Circular Half-Ring Monopole Antenna with DGS

**S. S. Mohan Reddy, P. Mallikarjuna Rao, B. Prudhvi Nadh, B. T. P. Madhav and K. Aruna Kumari**

**Abstract** A compact circular monopole shaped antenna is designed to operate at the dual band, and it is modified as half circle strip loaded circular monopole antenna with defected ground structure (DGS). The proposed antenna consisting of the combination of circular radiating element and half circle strip on the top side of the FR4 substrate material with permittivity 4.4. The bottom side a dumbbell-shaped DGS is incorporated to enhance the antenna bandwidth parameter. The proposed antenna is providing a bandwidth of 8 GHz and covering the UWB range. A peak realized a gain of more than 3 dB and directivity of 3.6 dB is achieved from the current design. The proposed antenna is providing excellent impedance and radiation characteristics in the operating band.

**Keywords** Circular monopole · Defected ground structure (DGS)
Half ring · Permittivity · Wideband

S. S. Mohan Reddy (✉) · B. Prudhvi Nadh
Department of ECE, SRKR College of Engineering, Bhimavaram, Andhra Pradesh, India
e-mail: rahulmohan720@gmail.com

B. Prudhvi Nadh
e-mail: prudhvi273835@gmail.com

P. Mallikarjuna Rao
Department of Electronics and Communication Engineering, Andhra University,
Visakhapatnam, Andhra Pradesh, India
e-mail: pmraoauece@yahoo.com

B. T. P. Madhav
Department of ECE, K L University, Guntur, Andhra Pradesh, India
e-mail: btpmadhav@kluniversity.in

K. Aruna Kumari
Department of CSE, SRKR College of Engineering, Bhimavaram, Andhra Pradesh, India
e-mail: arunasatti720@gmail.com

# 1 Introduction

At present UWB antennas generally, play an important role in the modern communication field. The monopole antenna has received increasing interest in UWB applications because they exhibit very attractive features like simple structure, wide impedance bandwidth [1] and Omni directional patterns, which cover the frequency defined by FCC (Federal Communication Commission) from 3.1 to 10.6 GHz for ultra wide band (UWB) applications [2, 3]. UWB supports a large amount of data and very narrow pulses in nanoseconds [4, 5]. It has a great capacity for transferring high-speed short-range data transferring in wireless environments. The microstrip antennas are portable devices having less weight, reduced size and low cost [6–8]. The microstrip antennas are embedded in hand held devices such as cellular phones, pagers and low profile antennas such as aircraft, satellites and missile applications.

The major disadvantage in the microstrip antennas is low gain and narrow bandwidth. Defected ground structure (DGS) is used to overcome the problem. In DGS technique ground plane is etched with a defect which causes the disturbance in the current distribution and provides the controlled propagation of electromagnetic waves [9]. The surface waves are present if the substrate permittivity is greater than 1. When the microstrip patch antennas are used, there are high losses which are because of surface waves excitation in the substrate layer. Generally, losses will occur, when the antenna is transmitting the signals. Due to surface waves antenna gain, efficiency and bandwidth are decreased. The surface waves propagate and reach the edges of the ground and reflected back with diffraction on the ground plane. When the microstrip patch antenna is fabricated on the high-dielectric substrate, it causes back radiation pattern increases the surface waves from ground plane edges.

In the DGS, various geometrical shapes are present; they are circles, rectangles, spirals, dumbbells shaped, H-shaped, L-shaped defects are used in the ground plane [10]. The microstrip patch antenna with Defected ground will provide little bit higher bandwidth and nominally less return loss. Whereas without DGS, it will provide narrow bandwidth and return loss will be high. For improving the radiation DGS is integrated on the ground plane.

Previously the circular ring and y-shape-strip with defected ground plane for WiMAX and WLAN application. The paper achieved three resonant frequencies at 2.61, 3.5, 5.4 GHz [11, 12]. The proposed antenna is having the UWB nature. The simulated results are presented with key structure parameters are analyzed (Fig. 1).

# 2 Antenna Parameters

Table 1 explains the antenna parameters of UWB proposed antenna. It consists of the radiating patch at the upper part of the substrate, which consists of a circular shaped ring and semicircle strip. The substrate has a length and width of

Fig. 1 Geometry of proposed antenna **a** front view, **b** back view, **c** inner dumbbell dimensions

**Table 1** Parameters of the designed antenna

| Parameter | Size (mm) | Parameter | Size |
|-----------|-----------|-----------|------|
| $L_s$ | 38 | $L_g$ | 11.5 mm |
| $W_s$ | 25 | $D_b$ | 0.2 mm * 1 mm |
| $L_f$ | 12.06 | $D_s$ | 0.7 mm * 0.5 mm |
| $W_f$ | 2.4 | $D_g$ | 1.3 mm |
| $R_1$ | 12 | H | 1.59 mm |
| $R_2$ | 9 | S | 1.5 mm |
| $R_3$ | 8 mm | D | 2.51 mm |
| $R_4$ | 5.5 mm | | |

$L_s$ = 38 mm, $W_s$ = 25 mm. The circular ring patch consists of outer radius $R_1$ and inner radius $R_2$, respectively. The semicircle consists of the outer radius of $R_3$ and inner radius of $R_4$. The circular ring patch antenna is fed by a 50 Ω microstrip feed line which is having a width of $W_f$ = 2.4 mm and length $L_f$ = 12.06 mm. The partial ground plane is placed on the back ground of the substrate with dumbbell-shaped defect. The gap between the circular ring and the semicircle is separated by distance 's', which affect the impedance performance.

The antenna prototyped on FR4 substrate material ($\varepsilon_r$ = 4.4) with thickness $h$ = 1.59 and dielectric loss tangent 0.02. The bottom side the DGS method is applied with half ground plane has $L_g$ = 11.5 mm. The dumbbell shape is etched on the ground plane and the parameters for dumbbell shapes are dumbbell size ($D_s$ = 0.7 mm * 0.5 mm), dumbbell gap ($D_g$ = 1.3 mm) and dumbbell bridge ($D_b$ = 0.2 mm * 1 mm).

The design model of the ultra wideband antenna and corresponding results are presented in the Fig. 2. The Antenna1 design begins with the circular ring and a rectangular ground is taken. The results obtained are Antenna1 operates at two

**Fig. 2** Design evolution method of antenna

different bands; one is 1.93–2.03 GHz which is centred at 2 GHz and the second band at 3.42–3.54 GHz. The first frequency bandwidth is 0.1 GHz, which is narrow and second frequency bandwidth has 0.12 GHz. The proposed Antenna2 consists of semicircle and dumbbell-shaped defected ground shape. It is observed that the impedance bandwidths of Antenna2 have the wide bandwidth at a frequency range of 3.4–9.7 GHz.

## 3 Results and Discussion

The design model of wide band antenna the conventional antenna is compared Fig. 3 shows the compression results of return loss of the Antenna1 and Antenna2. The return losses of the Antenna2 show improvement than Antenna1. The measured bandwidth of Antenna2 is 6.3 GHz in the band of 3.4–9.7 GHz, which meet the required bandwidth for WLAN/WiMAX applications.

The further study of wide band antenna the simulated E-filed current distributions at 2.5, 3.5, 5.5 and 5.8 GHz are shown in the Fig. 4.

Figure 5 demonstrates the gain of the antenna. It is observed that the antenna gain at 5.5 GHz, it shows the maximum gain. There is a gain reduction at 2.5 and 3.5 GHz; it shows 1.8 and 1.6 dBi. At 5.5 GHz, the Antenna shows the gain of 3.18 and 2.94 dBi at 5.8 GHz. The red colour in the diagram shows the maximum field intensity in particular direction.

The radiation patterns of the antenna describe how the energy is radiated. If the antenna is 100% efficient, it radiates the total energy for given same input power irrespective of pattern shape. Figure 6 demonstrates the simulated radiation patterns of Antenna 2. The E-plane, as well as H-plane radiation pattern at 2.5, 3.5, 5.5 and

**Fig. 3** Simulated return loss with and without DGS



**Fig. 4** Simulated surface current distribution E-field at **a** 2.5 GHz, **b** 3.5 GHz, **c** 5.5 GHz, **d** 5.8 GHz



**Fig. 5** Simulated 3D polar plots at different frequency **a** 2.5 GHz, **b** 3.5 GHz, **c** 5.5 GHz, **d** 5.8 GHz

**Fig. 6** Radiation patterns E-Plane and H-plane **a** 2.5 GHz, **b** 3.5 GHz, **c** 5.5 GHz, **d** 5.8 GHz



**Fig. 7** Surface current distributions at **a** 2.5 GHz, **b** 3.5 GHz, **c** 5.5 GHz, **d** 5.8 GHz

5.8 GHz, are given. The red line indicates $\phi = 0°$ and dotted blue line indicates $\phi = 90°$.

The simulated model surface current distributions at different frequencies are analyzed. The strong current distribution is shown in Fig. 7 flows in the circular ring and semicircle strip. The current flowing in the circular patch should travel the

longer distance. To get the correct radiation pattern, there should be a control on surface waves. The colours in the current distribution indicate the current magnitudes.

## 4　Prototyped Model and Measured Results

The prototype of the antenna has been fabricated, is shown Fig. 8. Through reflection coefficient curve, the bandwidth is measured by vector network analyzer. Good agreement between the measured and simulated results is observed. It is observed from the measured results that the designed antenna with the DGS shows the bandwidth of 3.6–9.8 GHz, while the simulated results show the bandwidth ranging from 3.6 to 9.7 GHz covering the entire UWB frequency band. The VSWR of the antenna is less than 2 in that entire operating band (Fig. 9).

## 5　Parametric Analysis

In the conventional antenna, only the full rectangular ground is taken without dumbbell shape. By adding the dumbbell size 0.7 * 0.5 as a DGS is implemented. For doing the parametric analysis, the size of the each dumbbell is varied from 0.3 to 0.9 * 0.5. At first, the dumbbell d1 is varied, respectively, the d2, d3, d4, d5 and d6. The Figs. 10, 11, 12, 13, 14, and 15 give the frequency verse return loss curves



**Fig. 8　a** Front view, **b** back view

**Fig. 9** Measured return loss from ZNB 20



**Fig. 10** Return losses for dumbbell size d1 is varied



**Fig. 11** Return losses for dumbbell size d2 is varied

**Fig. 12** Return losses for dumbbell size d3 is varied



**Fig. 13** Return losses for dumbbell size d4 is varied



**Fig. 14** Return losses for dumbbell size d5 is varied

**Fig. 15** Return losses for dumbbell size d6 are varied

for different dumbbell sizes. The optimized dimensional parameters are analyzed and accordingly the prototyped antenna is fabricated.

## 6 Conclusion

A simple circular monopole antenna for ultra wideband applications is proposed in this work. The simulation characteristics of the proposed antenna satisfy the FCC standard UWB band with considerable gain and directivity. An impedance bandwidth of 86% in the operating band with VSWR < 2 elevates this model for short range UWB applications. The Omni directional radiation pattern in the operating band and proper impedance matching makes this antenna suitable for the desired band of applications in the communication systems.

## References

1. Lee, C.P., Chakrabarty, C.K.: Ultra wideband microstrip diamond slotted patch antenna with enhanced bandwidth. Int. J. Commun. Netw. Syst. Sci. **4**(7), 468 (2011)
2. FCC First report and order on ultra wide band technology, Washington, DC (2002)
3. Sundar, P.S., Kotamraju, S.K., Ramakrishna, T.V., Madhav, B.T.P.: Novel miniatured wide band annular slot monopole antenna. Far East J. Electron. Commun. **14**(2), 149–159. ISSN: 0973-7006 (2015)
4. Ritu, K.S.: Microstrip antenna design for UWB applications. Int. J. Adv. Res. Comput. Commun. Eng. **2**(10), 3824–3828 (2013)
5. Singh, Y., Singh, D., Gill, G. S.: Design of wideband microstrip antenna for UWB applications. Int. Res. J. Eng. Technol. **2**(5), 995–998 (2015)
6. Madhav, B.T.P., Kaza, H.: Novel printed monopole trapezoidal notch antenna with S-band rejection. J. Theor. Appl. Inf. Technol. **76**(1), 42–49. ISSN: 1992-8645 (2015)

7. Singh, P., Tomar, R.: The use of defected ground structures in designing microstrip filters with enhanced performance characteristics. Proc. Technol. **17**, 58–64 (2014)
8. Mohan Reddy, S.S., Mallikarjunarao, P., Madhav, B.T.P.: Asymmetric defected ground structured monopole antenna for wideband communication systems. Int. J. Commun. Antenna Propag. **5**(5), 256–262. ISSN: 2039-5086 (2015)
9. Weng, L.H., Guo, Y.C., Shi, X.W., Chen, X.Q.: An overview on defected ground structure. Progr. Eletromagn. Res. **B7**, 173–189 (2008)
10. Pei, J., et al.: Miniaturized triple-band antenna with a defected ground plane for WLAN/WiMAX applications. IEEE Antennas Wirel. Propag. Lett. **10**, 298–301 (2011)
11. Ahuja, N., Khanna, R., Kaur, J.: Dual band defected ground microstrip patch antenna for WLAN/WiMax and satellite application. Int. J. Comput. Appl. **48**(22), 1–5 (2012)
12. Bhavani, K.V.L., Khan, H., Madhav, B.T.P.: Multiband slotted aperture antenna with defected ground structure for C and X-band communication applications. J. Theor. Appl. Inf. Technol. **82**(3), 454–461. ISSN: 1992-8645 (2015)

# Developing Higher Education Ontology Using Protégé Tool: Reasoning

**Ch. V. S. Satyamurty, J. V. R. Murthy and M. Raghava**

**Abstract** The high volume of data presents on the Web in semi-structured forms such as XML and HTML render the task of retrieval of information difficult for either for the search engines or information retrieval systems. Moreover, majority of the search engines are depending mostly on keywords that make automated Web search operation, about 70% of the times, to end up with retrieval of un-useful and irrelevant documents. To resolve the above issue, Tim Berners-Lee proposed intelligent and machine readable Web called as the Semantic Web in which metadata is associated with every entity on the Web. Resource Description Format (RDF) is offering a set of concepts to describe the metadata of the elements. RDF is a universal semantic model which lets the designer to describe the resources in his vocabulary without any regard to the domain of the problem. Hence, ontologies are introduced to formalize the association of metadata and semantics to each entity which becomes very much useful in Semantic Web for information processing, knowledge sharing across application systems and makes the data is machine readable. Web Ontology Language (OWL) is a markup language to implement Semantic Web and servers to share the ontologies and enables the agent-based communication while correctly interpreting the meaning of the Web document. However, there are fewer ontology development tools available for researchers that are generic to any domain. In our paper, we develop an ontology in the domain of higher education (engineering) domain using Protégé 4.3 tool. The associations among the resources are presented in the form of onto graphs/OWLViz, and the correctness of the class hierarchy is assessed by Fact++. Finally, the execution of semantic query for retrieval of the knowledge using SPARQL.

**Keywords** Ontology · Protégé editor · OWL · Semantic Web · SPARQL

Ch. V. S. Satyamurty (✉) · M. Raghava
CVR College of Engineering, Hyderabad, India
e-mail: satyamurtycvs@yahoo.co.in

J. V. R. Murthy
JNTUK, Kakinada, India

233

# 1   Introduction

The data present in the Web will be better understood by knowing the semantics of the data, and it is well recognized by business and research communities. This led W3C to extend the extended markup language (XML) to Web Ontology Language (OWL) by developing ontology inference layer (OIL) under the umbrella of Semantic Web.

## 1.1   Semantic Web

The Semantic Web [1, 2] is an intelligent Web. It is an advanced knowledge management system which arranges the knowledge in class hierarchy and uses description logics to create intelligent and agent-based application. It enables the Web content more machine readable, provides automatic integration and reuses of Web documents across multiple applications and enhances the capabilities of search engine. It also opens a port to artificial intelligence process on Web to gain a semantics-based automatic vision process which could link data, better process it and display the relations by the machines.

## 1.2   Ontology

Gruber [3] describes the term ontology as an explicit specification of conceptualization. In other words, ontologies describe the structure of a domain and the associated implicit knowledge of semantic conceptualizations. It formally describes an area of knowledge in terms of classes, subclasses, and their relationships. It is also helpful to compare and contrast information across two knowledge bases on the Web. Formally, ontology is a program which finds ways to associate with a common meaning across the knowledge bases and to extract the inferences at a common place without any concern to the existence of physical objects and their structural features. Overall, ontologies are helpful for knowledge sharing and reuse. In this context, an explicit specification of knowledge as classes is necessary to enable the agents with share and reuse capabilities. The construction of domain ontology by OWL is helpful and extends the power of Resource Description Framework (RDF) language in terms of reasoning, expressivity of Web data. It was also useful in automatic reasoning of tasks on Web data.

## 1.3  Protégé Ontology Editor

Many ontology editors are available such as Onto-Studio, TopBraid Composer free edition, SWOOP, Protégé, OntoEdit out of which Protégé is mostly used. Cardso [4] conducted a survey on different open source ontology editors to develop the domain specific ontology. He found that Stanford University's Protégé is largely used by many people. Figure 1a depicts the fact that nearly 70% of ontology editor share is grabbed by Protégé [5], followed by SWOOP with 14% and OntoEdit, TextEditor, Altova SW sharing the remaining. Figure 1b clearly shows that the creation of ontology is extensive in education domain [6–10] followed by computer science domain. In this paper, we used Protégé tool for creation of ontology in education domain.

## 2  Ontology Development in Education Domain

The survey on education domain reveals a fact that the information available in various data interpretations and their integration through SQL-based databases possess a structural dissimilarity. In other words, the basic disparity in the terminology utilized referring to the same subject by different educational institutions led



**Fig. 1  a**, **b** illustrating ontology development

to lot of confusion and uncertainty in extracting the information from the databases. Moreover, the database designer ponders a field appearing in many tables serving a context specific purpose.

The first step in Protégé based modeling is a creation of a new project. Then, in the new project, we can develop the ontology, using various tabs, e.g., entities, classes, object properties, and data properties of the instances. Figure 2 shows different entities in the education ontology with selected things/classes such as a module: referring to a set of related subjects, person; referring to teachers, students, and hostlers and sports comprising indoor and outdoor sports.

Figure 3 shows the different classes and the associated subclasses. Therein, the module class has basic science (BS), professional courses (PC) and Mathmodule (MM) as submodules.

Similarly, the person and sport classes are arranged into the corresponding subclasses. Subsequently, the identified semantic relationships among the subclasses are modeled as object properties. The object property maps the one class treated as the domain onto another class referred as range.



**Fig. 2** Entities in the education ontology



**Fig. 3** Various classes and sub classes

This aspect is demonstrated in Fig. 4, therein favorite sport, relates the person class with sports class, stays in associates with the assigned room number in the hostel with student class and studies link the student class with module class.

The instances of the class are created through the individual tab upon finalizing the data properties of each class by specifying an appropriate data type. For instance, Fig. 5 shows the data properties such as designation, facultyid, and age pertaining to lecturer subclass of class person.

Finally, the instances for the classes are created. Figure 6 shows the example instances of professional course are computer organization, java programming, and the instances of lecturer subclass are faculty1, faculty2, whereas student subclass has instances student1, student2. Here, the essential step is to define the assertions at object property level and data property level that constitutes the RDF. For instance, object property assertion, "teaches" corresponds the faculty2 with java programming, and the assertion "hasfovouritesport" maps faculty2 with shuttle. Coming to data property, assertions "last_name" and "first_name" are RAGHAVA and MORUSUPALLI, respectively declared as string variables. Similarly, the age, designation and faculty_id are specified appropriately.



**Fig. 4** Object properties

**Fig. 5** Data properties



**Fig. 6** Individuals or instances

**Fig. 7** Individuals or instances with usage tab

Figure 7 depicts the usage of the individual for faculty2, with the last_name: "RAGHAVA," has a type: lecturer, teaches: "javaprogramming," hasfavouritesport: "Shuttle," first_name: "MORUSUPALLI," age "40," designation "Associate Professor," teaches: "Computerorganization," facultyid: "CVRCSEF003," Individual: "faculty2." Thus, the entire set of usages can be created for the higher education system.

Figure 8 shows the ontology-based reasoning through tapping the OWLViz tab. It describes the "is a" relationship between classes and subclasses pictorially. For example, Hostler "is a" person. For the sake of evaluation of the ontology and its



**Fig. 8** Using OWLViz for reasoning

**Fig. 9** OntoGraph of classes drawn from Fig. 2

consistency, Protégé provides FACT++ reasoner which can highlight the class/subclass inconsistent relationship, if they prevail, through red mark. Figure 8 validates the proposed model to be correct as it is free form red marks.

Figure 9 depicts the output of OntoGraft tab, the fully grown classes without any symbol prefixed and expandable classes with + as the prefix to the label. The diamond symbol on a tab corresponds to the class instance. The presentation style of the hierarchy can be selected from the menu available in the Protégé visualization tab.

The graph can be radial, grid, tree vertical, and so on, and finally, we can save the graph, through different options such as save, export as an image, and ".dot" file.

## 3   SPARQL Query

For querying the ontology, we can use SPARQL. The output of the query can be shown as triplets, i.e., subject, object, and predicate. Fig. 10 shows output of the basic query in the subject and object form with Mathmodule as the subject, and the module as object which can be interpreted as "Mathmodule is a Module."

**Fig. 10** Output of SPARQL query

## 4 Conclusions

This paper presented the aspects of ontology from initiation to creation using Protégé 4.3 and generated OWL file. We have created the ontology for higher education (engineering) domain, based on various courses in an autonomous institution. We have created classes and subclasses and restrictions. Fact++ reasoner is utilized to check for consistency of the created ontology. We applied SPARQL query on the ontology to meaningfully retrieve the information.

## References

1. Berners Lee, T., Handler, J., Lassila, O.: The semantic web. Sci. Am. **284**, 28–37 (2001)
2. Passin, T.B.: Explorer's Guide to the Semantic Web. Manning, New York (2004)
3. Gruber, T.R.: A translation approach to portable ontology specifications. Tech. Rep. Knowl. Acquis. **5**(2), 199–220 (1993)
4. Cardoso, J.: The semantic web vision: where are we? IEEE Intell. Syst. **22**, 22–26 (2007)
5. Stanford University. http://Protege.stanford.edu
6. University of Maryland. http://www.cs.umd.edu/projects/plus/SHOE/onts/univ1.0.html (2000)
7. Malik, S.K.: Developing an university ontology in education domain using Protégé for semantic web. Int. J. Sci. Technol. **2**(9), 4673–4681 (2010)
8. Zeng, L.: Study on construction of university course ontology: content, method and process. In: International Conference on Computational Intelligence and Software Engineering. IEEE Press (2009)
9. Ameen, A., Rehman, K.U.R., Padmajarani, B.: Construction of university ontology. In: 2012 World Congress on Information and Communication Technologies. IEEE Xplore (2012)
10. Romeo, L.: Educational metadata ontology enrichment for E-assessment semantic description. In: 10th Iberian Conference on Information Systems and Technologies (CISTI). IEEE Press (2015)

# A New Approach to Interval-Valued Intuitionistic Hesitant Fuzzy Soft Sets and Their Application in Decision Making

**T. R. Sooraj, R. K. Mohanty and B. K. Tripathy**

**Abstract**  There are several models of uncertainty found in the literature like fuzzy set (FS), rough set, intuitionistic fuzzy set, soft set, and hesitant fuzzy set. Also, several hybrid models have come up as a combination of these models and have been found to be more useful than the individual models. In everyday life, we make many decisions. Making efficient decisions under uncertainty needs better techniques. Many such techniques have been developed in the recent past. These techniques involve soft sets (SS) and intuitionistic fuzzy sets. It is well known that intuitionistic hesitant fuzzy sets are more general than intuitionistic fuzzy sets. In this paper, we redefine interval-valued intuitionistic hesitant fuzzy soft sets (IVIHFSS) and also propose a decision-making technique which extends some of the recently developed algorithms. We also provide an application from real-life situations which illustrates the working of the algorithm and its efficiency over the other algorithms.

## 1  Introduction

Fuzzy set (FS) theory [1], rough set theory, probability theory, and interval mathematics are the theories that are used to deal with uncertainty and vagueness. The concept of FS was introduced by Zadeh in 1965 [2]. Even though it is one of the best models to handle uncertainty-based issues, it is having some drawbacks.

T. R. Sooraj (✉) · R. K. Mohanty · B. K. Tripathy
SCOPE, VIT University, Vellore, Tamilnadu, India
e-mail: soorajtr19@gmail.com

R. K. Mohanty
e-mail: rknmohanty@gmail.com

B. K. Tripathy
e-mail: tripathybk@vit.ac.in

Some of the drawbacks of fuzzy sets, probability theory, etc., are discussed by Molodtsov in [3]. The drawbacks are due to the lack of parameterization tool. So, Molodtsov introduced soft set theory in 1999, which is a parameterized family of subsets. Later, researchers extended the concept of soft sets (SS) to hybrid models like fuzzy soft sets (FSS) [4, 5], intuitionistic FSS, etc. [2–18]. In 2015, Tripathy et al. [9] introduced characteristic function for SS, and it generalizes most of the operations on SS. Similarly, they introduced membership function for FSS, IFSS, interval-valued FSS (IVFSS), etc. [1, 2, 14–17]. Jiyang et al. introduced the concept IVFSS. Intuitionistic fuzzy sets (IFS) [19], are extensions of the fuzzy set models and the membership, non-membership as well as hesitation values for each element is taken into consideration. Later Jiang et al. [20] extended intuitionistic fuzzy set model to define interval-valued IFS, where the membership, non-membership and hesitation values are sub-intervals of [0, 1]. The concept of hesitant fuzzy sets was introduced by Torra [8], where membership values of an element are a set of values. Later, many hybrid models of hesitant fuzzy sets like intuitionistic hesitant fuzzy sets, interval-valued intuitionistic hesitant fuzzy sets [21] are introduced. Peng et al. combined soft sets with interval-valued hesitant fuzzy sets to form interval-valued hesitant fuzzy soft sets [7], and they applied it in decision-making applications. In this paper, we introduce IVIHFSS and its operations.

## 2   Definitions and Notions

Let $U$ & $F$ be the set of universal and parameters, respectively.

**Definition 2.1** A soft set $(Z, F)$ is defined as follows:

$$Z : F \rightarrow \mathrm{PS}(U), \tag{2.1}$$

where $\mathrm{PS}(U)$ is the power set of U.

**Definition 2.2** The membership function (MF) of IVFSS $(Z, F)$ is as follows:

$$\mu_{(Z,F)}^a(x) = \left\{ \mu_{(Z,F)}^{a-}(x), \mu_{(Z,F)}^{a+}(x) \right\}, \tag{2.2}$$

where $\mu_{(Z,F)}^a(x) \in D[0, 1]$. $D[0, 1]$ is the set of all closed subintervals on [0,1].

**Definition 2.3** A hesitant FSS (HFSS) $(Z, F)$ is defined as follows:

$$Z : F \rightarrow \mathrm{HF}(U), \tag{2.3}$$

where $\mathrm{HF}(U)$ is the set of all hesitant fuzzy subsets of $U$.

**Definition 2.4** An intuitionistic HFSS, $(Z, F)$ is defined as follows:

$$Z : F \rightarrow \mathrm{IHFS}(U). \tag{2.4}$$

## 3 Interval-Valued Intuitionistic Hesitant Fuzzy Soft Sets (IVIHFSS)

In this section, we introduce the notion of IVIHFSS and establish some of its operations.

**Definition 3.1** An IVIHFSS $(Z, F)$ is defined as follows:

$$Z : F \rightarrow \mathrm{IVIHF}(U), \tag{3.1}$$

such that $\forall a \in F, \left\{ \tau_{(Z,F)}^a \right\}$ is a family of characteristic functions, where $\tau_{(Z,F)}^a(x) = \left( \mu_{(Z,F)}^a, v_{(Z,F)}^a \right) \in P(\mathrm{IVIHF})$, and $\mathrm{IVIHF}(U)$ is the set of all interval-valued intuitionistic hesitant fuzzy subsets of $U$. And $\forall a \in Z$ and $\forall x \in U$,

$$0 \leq \ \sup \mu_{(Z,F)}^a(x) + v_{(Z,F)}^a(x) \leq 1 \tag{3.2}$$

**Definition 3.2** The complement of an IVIHFSS $(Z, F)$ is $(Z, F)^c$, which is defined as follows:

$$\mu_{(Z,F)^c}^a(x) = v_{(Z,F)}^a(x), v_{(Z,F)^c}^a(x) = \mu_{(Z,F)}^a(x). \tag{3.3}$$

**Definition 3.3** The condition of null IVIHFSS is given below $\forall a \in Z$ and $\forall x \in U$, $\mu_{(Z,F)}^a(x) = [0,0], v_{(Z,F)}^a(x) = [1,1]$, and for absolute IVIHFSS, $\mu_{(Z,F)}^a(x) = [1,1], v_{(Z,F)}^a(x) = [0,0]$.

## 4 Application of IVIHFSS

Here, we discuss an application of DM in IVFSSs. Tripathy et al. [4] classified the parameters into positive and negative. We use Formula 4.1 to get a fuzzy value as score from an intuitionistic fuzzy value. It reduces the complexity and makes the comparison easier.

$$\begin{aligned} \mathrm{Score} &= \mu(1+h) \\ &= \mu + \mu h \end{aligned} \tag{4.1}$$

$\mu, \nu$, and $h$ are mutually dependent measures.

The normalized score is used to find the final rank in the decision making.

$$\text{Normalized Score} = \frac{2\left((|J| * |K| * |C|) - \sum_{i=1}^{|J|} \sum_{x \in K} R_{C_{ix}}\right)}{|J| * |K| * |C| - (|C| - 1)} \qquad (4.2)$$

Here, $|J|$ = No of Decision Makers, $|K|$ = No of approaches, and $|C|$ = No of objects.

Parameter data table: A parameter data table contains all the information of the parameters such as priority and parameter rank. Priority of a parameter is the degree of interest with respect to decision maker. The value of the priority of parameter lies in [0,1] or [−1,0] for a positive and negative parameter, respectively. Rank of parameters is calculated by taking the absolute value of the priorities given to a parameter.

**Algorithm**

1. Input parameter data table & IVIHFSS.
2. Select optimistic, pessimistic, and neutral values from IVIHFSS.
3. Procedure Deci_make(IFSS table)

   3.1. Construct priority table [5].
   3.2. Construct comparison table [5], and compute the score using Eq. 4.1.
   3.3. Rank the candidates based on the score value.
   3.4. Return (decision table).

4. Similarly, construct the decision tables for optimistic, pessimistic, and neutral values.
5. Construct the final rank matrix by using Eq. 4.2.

*A Real-Life Application*:

Nowadays, basketball team managers are allowed to buy players from the transfer market. Before a manager bids a player, he has to make an analysis about the quality of players. This can be achieved by checking the attributes or parameters of the player. Some of the parameters that he needs to consider are discipline, loyalty, overall attributes, nationality, transfer fee, etc. Here, the parameter "transfer fee" is a negative parameter because the manager's interest may decrease as the price of the player increases. In this paper, a synthetic dataset is taken to show the working principle of our algorithm.

The parameter data table is given in Table 1. It contains all the details about the parameters.

| Table 1 Parameter data table | $U$ | $e_1$ | $e_2$ | $e_3$ | $e_4$ | $e_5$ | $e_6$ |
|---|---|---|---|---|---|---|---|
| | Priority | 0.3 | 0.2 | −0.1 | 0.5 | 0.35 | 0 |
| | Parameter rank | 3 | 4 | 5 | 1 | 2 | 6 |

Let $U$ & $E$ be the set of players and parameters, respectively. For brevity, we can take

$U = \{p_1, p_2, p_3, p_4, p_5, p_6\}$ & $E = \{e_1, e_2, e_3, e_4, e_5, e_6\}$. Here, Table 2 denotes the representation of IVIHFSS (i.e., selection of player). Table 3 is IVIFSS which is extracted from Table 2.

From this table, we can extract the optimistic values, pessimistic values, and neutral values. For example, for the player $P_1$, the optimistic value is 0.5, pessimistic value is 0.1, and neutral value is 0.3 with respect to the parameter $e_1$. The priority assigned by the manager for parameters is shown in Table 1. The priority Tables 3 and 4 can be obtained by multiplying the user given parameter priority values in Table 1 with the respective membership, non-membership, and hesitation values in the optimistic, pessimistic, and neutral value tables. The priority table for optimistic values is shown in Table 4.

Likewise, PTs for pessimistic and neutral values can be calculated. Optimistic CT is constructed and shown in Table 5, and similarly, pessimistic CT and neutral CT can be calculated.

The score is obtained by using Formula 4.1. On the basis of the score value obtained, we rank the candidates. The decision tables are shown in Tables 6, 7, and 8.

Final rank table can be obtained by taking the rank columns from each of the decision tables as shown in Table 9. The normalized score is obtained by using Eq. 4.2. Then, the player who has got higher normalized score is the ideal choice to the manager. Here, $P_4$ is the best player, and the ranking of players is $P_4, P_2, P_3, P_6, P_1$, and $P_5$.

## 5 Conclusions

In this paper, we redefined IVIHFSS in a new way using membership functions and also redefined some operations on it. Furthermore, we proposed a decision-making algorithm using IVIHFSS with the help of an improved normalized score function which gives better results than the score function used in existing articles. An application of this algorithm in the real-life situations is discussed which shows realistic results. This approach can be further extended to the neutrosophic soft set theory and other hybrid neutrosophic models.

**Table 2** IVIHFSS

| U | e₁ μ | e₁ ν | e₂ μ | e₂ ν | e₃ μ | e₃ ν | e₄ μ | e₄ ν | e₅ μ | e₅ ν | e₆ μ | e₆ ν |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $P_1$ | 0.1–0.5, 0.1–0.3, 0.2–0.5 | 0.4–0.5, 0.3–0.5, 0.3–0.4 | 0.2–0.3, 0.2–0.4, 0.1–0.3 | 0.4–0.5, 0.5–0.6, 0.3–0.6 | 0.2–0.3, 0.3–0.5 | 0.4–0.5, 0.4–0.4 | 0.8–0.8, 0.8–0.9 | 0.0–0.1, 0.1–0.1 | 0.4–0.5, 0.5–0.7, 0.4–0.7 | 0.1–0.3, 0.1–0.2, 0.2–0.3 | 0.6–0.9 | 0.0–0.1 |
| $P_2$ | 0.8–1.0 | 0.0–0.0 | 0.4–0.8, 0.5–0.7 | 0.1–0.2, 0.2–0.2 | 0.7–0.9, 0.6–0.9, 0.6–0.8 | 0.0–0.0, 0.0–0.1, 0.1–0.1 | 0.2–0.5, 0.2–0.4, 0.3–0.4 | 0.1–0.3, 0.2–0.4, 0.1–0.4 | 0.7–1.0 | 0.0–0.0 | 0.5–0.5, 0.5–0.6 | 0.2–0.3, 0.3–0.3 |
| $P_3$ | 0.1–0.5, 0.1–0.3, 0.2–0.5 | 0.4–0.5, 0.3–0.5, 0.3–0.4 | 0.5–0.8 | 0.0–0.2 | 0.7–0.9, 0.7–0.8 | 0.0–0.1, 0.1–0.1 | 0.7–0.8 | 0.0–0.2 | 0.8–1.0 | 0.0–0.0 | 0.5–0.7, 0.6–0.7 | 0.2–0.2, 0.2–0.3 |
| $P_4$ | 0.5–0.9, 0.6–0.8 | 0.0–0.1, 0.1–0.1 | 0.6–0.8, 0.6–0.7 | 0.1–0.1, 0.1–0.2 | 0.5–0.9 | 0.1–0.1 | 0.8–1.0 | 0.0–0.0 | 0.5–0.9, 0.6–0.8 | 0.0–0.0, 0.0–0.1 | 0.7–0.8 | 0.0–0.1 |
| $P_5$ | 0.1–0.2 | 0.4–0.7 | 0.1–0.4 | 0.5–0.6 | 0.9–1 | 0.0–0.0 | 0.3–0.6 | 0.3–0.4 | 0.1–0.5 | 0.2–0.4 | 0.8–1.0 | 0.0–0.0 |
| $P_6$ | 0.9–1.0 | 0.0–0.0 | 0.7–0.9, 0.7–0.8 | 0.0–0.0, 0.0–0.1 | 0.6–0.7, 0.5–0.7 | 0.1–0.3, 0.2–0.3 | 0.1–0.3 | 0.4–0.7 | 0.2–0.4, 0.3–0.4 | 0.1–0.2, 0.1–0.4 | 0.3–0.7 | 0.2–0.3 |

**Table 3** IVIFSS

| $U$ | $e_1$ | | $e_2$ | | $e_3$ | | $e_4$ | | $e_5$ | | $e_6$ | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\mu$ | $v$ | $\mu$ | $v$ | $\mu$ | $v$ | $\mu$ | $v$ | $\mu$ | $v$ | $\mu$ | $v$ |
| $P_1$ | 0.1–0.5 | 0.3–0.5 | 0.2–0.4 | 0.3–0.6 | 0.3–0.5 | 0.4–0.5 | 0.8–0.9 | 0.0–0.1 | 0.4–0.7 | 0.1–0.3 | 0.6–0.9 | 0.0–0.1 |
| $P_2$ | 0.8–1.0 | 0.0–0.0 | 0.4–0.8 | 0.1–0.2 | 0.6–0.9 | 0.0–0.1 | 0.2–0.5 | 0.1–0.4 | 0.7–1.0 | 0.0–0.0 | 0.5–0.6 | 0.2–0.3 |
| $P_3$ | 0.1–0.5 | 0.3–0.5 | 0.5–0.8 | 0.0–0.2 | 0.7–0.9 | 0.0–0.1 | 0.7–0.8 | 0.0–0.2 | 0.8–1.0 | 0.0–0.0 | 0.5–0.7 | 0.2–0.3 |
| $P_4$ | 0.5–0.9 | 0.0–0.1 | 0.6–0.8 | 0.1–0.2 | 0.5–0.9 | 0.1–0.1 | 0.8–1.0 | 0.0–0.0 | 0.5–0.9 | 0.0–0.1 | 0.7–0.8 | 0.0–0.1 |
| $P_5$ | 0.1–0.2 | 0.4–0.7 | 0.1–0.4 | 0.5–0.6 | 0.9–1 | 0.0–0.0 | 0.3–0.6 | 0.3–0.4 | 0.1–0.5 | 0.2–0.4 | 0.8–1.0 | 0.0–0.0 |
| $P_6$ | 0.9–1.0 | 0.0–0.0 | 0.7–0.9 | 0.0–0.1 | 0.5–0.7 | 0.1–0.3 | 0.1–0.3 | 0.4–0.7 | 0.2–0.4 | 0.1–0.4 | 0.3–0.7 | 0.2–0.3 |

**Table 4** Priority table for optimistic values

| U | $e_1$ | | | $e_2$ | | | $e_3$ | | | $e_4$ | | | $e_5$ | | | $\Sigma\mu$ | $\Sigma v$ | $\Sigma h$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\mu$ | $v$ | $h$ | $\mu$ | $v$ | $h$ | $\mu$ | $v$ | $h$ | $\mu$ | $v$ | $h$ | $\mu$ | $v$ | $h$ | | | |
| $P_1$ | 0.15 | 0.09 | 0.06 | 0.08 | 0.06 | 0.06 | −0.05 | −0.04 | −0.01 | 0.135 | 0 | 0.015 | 0.245 | 0.035 | 0.07 | 0.56 | 0.145 | 0.195 |
| $P_2$ | 0.3 | 0 | 0 | 0.16 | 0.02 | 0.02 | −0.09 | 0 | −0.01 | 0.075 | 0.015 | 0.06 | 0.35 | 0 | 0 | 0.795 | 0.035 | 0.07 |
| $P_3$ | 0.15 | 0.09 | 0.06 | 0.16 | 0 | 0.04 | −0.09 | 0 | −0.01 | 0.12 | 0 | 0.03 | 0.35 | 0 | 0 | 0.69 | 0.09 | 0.12 |
| $P_4$ | 0.27 | 0 | 0.03 | 0.16 | 0.02 | 0.02 | −0.09 | −0.01 | 0 | 0.15 | 0 | 0 | 0.315 | 0 | 0.035 | 0.805 | 0.01 | 0.085 |
| $P_5$ | 0.06 | 0.12 | 0.12 | 0.08 | 0.1 | 0.02 | −0.1 | 0 | 0 | 0.09 | 0.045 | 0.015 | 0.175 | 0.07 | 0.105 | 0.305 | 0.335 | 0.26 |
| $P_6$ | 0.3 | 0 | 0 | 0.18 | 0 | 0.02 | −0.07 | −0.01 | −0.02 | 0.045 | 0.06 | 0.045 | 0.14 | 0.035 | 0.175 | 0.595 | 0.085 | 0.22 |

**Table 5** Comparison table for optimistic values

| U | P₁ | | | P₂ | | | P₃ | | | P₄ | | | P₅ | | | P₆ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\mu$ | $v$ | $h$ | $\mu$ | $v$ | $h$ | $\mu$ | $v$ | $h$ | $\mu$ | $v$ | $h$ | $\mu$ | $v$ | $h$ | M | V | h |
| $P_1$ | 0 | 0 | 0 | −0.235 | 0.11 | 0.125 | −0.13 | 0.055 | 0.075 | −0.245 | 0.135 | 0.11 | 0.255 | −0.19 | −0.065 | −0.035 | 0.06 | −0.025 |
| $P_2$ | 0.235 | −0.11 | −0.125 | 0 | 0 | 0 | 0.105 | −0.055 | −0.05 | −0.01 | 0.025 | −0.015 | 0.49 | −0.3 | −0.19 | 02 | −0.05 | −0.15 |
| $P_3$ | 0.13 | −0.055 | −0.075 | −0.105 | 0.055 | 0.05 | 0 | 0 | 0 | −0.115 | 0.08 | 0.035 | 0.385 | −0245 | −0.14 | 0.095 | 0.005 | −0.1 |
| $P_4$ | 0.245 | −0.135 | −0.11 | 0.01 | −0.025 | 0.015 | 0.115 | −0.08 | −0.035 | 0 | 0 | 0 | 0.5 | −0.325 | −0.175 | 0.21 | −0.075 | −0.135 |
| $P_5$ | −0.255 | 0.19 | 0.065 | −0.49 | 0.3 | 0.19 | −0.385 | 0.245 | 0.14 | −0.5 | 0.325 | 0.175 | 0 | 0 | 0 | −0.29 | 0.25 | 0.04 |
| $P_6$ | 0.035 | −0.06 | 0.025 | −02 | 0.05 | 0.15 | −0.095 | −0.005 | 0.1 | −0.21 | 0.075 | 0.135 | 0.29 | −0.25 | −0.04 | 0 | 0 | 0 |

**Table 6** Decision table by optimistic values

|        | $\mu$  | $v$   | $h$    | Score    | Rank |
|--------|--------|-------|--------|----------|------|
| $P_1$  | −0.39  | 0.17  | 0.22   | −0.4758  | 5    |
| $P_2$  | 1.02   | −0.49 | −0.53  | 0.4794   | 2    |
| $P_3$  | 0.39   | −0.16 | −0.23  | 0.3003   | 3    |
| $P_4$  | 1.08   | −0.64 | −0.44  | 0.6048   | 1    |
| $P_5$  | -1.92  | 1.31  | 0.61   | −3.0912  | 6    |
| $P_6$  | -0.18  | −0.19 | 0.37   | −0.2466  | 4    |

**Table 7** Decision table by pessimistic values

|        | $\mu$  | $v$   | $h$    | Score    | Rank |
|--------|--------|-------|--------|----------|------|
| $P_1$  | −0.48  | 0.54  | −0.06  | −0.4512  | 5    |
| $P_2$  | 0.93   | −0.96 | 0.03   | 0.9579   | 1    |
| $P_3$  | 0.39   | −0.24 | −0.15  | 0.3315   | 3    |
| $P_4$  | 0.81   | −0.93 | 0.12   | 0.9072   | 2    |
| $P_5$  | −2.04  | 1.68  | 0.36   | −2.7744  | 6    |
| $P_6$  | 0.39   | −0.09 | −0.3   | 0.273    | 4    |

**Table 8** Decision table by neutral values

|        | $\mu$   | $v$    | $h$   | Score    | Rank |
|--------|---------|--------|-------|----------|------|
| $P_1$  | −0.435  | 0.355  | 0.08  | −0.4698  | 5    |
| $P_2$  | 0.975   | −0.725 | −0.25 | 0.73125  | 2    |
| $P_3$  | 0.39    | −0.2   | −0.19 | 0.3159   | 3    |
| $P_4$  | 0.945   | −0.785 | −0.16 | 0.7938   | 1    |
| $P_5$  | −1.98   | 1.495  | 0.485 | −2.9403  | 6    |
| $P_6$  | 0.105   | −0.14  | 0.035 | 0.108675 | 4    |

**Table 9** Rank table

| Players | Optimistic | Pessimistic | Neutral | Normalized score | Final rank |
|---------|------------|-------------|---------|------------------|------------|
| $P_1$   | 5          | 5           | 5       | 0.066667         | 5          |
| $P_2$   | 2          | 1           | 2       | 0.288889         | 2          |
| $P_3$   | 3          | 3           | 3       | 0.2              | 3          |
| $P_4$   | 1          | 2           | 1       | 0.311111         | 1          |
| $P_5$   | 6          | 6           | 6       | 0                | 6          |
| $P_6$   | 4          | 4           | 4       | 0.133333         | 4          |

# References

1. Tripathy, B.K., Sooraj, T.R., Mohanty, R.K.: Advances decision making using hybrid soft set models. Int. J. Pharm. Technol. **8**(3), 17694–17721 (2016)
2. Zadeh, L.A.: Fuzzy sets. Inf. Control **8**, 338–353 (1965)
3. Molodtsov, D.: Soft set theory—first results. Comput. Math Appl. **37**, 19–31 (1999)
4. Sooraj, T.R., Mohanty, R.K., Tripathy, B.K.: Fuzzy soft set theory and its application in group decision making. Adv. Intell. Syst. Comput. **452**, 171–178 (2016)
5. Tripathy, B.K., Sooraj, T.R., Mohanty, R.K.: A new approach to fuzzy soft set theory and its application in decision making. Adv. Intell. Syst. Comput. **411**, 305–313 (2016)
6. Mohanty, R.K., Sooraj, T.R., Tripathy, B.K.: IVIFS and decision-making. Adv. Intell. Syst. Comput. **468**, 319–330 (2017).
7. Peng, X., Yang, Y.: Interval-valued hesitant fuzzy soft sets and their application in decision making. Fundam. Informatica (2015)
8. Torra, V.: Hesitant fuzzy sets and decision. Int. J. Intell. Syst. **25**(6), 395–407 (2010)
9. Tripathy, B.K., Arun, K.R.: A new approach to soft sets, soft multisets and their properties. Int. J. Reasoning-based Intell. Syst. **7**(3/4), 244–253 (2015)
10. Tripathy, B.K., Sooraj, T.R., Mohanty, R.K., Arun, K.R.: A new approach to intuitionistic fuzzy soft sets and its application in decision-making. Adv. Intell. Syst. Comput. **439**, 93–100 (2016)
11. Tripathy, B.K., Sooraj, T.R., Mohanty, R.K., Arun, K.R.: Parameter reduction in soft set models and application in decision making. Handbook of research on fuzzy and rough set theory in organizational decision making, pp. 331–354 (2016). doi:10.4018/978-1-5225-1008-6.ch015
12. Tripathy, B.K., Mohanty, R.K., Sooraj, T.R.: On intuitionistic fuzzy soft sets and their application in decision-making. Lect. Notes. Electr. Eng. **396**, 67–73 (2016)
13. Tripathy, B.K., Mohanty, R.K., Sooraj, T.R.: On intuitionistic fuzzy soft set and its application in group decision making. Paper presented at the 1st international conference on emerging trends in engineering, technology and science, ICETETS 2016—Proceedings, doi:10.1109/ICETETS.2016.7603002 (2016)
14. Tripathy, B.K., Mohanty, R.K., Sooraj, T.R., Tripathy, A.: A modified representation of IFSS and its usage in GDM. Smart. Innov. Syst. Technol. **50**, 365–375 (2016)
15. Tripathy, B.K, Sooraj, T.R., Mohanty, R.K.: A new approach to interval-valued fuzzy soft sets and its application in decision-making. Adv. Intell. Syst. Comput. **509**, 3–10 (2017)
16. Mohanty, R.K., Tripathy, B.K.: Hesitant fuzzy soft set theory and its application in decision making. Adv. Intell. Syst. Comput. **517**, 221–233 (2017)
17. Sooraj, T.R., Tripathy, B.K.: Interval valued hesitant fuzzy soft sets and its application in stock market analysis. Adv. Intell. Syst. Comput. **517**, 755–764 (2017)
18. Sooraj, T.R., Mohanty, R.K., Tripathy, B.K.: Hesitant fuzzy soft set theory and its application in decision making. Adv. Intell. Syst. Comput. **517**, 315–322 (2017)
19. Atanassov, K.: Intuitionistic fuzzy sets. Fuzzy Set Syst. **20**, 87–96 (1986)
20. Atanassov, K., Gargov, G.: Interval valued intuitionistic fuzzy sets. Fuzzy Sets Syst. **31**, 343–349 (1989)
21. Broumi, S., Smarandache, F.: New operations over interval valued intuitionistic hesitant fuzzy sets. Math. Stat. **2**(2), 62–71 (2014)

# Design of Rivalize and Software Development to Convert RDB to RDF

**Syed Umar, P. Gayathri, Ch. Anil, N. Priya and R. Srikanth**

**Abstract** In the storage of databases, the relational databases are used utmost for storing data in very large amount which should be with good and high integrated security to safe the data semantics. So, the Semantic Web services facilitate a good solution which allows the data to be used safely and restage with various applications. But the limitation is to transform relational data to Resource Description Infrastructure/Framework for the business applications which will mostly depend on the relational model and RDF. So in this paper, we are concerning on to develop a new methodology R2R which is a RDB–RDF with integration approach contains architecture, algorithms, and mapping which will enable the RDF-based RW access to enrich the various characteristics of RDB. So the R2R ingression technology helps to align the RDM and RDF to handle the blank nodes, a clear description is shown in below topics.

**Keywords** RDF · RDM · Ingression methodology · PoC

## 1 Introduction

Software engineering field is constantly changing. New paradigms, programming languages and tools to obtain a huge wave suddenly and then they are as fast as they were unconscious. A little time to sacrifice, because it is important for the success

S. Umar (✉) · Ch. Anil · N. Priya
Department of Computer Science Engineering, MLR Institute of Technology,
Hyderabad, India
e-mail: umar332@gmail.com

P. Gayathri
Department of Information Technology, Gokaraju Rangaraju Institute of Engineering
and Technology, Hyderabad, India

R. Srikanth
Department of Information Technology, Institute of Aeronautical Engineering,
Hyderabad, India

of the new instruments, in both industry and research environment. Unlike short-term means, a body of knowledge provided by the companies usually still proceeds slowly and sometimes decades record for the largest database (RDB). This will inevitably lead to a conflict that must work as another generation of applications in the data. Old systems are reliable to easily adjust the views of their database applications or online transactions, and only when needed to change. Older systems are often important to the sun and should be too. They value for many organizations, but the material, maintenance, and further development are particularly difficult and costly.

This paper describes how RDB–RDF access [1] promotes the current transitional system, the Semantic Web application. Access offers a top semantics of relational database that exists. This allows RDF-based read and write relational databases. Based on the concept image that fills the gap between RDF and the relational model, the intermediary translates the Semantic Web applications, SQL flights. This enables the applications based on relational RDF and works with the same data and makes better use of the benefits of proven technologies such as database application's flexibility, security, and trade. The contribution of this thesis is to analyze how they are used successfully in access after the analysis of the development environment based on Eclipse software animation Rivalize [2] SOFAS [3] validated software platform to analyze cases and service-oriented distributed rations. We present a case where the focus of RDF and RDB is unsatisfactory and vision access.

## 2   Related Survey: Rivalize and SOFAS

### 2.1   Rivalize

In this case, a description of the two platforms to work together with the diagnostic program provides access information. The first option, Rivalize animated mannequin, considered legacy systems, although the bank scalable way of the past, which is deployed to analyze the development of software systems. Rivalize animation is based on the concept of historical license database (RHDB) [4]. It is used as a set of Eclipse plug-ins and combines data from multiple software repositories, such as version control, issue tracking, mailing list, and the combination of different databases, but listen and touch a variety of development components and software analysis. Animation is typical for old system's Rivalize: Even if the platform is still in active use, it is difficult to adapt to changing needs and the latest designs. Closely associated with the Eclipse makes it difficult to adapt and reuse tools and animation Rivalize algorithms in a new environment, such as within a service-oriented. The RHDB technology is also based on a traditional knowledge database. Therefore, it is difficult to exchange information with other sources of external data RHDB for communication with the local market, but not in general, and our

members do not have a special event of the center, which can be dereferenced online.

## 2.2 SOFAS

Evolyzer allows us to integrate and analyze various aspects of software development and success. It recognizes, however, that the great potential that is easily accessible and modular learning is not limited to the platforms and languages and requires no special tools to install and configure. On this basis, this place is what led to the "Evaluation Software as a service" [5]. Get easy access to a wide range of analytical tools and Web services that use a number of suppliers to implement this concept easier and flexible platform bed (Software Analysis Services) [3]. SOFAS-based RESTful architecture [6] enables the determination which is based on the principle of a simple but powerful analysis software and the use of state resources; the transfer of the Web architecture consists of three main parts: malformation SERVITI software analysis, analysis of broker software, and software analysis and encyclopedias. The service exposes the functionality and data from Web services interface standard is silent. Broker analysis software features the interface between the provider and the services and users. [7] This directory contains all of the analytical services that are related to the specific study of the taxonomy of software.

## 3 Case Studies of Bridge Data Analysis Software

Stimulation of the incident, we will introduce the case when it comes to communication Rivalize tween animation and use banks. They must be removed from the case in two directions, that is, the bank Rivalize animation and vice versa. First Rivalize animation data for hundreds of software and systems, introduced during the last year [8]. Banks should be able to access the platform to get without having to re-enter. To do this, read RDF-based animation Rivalize access to the database. Importers Rivalize animation equipment second raw data and version control systems like CVS and SVN core history [8]. The data is important because Rivalize animation can be existing tools used to build it. However, it requires write access to the database Rivalize RDF-based animation. Banks finally implemented and extensible software metrics to calculate the data. Again, the RDF data model, but the appropriate ratio to achieve animation Rivalize database schema. KA by subscription rights RHDB necessary information, the animation Rivalize indicators [9]. Efforts must use reflecting bridge and Bank animation Rivalize RDB mapping RDF data access read and write relative. Because the approach is currently limited to read-only services issues, we have developed without access to read and write. The rest of this section, we present the ontology model is relational data model

animation Rivalize banks and many entries using the conceptual difference between two data models. We will continue to the issues raised in this discussion [10].

## 3.1  Data Analysis Software System for Animation Rivalize

Animation program consists of several parts Details Rivalize, covering many aspects of the software engineering field. In this study, we will focus on the elements of the historical source, one of the most important audits. Review prepared a special version of the file. Personally, it means that software developers to edit the file, and must, version control system. The latest news of the change (i.e., the commit message from the client), as well as additional information such as the number of lines involved. Release is an important step in the life cycle of a software system [11]. It is often referred to the name and the code of the latest versions of files and data to the picture. New or experimental features and bug fixes often in branch. If the code is stable and is connected to the housing. To be a significant source model must meet state assessment system. The release is usually only guaranteed, so we can improve the initial part of the code and model analysis based initialize [12]. Redesign class if it meets the highest levels of Java, C #, etc. Class members, including the properties and methods of the numbers. Classes can, attributes and methods summarized in the machine.

Relations between organizations such as the origin of the calls between methods' access properties and methods' class heritage, they represent a class association or connection to the desktop. While it is difficult to distinguish the real ones, this is the only way to make the relational model; we would like to ask you to express contacts. Units can be measured.

## 3.2  Ontology Analysis Software

Describing the data and software used in banks, we have developed to the evolution of the lexicon (John) family. They describe a variety of software, such as version control, issue tracking, change each static source code structure of software engineering metrics [13], and so on. Sen was some ritual pyramid ontology. All major sub-domain with more vocabulary to define common concepts. The system-specific or language-specific concepts we have developed a special glossary. We must build on the advanced versions, and some concrete ontology ontology small systems such as control version CVS, SVN, Git, and high-level ontology version. In this article, we will discuss border basic terms and conditions code version control systems ontology and ontology. The source ontology model based on only static structures FAMIX source meta-model. Therefore, like animated Rivalize data in Table 1 summarizes the main features of the classroom, and to compel the source ontology.

**Table 1** Overview of source code ontology

| Class:Class | Class:Method |
|---|---|
| → declaresMethod:Method | → accessesField:Field |
| → declaresField:Field | → hasParameter:Parameter |
| → isReturnTypeOf:Method | → invokesMethod:Method |
| → isSubclassOf:Class | → hasReturnClass:Class |
| → isSuperclassOf:Class | → isInvokedByMethod:Method |
| → hasName:xsd:string | → isMethodOf:Class |
| Class:Field | → hasName:xsd:string |
| → isDeclaredFieldOf:Class | Class:Parameter |
| → isAccessedByMethod:Method | → isParameterOf:Method |
| → hasName:xsd:string | → hasName:xsd:string |

**Table 2** Overview of version ontology

| Class:Version | Class:ChangeSet |
|---|---|
| → hasID:xsd:string | → hasCommitDate:xsd:date |
| → follows:Version | Class:Branch |
| → precedes:Version | → hasTag:xsd:string |
| → hasCreationDate:xsd:date | Class:Release |
| → linesAdded:xsd:int | → hasReleaseDate:xsd:string |
| → linesDeleted:xsd:int | → hasTag:xsd:string |
| → hasMessage:xsd:string | |

Full ontology many other concepts, such as the surface, local variables, and exceptions.

Ontology Versioning Version control structure models and models based on the data. Table 2 summarizes the most important classes and attributes overview of SEON ontology version.

## 3.3 Access and Town Bridge Analysis Software

They live in a conceptual gap between animations and AVC Rivalize Semantic Web-connected banks. RDB mapping RDF translation and fly place to read RDF-based applications requires Rivalize RHDB animation. Tables 3 and 4 provide a systematic overview of the interpretation of the survey. Again, we focus on Map RHDB matter Rivalize animation [14]. Mapping uses the namespace. SEON versioned from the ontology. Table 3 listed in Tables Fig. 1 shows a map of the domain concept and features. The table is a pile of ontology mapping, but the honors of office-named properties. Company and certain other concepts only (below) represent the concepts of the ontology.

Table 4 lists the mapping table's connectivity combined with M N AVC. Since RDF provides several ways to represent M N Bond, and not necessarily as an

**Table 3**  Mapping overview of Rivalize source code

| Link table | → property | : inverse property |
|---|---|---|
| *Release_Revision* | → ver:comprises | : ver:appearsIn |
| *Branch_Revision* | → ver:comprises | : ver:isOn |
| *Transaction_Revision* | → ver:comprises | : ver:commitedIn |
| *File_Revision* | → ver:hasVersion | : ver:belongsTo |
| Person_Revision | → – | : ver:committedBy |
| *Class_Revision* | → ver:hasSource | : – |
| *Method_Class* | → java:isDeclaredMethodOf | : java:declaresMethod |
| *Attribute_Class* | → java:isDeclaredFieldOf | : java:declaresField |
| *Measurements_Entity* | → met:isMetricOf | : met:hasMetric |
| *Inheritance* | → java:hasSubClass | : java:hasSuperClass |
| *Invocation* | → java:invokesMethod | : java:isInvokedByMethod |
| *Access* | → java:accessField | : java:isAccessedByMethod |

**Table 4**  Mapping overview of version source code

| Table | Class | Attribute | Property |
|---|---|---|---|
| *Revision* | → ver:Version | number | → ver:hasID |
| | | previousRevision | → ver:follows |
| | | nextRevision | → ver:precedes |
| | | date | → ver:hasCreationDate |
| | | linesAdded | → ver:linesAdded |
| | | linesDeleted | → ver:linesDeleted |
| | | message | → ver:hasMessage |
| *Transaction* | → ver:ChangeSet | start | → – |
| | | end | → ver:hasCommitDate |
| *Branch* | → ver:Branch | name | → ver:hasTag |
| *File* | → top:File | path | → top:filePath |
| *Release* | → ver:Release | name | → ver:hasTag |
| | | date | → ver:hasReleaseDate |
| *Person* | → foaf:Person | name | → foaf:name |
| | | email | → foaf:mbox |
| *Entity* | → – | isAbstract | → Java:isAbstract |
| | | isStatic | → Java:isStatic |
| *Class* | → java:Class | | |
| *Method* | → java:Method | returnType | → java:hasReturnType |
| *Attribute* | → java.Field | | |
| *Measurement* | → met: SoftwareMetric | metric | → met:hasName |
| | | value | → met:hasValue |

**Fig. 1** Rivalize data representation and fetching of source code and historical analysis

assistant construction, ontology mapping tables shelves of the time occurs. The table consists of three columns. First name combined with table. Fig. 1 shows how the line connects the two concepts or concepts of the course.

## 4   Challenges and Limitations

We showed how the lacks of a bridge to successfully implement the new city case study. It provides a gradual transition from legacy systems, such as animation Rivalize, and new programs, such as our SOFAS. We showed how the conceptual gap between the relational data model Access Rivalize Banks and RDF-based animation creation. We also discussed the map; read-only event RDB OST is not

suitable for this application because they limit access to read-only service request which is based on RDF. During the case study, we faced a number of challenges with respect to Access Map. In addition, we added two majors and develop solutions to overcome them.

The first challenge is the concept of legacy relational database systems. Heritage core concept of object-oriented approach, etc., is often used in object-oriented systems, including animated Rivalize. Family members, as opposed to object-oriented databases, are direct legacy supporters. It is, however, to implement a major strategic legacy relational database system. Table class inheritance hierarchy represents the hierarchy in the table. Table columns for the category and characteristics of each type are called a separator (i.e., class) one for every occasion. As case studies, we need to add explicit support for inheritance hierarchy which does not have access to the inventory. The strategy at a table dais-of-the-box state as it applies to each category of the table and independent tables. Mapping two different strategies to support features such as columns and tables in discriminatory relations between parents and children is required. We bring this product to build a simple mapping that is not the possibility of mapping.

Another challenge is to define the descriptions of RDB to RDF. Mapping a cryptographically protected RDF makes them very suitable for automatic processing, but prevents access to human users. In fact, set a time for such reflection, and error-prone task, which consists mainly of repetitions. Therefore, it is an indispensable tool to determine the amount of the map for more complex applications, database tables, and columns of interest. We have developed a tool [6] to reduce the equations defined in the descriptions' access. It causes the semi-automatic mapping of RDB two-phase system. First of all, it automatically creates a mapping that is necessary, according to the list of database schema. Terms of the target ontology, already at this stage, are including the name and category column BLE database schema. Then, to show the graphic processing of mapping tools.

# 5   Conclusion

In theory, the Semantic Web provides a common framework that several major and reuse of data in the application borders' ciliates installation and the company. In practice, it is complicated by the wide acceptance in addition to the fact that many companies closed their data in relational databases. The business-critical legacy applications rely on databases to maintain daily operations and develop new systems are often implemented together with its predecessor, so it should be withdrawn gradually. Same as legacy systems, as well as their descendants, mainly active cooperation existing data. In this article, we present business access deVere-to-RDB platform to allow RDF-based relational database to read and write access. This will greatly facilitate the transition from old systems to Semantic Web-enabled application in practice of the semantic layer on top of the relational databases. Updating

applications for the implementation of Semantic Web query and on-the-fly SQL database system. So therefore there is no need to access and view the data that shows a relative synchronization RDF and also allows for a more suitable technique exploiting the benefits of the database. Slow migration path to our study, we explained access onto our implementation of two major programs that Rivalize animation legacy application and its successor, the platform for banks. We met a challenge when it comes to mapping genetic structure clearly inadequate approach and expanded the inventory of heritage strategy. In addition, we have created a map of a semi-automatic processing, ontology using a relational database schema.

# References

1. Hert, M.: Relational databases as semantic web endpoints. In: Proceedings of European Semantic Web Conference, June 2009
2. Gall, H.C., Fluri, B., Pinzger, M.: Change analysis with evolizer and changedistiller. IEEE Softw. (January/February 2009)
3. Ghezzi, G., Gall, H.C.: SOFAS: a lightweight architecture for software analysis as a service. In: Working IEEE/IFIP Confernce on Software Architecture, June 2011
4. Fischer, M., Pinzger, M., Gall, H.: Populating a release history database from version control and bug tracking systems. In: Proceedings of International Conference on Software Maintenance, September 2003
5. Ghezzi, G., Gall, H.C.: Towards software analysis as a service. In: Proceedings of International ERCIM Workshop on Software Evolution and Evolvability, September 2008
6. Fielding, R.T.: Architectural styles and the design of network-based software architectures. Ph.D. thesis, University of California, Irvine (2000)
7. Barrasa, J., Corcho, O., Perez-Gomez, A.: R2O expandable and semantically Home based ontology language mapping. In: Proceedings of Workshop Seminar Web and Resources, August 2004
8. Berners-Lee, T.: Links. http://www.w3.org/DesignIssues/LinkedData. HTML (2009), the most recent trip in June 2011
9. Berners-Lee, T.. Relational databases in the semantic web. Design issues. http://www.w3.org//AVCRDF.html. (2009), the most recent trip in June 2011
10. BIZERN, C., Cyganiak, R.: D2R server database publishing relational semantic web. In: Proceedings of International Seminar Web Conferences, November 2006
11. BIZERN, C.: sea, D2RQ—smoking databases virtual RDF RDF graph. In: Proceedings of International Seminar Web Conferences, November 2004
12. Brugger, N.: RDB mapping produce RDF schema relational database. Master's thesis, University of Zurich (2009)
13. Das, S.A., Sundara, S.A., Cyganiak, R.: R2RML: RDB mapping RDF language. W3C working draft. http://www.w3.org/TR/2010/WD-r2rml-20101028/. (Stone October 2010)
14. Demeyer, S., Ducasse, S., Nierstrasz, O.: Object Models Nierstrasz Reconstruction. Morgan Kaufmann Publishers, Inc., San Francisco, CA, USA (2002)

# An IoT-Based Low-Cost Weather Monitoring and Alert System Using Node MCU

P. Siva Nagendra Reddy, D. Vishnu Vardhan,
K. Tharun Kumar Reddy and P. Ajay Kumar Reddy

**Abstract** In environment, due to many factors, day by day pollution levels are increasing. So many natures' gifts like air, water, etc. Air is polluted because of many factors like landfill sites, toxic gas release from factories, etc. Planting trees to reduce air pollution but for many reasons like roads expansion, building homes, and forming, humans are cutting trees. So air pollution is increasing. It will cause harmful effects to the living things in the society. The motto behind this is to intimate people in the society about environmental conditions like temperature, humidity, and air quality in the surrounding area. An IoT is a new era in the present world. With the help of IoT, this system has developed. It is a low-cost and high-efficient system. It can monitor environmental conditions for every 60 s and upload to the web server. If any harmful gases like Hydrocarbons (HC), Carbon monoxide (CO), Sulfur dioxide (SO2) and Nitrogen oxides (NOx) etc., are released into the environment from various factors, this system will give an alert message to the higher authorities and people living in nearest areas, and it will give the alarm in the surrounding areas.

**Keywords** IoT · Air quality · Node MCU · DHT11

P. Siva Nagendra Reddy (✉) · K. Tharun Kumar Reddy · P. Ajay Kumar Reddy
Kuppam Engineering College, Kuppam, Chittoor, Andhra Pradesh, India
e-mail: snreddy715@gmail.com

K. Tharun Kumar Reddy
e-mail: kethireddy.tharun@gmail.com

P. Ajay Kumar Reddy
e-mail: ajaypedamalli@gmail.com

D. Vishnu Vardhan
JNTU Kalikiri, Chittoor, Andhra Pradesh, India
e-mail: vishnu.ece@jntua.ac.in

# 1 Introduction

Global Health Observatory (GHO) body of the government of India reported that the major sources for contribution of release of harmful gases are industries, cars, and trucks. Major contribution for a fine particulate matter comes from fuel combustion, which includes both mobile sources such as vehicles and immobile sources such as power plants, industry, households or biomass burning [1]. Due to this, the human health is majorly effected which can be given by the causes for deaths are as follows about 9% of COPD deaths, 26% of lung cancer deaths and about 14% of ischemic heart disease and stroke. This particulate matter pollution is mostly environmental-related health problem which is affecting the people living around the world, but the majority of people who are living in developing countries like India inexplicably experiencing this burden [2]. According to the World Health Organization (WHO), Hospitals in major cities are having more number of admissions of people who are suffering from respiratory diseases out of which our country placed top in the world, with 205 deaths per 100,000 people in 2015. Recent survey reported that children are predominantly weak pronouncement about 60% of the city's 5.3 million school children had diminutive lung growth [3].

Recently, in India's capital, Delhi also faced the problem of air pollution. Indian government declared an emergency situation in Delhi and temporarily closed construction sites, coal-fired power stations and other industries which produce harmful gases into the atmosphere. Not only Delhi, so many other places are also facing many problems because of air pollution. So necessary precautions are required [4].

Sunstroke is also one of the major problems in south India. So many peoples were dead or injured during summer season because of sunstroke. The Fig. 1 shows number of deaths in India from 1992 to 2015 [taken from ndma.gov.in].

In summer, temperature level reaches up to 48 °C in south India. Up to 37 °C no effect for human body. If the temperature increases more than 37 °C, human body starts observing heat from the atmosphere. Humidity is also a major concern in the environment. The Fig. 2 shows the temperature versus humidity chart, and temperature actually felt with an increase in humidity [5].

So monitoring of temperature levels is also a major concern in the society. In this proposed method, whenever temperature levels reached threshold value it will give an alert message through the SMS. So that this system can reduce the number of deaths because of sunstroke.

# 2 Proposed Method

This proposed method consists of MQ135 air quality sensor, Node MCU, DHT 11, and Buzzer. The block diagram of proposed system is shown in Fig. 3.

**Fig. 1** Number of deaths in different years [ndma.gov.in]

| Year | No. of Deaths |
|------|---------------|
| 1992 | 612 |
| 1993 | 631 |
| 1994 | 773 |
| 1995 | 1677 |
| 1996 | 434 |
| 1997 | 393 |
| 1998 | 1016 |
| 1999 | 628 |
| 2000 | 534 |
| 2001 | 505 |
| 2002 | 720 |
| 2003 | 807 |
| 2004 | 756 |
| 2005 | 1075 |
| 2006 | 754 |
| 2007 | 932 |
| 2008 | 616 |
| 2009 | 1071 |
| 2010 | 1274 |
| 2011 | 793 |
| 2012 | 1247 |
| 2013 | 1216 |
| 2014 | 1677 |
| 2015 | 2422 |

Node MCU has inbuilt Wi-Fi module. It can process both analog and digital signals. MQ135 is used to measure the air quality levels. The output of MQ135 is in analog form. It is connected to the analog pin(A0) of Node MCU. MQ135 is capable to detect harmful gases [5]. DHT 11 gives the temperature and humidity values in digital form. This output is connected to the digital pin (D0) of Node MCU.

The measured values of MQ135 and DHT11 sensors are continuously uploaded into Xively or Thing Speak IoT platforms. In this method, sensor values are uploaded into the Thing Speak continuously for every minute [6]. In the summer season, many people are injured because of sunstroke. This sensor measures temperature levels, and if the temperature levels are above the threshold levels it will give alert messages to the surrounding people.

Whenever the pollution levels increase in the air, SMS automatically will be sent to higher authorities and nearest people using SMS gate way. It also alerts the surrounding people by using the alarm system. SMS gate way is used to send SMS from the web server. No need of GSM to send SMS and it is also a costly process to send bulk SMS from GSM. As per TROY rules, we cannot send more than 100 SMS per day from each SIM. So SMS gate ways are used to send the messages [7]. All the people can access or observe the environment conditions in their nearest place from the Thing Speak database.

## Temperature/ Humidity Index

| Relative Humidity % | Temperature °C | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 27 | 28 | 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 | 37 | 38 | 39 | 40 | 41 | 42 | 43 |
| 40 | 27 | 28 | 29 | 30 | 31 | 32 | 34 | 35 | 37 | 39 | 41 | 43 | 46 | 48 | 51 | 54 | 57 |
| 45 | 27 | 28 | 29 | 30 | 32 | 33 | 35 | 37 | 39 | 41 | 43 | 46 | 49 | 51 | 54 | 57 | |
| 50 | 27 | 28 | 30 | 31 | 33 | 35 | 36 | 38 | 41 | 43 | 46 | 49 | 52 | 55 | 58 | | |
| 55 | 28 | 29 | 30 | 32 | 34 | 36 | 38 | 40 | 43 | 46 | 48 | 52 | 54 | 58 | | | |
| 60 | 28 | 29 | 31 | 33 | 35 | 37 | 40 | 42 | 45 | 48 | 51 | 55 | 59 | | | | |
| 65 | 28 | 30 | 32 | 34 | 36 | 39 | 41 | 44 | 48 | 51 | 55 | 59 | | | | | |
| 70 | 29 | 31 | 33 | 35 | 38 | 40 | 43 | 47 | 50 | 54 | 58 | | | | | | |
| 75 | 29 | 31 | 34 | 36 | 39 | 42 | 46 | 49 | 53 | 58 | | | | | | | |
| 80 | 30 | 32 | 35 | 38 | 41 | 44 | 48 | 52 | 57 | | | | | | | | |
| 85 | 30 | 33 | 36 | 39 | 43 | 47 | 51 | 55 | | | | | | | | | |
| 90 | 31 | 34 | 37 | 41 | 45 | 49 | 54 | | | | | | | | | | |
| 95 | 31 | 35 | 38 | 42 | 47 | 51 | 57 | | | | | | | | | | |
| 100 | 32 | 36 | 40 | 44 | 49 | 56 | | | | | | | | | | | |

Caution    Extreme Caution    Danger    Extreme Danger

Source: Calculated °F to °C from NOAA's National Weather Service

**Fig. 2** Temperature versus humidity chart and temperature actually felt with increasing in humidity

## 3 Node MCU

It is an open-source IoT platform. It has firmware to run ESP8266 Wi-Fi System on Chip Module. It is the best platform to develop IoT Applications and Prototype projects. It has ESP8266 Wi-Fi Module, integrated GPIO, PWM, IIC, and serial communication pins. It can process both analog and digital signals.

Node MCU 1st and 2nd generation boards have 4 MB flash memory. The firmware uses the LUA Scripting Language. Using Arduino IDE, Node MCU can be programmed. By installing ESP8266 board drivers from board manager option or web sources, Node MCU can be programmed from Arduino IDE. Node MCU has 128 Kbytes of Memory and uses eXtendable Test Operating System (XTOS). Node MCU is very low-cost IoT development board compared with the Intel Galileo, Raspberry Pi, UDOO, and other IoT development boards. So, in this design, Node MCU is taken as a controller [8] and it is shown in Fig. 4.

**Fig. 3** Proposed low-cost weather monitoring and alert system

**Fig. 4** Node MCU board



## 3.1 Air Quality Sensor

MQ 135 has high sensitivity to Sulfide, Ammonia, Benze, smoke, and other harmful gases. It is a low-cost air quality sensor. MQ135 conductive material gives low output levels in clean air. Conductivity level increases with increase in

pollution in air (increase in harmful gas concentration in the air). It has 3 pins. It gives analog output at pin 2, others pins are VCC and Ground. The lifetime of the MQ 135 sensor is also more (Fig. 5).

## 3.2  DHT 11

DHT11 is a humidity and temperature sensor. It has 3 pins. DHT11 gives digital form of output (pin2), and others pins are VCC and ground. It works with supply voltage of 3.3–5 V. It can measure humidity range from 20 to 90% RH. It gives an accuracy of ±5 RH. It can measure temperature range from 0° to 60 °C with an accuracy of ±2 °C (Fig. 6).

## 4  Results

This proposed method gives exact values of air quality, humidity, and temperature. In this design, MQ135 is used to measure air quality and DHT 11 is used to measure humidity and temperature. All the measured values measured by sensors are processed by Node MCU and uploaded into the ThingSpeak using ESP8266. The uploaded results are shown in below. Figure 7 shows the uploaded values of MQ135 sensor. This system will give an alert message when air quality levels below the value 2.5.

The Figs. 8 and 9 show uploaded values of temperature and humidity. Whenever temperature and humidity level exceeds above threshold level it will send a message to people through SMS gate way.

The Figs. 10 and 11 show SMS received from the SMS gate way when the MQ 135 sensor values and DHT 11 values crossed certain threshold levels. The Fig. 12 shows the experimental setup of An Iot-Based Low-Cost Weather Monitoring and Alert System Using Node MCU.



Fig. 5  MQ135 gas sensor

**Fig. 6** DHT 11 sensor





**Fig. 7** Air quality sensor values uploaded into ThingSpeak



**Fig. 8** Temperature values uploaded into ThingSpeak

**Fig. 9** Humidity values uploaded into ThingSpeak

**Fig. 10** Alert message when temperature exceeds threshold level



**Fig. 11** Alert message when pollution exceeds threshold level

**Fig. 12** Experimental setup of proposed design

## 5 Conclusion

The brisk growth in transportation and industrial plants creating ecological issues like atmosphere change, flawed and pollution has greatly inclined for the need of efficient operationally malleable and smart monitoring systems. As a developing country like India which is second most populous in the world, it is necessary to monitor the quality of natural resources like air. There are so many industries in our country which are releasing polluted air without implementation of the proper cleaning process. In this paper, we developed a module, which can continuously monitor the air quality in an area with humidity sensor by keeping threshold levels accepted by governments. When the limit exceeds certain levels we can immediately send an alert message to the concerned authorities so that they can take proper action to avoid the major human loses. We also monitor the temperature in that area and send those values to the concerned authorities so that they can warn the people go outside to avoid sunstroke which is another major problem faced by people in summer. From this paper, we can conclude that this project will help the government to protect the people from diseases and also possible to take necessary action on the industries which are releasing the polluted air.

## References

1. Tharun Kumar Reddy, K., Ajay Kumar Reddy, P., Siva Nagendra Reddy, P.: An IoT based remote monitoring of landfill sites using raspberry Pi2. Emerging Trends in Electrical, Communications and Information Technologies, p. 219 (2017)
2. Abraham, S., Li, X.: A cost-effective wireless sensor network system for indoor air quality monitoring applications. Proc. Comput. Sci. **34**, 165–171 (2014)
3. Mendez, D., et al.: P-sense: a participatory sensing system for air pollution monitoring and control. In: 2011 IEEE International Conference on Pervasive Computing and Communications Workshops (PERCOM Workshops), IEEE (2011)

4. Simić, M.: Design and development of air temperature and relative humidity monitoring system with AVR processor based web server. In: 2014 International Conference and Exposition on Electrical and Power Engineering (EPE), Iasi, pp. 38–41 (2014)
5. Naresh Naik, R., Siva Nagendra Reddy, P., Nanda Kishore, S., Tharun Kumar Reddy, K.: Arduino based LPG gas monitoring & automatic cylinder booking with alert System IOSR-JECE, **11**(4), Ver. I, pp. 6–12 (Jul–Aug 2016), e-ISSN: 2278-2834, p-ISSN: 2278-8735
6. Noordin, K.A., Onn, C.C., Ismail, M.F.: A low-cost microcontroller-based weather monitoring system. CMU J. **5**(1), 33–39 (2006)
7. Busari Sherif, A., Dunmoye Abibat, F., Akingbade Kayode, F.: Development of Arduino-based data acquisition system for environmental monitoring using Zigbee communication protocol (2016)
8. Devikar, P., Krishnamoorthy, A., Bhanage, A., Chauhan, M.S.: IoT based biometric attendance system. Int. J. Adv. Res. Comput. Commun. **5**, Special Issue 2 (2016)

# Energy Efficient Secure Data Transmission in Wireless Sensor Network

**Venu Madhav Kuthadi, Rajalakshmi Selvaraj and Tshilidzi Marwala**

**Abstract** Wireless sensor network (WSN) turns out to be a vital innovation identified with systems as well as it is by and large broadly connected and received in the present day control and checking functions. In any case, problems with respect to vitality proficiency create huge worries during the time spent in remote sensor system. For problems concerning vitality, power efficient adaptive clustering hierarchy (PEACH) tradition has been suggested essentially for WSNs so as to criticize vitality use of hubs and expanding life range of the system. The hubs in the remote sensor system structure as a gathering. The most noteworthy vitality hub in the gathering has been chosen as a bunch head. The bunch head determination has a few criteria in choosing: (1) based on the vitality levels of sensors in group (i.e., the hub with higher vitality would be chosen as a bunch head). (2) based on hub degree (total number of hubs associated with it) (i.e., the hub with advanced hub degree would be chosen as a bunch head). The time differential of arrival (TDA) has been proposed for vitality consumption of sensor hubs in the system. TDA decreases the vitality utilization of hubs and expanding the system life duration. The node connectivity algorithm (NCA) has been proposed to discover the availability with neighbor's hub in the remote sensor system. The bunch head confirms the gathering hubs and transfers the confirmation data to server. On the off chance that any noxious hub available in the system, the server would recognize through self-mending key appropriation procedure and the bunch head would shut the hub and toss out the hub from the bunch. The noxious hub assets would be stored into the server. It will make the system more ensured and safe one. To broadcast the parcels from resource to goal, figure the most limited way with the assistance of

V. M. Kuthadi (✉) · R. Selvaraj
Department of Computer Science, BIUST, Palapye, Botswana
e-mail: kuthadiv@biust.ac.bw

R. Selvaraj
e-mail: selvarajr@biust.ac.bw

V. M. Kuthadi · T. Marwala
Faculty of Engineering and the Built Environment, University of Johannesburg,
Johannesburg, South Africa
e-mail: tmarwala@uj.ac.za

Dijkstra algorithm. The bundle encodes at source and unscrambled once it gets to the goal by means of Advanced Encryption Standard (AES) encryption calculation.

# 1   Introduction

The development in innovation implies that numerous remote sensor hubs have been presently in general minimal at price; in any case, the expense of conveying them could stay elevated. There has been a necessity to obtain the best life out of a system of sensors, and the life has been for the most part constrained by battery supply utilization.

Remote sensor systems (WSN's), the change in sensor innovation has prepared it conceivable to contain to a great degree; little, low-fueled detecting gadgets furnished with programmable figuring, numerous factor detecting and remote correspondence capacity. Likewise, the minimal price ensures it is conceivable to have a system of hundreds or more number in thousands of these sensors, in this manner upgrading the unwavering quality and exactness of information and the zone scope. Remote sensor systems present data regarding remote structures, across the board natural changes, and so forth in obscure and unwelcoming territories. There have been various favorable circumstances of remote sensor networks with regard to wired ones, for example, simplicity of sending (lessening establishment price), expanded reach (system of small sensors can be disseminated over a more extensive locale. Remote sensor system has been a developing field at which loads of exploration work is being done including equipment and framework plan, organizing, security consider, and appropriated algorithm [1–3]. Sensor hubs regularly read the information parcel and exchange it to the base station by means of a little middle-of-the-road hubs. The sensor hubs have been at ease, little influence, and diminutive in sort of transmission [4]. Hubs employ to transfer information bundle nearby to its solitary jump nearby hubs; thus on lastly, it compasses to its base station. At first, hubs have been sent soaring from flying machines or arbitrarily, and a number of times hub alters its underlying location (the season of arrangement) and shifts over the area taking into account the prerequisite; so this sort of hubs has been known as portable hubs. Hence, there have been two sorts of information transmission in remote sensor organize; they are direct transmission and multi-jump information transmission. In the information transfer which is straight, the information have been sent straightforwardly to the sink at which as multi-jump transmission information transfer takes place by means of number of middle-of-the-road hubs lies between source hub and base station. In the set up of sensor arrangement, the stream of information has been deemed an essential angle on the grounds that every information bundle includes the occasion which might be imperative for certain functions. Hence, transmission of information should be protected. In any

case, sensor hub has constrained vitality and restricted memory limit so keeping up safety has been troublesome for them [5]. It ought to be ensured that the accounts from the 'sensors in real life' are to be valid and achieve the base station (BS) with no creation or alteration. The errand of safeguarding the remote sensor systems has been in any case, entangled in light of the fact that sensors have been profoundly mysterious gadgets with a restricted vitality and reminiscence limit, and at first, they have no information of their areas in the arrangement situation. The single hubs have been skilled at detecting the surroundings around them, preparing the data insights in the region, as well as transferring information to one or more arrangement focuses in a WSN. Proficient transmission of information has been a standout among the most huge issues for WSNs. Typically numerous WSNs have been introduced in secret, unforgiving, and frequently ill-disposed bodily situations for particular functions, for example, military spaces and detecting assignments with problematic environment. Effective and protected transmission of information plays a key role in sensible WSNs. One of most imperative examination issue in WSN has been the means by which to safeguard the sensor system topology and its correspondence methodology from possible modifications which could be completed through malignant hubs within the current system. There have been a few strategies which are utilized by numerous analysts; in the wake of checking on a few papers, it has been found that assortment of system connected with sensor system to protect information transmission like cryptographic calculation, network layer convention, physical MAC control and steering system [6]. In the study, our primary methodology avoids malignant hubs passage in sensor system. By this system, information could be protected in system, furthermore, vitality of hubs in the system likewise kept up. At this point, bunching strategy has been utilized for system versatility and regulation that amplifies hub life traverse and diminishes transfer speed usage by utilizing neighborhood collaboration amid sensor hubs. In a bunch-directed WSN (CWSN), every group is accommodated with a pioneer sensor hub, called as group head (CH). A CH gathers the information assembled by the leaf hubs (non-CH sensor hubs) in its bunch and transfers the pooled information to the base station (BS). At long last, cryptography system has been additionally utilized to give classification of information in system. Whatever is left of the paper has been sorted out as takes after. Segment 2 gives a review of background and related research carried out in range of WSN safety with existing strategy. Segment 3 presents suggested work plan and execution technique for accomplishing fancied attributes of WSN framework. Area 4 gives outcome and execution examination procedure accomplished for recommended work plan utilizing standard factors assessment. The research then closes in segment 5.

## 2 Related Work

In this area, the researchers have exhibited the short presentation concerning remote sensor system and their safety limitations to comprehend their arrangement. The point of venture effort has been to give a productive and end-to-end client security instrument for remote sensor system and that ought not influence the execution, dependability, and calculation assignment of sensor system. The entire sensor system spoke to utilizing layered design to speak to its diverse level safety viewpoint. The essential practical security prerequisite for several WSN function could be arranged as far as integrity, availability, confidentiality, robustness, authenticity, self-synchronization, survivability, and non-denial. Group-based information transmission in WSNs has boosted the sensor hubs span of life and their data transfer capacity utilization by coordinated effort with nearby sensor hubs. In CWSN [7], each group possesses a pioneer hub in particular, cluster head that totals the information gathered from the leaf hubs in its bunch and sends accumulated information to the base station (BS). The small vitality versatile grouping hierarchy (LEACH) protocol [8] has been one recognized to lessen and adjust the vitality utilization for CWSNs and furthermore accomplishes an extensive hubs span of life. Versatile periodic threshold sensitive energy efficient sensor network (APTEEN) [9] and power efficient adaptive clustering hierarchy convention have been two conventions ideas of LEACH. In addition the safety to LEACH is exceptionally troublesome as they pivot haphazardly, intermittently, and they doesn't have durable association with hubs as they turn in encircling furthermore, experiences orphan node problem [10], implies where a hub does not impart a pairwise key to others in its preloaded key ring, and while it is not adequate to impart its pairwise symmetric keys to all of the hubs in its ring, it can't take an interest in any group, and along these lines, needs to choose itself as a cluster head (CH). Cryptographic calculations have been a vital part of the safety engineering of WSNs, utilizing the most proficient and adequately safe calculation that is a successful method for monitoring assets of small sensors. While submitting an application any encryption plan, it needs transmission of additional bits. It has been an extremely valuable to give information safety utilizing information encoding copyright and similarity with existing system layer conventions. The cryptographic calculations utilized as a part of WSNs have been for the most part classified into two sections: symmetric key calculations and Asymmetric key calculations. Symmetric key cryptographic instruments utilize a solitary divided key among the two imparting host that has been utilized both for encryption and unscrambling. Symmetric key calculations could be additionally partitioned into piece figures for settled changes on plain-message information, and rivulet figures for time differing changes. It would provide a correlation for certain encryption calculations at various settings, for example, diverse sizes of information squares, distinctive information sorts, battery power utilization, distinctive key size lastly encryption/unscrambling speed. Symmetric key cryptosystems, for example, the RC4, RC5, AES, CAST, DES calculation has been utilized as a part of WSN. In

awry open key cryptosystems, every hub has an open key and a private key. General society key is distributed, while the private key has been kept mystery. This lopsided key cryptography has been giving improved safety regarding many-sided quality in excess of larger amount of assaults. Topsy-turvy open key cryptosystems, for example, the Diffie–Hellman key ascension, ECC, or RSA marks have been ordinarily excessively moderate in their efforts to establish safety; however, including an excessive amount of many-sided quality and convention overhead to be usable in WSN arrangements [11]. The primary disadvantage of this cryptography procedure is that it experiences high computational unpredictability and correspondence overhead with requirements of sensor hubs. Atique et al. [12] suggested a protected directing method to recognize fake records and dim opening assault by the utilization of measurable on the way separating (SEF) to expand the safety stage in WSN [13]. To advance enhance the safety and diminish the vitality utilization amid transmission of detected information, the creators executed circular bend cryptography (ECC) [14]. The suggested strategy gave assured results as far as safety is considered and generates a few hindrances in the method for an aggressor.

## 3   Proposed Mechanisms

### 3.1   Overview

The vital contemplations of outlining remote sensor system have been safe course in transmission of information. A number of sensor hubs have been introduced in unfasten spot, and it has least vitality that might prompt deprived execution and safety of sensor hubs information. This segment shows the suggested plan to distinguish pernicious hub action utilizing idea of self-recuperating, and hub confirmation strategy. At this point, bunching technique connected with hubs gathering and keeps up data of all hubs in WSN and cryptography procedure has been utilized for giving information privacy on system. By utilizing these techniques, information could be protected and hubs vitality has been looked after dependably.

### 3.2   Cluster Network-Based Data Transmission

In group systems, sensors hub transmit information to the bunch head at which information transmission has been executed. In vitality-obliged sensor systems of vast size, it has been wasteful for sensors to broadcast the information straightforwardly to the sink. The group heads could speak with the sink straightforwardly by means of long range transmissions or multi trusting during other bunch heads. Hence, in this suggested work, grouping strategy has been utilized to keep up hubs

data, transmission of information, and find noxious hub action in WSN. In this framework, power efficient and adaptive clustering hierarchy (PEACH) convention for WSNs is to reduce the vitality utilization of every hub, and boost the systems life span. In this study, bunch development has been executed through the utilization of vitality and hubs availability with different hubs in remote correspondence to bolster versatile multi-level grouping and keep away from extra overheads. In WSNs, catching a hub could perceive the resource and the goal of bundles transmitted by the nearby hubs. PEACH has been relevant in both areas: unconscious and area mindful remote sensor systems. PEACH has been intended to work on probabilistic steering conventions, so as to give a versatile multi-level grouping as shown in Fig. 1.

PEACH has been for the most part more versatile and proficient to the different conditions than the current bunching conventions of the remote sensor networks. The researchers consider a system that comprises of groups at which every bunch consists of a head hub. The group head hub has the most astounding vitality in the bunch and is progressively appointed in view of the vitality contemplations and network with of the considerable number of hubs in the group. To discover high vitality hub in bunch, it utilizes TDA strategy. TDA technique distinguishes most noteworthy vitality hub in view of the extents from one hub to other. At this time, the hub that has greatest network territory to other hub has been chosen as group head. In the meantime, hub network by means of additional hub furthermore reflect on to choose bunch head. For example if two or more hubs are likely to have similar vitality then bunch header will be considered. At this time once bunch has been shaped, group head send recognized bundle to all hubs in that gathering, then hub in bunches propels answer affirmation to group head. In the event that any hubs in bunch did not answer to group head that hub would be considered as vindictive hub. The similar procedure would be accomplished for all bunch bunches. The bunch head forwards the hubs noxious data to server for prospect utilization. Server has been utilized to put away noxious hubs data.

**Algorithm: Functions of cluster formation and joining**

Function $NO_i$ createCluster($T_{wait}$)
$NOi.setTimer(T_{wait})$
While $No_i.$ timeTerminated() – false do
$N_{des} \leftarrow No_i.acceptPacket()$
If $NO_{des} - No_i$ then
$NO_i.aggrPacket()$
elseif$NO_i.distFrom(NOdes) \leq NO_i.distFrom(NO_i.nextHop)$then
$NO_i.joinCluster(NO_{des})$
end if
end while
end function
function
$NO_i.joinCluster(NO_{des})$

**Fig. 1** Secure and energy-efficient data transmission architecture

$NO_i.nextHop \leftarrow NO_{des}$
end function

## *3.3  Node Authentication and Self-Healing*

Information respectability has been the nature of accuracy, fulfillment, wholeness, soundness, and consistence with the expectation of the makers of the information. It is accomplished through averting unapproved addition, adjustment, or obliteration of information. In WSNs, a pernicious hub might alter communications to bother the system usefulness. Additionally, because of questionable correspondence channels, it is anything but difficult to infuse tainted bundles or alarmed information. For giving information uprightness in WSN hub, verification procedure has been utilized to recognize malevolent hubs movement. Verification systems have been utilized to distinguish vindictively or ridiculed hubs. They have been particularly critical in WSNs that utilize a mutual remote medium. Hubs confirmation has been essential errand to keep up protected system. In this framework at whatever point, any novel hub goes into any bunch, primary group head checks if it is malignant or not. For that, it transfers affirmation bundle to novel hub, if novel hub answers affirmation parcel to bunch head after that group head gives consent to novel hub to go into bunch. This validation procedure provides additional accommodation to keep up the data trustworthiness of WSN. The cause has been if any hub in bunch gathering has been bargained hub, it would ruin information respectability of WSN. When vindicive hub is recognized, self-recuperating procedure is used to control the hubs from forming the group in WSN. In hub verification procedure, it distinguishes noxious hub in the system. On the off chance that any hub has distinguished as noxious hub in hub confirmation prepare this framework would store pernicious hub data in like manner server of WSN. This data has been utilized as a part of self-mending procedure duration. Essentially self mending procedure is used for wedging passage of noxious hub. The similar procedure has been implemented in the framework to check whether the vindictive hub data from the server is moving towards the novel hub to form grouping.

**Algorithm: Node connectivity algorithm**

**//Find Node Position**
```
nn=number of node;
for (int i=0; i<nn; i++)
{
int X1[i] = nx[i]
int Y1 [i]= ny[i]
}
```
**//Find neighbor node**
```
initilize m=0;
for(int i=0; i<nn; i++)
{
initilize k=0;
initilize p=0;
initilize a,b,cx,dx,e,fx,g;
```

```
for(int j=0; j<nn; j++)
{
a = X1[j]-X1[i];
b = a*a;
cx = Y1[j] - Y1[i];
dx = cx*cx;
e = b+dx;
fx = 0.5;
g = pow[e,fx];
if {g<=250 && i!=j}
{
k=k+1;
initilize ip_n(p) = j;// node i is neighbour to node j;
initilize nei (m) = j;// node i is neighbour to ip_n(p);
m = m+1;
}
}
count[i] = k; //number of neighbour node for node i is ———————————> count
[i];
}
initilize start = 0, snd = count[0], temp = 0;
for(int m=0;m<nn; m++)
{
snd = count[m]+temp;
for(int i= start; i<snd; i++)
{
initilize neighbour[i] = nei(i); // Node m is neighbour to neighbour i
start = snd;
temp = snd;
}
}
```

## 3.4 Finding Shortest Path and Forward Packet to Destination

During the bunch arrangement, researchers could be without much of a stretch enhance the lifespan of system and effortlessly diminish activity of the system. Hubs having a place with the DS have been completing all correspondence coming up short on vitality rapidly that is being put up on the lingering vitality of every hub. In this task, Dijkstra calculation has been utilized to discover briefest way frame one hub to further. Dijkstra's calculation has been connected to naturally discover bearings among material areas, for example, pouring headings on sites

such as Mapquest or Google Maps. It is additionally utilized for explaining an assortment of most limited way issues emerging in plant and office format, mechanical autonomy, transportation, and VLSI outline. In this undertaking, it discovers most brief way from a specified hub s to every additional hub in the WSN. Hub s is known as a beginning hub or an underlying hub. Dijkstra's calculation begins through allotting a number of underlying qualities for the separations from hub s and to each other hub in the system. It works in stepwise, at which at every progression the calculation enhances the separation charges. At every progression, the most limited separation from hub s to another hub is resolved. By applying Dijkstra's procedure, briefest way from one to other has been ascertained at which resource and goal have been distinguished. When resource and goal have been settled, parcel would be forwarded in system with encryption position. While pertaining any encryption plan, it needs transmission of additional bits. Its extremely valuable to give information safety utilizing information encoding patent and similarity with obtainable system layer conventions. At this time, AES encryption procedure has been utilized for encryption. This encryption procedure gives privacy of sent bundle. In the event that any pernicious hub attempts to get to sent parcel, it would not get genuine information data.

## 4 Results and Discussion

This area included result investigation for suggested model of energy-efficient secure data transmission method utilizing malicious hub ID and encryption procedure. This examination is useful to satisfy necessity of remote sensor system imperatives. At this time, this execution of suggested plan of vitality productive has been accomplished with node validation, self-healing, and AES encryption systems. The usage of vitality for remote sensor system relies upon various components, for example, bunching and vitality utilization strategies. PEACH and grouping techniques have been utilized for vitality utilization as a part of WSN. In existing framework, LEACH convention has been utilized to lessen and adjust the vitality utilization in WSN. At this time, the researchers demonstrate a number of graphical examinations of prevailing and suggested business connected with vitality, pernicious, and safety.

The Fig. 2 indicates comparison procedure of PEACH and LEACH convention. PEACH has been suggested convention. Execution of proposed convention improved the zone of vitality utilization and systems lifespan analyze than existing convention. In the suggested framework, the researchers utilize hub confirmation and self-mending procedure to distinguish vindictive hub in the system. This is one of the improvements of the work in suggested framework. Noxious hub recognizable proof likewise assumes imperative part for vitality utilization in WSN by maintaining a strategic distance from parcel retransmission and way choice.

Figure 3 demonstrates pernicious hub recognizable proof in suggested framework. At this time, malevolent hub recognizable proof has been accomplished for

**Fig. 2** Comparison of proposed and existing system



**Fig. 3** Energy consumption by malicious node identification



safety and vitality utilization. Blocking malevolent hub aids for information trustworthiness in WSN. Classification of information has been essential perspective in WSN. In the proposed system AES encryption has been utilized for getting the information classification and its moving time in WSN. The encryption instrument changes over plaintext into figure content to give information safety.

Figure 4 clarifies examination procedure of AES method further encryption system such as RSA, DES. The primary preferred standpoint of AES has been at which it invests least energy for encryption and decoding procedure. It additionally utilizes least support size for bundles. On the off chance that parcel size is increment the other encryption procedure of RSA and DES gets additional time encryption and decryption procedure yet AES gets not exactly RSA and DES.

## 5   Conclusion

In the given framework, it employs grouping convention for empowering sensor hubs to decrease information bundles by information total in remote sensor systems. The remote correspondence expense has been diminished by decreasing the information bundles, and the bunching conventions enhance the duration of life and

the vitality utilization of the systems. In the given study, it utilizes a vitality-effective bunch-based directing conventions, where PEACH has no overhead on group head determination and structures versatile multi-level grouping PEACH altogether enhances the duration of lifespan and the vitality utilization of the remote sensor systems contrasted and other bunching conventions. Malignant hub additionally assumes vital part for vitality utilization of hubs in system. In the given study, blocking and ID malignant hub procedure has been accomplished by hub verification and self-recuperating procedure. By chunking passage of malevolent hub diminish undesirable retransmission parcel transferring. This procedure has been exceptionally useful for recovering the vitality of hubs in WSN.

# References

1. Akyildiz, F., Su, W., Sankarasubramaniam, Y., Cayirci, E.: Wireless sensor networks: a survey. Comput. Netw. **38**, 393–422 (2002)
2. Dai, S., Jing, X., Li, L.: Research and analysis on routing protocols for wireless sensor networks. In: International Conference on Communications, Circuits and Systems, vol. 1, pp. 407–411 (2005)
3. Kuthadi, V.M., Rajendra, C., Selvaraj, R.: A study of security challenges in wireless sensor networks. JATIT **20**(1), 39–44 (2010)
4. Culler, D.E., Hong, E.: Wireless sensor networks. Commun. ACM **47**(6), 30–37 (2007)
5. Kim, S.C., Jeon, J.H., Park, H.J.: Energy Efficient Data transmission Mechanism in Wireless Sensor Networks, Lecture Notes in Electrical Engineering, vol. 235, pp. 845–852. Springer, Berlin (2013)
6. Bok, K., Lee, Y., Park, J., Yoo, J.: An energy-efficient secure scheme in wireless sensor networks. J. Sens. **16**, 1–11 (2016)
7. Kuthadi, V.M., Selvaraj, R., Marwala, T.: An efficient web services framework for secure data collection in wireless sensor network. Br. J. Sci. **12**(1), 18–31 (2015)
8. Abbasi, A.A., Younis, M.: A survey on clustering algorithms for wireless sensor networks. Comput. Commun. **30**(14), 2826–2841 (2007)
9. Heinzelmann, W., Chandrakasan, A., Balakrishnan, H.: An application-specific protocol architecture for wireless micro sensor networks. IEEE Trans. Wirel. Commun. **1**(4), 660–670 (2002)

10. Manjeshwar, A., Zeng, Q.-A., Agrawal, D.P.: An analytical model for information retrieval in wireless sensor networks using enhanced APTEEN protocol. IEEE Trans. Parallel Distrib. Syst. **13**(12), 1290–1302 (2002)
11. Sharma, S., Jena, S.K.: A survey on secure hierarchical routing protocols in wireless sensor networks. In: Proceedings International conference on Communication, Computing & Security (ICCCS), pp. 146–151 (2011)
12. Kuthadi, V.M., Selvaraj, R., Marwala, T.: An enhanced security pattern for wireless sensor network. In: Proceedings of the Second International Conference on Computer and Communication Technologies: IC3T 2015, vol. 2, p. 61. Springer (2015)
13. Kukkurainen, J., Soini, M., Sydanheimo, L.: RC5-based security in wireless sensor networks: utilization and performance. WSEAS Trans. Comput. **9**(10), 1191–1200 (2010)
14. Sakharkar, S.M., Mangrulkar, R.S., Atique, M.: A survey: a secure routing method for detecting false reports and gray-hole attacks along with Elliptic Curve Cryptography in wireless sensor networks. In: IEEE Students conference on Electrical, Electronics and Computer Science (SCEECS) 2014, India, pp. 1–5

# Significant Rule Power Factor: An Algorithm for New Interest Measure

**Ochin Sharma, Suresh Kumar and Nisheeth Joshi**

**Abstract** In the process of knowledge discovery, a large number of patterns are generated. It is therefore infeasible for an expert to provide his opinion considering this vast amount of patterns. One of the methods that are useful to overcome this problem is association rule interestingness measures. Interest measure essentially helps to fetch the data of interest in data mining. The data of interest completely depends on the competitive business requirements. Therefore, many interest measures exist; a few are support, confidence, lift, leverage. In this paper, we are proposing a new interest measure which is more informative and accurate when compared with many existing interest measures. We have used WEKA, JEDIT, ANT open source tools to conduct experiments. Through experiments, it has been observed that our proposed interest measure can fetch more important rules than many of the existing interest measures.

## 1 Introduction

In association rule mining, the aim is to extract interesting associations, correlations and frequent patterns from set of items from transaction data sets. Association rule mining is applicable to many fields such as market basket analysis, genetic analysis, social effectiveness analysis, infection prediction, big data analysis, chemical

O. Sharma (✉) · S. Kumar
Manav Rachna International University, Faridabad, Haryana, India
e-mail: Ochinsharma2@gmail.com

O. Sharma · N. Joshi
Banasthali Vidyapith, Nawai, Rajasthan, India

analysis, internet of things [1, 2]. Related to association rule mining, some important definitions are given below [2–6]:

**Definition 1** A rule is a representation of the association among different data variables. It can be represented as X->Y, where X is called antecedent part of a rule and Y is called consequent of a rule. X can be a single item or an itemset.

**Definition 2** An itemset is a set of attributes and value denoted as $\{(a_i, b_i), vi\}, \{(c_j, d_j), v_j)$, where $a_i, b_i, c_j, d_j$ are related attributes and v is the value of itemset.

**Definition 3** 'Support' interest measure is to find out how many times a set of items occur in the database.

In a rule X -> Y, Support (X->Y) = P (X) / TN, where TN is the total number of transactions in the database.

**Definition 4** 'Confidence' interest measure is to find out how many times consequent of a rule exists in comparable to rule's antecedent. Confidence (X->Y) = P (AB) / P (A).

**Definition 5** In a rule X->Y, 'Conviction' matches the possibility that X appears without Y.

Conviction (X->Y) = 1−P(Y) / 1−Confidence (XY).

**Definition 6** Leverage measures the difference of X and Y together appears in the data set and what is expected if X and Y are statistically dependent. Leverage (X->Y) = P (XY) – P (X) P (Y)

Consider a dataset of supermarket where thousands of items are listed may result into generation of thousands of thousands association rules. Fetching decision-making knowledge from these vast amounts of rules is a big challenge. To fetch the interesting rules, there must be some desirable constraint and this is actually known as interest measure. Interest measure plays a very imperative role to not only finding interesting rules but also ranking them [3, 7–10]. There exist various interest measures because of different business needs. We are briefing defining here a few well-accepted interest measures.

According to Lenca et al. [7], interest measures can be used to filter and sort the discovered rules. Interest measures play a crucial role in decreasing the number of generated rules and in collecting only the best ones. Since in the existed measures there is no optimal measure, so need either to compromise certain aspects or to develop new interest measures as per business needs. Azevedo and Alípio [4, 5, 11] emphasised that it would be worthwhile to combine different interest measures to produce ensembles of classifiers. In [12], Geng and Hamilton emphasised that interest measure plays a vital role to reduce the number of mined rules. In this paper, they have investigated more than 30 interest measures based on nine different principles given by Piatetsky-Shapiro [13], Tan et al. [11], Lenca et al. [7] and Hilderman and Howard [14] concluded that interest measure plays an important role in association rule mining, rule pruning and classifying important rules.

Kannan and Bhaskaran in paper [10] admitted that interest measure is quite useful in classification, prediction and clustering of data sets. In [5], Ochin et al. discussed the drawbacks of support, confidence, chi-square and lift measures to judge the significant rules and emphasised to develop new interest measures.

Different interest measures based on different business requirements have been proposed to fetch desired patterns. As discussed, many researchers have emphasised on the importance of interest measures and focussed to give continuous attention to develop new interest measures to fetch the important patterns for current business needs. Interest measures broadly have two categories, objective and subjective interest measures. Subjective interest measure can be applied where a prior knowledge and goals are well known. Objective measure is based on the probability and much suitable in many real-time scenarios where a priori knowledge is not available or not adequate [7, 8, 15]. In this paper, we are proposing a new objective interest measure named Significant Rule Power Factor.

## 2   Significant Rule Power Factor

The algorithm Significant Rule Power factor (SRPF) is being proposed in this section. Association rule generation is based on Apriori algorithm and is well explained in [3]. Further experiments conducted have proven the worth of the rules selected by SRPF, which are better when compared with many existed interest measures such as support, confidence, lift, coverage, leverage.

### 2.1   Algorithm_Calculate_SRPF ($L_k$, Minsrpf)

Input: Frequent Itemset $L_k$, SRPF minimum threshold Minsrpf
Output: Refined rules based on SRPF value
Description: This algorithm returns the set of rules that are frequent and has SRPF value >greater than the specified value.

1) Begin
2) For each (r in $L_k$) // frequent rules r with k length
3) r.SRPFvalue = (((P (AB) / P (A)) * P(AB)) / P(B)) // where P(AB) is the occurrence of item A and B both simultaneously in total transactions, P (A) is the occurrence of item A A in the total transactions and P(B) is the occurrence of item A in the total transactions. The range of SRPF is between 0 and 1.
4) Sk = {r ∈ $L_k$ | r.SRPFvalue > minSRPF}
5) End of For loop
6) return $S_k$
7) End of Calculate_SRPF algorithm

# 3    Analysis of Significant Rule Power Factor

Piatetsky-Shapiro has given three principles to judge the importance of an interest measure [13]. Table 1 gives an analysis of these three principles with context of SRPF.

By principle 1, if we choose values for P(A), P(Z), P(YZ) in such a manner that P (YZ) = P (Y) P (Z) then for this chosen P(Y), P(Z), P(YZ), interest measure value should be equal to 0. Please refer first column of Table 1, SRPF = 0. Hence, SRPF follows first principle of Piatetsky-Shapiro. Second principle says that interest measure monotonically increases with the value of P (YZ) when P (Y), P (Z) is kept same. So, we kept P (Y) =.2, P (Z) = .1 and increased P (YZ) from 0.05 to .1. SRPF also has shown increase in value from .11 to .50. Hence, SRPF follows Piatetsky's principle 2. Third principle says that if support of Y → Z and Z (or Y) are fixed, the smaller the support for Y (or Z), the more interesting is the pattern. Keeping P (YZ) and P (Z) same, by decreasing support of P (Y) from .2 to .1, the value of SRPF increases from .12 to .25. So, SRPF obeys principle 3 as well. Hence as SRPF obeys all the three principles given by Piatetsky-Shapiro, we can say according to him SRPF is an effective interest measure.

Further, Tan et al. proposed five properties to verify the importance of an interest measure [11].

O1: Under variable permutation, interest measure value is symmetric.
O2: When we scale any row or column by a positive factor, interest measure value is the same.

**Table 1**   Analysis of Piatetsky principle's in context of SRPF

| Principle 1 | Principle 2 | Principle 3 |
|---|---|---|
| The value of interest measure should be 0 provided both Y and Z are statistically independent and P (YZ) = P (Y) P Z) This principle depicts that the value of measure is nil when association among variables is merely a coincidental association. Take an example P (Y) = 20 / 100 = 0.2, P(Z) = 10 / 100 = 0.1, P (YZ) = 2 / 100 = 0.02, Conf (Y->Z) = .02 / .2 = 0.1, P (YZ) = P (Y) P (Z). 0.02 = 0.2 * .1 = 0.02, SRPF ((P (YZ) / P (Y)) * P (YZ)) / P(Z) = ((0.1 * 0.02) / .1 = 0.00 = 0 | Interest measure increases with P (YZ) when P (Y) and P (Z) are kept same P (Y) = 20 / 100 = 0.2, P (Z) = 10 / 100 = 0.1, P (YZ) = 5 / 100 = 0.05, conf = P (YZ) / P(Y) = 0.05 / .2 = 0.25 SRPF = (.25 * 0.05) / .1 = 0.12 Now, if P (YZ) = 10/100 = 0.1, P (Z) = 10/100 = 0.1, P (Y) = 20/100 = 0.2, P (YZ) / P (Y) = 0.1 / 0.2 = 0.50, SRPF = ((P (YZ) / P (Y)) * P (YZ)) / P (Z) = 0.5 | The third principle is about the case when support of Y → Z and Z (or Y) are fixed, the lesser the support for Y (or Z), the pattern is more interesting the [7] P (Y) = 20 / 100 = 0.2, P (Z) = 10 / 100 = 0.1, P (YZ) = 5 / 100 =0 .05, P(YZ) / P(Y) =0.05/0.2 = 0.25 SRPF = ((P (YZ) / P(Y)) * P (YZ)) / P (Z) = ((.25 * 0.05) / 0.1) = 0.12 If P (YZ) = 5 / 100 = 0.05, P (Z) = 10 / 100 = 0.1, P (Y) = 10 / 100 = 0.1, SRPF = ((0.05/0.1)* 0.05) / 0.1) = 0.25 |

**Table 2** Overview of different interest measures on the basis of Piatetsky-Shapiro's and Tan's principles

| Principle | Support | Confidence | Lift | Coverage | Leverage | Conviction | SRPF |
|---|---|---|---|---|---|---|---|
| P1 | X | X | X | X | X | X | $\sqrt{}$ |
| P2 | $\sqrt{}$ | $\sqrt{}$ | $\sqrt{}$ | X | $\sqrt{}$ | $\sqrt{}$ | $\sqrt{}$ |
| P3 | X | X | $\sqrt{}$ | X | $\sqrt{}$ | X | $\sqrt{}$ |
| O1 | $\sqrt{}$ | X | $\sqrt{}$ | X | X | X | $\sqrt{}$ |
| O2 | X | X | X | X | X | $\sqrt{}$ | X |
| O3 | X | X | X | X | X | X | X |
| O4 | X | X | X | X | X | X | X |
| O5 | X | X | X | X | $\sqrt{}$ | $\sqrt{}$ | $\sqrt{}$ |

O3: Under row/column permutation, interest measure value becomes equal to negation of its value.
O4: Under both row and column permutations, interest measure value remains the same.
O5: With the records that do not contain A and B, interest measure should has no relationship.

In [11], Tan et al. suggested that a measure should obey five principles. Unlike Piatetsky-Shapiro's principles, these principles should not be considered as essential for interest measures. They are useful to classify the measures into different groups [12]. Many principles are similar to Piatetsky-Shapiro principles and others are not directly applicable to association rule mining because these consider that the count attribute values are in sorted order and such accumulating is less instinctive for contingency tables. Hence, no rule interest measure could obey completely the five principles.

In Table 2, we could observe the importance of various interest measures on the basis of Shapiro's and Tan's principles. Support measure follows Shapiro's principle P1 and Tan's principle O1. Confidence measure obeys only P2 principle. Lift measure follows P2, P3, O1. Leverage measure obeys P2, P3, O5. Conviction follows P2, O2, O5, and SRPF measure follows P1, P2, P3, P4 and O5. So, it is clear that SRPF has qualified with more edge to many of the existing well-accepted interest measures.

# 4 Experimental Observation and Comparison with Existing Measures

Weka is a well-known data mining tool. Weka (version 3.7.13) is used to conduct experiments. We have used various data sets on Windows 7 OS, Intel core i3, 4 GB RAM, processor 2.40 GHz. To edit WEKA source code, JEDIT software is being used and to recompile the Weka source code .jar file, ANT tool is used. All of these

tools are open source tools. Figure 1 shows recompiling of WEKA .jar files using ANT, and Fig. 2 displays successful building of .jar file again.

In Fig. 3, an experiment is shown with 10 best instances with weather data set with >90% confidence.

Consider Table 3, are rules 1 to 10, confidence is same for all the rules. So confidence is not helpful to identify important rules. Lift is the superset of confidence but observe rule 2 and 6, lift is actually misguiding. In rule 2,



**Fig. 1** Executing Weka .Jar file through ant



**Fig. 2** Successful build compiled



**Fig. 3** Weka experiment

**Table 3** Analysis on the basis of values of different interest measures for 10 best rules

| Rules | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| Confidence (Conf) | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| Lift | 2.8 | 1.56 | 1.17 | 2.8 | 2.33 | 1.56 | 2.8 | 1.17 | 1.75 | 3.5 |
| SRPF | 0.6 | 0.33 | 0.25 | 0.4 | 0.33 | 0.22 | 0.4 | 0.17 | 0.25 | 0.5 |
| Leverage (Lev) | 0.14 | 0.08 | 0.03 | 0.09 | 0.08 | 0.05 | 0.09 | 0.02 | 0.06 | 0.10 |
| Cosine (Cos) | 0.80 | 0.60 | 0.52 | 0.65 | 0.60 | 0.49 | 0.65 | 0.42 | 0.52 | 0.73 |
| Coverage (Cov) | 0.21 | 0.21 | 0.29 | 0.14 | 0.21 | 0.14 | 0.14 | 0.21 | 0.14 | 0.14 |
| Added value (AV) | 0.64 | 0.36 | 0.36 | 0.64 | 0.71 | 0.36 | 0.64 | 0.43 | 0.43 | 0.71 |



**Fig. 4** Rule ranking comparison through different measures on weather dataset

P (AB) = 3/14=.21 (this data set contains 14 total row instances). In this, 21% of both antecedent and consequent is true, and in rule 6, P (AB) =.14, so 14% instances in data set are following this rule, but surprisingly lift says both rules are equally good.

Whereas SRPF properly projected rule 2 is high ranked as compared to rule 6. Similar case is being observed with rules 1, 4, 7 and rules 3, 9.

'Support' and 'Coverage' measures are good measures to find out the frequent itemsets. Confidence is a measure that focuses to fetch important rules. Leverage focuses to find out the difference between the occurrence of item A and B together in actual and that of expected. So, leverage prospective is different from SRPF. Same

with Conviction, it compares the probability that item A appears without item B. Lift is a measure that try to fetch important rules, by giving adequately importance to all the components of a rule.

Figure 4 represents that SRPF is able to fetch more important rules than lift, leverage, cosine, coverage and added value.

## 5    Conclusion and Future Work

Association rule mining is an important technique for pattern recognition. It helps not only in establishing relationship among attributes but also while taking effective decisions based on these associations. Through our work, we tried to contribute in this field to extract more valuable patterns. These patterns will be more helpful to the experts to come up with better decisions and with no overhead. Regarding future work, SRPF can be used in associative classification technique and can classify the attributes to contribute in direct decision-making process.

## References

1. Bayardo Jr, R.J., Agrawal, R.: Mining the most interesting rules. In: Proceedings of the Fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. ACM (1999)
2. Chen, Y., Zhang, A.H.C.: Research on data mining model in the internet of things. In: International Conference on Automation, Mechanical Control and Computational Engineering, Atlantis (2015)
3. Agrawal, R., Srikant, R.: Fast algorithms for mining association rules. In: Proceedings of 20th International Conference on Very Large Data Bases, VLDB, vol. 1215 (1994)
4. Azevedo, P.J., Jorge, A.M.: Comparing rule measures for predictive association rules. In: European Conference on Machine Learning. Springer, Berlin (2007)
5. Ochin, K.S., Nisheeth, J.: Rule power factor: a new interest measure in associative classification. Proc. Comput. Sci. **93**, 12–18 (2016)
6. Merceron, A., Yacef, K.: Interestingness measures for association rules in educational data. In Educational Data Mining (2008)
7. Lenca, P., et al.: Association Rule Interestingness Measures: Experimental and Theoretical Studies. Quality Measures in Data Mining, pp. 51–76. Springer, Berlin (2007)
8. Shaharanee, I.N.M., Hadzic, F., Dillon, T.S.: Interestingness measures for association rules based on statistical validity. Knowl. Based Syst. **24**(3), 386–392 (2011)
9. Brin, S., et al.: Dynamic itemset counting and implication rules for market basket data. ACM SIGMOD Record, vol. 26, no. 2. ACM (1997)
10. Kannan, S., Bhaskaran, R.: Association rule pruning based on interestingness measures with clustering. arXiv preprint arXiv:0912.1822 (2009)
11. Tan, P.-N., Kumar, V., Srivastava, J.: Selecting the right interestingness measure for association patterns. In: Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, ACM (2002)
12. Geng, L., Hamilton, H.J.: Interestingness measures for data mining: a survey. ACM Comput. Surv. (CSUR) **38**(3), 9 (2006)

13. Piatetsky-Shapiro, G. Discovery, analysis, and presentation of strong rules. Knowl. Discov. Databases, 229–238 (1991)
14. Hilderman, R., Hamilton, H.J.: Knowledge Discovery and Measures of Interest, vol. 638. Springer, Berlin (2013)
15. Freitas, A.A.: On rule interestingness measures. Knowl. Based Syst. **12**(5), 309–315 (1999)

# Covering-Based Pessimistic Multigranular Approximate Rough Equivalences and Approximate Reasoning

**B. K. Tripathy and Suvendu K. Parida**

**Abstract** The multigranular rough set (MGRS) models of Qian et al. were extended to put forth covering-based multigranular rough sets (CBMGRS) by Liu et al. in 2012. The equality of sets, which is restrictive and redundant, was extended first in Pawlak (Rough sets: theoretical aspects of reasoning about data. Kluwer, London, 1991) and subsequently in Tripathy (Int J Adv Sci Technol 31:23–36, 2011) to propose four types of rough set-based approximate equalities. These basic concepts of rough equalities have been extended to several generalized rough set models. In this paper, covering-based pessimistic multigranular (CBPMG) approximate rough equivalence is introduced and several of their properties are established. Real life examples are taken for constructing counter examples and also for illustration. We have also discussed how these equalities can be applied in approximate reasoning and our latest proposal is no exception.

## 1 Introduction

The rough set model introduced [4, 5] has been extended in many ways [19] and there are 12 types of CBRSs in the literature [8]. The basic uni-granular rough sets [25] were extended for multiple granulations in the form of optimistic and pessimistic MGRS [6, 7] in 2006 and 2010, respectively. This has further been extended to the setting of CBRS and four types of CBMGRS are defined and

---

B. K. Tripathy (✉)
SCOPE, VIT University, Vellore, Tamilnadu, India
e-mail: tripathybk@vit.ac.in

S. K. Parida
SCS Autonomous College, Puri, Odisha 752001, India
e-mail: paridasuvendukumar@gmail.com

studied in [1, 2]. Approximate equality of concepts using rough sets was introduced in [5] when they defined three types of equalities through rough sets, which automatically take care of user knowledge. It was generalized by Tripathy et al. [11] by introducing rough equivalence of sets [11, 20]. Two more equalities were proposed in [9, 10, 17] for completion of the family of these approximate equalities. In this paper, we use pessimistic multigranular rough sets to define approximate rough equivalence, establish their properties, and provide examples from real life for illustration and counter examples in proving properties. Also, we show how approximate reasoning can be handled through this model following its properties. It may be noted that the corresponding notion for optimistic case has already been studied [12–15, 21–24]. As far as organization of the paper is concerned, the concepts and notations to be used are presented in the next section. Section 3 presents MGRS followed by CBMGRS in Sect. 4. In Sect. 5, properties of pessimistic multigranular rough equivalence are established. We outline some ideas for application of the results to approximate reasoning in Sect. 6. This is followed by conclusion in Sect. 7 and references.

## 2  Definitions and Notations

Some of the definitions used here are as follows. The definition of rough set with standard notations can be found in [5].

The CBRS definition is similar to that of basic rough sets by replacing the equivalence classes with elements of a cover. The following two notions are to be used in defining CBMGRS.

**Definition 2** Let $(V, E)$ be a covering approximation space. Then for any $v \in V$, sets $md_E(v)$ and $MD_E(v)$ are respectively min and max descriptors of vin $E$, defined as

$$md_E(v) = \{K \in E | v \in K \text{ and } \forall L \in E : (v \in L \text{ and } L \subseteq K) \text{ we have } K = L\} \quad (1)$$

$$MD_E(v) = \{K \in E | v \in K \text{ and } \forall L \in E : (v \in L \text{ and } L \supseteq K) \text{ we have } K = L\} \quad (2)$$

## 3  MGRS

Here, we introduce the pessimistic MGRS which we used in defining CBMGRS. Qian et al. introduced this concept in [20]. We use two equivalence relations at a time.

**Definition 3** Let $S$ be a family of equivalence relations, $K$, $L$ be in $S$. Then the pessimistic MG lower and upper approximation of $Z$ in $V$ as

$$\underline{K * L}(Z) = \cup \left\{ v/[v]_K \subseteq Z \,\&\, [v]_L \subseteq Z \right\} \tag{3}$$

$$\overline{K * L}(Z) = \left( \underline{K * L}(Z^C) \right)^C \tag{4}$$

If $\underline{K * L}(Z) \neq \overline{K * L}(Z)$ then $Z$ is said to be pessimistic MG rough for $K$ and $L$. Otherwise, $Z$ is said to be pessimistic MG definable for $K$ and $L$.

## 4 CBMGRS

Extending the notion of MGRS, 4 types of optimistic and pessimistic CBMGRSs were studied in [1–3] using (1) and (2) as follows.

Let $E_1$ and $E_2$ be any two covers defined over $V$. The lower and upper approximations for CBPMGRS of different types with respect to $E_1$ and $E_2$ are defined below.

**Definition 4** The first type is given by (5) and (6).

$$\underline{FR_{E_{1*}E_2}}(Z) = \{ v \in V | \cap md_{E_1}(v) \subseteq Z \text{ and } \cap md_{E_2}(v) \subseteq Z \}; \tag{5}$$

$$\overline{FR_{E_{1*}E_2}}(Z) = \{ v \in V | \cap md_{E_1}(v) \cap \neq \phi \text{ or } (\cap md_{E_2}(v)) \cap Z \neq \phi \} \tag{6}$$

**Definition 5** The second type is given by (7) and (8).

$$\underline{SR_{E_{1*}E_2}}(Z) = \{ v \in V | \cup md_{E_1}(v) \subseteq Z \text{ and } \cup md_{E_2}(v) \subseteq Z \} \tag{7}$$

$$\overline{SR_{E_{1*}E_2}}(Z) = \{ v \in V | \cup md_{E_1}(v) \cap Z \text{ or} (\cup md_{E_2}(v)) \cap Z \neq \phi \} \tag{8}$$

**Definition 6** The third type is given by (9) and (10).

$$\underline{TR_{E_{1*}E_2}}(Z) = \{ v \in V | \cap MD_{E_1}(v) \subseteq Z \text{ and } \cap MD_{E_2}(v) \subseteq Z \}; \tag{9}$$

$$\overline{TR_{E_{1*}E_2}}(Z) = \{ v \in V | \cap MD_{E_1}(v) \cap Z \neq \phi \text{ or} (\cap MD_{E_2}(v)) \cap Z \neq \phi \} \tag{10}$$

**Definition 7** The fourth type is given by (11) and (12).

$$\underline{LR_{E_{1*}E_2}}(Z) = \{ v \in V | \cap MD_{E_1}(v) \subseteq Z \text{ and } \cup MD_{E_2}(v) \subseteq Z \}; \tag{11}$$

$$\overline{LR_{E_{1*}E_2}}(Z) = \{ v \in V | (\cup MD_{E_1}(v)) \cap Z \neq \phi \text{ or } (\cup MD_{E_2}(v)) \cap Z \neq \phi \} \tag{12}$$

CBPMGRSs satisfy the following properties with respect to union and inter-section of sets. All the four types have same properties. Hence, we reproduce these properties here only for the first type.

$$\underline{FR_{E_{1*}E_2}}(Z \cap W) = \underline{FR_{E_{1*}E_2}}(Z) \cap \underline{FR_{E_{1*}E_2}}(W) \tag{13}$$

$$\underline{FR_{E_{1*}E_2}}(Z \cup W) \supseteq \underline{FR_{E_{1*}E_2}}(Z) \cup \underline{FR_{E_{1*}E_2}}(W) \tag{14}$$

$$\overline{FR_{E_{1*}E_2}}(Z \cup W) \supseteq \overline{FR_{E_{1*}E_2}}(Z) \cup \overline{FR_{E_{1*}E_2}}(W) \tag{15}$$

$$\overline{FR_{E_{1*}E_2}}(Z \cap W) \supseteq \overline{FR_{E_{1*}E_2}}(Z) \cap \overline{FR_{E_{1*}E_2}}(W) \tag{16}$$

**A Real life Example**

Table 1 contains the information of faculty in a department of a college out of which an examination committee is to be formed.

$E_1$ = It is the cover formed by the relation that the group should have 2 years mean experience for collection and the size of male faculty is more than that of female faculty. Then,

$$V/E_1 = \{\{v_1, v_2, v_3\}, \{v_3, v_4\}, \{v_3, v_8\}, \{v_4, v_5\}, \{v_6, v_7\}\}.$$

$E_2$ = It is the cover formed by the relation that the group should have 2 years mean experience for distribution and the size of male faculty should be at least equal to that of female faculty. Then

$$V/E_2 = \{\{v_1, v_2, v_3\}, \{v_3, v_4\}, \{v_3, v_5\}, \{v_4, v_8\}, \{v_6, v_7\}\}.$$

The CBPMGRSs with respect to $E_1$ and $E_2$ can be defined through their minimal descriptors as follows:

For example, there are two minimum descriptors for the element $v_3$ in $V/E_2$, namely $\{v_3, v_4\}$ and $\{v_3, v_5\}$. So, the entry in place (2, 4) is their intersection $\{v_3\}$. From Table 2, the lower and upper covering-based pessimistic multigranular approximations can be obtained using (5)–(12) depending upon the type of such approximations required.

**Table 1** Faculty experience details

| S. no. | Faculty name | Collection experience (years) | Distribution experience (years) | Sex |
|---|---|---|---|---|
| 1 | Allen-vl | 1 | 2 | Male |
| 2 | Brinda-v2 | 2 | 1 | Female |
| 3 | Celina-v3 | 4 | 4 | Male |
| 4 | Danya-v4 | 2 | 3 | Female |
| 5 | Ershad-v5 | 2 | 2 | Male |
| 6 | Feroz-v6 | 3 | 2 | Male |
| 7 | Geeta-v7 | 2 | 3 | Male |
| 8 | Harsha-v8 | 1 | 2 | Male |

**Table 2** Minimum descriptions of elements

| Elements $(v)$ | $v_1$ | $v_2$ | $v_3$ | $v_4$ | $v_5$ | $v_6$ | $v_7$ | $v_8$ |
|---|---|---|---|---|---|---|---|---|
| $\cap md_{E_1}(v)$ | $\{v_1, v_2, v_3\}$ | $\{v_1, v_2, v_3\}$ | $\{v_3\}$ | $\{v_4, v_5\}$ | $\{v_4, v_5\}$ | $\{v_6, v_7\}$ | $\{v_6, v_7\}$ | $\{v_3, v_8\}$ |
| $\cap md_{E_2}(v)$ | $\{v_1, v_2, v_3\}$ | $\{v_1, v_2, v_3\}$ | $\{v_3\}$ | $\{v_4\}$ | $\{v_3, v_5\}$ | $\{v_6, v_7\}$ | $\{v_6, v_7\}$ | $\{v_4, v_8\}$ |

## 5 CVPMG Approximate Equivalences

We now introduce in the following different rough equalities for the first type $(F)$ of CBPMGRS and study their properties. The definitions for other 3 types are similar.

Let $E_1$ and $E_2$ be two covers on $V$, $E_1, E_2 \in E$ and $Z, W \subseteq V$.

**Definition 8** We say that $Z$ and $W$ are CB Pessimistic bottom approximate rough equivalent for $E_1$ and $E_2$ $(Zb\_E_1 * E_2\_aeqv \ W)$ if and only if

$$\underline{FR_{c_{1*}c_2}}(X) \quad \text{and} \quad \underline{FR_{E_{1*}E_2}}(W) \quad \text{are } \phi \text{ or not } \phi \text{ together.} \tag{17}$$

$Z$ and $W$ are CB Pessimistic top approximate rough equivalent for $E_1$ and $E_2$ $(Z \ t\_E_1 * E_2\_aeqv \ W)$ iff

$$\overline{FR_{E_1 * E_2}}(Z) = \overline{FR_{E_1 * E_2}}(W) \tag{18}$$

$Z$ and $W$ are CB Pessimistic total approximate rough equivalent for $E_1$ and $E_2$ $(Z \ E_1 * E_2\_aeqv \ W$ if and only if $\underline{FR_{E_1 * E_2}}(Z) \ and \ \underline{FR_{E_1 * E_2}}(W)$ are $\phi$ or not $\phi$ together and

$$\overline{FR_{E_1 * E_2}}(Z) = \overline{FR_{E_1 * E_2}}(W) \tag{19}$$

### 5.1 Properties of CBPMGRS

As stated above, the general properties of first type of covering-based rough equivalences are established. Similar properties hold for other types also.

Example: We shall use the following example to illustrate or provide counter examples in proving the properties.

Let $V = \{v_1, v_2, v_3, v_4, v_5, v_6, v_7, v_8\}$. We define two covers $E_1$ and $E_2$ as follows.

$$V/E_1 = \{\{v_1, v_5\}, \{v_3, v_4, v_5, v_6\}, \{v_2, v_7, v_8\}\}.$$

and

$$V/E_1 = \{\{v_1, v_6\}, \{v_2, v_3, v_4\}, \{v_5, v_6\}, \{v_7, v_8\}\}.$$

**Table 3** The minimal descriptions

| Elements (V) | $v_1$ | $v_2$ | $v_3$ | $v_4$ | $v_5$ | $v_6$ | $v_7$ | $v_8$ |
|---|---|---|---|---|---|---|---|---|
| $\cap md_{E_1}(v)$ | $\{v_1, v_5\}$ | $\{v_2, v_7, v_8\}$ | $\{v_3, v_4, v_5, v_6\}$ | $\{v_3, v_4, v_5, v_6\}$ | $\{v_5\}$ | $\{v_3, v_4, v_5, v_6\}$ | $\{v_2, v_7, v_8\}$ | $\{v_2, v_7, v_8\}$ |
| $\cap md_{E_2}(v)$ | $\{v_1, v_6\}$ | $\{v_2, v_3, v_4\}$ | $\{v_2, v_3, v_4\}$ | $\{v_2, v_3, v_4\}$ | $\{v_5, v_6\}$ | $\{v_6\}$ | $\{v_7, v_8\}$ | $\{v_7, v_8\}$ |

The minimum descriptions for the different elements are tabulated below. (Table 3)

Let $E_1$ and $E_2$ be two covers on $V$ and $E_1, E_2 \in Z$ and $Z, W \subseteq V$. Let $F$ denote the first type of CBMGRS. Then,

**Property 1**

$$Zb\_E_1 * E_2\_aeqv \, W \text{ if } Z \cap W \, b \_E_1 * E_2\_aeqvZ \text{ and Wboth.}$$

**Proof: (If part)**

$Z\_bE_1 * E_2\_aeqv \, W \Rightarrow \underline{F_{E_1*E_2}}(Z \cap W) \text{ and } \underline{F_{E_1*E_2}}(Z) \text{ are } \phi \text{ or not } \phi \text{ together and } Z \cap W\_bE_1 * E_2\_aeqv \, \overline{W}.$

$$\Rightarrow \underline{F_{E_1*E_2}}(Z \cap W) \text{ and } \underline{F_{E_1*E_2}}(W) \text{ are } \phi \text{ or not } \phi \text{ together.}$$

So, when $Z \cap W \, b\_E_1 * E_2\_$aeqv $Z$ and $W$ both $\underline{F_{E_1*E_2}}(Z \cap W)$ being a common factor it implies that $\underline{F_{E_1*E_2}}(Z)$ and $\underline{F_{E_1*E_2}}(W)$ are $\phi$ or not $\phi$ togetherand so $Zb\_E_1 * E_2\_aeqv \, W$.

(Only if part) From (15) we get $\underline{F_{E_1*E_2}}(Z \cap W) = \underline{F_{E_1*E_2}}(Z) \cap \underline{F_{E_1*E_2}}(W)$

So, $Z \, b\_E_1 * E_2\_aeqv \, W \Leftrightarrow \underline{F_{E_1*E_2}}(Z)$ and $\underline{F_{E_1*E_2}}(W)$ are $\phi$ or not $\phi$ together. If both are $\phi$ then the conclusion holds. But if both are not $\phi$, then we cannot be sure that $\underline{F_{E_1*E_2}}(Z \cap W)$ is not $\phi$. So, the conclusion may not be true. We provide the following example.

In the example, let us take $Z = \{v_1, v_5, v_6\}$ and $W = \{v_2, v_7, v_8\}$. Then $\underline{F_{E_1*E_2}}(Z) = \{v_1, v_5\} \neq \phi$, $\underline{F_{E_1*E_2}}(W) = \{v_7, v_8\} \neq \phi$. But, $Z \cap W = \phi$. So, $\underline{F_{E_1*E_2}}(Z \cap W) = \phi$.

**Property 2** $Z \, t\_E_1 * E_2\_aeqv \, W$ if and only if $Z \cup W \, t\_E_1 * E_2\_aeqv \, W$ and $Z$ both.

***Proof*: (If part)**

$Z\,t\_E_1 * E_2\_aeqv\,W \Rightarrow \overline{F_{E_1*E_2}}(Z \cup W) = \overline{F_{E_1*E_2}}(Z)$
$and\,Z \cup W\,t\_E_1 * E_2\_aeqv\,W \Rightarrow \overline{F_{E_1*E_2}}(Z \cup W) = \overline{F_{E_1*E_2}}(W).$ So, $Z \cup W\,t\_E_1 * E_2\_aeqv\,Z$
$and\,W$ both $\Rightarrow \overline{F_{E_1*E_2}}(Z) = \overline{F_{E_1*E_2}}(W) \Rightarrow Z\,t\_E_1 * E_2\_aeqv\,W.$

(Only if part)

From (15) we have $\overline{F_{E_1*E_2}}(Z \cup W) = \overline{F_{E_1*E_2}}(Z) \cup \overline{F_{E_1*E_2}}(W).$ So, we have

$$\overline{F_{E_1*E_2}}(Z \cup W) = \overline{F_{E_1*E_2}}(Z) \cup \overline{F_{E_1*E_2}}(W). \quad So, Z\,t\_E_1 * E_2\_aeqv\,W$$
$$\Leftrightarrow \overline{F_{E_1*E_2}}(Z \cup W) = \overline{F_{E_1*E_2}}(Z) = \overline{F_{E_1*E_2}}(W).$$

This completes proof.

**Property** **3** $Z\,t\_E_1 * E_2\_aeqv\,Z'$ and $Z\,t\_E_1 * E_2\_aeqv\,W' \Rightarrow (Z \cup W)\,t\_E_1 * E_2\_aeqv\,(Z' \cup W')$

***Proof*** Using the proof of the above property and the definition of top covering-based multigranular equality, we have

$$\overline{F_{E_1*E_2}}(Z \cup W) = \overline{F_{E_1*E_2}}(Z) \cup \overline{F_{E_1*E_2}}(W) = \overline{F_{E_1*E_2}}(Z') \cup \overline{F_{C_1*C_2}}(W') = \overline{F_{E_1*E_2}}(Z' \cup W')$$
$$So,\ Z \cup W\,t\_E_1 * E_2\_aeqv\,Z' \cap W'$$

**Property 4** $Z\,b\_E_1 * E_2\_aeqv\,Z'$ and $W\,b\_E_1 * E_2\_aeqv\,Z'$ may not imply that $Z \cap W\,b\_W\_aeqv\,Z' \cap W'.$

***Proof*** From hypothesis, both $\underline{F_{E_1*E_2}}(Z)$ and $\underline{F_{E_1*E_2}}(Z')$ are $\phi$ or not $\phi$ together. Also, both $\underline{F_{E_1*E_2}}(W)$ and $\underline{F_{E_1*E_2}}(W')$ are $\phi$ or not $\phi$ together. It is known that

$$\underline{F_{E_1*E_2}}(W \cap Z) = \underline{F_{E_1*E_2}}(Z) \cap \underline{F_{E_1*E_2}}(W) \text{ and } \underline{F_{E_1*E_2}}(Z') \cap \underline{F_{E_1*E_2}}(W')$$
$$= \underline{F_{E_1*E_2}}(Z' \cap W').$$

So, when at least one of the pairs $\underline{F_{E_1*E_2}}(Z)$ and $\underline{F_{E_1*E_2}}(Z')$ and $\underline{F_{E_1*E_2}}(W)$ and $\underline{F_{E_1*E_2}}(W')$ are $\phi$ together, it is clear that

$\underline{F_{E_1*E_2}}(Z \cap W)$ and $\underline{F_{E_1*E_2}}(Z' \cap W')$ are both $\phi$ and hence $Z \cap W\,b\_E_1 * E_2\_aeqv\,Z' \cap W'.$

But when both the pairs are not $\phi$, then it may happen that one of $\underline{F_{E_1*E_2}}(Z \cap W)$ and $\underline{F_{E_1*E_2}}(Z' \cap W')$ is $\phi$ whereas the other one is not so. So, in that case we cannot say that $Z \cap W\,b\_E_1 * E_2\_aeqv\,Z' \cap W'.$ For this, we provide an example. In the example, let us take $Z = \{v_1, v_5, v_6\}$ and $W = \{v_2, v_7, v_8\}.$ Let us take $Z' = \{v_2, v_5, v_6, v_7, v_8\}$ and $W' = \{v_2, v_7, v_8\}.$ Then, $\underline{F_{E_1*E_2}}(Z) = \{v_1, v_5\} \neq \phi, \underline{F_{E_1*E_2}}(W) = \{v_7, v_8\} \neq \phi.$

$\underline{F_{E_1*E_2}}(Z') = \{v_5, v_7, v_8\} \neq \phi$ and $\underline{F_{E_1*E_2}}(W') = \{v_7, v_8\} \neq \phi.$ Now $Z \cap W = \phi,$

So, $\underline{F_{E_1*E_2}}(Z \cap W) = \phi$. Again $Z \cap W = \{v_2, v_7, v_8\}$. So, $\underline{F_{E_1*E_2}}(Z' \cap W') = \{v_7, v_8\} \neq \phi$. Hence, $Z\,b\_E_1*E_2\_aeqv\,Z'$ and $W\,b\_E_1*E_2\_aeqv\,W'$. But, $Z \cap W$ is not $b\_E_1*E_2\_aeqv\,Z' \cap W'$.

**Property 5** $Z\,t\_E_1*E_2\_aeqv\,Z \cup W^c\,t\_E_1*E_2\_aeqv\,V$

**Proof** Given $Z\,t\_E_1*E_2\_aeqv\,W \Rightarrow \overline{F_{E_1*E_2}}(Z) = \overline{F_{E_1*E_2}}(W)$.
Using (15), we get

$$\overline{F_{E_1*E_2}}(Z \cup W^c) = \overline{F_{E_1*E_2}}(Z) \cup \overline{F_{E_1*E_2}}(W^c) = \overline{F_{E_1*E_2}}(W) \cup \left(\underline{F_{E_1*E_2}}(Z)\right)^c$$
$$= \overline{F_{E_1*E_2}}(W) \cup \left(\overline{F_{E_1*E_2}}W \cap (BN_{E_1*E_2}W)^c\right)^C$$
$$= \overline{F_{E_1*E_2}}(W) \cup \left(\overline{F_{E_1*E_2}}(W)\right)^C \cup (BN_{E_1*E_2}W) = V \Rightarrow Z \cup W^c\,t\_E_1*E_2\_aeqv\,V$$

**Property 6** $Z\,b\_E_1*E_2\_aeqv\,W$ may not imply that $Z \cap W^c\,t\_E_1*E_2\_aeqv\,\phi$

**Proof** To establish this, we provide an example.
Let $Z = \{v_1, v_5, v_6\}$ and $W = \{v_2, v_7, v_8\}$. Then $\underline{F_{E_1*E_2}}(Z) = \{v_1, v_5\} \neq \phi$, and $\underline{F_{E_1*E_2}}(W) = \{v_7, v_8\} \neq \phi$. So, $Z\,b\_E_1*E_2\_aeqv\,W$. But,

$$Z \cap W^c = \{v_1, v_5, v_6\} \cap \{v_1, v_3, v_4, v_5, v_6,\} = \{v_1, v_5, v_6\}. \text{ Hence,}$$
$$\overline{F_{E_1*E_2}}(Z \cap W^c) = \{v_1, v_3, v_4, v_5, v_6,\} \neq \phi.$$

So, $Z \cap W^c$ is not $t\_E_1*E_2\_aeqv\,\phi$.

**Property 7** If $Z \subseteq W$ and $W\,t\_E_1*E_2\_aeqv\,\phi$ then $Z\,t\_E_1*E_2\_aeqv\,\phi$

**Proof** Given $Z \subseteq W$ and $W\,t\_E_1*E_2\_aeqv\,\phi$. So we have $\overline{F_{E_1*E_2}}(Y) = \phi$. As $Z \subseteq W \Rightarrow Z = \phi \Rightarrow \overline{F_{E_1*E_2}}(Z) = \phi \Rightarrow Z\,t\_E_1*E_2\_aeqv\,\phi$

The following properties follow directly from the definitions and hence we only state them.

**Property 8** If $Z \subseteq W$ and $Z\,t\_E_1*E_2\_aeqv\,V$ then $W\,t\_E_1*E_2\_aeqv\,V$

**Property 9** $Z\,t\_E_1*E_2\_aeqv\,W$ iff $Z^C\,b\_E_1*E_2\_aeqv\,W^C$

**Property 10** If $Z\,b\_E_1*E_2\_aeqv\,\phi$ or $W\,b\_E_1*E_2\_aeqv\,\phi$ then $Z \cap W\,b\_E_1*E_2\_aeqv\,\phi$

**Property 11** If $Z\,t\_E_1*E_2\_aeqv\,V$ or $W\,t\_E_1*E_2\_aeqv\,V$ then $Z \cup W\,t\_E_1*E_2\_aeqv\,V$

# 6 Approximate Reasoning

Human reasoning process is mostly approximate by nature. In a world of uncertainties, it is foolish to think of exact reasoning being applicable. This varies from day-to-day activities such as characterizing two images to be identical to intellectuals being considered as having same level of intelligence. We allow uncertainty to creep into our activities to make it more feasible and robust. However, due to relative scarcity of equivalence, relations covers are used instead of partitions. Similarly, multiple granularity is sometimes more useful than individual ones. So, covering-based multigranulation is a more general and common process in human cognition. Approximate rough equivalence was introduced as an extension of mathematical equality and is a tool to deal with equality of concepts using human knowledge. As a consequence, covering-based multigranular rough sets can be used for approximate equality at a higher level.

The properties of approximate rough equivalence established here are helpful in partial or complete equality of objects/concepts with respect to multiple granulations of the universe. The pessimistic version considered here is a restrictive one than the optimistic one, but provides less ambiguity than it.

To illustrate this further, let us consider Property 10. It says that if there is no minimal covering element $v$ in $Z$ for covers $E_1$ and $E_2$ separately and there is no minimal covering element w in $W$ for covers $E_1$ and $E_2$ separately, then there exists no minimal covering element z in $Z \cap W$ with respect to covers $E_1$ and $E_2$ separately.

# 7 Conclusion

We introduced two new notions, CBPMG approximate rough equivalence of sets, which a generalization of the notion of CB approximate rough equivalence of sets, which in turn is an extension of approximate rough equivalence of sets further generalizing the notions of normal equality of sets. We established several of the properties of the lower and upper approximations of these approximate rough equivalences. These properties are useful in their application on real data sets as is evident from the example above. Some examples are provided to illustrate the computation and also as counter examples to establish negative results. It is evident that these notions being approximate by nature and dealing with equality of sets/concepts can be a useful tool in approximate reasoning where multiple granular structures are prevalent.

# References

1. Lin, G.P., Qian, Y.H, Li, J.J.: A covering-based pessimistic multi-granulation rough set. In: Proceedings of International Conference on Intelligent Computing, 11–14 August, Zhengzhou, China (2011)
2. Liu, C.H., Miao, D.Q.: Covering rough set model based on multi-granulations. In: Proceedings of Thirteenth International Conference on Rough Sets, Fuzzy Set, Data Mining and Granular Computing, LNCS (LNAI) 6743, pp. 87–90 (2011)
3. Liu, C.L., Miao, D., Quain, J.: On multi-granulation covering rough sets. Int. J. Approx. Reason **55**, 1404–1418 (2014)
4. Pawlak, Z.: Rough sets. Int. J. Comput. Inf. Sci. **11**, 341–356 (1982)
5. Pawlak, Z.: Rough Sets: Theoretical Aspects of Reasoning About Data. Kluwer, London (1991)
6. Qian, Y.H., Liang, J.Y., Dang, C.Y.: Pessimistic rough decision. In: Proceedings of RST 2010, Zhou Shan, China, pp. 440–449 (2010)
7. Qian, Y.H., Liang, J.Y: Rough set method based on multi-granulations. In: Proceedings of the 5th IEEE Conference on Cognitive Informatics, vol. 1, pp. 297–304 (2006)
8. Safari, S., Hooshmandasl, M.R.: On twelve types of covering-based rough sets. Springer Plus **5**(1003), 1–18 (2016)
9. Tripathy, B.K.: An analysis of approximate equalities based on rough set theory. Int. J. Adv. Sci. Technol. **31**, 23–36 (2011)
10. Tripathy, B.K., Jhawar, A., Vats, E.: An analysis of generalised approximate equalities based on rough fuzzy sets. In: Proceedings of the International Conferences on SocPros 2011, 22–24 December, IIT-Roorkee, AISC 130, Springer India, vol. 1, pp. 333–346 (2012)
11. Tripathy, B.K., Mitra, A., Ojha, J.: On Rough Equalities and Rough Equivalence of Sets, RSCTC 2008-Akron, U.S.A., Springer, LNAI 5306, pp. 92–102 (2008)
12. Tripathy, B.K., Mitra, A.: On the approximate equalities of multigranular rough sets and approximate reasoning. In: 4th IEEE Conference on Computing, Communication and Network Technologies, ICCCNT 2013, Elayapalayam (2013). doi:10.1109/ICCCNT.2013. 6726771
13. Tripathy, B.K., Mohanty, R.R.: Covering based pessimistic multigranular approximate rough equalities and their properties. Accepted for publication in the Int. J. Rough Sets Data Anal. (2016)
14. Tripathy, B.K., Panda, G.K.: Approximate equalities on rough intuitionistic fuzzy sets and an analysis of approximate equalities. IJCSI **9**(2), 3, 371–380 (2012)
15. Tripathy, B.K., Saraf, P., Parida, S.C.: On multigranular approximate rough equivalence of sets and approximate reasoning. In: Proceedings of the International Conference on CIDM, 20–21, December 2014, Computational Intelligence in Data Mining, vol. 2, Smart Innovation, Systems and technologies, pp. 605–616 (2015)
16. Tripathy, B.K., Tripathy, H.K.: Covering based rough equivalence of sets and comparison of knowledge. In: Proceedings of the IACSIT-SC 2009, Singapore, pp. 303–307 (2009). doi:10. 1109/IACSIT-SC.2009.88
17. Jhawar, A., Vats, E., Tripathy, B., & Chan, C.S.: Generalised approximate equalities based on rough fuzzy sets & rough measures of fuzzy sets. Paper presented at the IEEE International Conference on Fuzzy Systems, (2013). doi:10.1109/FUZZ-IEEE.2013.6622541
18. Tripathy, B.K.: On Approximation of Classifications, Rough Equalities and Rough Equivalences, Springer International Studies in Computational Intelligence, vol. 174, Rough Set Theory: A True Landmark in Data Analysis, pp. 85–133 (2009)
19. Tripathy, B.K.: Rough sets on intuitionistic fuzzy approximation spaces. IEEE Intelligent Systems, pp. 776–779 (2006)
20. Tripathy, B.K., Mitra, A.: On approximate equivalences of multigranular rough sets and approximate reasoning. Int. J. Inf. Technol. Comput. Sci. **10**, 103–113 (2013). doi:10.5815/ ijitcs.2013.10.11

21. Nagaraju, M., Tripathy, B.K.: Approximate equalities for covering based optimistic multi granular rough sets and their properties. IIOAB J. **6**(4), 77–97 (2015)
22. Tripathy, B.K., Parida, S.C.: Covering based pessimistic multigranular rough equalities and their properties. Int. J. Inf. Technol. Comput. Sci. **8**(4), 52–62 (2016)
23. Tripathy, B.K., Parida, S.C.: Covering based optimistic multigranular rough equalities and their properties. Int. J. Intell. Syst. Appl. **8**(6), 70–79 (2016)
24. Tripathy, B.K., Rawat, R., Divya Rani, V., Parida, S.C.: Approximate reasoning through multigranular approximate rough equalities. IJISA **6**(8), 69–76 (2014)
25. Yao, Y.Y.: Perspectives of granular computing. In: Proceedings of 2005 IEEE International Conference on Granular Computing, I, pp. 85–90 (2005)

# RETRACTED CHAPTER: Real-Life Facial Expression Recognition Systems: A Review

**Samta Jain Goyal, Arvind K. Upadhyay, R. S. Jadon and Rajeev Goyal**

The editor has retracted this chapter [1] because of significant overlap with an earlier published article by Benta and Vaida [2]. Authors Samta Jain Goyal, Arvind K. Upadhyay and Rajeev Goyal agree with this retraction. Author R. S. Jadon has not responded to any correspondence from the publisher about this retraction.

[1] Goyal S. J., Upadhyay A. K., Jadon R. S., Goyal R. (2018) Real-Life Facial Expression Recognition Systems: A Review. In: Satapathy S., Bhateja V., Das S. (eds) Smart Computing and Informatics. Smart Innovation, Systems and Technologies, vol 77. Springer, Singapore
https://doi.org/10.1007/978-981-10-5544-7_31.

[2] K.-I. Benta, M.-F. Vaida, "Towards Real-Life Facial Expression Recognition Systems," Advances in Electrical and Computer Engineering vol. 15, no. 2, pp. 93–102, 2015, https://doi.org/l0.4316/AECE.2015.02012.

# Graphical Structure of Bayesian Networks by Eliciting Mental Models of Experts

**Udai Kumar Kudikyala, Mounika Bugudapu and Manasa Jakkula**

**Abstract** In knowledge-driven approaches to construct Bayesian networks (BN), a knowledge engineer consults with a domain expert to elicit and represent the graphical structure of a BN. The directed graph along with the node probabilities are then used for prediction or diagnosis. In this paper, we present a formal approach to learning the graphical structure of BN using domain expert(s). The proposed PFNetBN technique elicits and represents the mental model of an expert as a directed Pathfinder network. This technique uses the Target method to capture causal/influence relationships among the probability nodes from experts. It then generates a directed graph by applying the Pathfinder algorithm. Consensus Pathfinder network may be generated if multiple experts are involved. This technique generated graphs that are similar to some academic examples from the BN literature. This technique may save time in eliciting and constructing the graphical structure of a BNfrom experts.

**Keywords** Bayesian network · Pathfinder network · Mental model
Directed graph · Domain expert · Target method

## 1 Introduction

Bayesian networks are directed acyclic graphs (DAG) where a directed link or dependency connects a cause node to an effect node—causal link or influence [1, 2, 5]. Bayesian networks have been successfully used for prediction, risk assessment, and diagnosis in many domains such as medicine, legal, transportation, financial sector, and software reliability. [1, 3, 4, 7–10]. The networks and the node probabilities may be generated from large data (data-driven) [6] or from the knowledge of experts (knowledge-driven) [2] or by the combination of both [7, 8].

U. K. Kudikyala (✉) · M. Bugudapu · M. Jakkula
Computer Science and Engineering Department, Keshav Memorial Institute of Technology,
Narayanguda, Hyderabad 500029, India
e-mail: udaikudikyala@gmail.com

Constructing Bayesian networks involves first, learning/eliciting the structure of the network and then learning/eliciting the probabilities of the nodes [5, 10]. The performance of the BN for prediction depends on how well it captures the knowledge and the relationships from the data or the experts [8, 11, 12].

The Domain Experts (DEs) are critical in determining the structure of the network as a starting point [7–9]. This is especially true in solving problems involving risk assessment or diagnosis where data is either not available or if available, is not sufficient to generate the structure of the BN [6, 9]. We propose that the structure of the BN is the mental model of an expert. It consists of probability nodes as concepts of domain and directed links as influence/cause-effect relationships among the nodes. We elicit and represent the mental model of an expert as a Pathfinder network (PFNet). This graph, which has the most salient dependencies, can then be considered for determining the structure of the BN in consultation with the experts.

Section 2 presents the literature survey on application of PFNets and construction of graphical structure of BNs using experts. Section 3 introduces the Pathfinder paradigm and methods to generate directed and consensus PFNets from experts. Section 4 proposes the method and tools to determine probability nodes, proximity data, and generate structure of BN as directed PFNets. Section 5 discusses the results obtained on two academic examples from BN literature. Section 6 presents our conclusions and future work.

## 2 Literature Review

The Pathfinder paradigm enables one to model certain aspects of human semantic memory like memory organization and recall as a network. The PFNets contain nodes that represent concepts of the domain and the links represent similarity or relatedness among the concepts [14]. Undirected PFNets have been used extensively to assess knowledge structures of different groups like students and experts especially in the field of education [13, 17]. They have been used to assess requirements understanding between different groups in the field of software engineering [16]. The Pathfinder paradigm has the ability to elicit and capture the knowledge structure as both directed and undirected networks [13]. Directed PFNets were used to model the knowledge structure of rhyming words of human subjects. The directed PFNets showed that based on its structure one could explain how human subjects organized and recalled the rhyming words [15].

Constructing a BN, first, consists of determining the graphical structure of the network (DAG) and then eliciting the node probability tables (NPT) [5, 6, 10, 12]. A good network structure should be elicited first before determining the NPTs. Similarity networks containing a similarity graph and a collection of local belief networks and partitions were used to construct the structure of the BN [10]. There is a structured approach to automatically elicit and construct causal map which are then converted to DAG for BN [2]. Fenton et al. proposed "idioms" that are abstract patterns that match the metal model of experts based on problem description by

experts. The instantiation of "idioms" form objects. Then object-oriented BN approach can be used to construct the structure of network [5]. There is an approach that uses table-like questionnaire to elicit causal relationships from DE [11]. Another approach uses diagnostic questions to elicit causal relationship from DEs [12].

We could not find structured techniques to elicit mental model from experts as directed graph using techniques from psychology. The survey could not find any evidence of a method on how consensus network may be constructed if multiple experts are involved in BN. The proposed method may enable KE to use the DE's time more effectively in eliciting the nodes and the preliminary structure of BN.

## 3 Pathfinder Networks

The PFNet is a directed or undirected graph that was originally used to represent the structural knowledge of human subject using psychological proximity data. The link weight represents the strength of the relationship between the two concepts (nodes). The lower the link weight, the closer they are related to one another as perceived by a subject [14, 16, 17]. Figure 1a shows all the links in the original network. The Pathfinder algorithm eliminates the link shown by dashed lines because they violate triangle inequality as shown in Fig. 1b.

The dashed line from A to B in Fig. 1b violates the triangle inequality because there is a shorter alternate two-edge path from nodes A to C and then from nodes C to B. The edge weights are dissimilarity measures. The higher the edge weight, the more dissimilar the pair of concepts were perceived to be by the subjects. The dashed links are not present in the PFNet. The path length is assumed to be sum of the edge weights. The Pathfinder algorithm generates a class of PFNets based on two parameters, $r$ and $q$. The *Minkowski r*-metric determines how the path length is computed between two nodes. It has a range from 1 through $\infty$. When $r = 1$, the path length is sum of the edge weights and when $r = \infty$, the path length is computed as maximum of the edge weight along that path. The second parameter $q$ determines the upper limit on the number of edges considered as alternate paths between two nodes in order to eliminate a direct link from the network. It takes values from 2 to $n - 1$, where $n$ is the number of nodes. When $q = 3$, all alternate



**Fig. 1** Generalized triangle inequality

paths with two and three edges are considered. The direct link is eliminated only if there is a shorter alternate path between them. When $q = n - 1$, then all alternate paths between two nodes are evaluated before considering eliminating the direct link from the PFNet. If the matrix is symmetrical, then undirected PFNets are generated otherwise a directed PFNet is generated [14].

Figure 2a shows a fully-connected directed graph with edge weights. Figure 2b shows the PFNet ($r = 1$, $q = 2$). The arc from node C to B is eliminated because there is an alternate shorter path C to A and then from A to B that has a path length of two. When $q = 2$ a maximum of two-edge paths are considered to check if there is shorter alternate path in order to keep a direct arc in the fully-connected network. Figure 2c shows PFNet where, path length is calculated as the maximum of the edge weights along that path.

The arc from node D to nodes A with an edge weight of five is eliminated because there is shorter alternate path from node D to node B and then from node B to node A. This path length is calculated as maximum of the edge weights three and two i.e., three. This class of PFNets have been found to be very useful in modeling knowledge structures of human subjects [13, 15, 16].

This property of the directed PFNet makes it suitable as a starting point to develop a DAG for the BN since links in BN represents influence/dependencies. Only the highly dependent/influential links as perceived by the experts are revealed in such class of PFNet.

A causal/influence matrix is elicited from the expert's ratings. The details of the elicitation method are discussed in Sect. 4. This matrix is used by a tool that implements the Pathfinder algorithm to generate a graphical PFNet. In this paper, all PFNets are directed and have values of $r = \infty$ and $q = n - 1$ unless otherwise stated. When multiples DEs are involved, the KE may want to generate the consensus network to consolidate the mental models of all experts [16]. Assume that four nodes A, B, C, and D were identified by experts. The size of each matrix would then be 4 by 4. The influence of a node with itself is assumed to be 0. The numeric rating scale for the influence/dependency would range from 1 to 4. In matrix $A_1$, node A has the most influence on node B but has no influence/dependency on nodes C and D.



**Fig. 2** **a** Fully-connected graph, **b** PFNet ($r = 1$, $q = 2$), **c** PFNet ($r = \infty$, $q = n - 1$)

Therefore, the numeric rating for row A and column B is 1 and that of row A and columns C and D are 4. The resultant matrices generated are $A_1$, $A_2$, $A_3$, and $A_4$ corresponding to the four experts. In order to obtain a consensus network of the experts, the consensus matrix is derived by adding the matrices $A_1$, $A_2$, $A_3$, and $A_4$. The resultant consensus matrix is denoted by the symbol $A_{1-4}$ as shown in Fig. 3. For the purpose of visual comparison, the individual and the consensus PFNETs of DEs are shown in Fig. 4. The PFNet of the first expert is denoted by $PFNet_{A1}$. The consensus PFNet of the four DEs is denoted by $PFNet_{A1-A4}$. This consensus matrix is then used by the Pathfinder tool to generate the consensus PFNet as shown in Fig. 4e. The common arcs of the four individual PFNets shown in Fig. 4a–d are preserved in the consensus network shown in Fig. 4e. It may have some additional arcs that are not common to all the four networks.

## 4   Elicitation of Mental Model

We propose a technique called PFNetBN technique for elicitation of the mental model of experts as PFNets as shown in Fig. 5. In the PFNetBN technique, first, the probability nodes or concepts are identified by the KE from information available from various sources. In the second step, the nodes are rated on a scale of 1–5 by the experts based on centrality of node to solve the problem. The nodes that are considered to be most relevant by the experts should be selected for the next stage. At this stage, feedback is taken from the experts regarding any missing nodes or any nodes that were excluded by KE. In the third stage, the finalized probability nodes are then used to collect the causal/influence data using Target method. Finally, the PFNet is generated by applying the Pathfinder algorithm. *JTarget* and *JPathfinder* tools were applied to implement the third and the fourth steps of the technique, respectively. Both the tools were already implemented in Java and their executables are available for free [18].

$$A_1 = \begin{bmatrix} 0 & 1 & 4 & 4 \\ 1 & 0 & 1 & 4 \\ 1 & 4 & 0 & 1 \\ 2 & 4 & 1 & 0 \end{bmatrix} A_2 = \begin{bmatrix} 0 & 1 & 4 & 1 \\ 4 & 0 & 4 & 4 \\ 1 & 4 & 0 & 1 \\ 2 & 4 & 1 & 0 \end{bmatrix} A_3 = \begin{bmatrix} 0 & 1 & 4 & 1 \\ 4 & 0 & 4 & 4 \\ 1 & 4 & 0 & 1 \\ 1 & 4 & 1 & 0 \end{bmatrix} A_4 = \begin{bmatrix} 0 & 1 & 4 & 1 \\ 4 & 0 & 4 & 4 \\ 1 & 4 & 0 & 1 \\ 2 & 1 & 1 & 0 \end{bmatrix}$$

$$A_{1-4} = A_1 + A_2 + A_3 + A_4 = \begin{bmatrix} 0 & 4 & 16 & 7 \\ 13 & 0 & 13 & 16 \\ 4 & 16 & 0 & 4 \\ 7 & 13 & 4 & 0 \end{bmatrix}$$

**Fig. 3** Co-occurance matrices for individuals and consensus matrix

Fig. 4  a PFNet$_{A1}$, b PFNet$_{A2}$, c PFNet $_{A3}$, d PFNet$_{A4}$, e PFNet$_{A1-A4}$



Fig. 5  Mental model elicitation: PFNetBN technique

A method to collect proximity data called the Target method was developed by Tossell, Schvaneveldt, and Branaghan [13]. In this method, the central concept is placed in the center of three concentric circles and the related concepts are dragged-and-dropped by subjects in appropriate concentric circles based on categories like *Most influential or Cause*, *Moderately Related/Influential, Somewhat Related/Influential* and *Not Related* or *Slightly influential* to the central node.

The PFNetBN technique was applied in the following manner for this study:

- Two academic examples, one for cancer diagnosis and another for flood risk assessment were selected from the BN literature [1]. A node from which an arc emanates is called the parent node and node to which the arrow is pointing is referred to as the child node.

- We directly applied the third step of the technique as the nodes are already finalized for each example. For an example problem, each node was considered as the central node. Its child nodes(s), if any, were categorized as *Most Influential/Cause*. All other nodes that are not directly linked and that are not the children of central node were categorized as *Somewhat Related/Influential*.
- The *JPathfinder* tool is used to directly load the matrix generated by the *JTarget* tool to generate the graphical PFNet.

## 5 Results and Discussions

Table 1 shows the proximity matrix for *Visit to Asia?* example. The column headers (nodes) follow the same order as the row headers. In Table 1, if V*isit to Asia?* is the central node, then *HasTuberculosis?* was the only node that was categorized as being *Most influential or Cause* with an influence/causal value of 1. The remaining nodes were categorized as *Somewhat Related/Influential* to the central node with an influence/causal value of 4. The influence of a node with itself is assumed to be 0. The proposed technique produced exact network structure as academic examples. The visual comparison between the academic examples and the graphical PFNets generated by the tool are shown in Figs. 6 and Fig. 7.

## 6 Conclusion and Future Work

We proposed a technique that contributes toward elicitation of mental model of DEs that can be used as starting point to derive the graphical structure of the BN. Consensus PFNet may be automatically generated if multiple experts are involved. This technique may reduce the time spent with expert(s) to elicit the initial graph structure for the BN. Further research needs to be carried out to see if this technique works with real-world problems. Further research is required to see if it scales with size of the network and how it compares with other techniques.

**Table 1** Influence/causal proximity data of visit to Asia? example

| Probability nodes | Influence/causal values | | | | | | |
|---|---|---|---|---|---|---|---|
| Visit to Asia? | 0 | 4 | 4 | 1 | 4 | 4 | 4 |
| Smoker? | 4 | 0 | 4 | 4 | 4 | 1 | 1 |
| Positive X-ray? | 4 | 4 | 0 | 4 | 4 | 4 | 4 |
| Has tuberculosis? | 4 | 4 | 1 | 0 | 1 | 4 | 4 |
| Dyspnoea? | 4 | 4 | 4 | 4 | 0 | 4 | 4 |
| Has bronchitis? | 4 | 4 | 4 | 4 | 1 | 0 | 4 |
| Has lung cancer? | 4 | 4 | 1 | 4 | 1 | 4 | 0 |

**Fig. 6** Detecting cancer, **a** example from literature [1], **b** directed PFNET



**Fig. 7** Assess flood risk, **a** example from literature [1], **b** directed PFNET

# References

1. Fenton, N.E., Neil, M.: The use of Bayes and causal modelling in decision making, uncertainty and risk. J. Council Eur. Prof. Inf. Soc. (CEPIS) **12**(5), 10–21 (2011)
2. Nadkarni, S., Prakash, P.P.: A causal mapping approach to constructing Bayesian networks. Decis. Support Syst. **38**(2), 259–281 (2004)
3. Yet, B., Bastani, K., Raharjo, H., Lifvergren, S., William, Marsh W., Bergman, B.: Decision support system for Warfarin therapy management using Bayesian networks. Decis. Support Syst. **55**, 488–498 (2013)
4. Fenton, N.E., Neil, M.: Decision support software for probabilistic risk assessment using Bayesian networks. IEEE Softw. **31**(2), 21–26 (2014)
5. Neil, M., Fenton, N.E., Nielson, L.: Buliding large-scale bayesian networks. Knowl. Eng. Rev. **15**(3), 257–284 (2000)
6. Helsper, E.M., Van der Gaag, L.C., Feelders, A.J., Loeffen, W.L.A., Geenen, P.L., Elbers, A. R.W.: Bringing order into Bayesian-network construction. In: The Proceedings of the 3rd International Conference on Knowledge Capture, K-CAP'05, Banff, Alberta, Canada, pp. 121–128, Oct 2–5 (2005)
7. Velikovaa, M., Lucasa, P.J.F., Samulskic, M., Karssemeijerd, N.: On the interplay of machine learning and background knowledge in image interpretation by Bayesian networks. Artif. Intell. Med. **57**, 73–86 (2013)
8. Yet, B., Perkins, Z., Fenton, N., Tai, N., Marsh, W.: Not just data: a method for improving prediction with knowledge. J. Biomed. Inf. **48**, 28–37 (2014)

9. Constantinou, A.C., Fenton, N., Marsh, W., Radlinski, L.: From complex questionnaire and interviewing data to intelligent Bayesian network models for medical decision support. Artif. Intell. Med. **67**, 75–93 (2016)
10. Heckerman, D.E., Nathwani, B.N.: Toward normative expert systems: part II probability based representations for efficient knowledge acquisition and inference. Methods Inf. Med. **31**, 106–116 (1992)
11. Xio-xuan, H., Hui, W., Shuo, W.: Using expert's knowledge to build bayesian network. In: IEEE International Conference on Computational Intelligence and Security Workshops, pp. 220–223 (2007)
12. Henrion, M.: Practical issues constructing a Bayes' belief network. In: The Proceeding of Uncertainty in AI '87, Seattle, WA, pp. 132–140 (1987)
13. Tossell, C.C., Schvaneveldt, R.W., Branaghan, R.J.: Targeting knowledge structures: a new method to elicit the relatedness of concepts. Cogn. Technol. **1**(2), 11–19 (2010)
14. Dearholt, D.W., Schvaneveldt, R.W.: Properties of pathfnder networks. In: Schvaneveldt, R. W. (ed.) Pathfinder Associative Networks: Studies in Knowledge Organization, Ablex Publishing Corp. Norwood, NJ, USA. pp. 1–30 (1990)
15. Rubin, D.C.: Directed graphs as memory representations: the case of rhyme. In: Schvaneveldt, R.W. (ed.) Pathfinder associative networks: studies in knowledge organization, pp. 1–30. Ablex Publishing Corp, Norwood (1990)
16. Kudikyala, U.K., Vaughn, R.B.: Software requirements understanding pathfinder networks: discovering and evaluating mental models. J. Syst. Softw. **74**(1), 101–108 (2005)
17. Dearholt, D.W., Alt, K.J., Halpin, R.F., Oliver, R.L.: Foundational aspects of student-controlled learning: a paradigm for design, development, and assessment appropriate for web-based instruction. J. Eng. Edu. **93**, 129–138 (2004)
18. Interlink Inc.:http://interlinkinc.net/

# Signal Monitoring in a Telemedicine System for Emergency Medical Services

**S. T. Aarthy and S. Kolangiammal**

**Abstract** Wireless gadgets have invaded the medical field with a wide range of capability. It is necessary to monitor the condition of the patient during regular intervals of time. The field of contemporary electronics design and communication has the greatest potential for the advancement and dissemination of medical knowledge. There are basically six different types of sensor to gather patient's medical information without injecting those inside the patient's body and from this, we are achieving remote monitoring and data gathering of patients. Hence, there are advantages of mobility. There is no need for a periodical visit to a doctor to the patient. The proposed system provides emergency treatment after analysis of recorded values using mechanical setup. The collected data is continuously shared and monitored among the medical professionals using IoT.

**Keywords** Wireless sensor technologies · Sensors · Remote monitoring
Data gathering

## 1 Introduction

The continuous development in mobile communication acts as a linkage between a mobile and telemedicine as well as emergency medical services (EMS). The Internet of things (IoT) is the inter-networking of physical devices and other items embedded with electronics, software, and network connectivity which enable these objects to collect and exchange data. These objects can automatically transfer data over a network without requiring any computer or human interaction. Its appeals the ubiquitous generalized access to the status and location of anything. We may be interested in the online monitoring system for continuous casting; equipment is

S. T. Aarthy (✉) · S. Kolangiammal
SRM University, Chennai, India
e-mail: aarthy.s@ktr.srmuniv.ac.in

S. Kolangiammal
e-mail: kolangiammal.s@ktr.srmuniv.ac.in

established based on IoT sensing technology and communication technology [1]. As this particular system contains different kinds of sensors like temperature sensor, blood pressure sensor, heartbeat sensor, MEMS motion sensor, eye blink sensor, and other data transmission protocols, it will lead to a collection of a large amount of heterogeneous data of different sensors current and voltage values.

## 2   Related Work

A. Background and Motivation
   The existing medicinal human monitoring process includes the constant presence of labor, and this led to the achievement of overall system architecture.
B. Limitations in the present system
   The present system needs more manpower. Gross errors are noticeable in the current scenario. This lacks in providing automated service during the period of emergency, and it lacks in the technological advancement of not being able to access the patient's report in absence of a doctor. The noticeable Gross errors in the current system can be reduced with technological advancement.

## 3   Requirements of a Biomedical System

The functional requirements of the overall system are expected to, minimizing the need for human labor and additional system interaction on-screen.

The patient's privacy must be kept as safe as possible i.e., no data is transmitted when it is not supposed to transmit.

Periodic transmission of monitored values is sent to the concerned doctor from anywhere in this world [2]. The data about the particular patient must be shared with the patient's physician and family. There will be proper indications whenever the basic measurements like temperature, pressure, heartbeat, motion, and Eye-blink sensors exceed the threshold level given to it, using proper simulation.

## 4   Flow Chart and Health Monitoring System

The flowchart is shown in Fig. 1 explains the entire working algorithm of the prototype. It begins with the five different sensors that include temperature sensor, heartbeat sensor, blood pressure sensor, IR based Eye-blink sensor, and MEMS Motion sensor. These sensors are interfaced into the 4 × 16 LCD display. On connecting to the patient's body, we sense all the medical vitals which are

constantly displayed on the LCD [2]. The program is written in such a way that if these sensed values increase above the given threshold then an automatic indication is provided. This indication is in the form of a buzzer in the prototype. When the values cross the threshold then the buzzer goes ON. When the values cross the threshold, it acts as an interrupt for the code which makes on the rotor in the mechanical set up to rotate. The servomotors in the mechanical setup begin functioning when the value increases the desired threshold loaded into the code.

The entire proposed work is shown in Fig. 2. The system starts the monitoring and sensing process usually with the help of all the sensors which are connected to the comatose patient [3]. Each sensor has a particular threshold level. The



**Fig. 1** The general flow chart

monitored values are displayed in real time via an LCD display which is fitted with the ARM microcontroller. Whenever there is a change in the value of the sensor, it will be displayed in the LCD display and whenever the value of the sensor increases the particular threshold level, the Indication process takes place and also these observations are sent to the Partial Cloud from where the doctors can access it in real time. For temperature sensors, the threshold value is 39 °C, for pressure sensor the threshold value is greater than 140/90 mmHg. For MEMS, motion sensor any change in its axis will be indicated through a buzzer prototype. The heartbeat sensor the threshold value is greater than 83 beats per minute. The infrared eye blink sensor indicates whenever there is a change in the eye blink which is considered as recovery from coma. There are slots for other possible sensors which depend on future development in telemedicine and Automation field [3, 4]. Apart from this, there exists a mechanical setup for treatment during emergency situations (i.e., medicines stored in the injection can be automatically injected into patients body). The automated treatment is used only to the point of emergency. The labor can induce this treatment manually.

ARM controller: it is Advanced RISC Machine (ARM) is a 32-bit microcontroller which belongs to the family of Reduced Instruction Set Computing (RISC) machines. LPC2148 microcontrollers are ideal for applications in which miniature is the key requirement. Its pipelining of executing more instructions per clock cycle makes it easier and simpler in this system.

LPC2148 is a 16 or 32-bit microcontroller in a small package (LQFP64). It has 8–40 kB on-chip static RAM, 32–512 kB on-chip flash memory, and 128-bit wide interface/accelerator. 6/14 analog inputs are provided by 10-bit ADCs with about 2.44 ms of conversion time per channel. It has two 32-bit timers or external event counters and PWM unit (six outputs).



**Fig. 2** The black diagram of health monitoring system

The interfacing is shown in the Fig. 3. In an SPI connection, a master device (usually a microcontroller) always controls the peripheral devices. Three lines common to all the devices are specified as master in slave out (MISO), master out slave in (MOSI), and serial clock (SCK).

Slave Select pin is the pin embossed on all the devices that the master can command to enable and disable specific peripherals. When a device's slave select pin is low, it follows the master. If it is high, it disowns the master. This allows one to have multiple devices on SPI to share the same MISO, MOSI, and CLK lines.

In SPI, the clock signal is controlled by LPC2148 Slicker Board. All data is transmitted in and out using this pin. These lines need to be connected to the pins on the LPC2148 Slicker Board. The ENC28J60 requires a single packet control byte to precede the packet for transmission to Microcontroller. The IP address is used to access the Ethernet control. The ENC28J60 SPI connections with LPC2148 have four I/O lines required.

UART It is the acronym of universal asynchronous receiver transmitter, a serial to parallel and parallel to a serial conversion device which has an inbuilt baud rate generator of nearly 1200 bps. It provides a better buffering of data so that the computer and other serial device's data stream remains coordinated properly. The reason why we chose Asynchronous over Synchronous is the Asynchronous device provides parity bit which is used in error checking.

Sensors:

*Temperature sensor* It is a wearable sensor which continuously monitors the body temperature. LM35 is used which has the range of −55–150 °C. This has three pins for input, output, and ground [3, 4].

*Blood pressure sensor* It is a non-invasive sensor designed to measure the systolic and diastolic pressure in blood continuously. The systolic blood pressure is the pressure between the veins when the heart muscles are completely contracted. The diastolic blood pressure is the pressure in veins between the beats of the heart. These systolic and diastolic blood pressures must be maintained between 120/80 mmHg [3, 4].

*Motion sensor* This motion sensor is based on micro electro mechanical systems (MEMS) sensor. This MEMS sensor intelligently senses the movement of the



Fig. 3 Ethernet interfacing with controller

human body with respect to change in its XYZ axis. Accelerometers and gyroscopes are the two main design components in MEMS motion sensors [3, 4].

*Eye blink sensor* This eye blink sensor is also the infrared sensor which senses the blink of an eye. Blinking of the eyes is the first symptom of recovering from Coma.

*Heart beat sensor* This is an infrared based sensor which measures the heart beat by sensing the change in the intensity of blood in forefinger's artery while the heart is pumping the blood. The normal heart beat for a human being is 72 beats per minute [3, 4].

A servomotor is an actuator that controls linear or angular position, velocity and acceleration. It has a motor that is coupled to a sensor for position feedback. The controller is relatively sophisticated and often a dedicated module designed specifically for servomotors is used.

The motor is in neutral position if the potential rotation of the servo in the both the clockwise or counter-clockwise direction is equal. The PWM sent to the motor determines position of the shaft. The rotor turns to the desired position based on the duration of the pulse sent via the control wire. The servo motor expects a pulse every 20 ms and the length of the pulse will calculate how far the motor turns. When these servos are commanded to turn, they will move to the position and hold that position. If an external force opposes the servo while the servo is holding a position, it should resist and hold its position. The position pulse has to be repeated so that the servo stays in position.

The Fig. 4 shows the prototype of Emergency Medical Service with sensors and Microcontroller configuration

If a patient's medical profile and inputs required for dialysis are determined using medical devices that are attached to the body then the patient can receive treatment just with the help of portable/home machines designed for such purpose. Data collected from this device is analyzed and stored. Data received from multiple sensors and medical devices together aids in taking precise and timely decisions.



**Fig. 4** The prototype of emergency medical service with sensors and micro controller configuration

Thus, patients can be monitored from any location. Based on the alert received, appropriate response can be given. Such advanced treatment can help improve a patient's quality of life. IoT can be considered as one of the most sophisticated technologies that have the potential to affect the health and safety of billions of people as well as a major impact on the economy [1, 5]. The primary constituents of IoT are linked to networks for data transportation. The output can be viewed in the hyper terminal system which is achieved using IoT concept [2, 5]. The Wi-Fi is configured using the Mikro-C which is an interface between the controller and the software.

The graph shown in Fig. 5 explains the various sensors and their corresponding threshold values for particular time intervals. The graph represents the relationship between the various sensor values for their corresponding time intervals in a comatose patient. When the sensors value cross the threshold, the indication is provided using buzzer. This graph gives a better illustration of the behavior of all the sensors used with respect to time.

## 5 Performance and Evaluation of the System

The Table 1 shows the threshold values of various sensors.

The standard display operates at 3.3 V. 5 V operation for the character display is also possible. The output of LCD, when any value exceeds the threshold value is shown in the Fig. 6.

Compiler is designed to be smart and efficient so that it can be relied on for doing the hard work. It features four levels of optimizations that can reduce code size up to 20%.

To ensure that the modem is properly connected or to view the modem's settings, the PC can send commands through HyperTerminal and check the results. HyperTerminal also has scroll functionality. The PC can use HyperTerminal to transfer large files from a computer to the portable computer using a serial port



**Fig. 5** Graph between sensor values and time intervals

**Table 1** xxx

| S. no | The basic sensors | | |
|---|---|---|---|
| | Name of the sensor | Normal value | Threshold value |
| 1 | Temperature sensor (LM35) | 37 °C | 39 °C |
| 2 | Blood pressure sensor | 120/80 mmHg | 140/90 mmHg |
| 3 | Heart beat sensor | 72 beats/min | 85 beats/min |
| 4 | MEMS motion sensor | $XYZ$ axis | Change in axis |
| 5 | Eye blink sensor | 0,1 | 1 |

**Fig. 6** Output of LCD



rather than going through the process of setting up a portable computer on a network [1, 6].

The paper makes use of biosensors to read signals from the human body, which is an analytical device that converts biological responses into electrical signals. A successful biosensor must provide accurate and precise response that is linear, noise-free and reproducible over the analytical range. The biocatalyst used in analysis should be highly specific. Biosensors used for invasive monitoring in clinical purposes should be biocompatible, antitoxic, and free from antigenic effects. It should be sterilizable and not be prone to fouling or proteolysis.

## 6 Conclusion

This advanced system will help doctors to communicate better with their patients and keep a constant tab on their vitals. The mechanical setup is a state-of-the-art system that immediately operates under critical conditions and helps an assistant or a nurse in choosing the right medication, both in prescription as well as quantity. If this proposed protocol is further evolved using real time equipment along with the correction of minor glitches and constraints, this system will change the world of medicine and evolve into a much faster and efficient medical environment.

# References

1. Kim, J.: Energy-efficient dynamic packet downloading for medical IoT platforms. IEEE Trans. Ind. Inform. **11**(6), 1653–1659 (2015)
2. Thelen, S., Czaplik, M., Meisen, P., Schilberg, D., Jeschke, S.: Using off-the-self medical devices for biomedical signal monitoring in a telemedicine system for emergency medical services. IEEE J. Biomed. Health Inform. **19**(1), 117–123 (2016)
3. Konganti, S.C., Suma, H.N., Abhishek, A.M.: Analysis and monitoring of coma patients using wearable motion sensor system. Int. J. Sci. Res. **4**(9), 1154–1158 (2015)
4. Fanucci, L., Seryiosaponala, S.: Sensing devices and sensor signal processing for remote monitoring of vital signs in CHF patient. IEEE Trans. Instrum. Measurem. **62**(3), 553–569 (2013)
5. Ko, Y.J., Huang, H.M.: A patient- centered medical environment with wearable sensors and cloud monitoring. In: IEEE 2nd World Forum On Internet of Things, Dec 2015, pp. 628–633
6. Basak, A., Narashimhan, S., Bhunia, S.: Kims: kids health monitoring system at day-care centers using wearable sensors and vacabulary-based acoustic signal processing. In: IEEE 13th International conference on E-Health Networking, Applications and Services, September 2011, pp. 1–8

# Evolving the Efficiency of Searching Technique Using Map-Reduce Hashing Technique

**Shivendra Kumar Pandey and Priyanka Tripathi**

**Abstract** Nowadays as data volume is increasing, it is becoming difficult to access the data within span of time. So, our aim is to process required data as fast as possible. Though we have variety of algorithms but none of them are specially designed to manage the large data (e.g. peta byte size of data). In this research paper, authors have proposed an algorithm based on hashing technique which uses Hadoop framework to reduce search time.

## 1 Introduction

String matching is a technique of searching the appearance of a pattern 'p' in text 'T'. There are two categories of string matching, first is Exact String Matching (KMP, BMH, etc.), and second is Approximate String Matching (Brute Force, Rabin–Karp algorithm, fuzzy string matching, etc.). String matching is an influential part of many areas including text retrieval, Text editors, Database queries, Bioinformatics, n-D mesh problem, intrusion detections in network, large string matching, music content retrievals, DNA sequence matching, MS word spelling checker, search engines and many more applications. The size of the enterprise database used nowadays has been growing at exponential rates. Simultaneously, we need to analyse and process Terabytes of data in efficient way on daily basis [1]. This is contributing big data problem to process and manage

S. K. Pandey (✉) · P. Tripathi
National Institute of Technical Teachers Training and Research, Bhopal
Madhya Pradesh, India
e-mail: shivendrapandey786@gmail.com

P. Tripathi
e-mail: ptripathi@nitttrbpl.ac.in

within the span of time. Processing of data includes various operations depending on usage like highlighting, tagging, searching, indexing, matching etc. [1, 2].

The heart of Big Data analytics is the Map-Reduce programming model. As a distributed computing framework, Map-Reduce programming model uses a divide and conquer approach to allow large scale parallel processing of Big Data [3, 4]. As the name suggests, Map-Reduce model has of two basic functions, first is Map function which splits the data, and second is Reduce function which carries out the final processing of the Mapper outputs in a reduced size.

The remainder of the work, we arranged as follows, Part II we have explained survey of the Hashing technique and string searching algorithms, Part III author has presented proposed algorithm with an example and in Part IV experiment specifications and result analysis with Brute Force algorithm and at the last Part V is conclusion of the paper.

## 2 Related Work

For this research work, we have studied several literature works in the both areas that illustrate related work on string matching algorithms and uses of Hash function that works on Big Data for text retrieval. Following is the important string matching and hashing technique on big data.

### 2.1 Simulation of Map Reduce with the Hash of Hashes Technique

Hinson et al. [5], research work illustrates simulation of Map Reduce in three simple examples with Hash of Hashes technique. First, text analytics, word count by conventional Map Reduce, second Web analytics, counting unique visitor and click-through and third finance, a frequency of stock market changes object to search a string in a Hash table. Author has described that a Hash table can be created within another Hash table to retrieve and to store data in the Hash table. Hash table provides three different types of component objects Java objects, Hash and Hash iterator objects and appended objects. Objects are simply data elements consisting of attributes, method and operators. Hash objects are RAM memory resident tables with each record on the table made of different lookup keys and its data sets.

In this paper, author also explained that objects can be created dynamically rather than declaring in advance. HOH technique is very suitable for data processing involving clustering or classification.

## 2.2 Brute Force String Matching

It is also called as proof by exhaustion or naive string matching. It is mathematical proof where the text to be proved by dividing into a finite number of sub-categorical cases and each sub-categorical case is matched. In this technique, pattern matching is done with the help of characters in text and slides one character at a time. BFS Matching technique does not need any pre-processing time and takes constant extra space. So, total matching time of Brute Force technique equals to its running.

## 2.3 Rabin–Karp String Matching Algorithm

Rabin–Karp string matching algorithm uses two phase to match a string (pattern) in text: Pre-processing phase and matching phase. String matching algorithm uses Hash function to match the pattern 'p' in a text 'T' [6]. In Rabin–Karp string matching algorithm, first determines the integer values in the given text. String Matching Algorithm then separates pattern by a user defined prime number 'q' and determines the residue of the pattern 'p'. Afterwards, determines the residue of initial 't' character of text 'T' [6]. If the residue of initial 't' characters of text and residue of pattern 'P' is identical, only then it performs matching else need not to do a comparison of residue. The basic principal of Rabin–Karp string searching technique uses Hash function to match the pattern in the given string. It is a string searching algorithm that uses Hash value to match string from the Hash table. Rabin–Karp string searching time complexity for best case is $O(n + m)$ and for worse case is $(nm)$.

## 2.4 Knuth Morris Pratt String Matching Algorithm

KMP String matching algorithm uses two phases to match string. A Knuth Morris Pratt algorithm takes linear time to find a particular pattern. It is a very efficient algorithm for string matching [6]. In the pre-processing phase, it makes the prefixes and in the matching phase, it finds the occurrence of a pattern in string and returns number of shifts of the pattern after which occurrence is found [7, 8]. In other words, KMP string matching algorithm matches occurrences of a pattern 'p' within the text 'T' by providing the remark whenever a string do not match.

## 2.5    Boyer Moore String Matching Algorithm

In the field of String machining, Boyer and Moore have proposed is a very effective string matching algorithm [8, 9]. It is proposed by J. Strother Moore and Robart S. Boyer in 1977. This matching algorithm works with two stages and main objective of this string matching algorithm is to match on the following part of the pattern with text rather than the head part of text. If it does not match, skip and jump the word of multiple characters rather than matching every single character available in the text. It is suitable for applications in which the pattern is much shorter than the text. Boyer–Moore algorithm uses bad character shifts. The many sided quality of the Boyer–Moore String Searching Algorithm can be demonstrated to straight time conduct, just if patent 'P' does not comes in text 'T'. Time complexity of Boyer–Moore String matching is $\Theta(nm)$.

## 3    Proposed Algorithm

Proposed algorithm developed with an aim to search index value of pattern from the text file. Working principle of the algorithm is similar to the standard word count algorithm in Map Reduce. In this research paper, author has applied string matching technique on Hadoop to search index of a pattern 'p' in a text file 'T'. The proposed algorithm uses hashing technique to match the string from the text [10–13]. It has two working phases: pre-processing phase and searching phase. Hear in the Hadoop architecture, pre-processing takes place in the Mapper where Mapper breaks the text file and converts it to substrings and assigns an integral value (that is an index value of a word) to all the substrings. Next phase is the searching phase that is done with help of Reducer part of Hadoop where Reducer uses Hash function to store the substrings integral value in the Hash table. Hash function could be user defined or built-in. Authors have experimented text data with both Hash functions. But for this research work, we have used a built-in Hash function. This technique uses continuous memory allocation to store the substrings with the help of Hashing which makes it very efficient to store and search a string.

The working principle of proposed algorithm uses Hadoop architecture and string matching algorithm. The Pre-processing phase of the algorithm is divided into several stages, in first stage, input split takes place, and a text file is divided into a number of chunks of size 60 BM or 128 MB sub-files and transfers text data to the Record Reader as shown in Fig. 1. The Record Reader will read the file line by line and convert it into a pair of offset value and text format (byte offset, line). Where offset value is starting address of line, and text is a one line text (characters) of a file [3]. In the next step, Mapper reads text, line by line and divides each line into key-value pairs. Mapper key value a pair is a substring and index value (keyword, index-value) of a word, as explained in Algorithm 3.1.

**Fig. 1** Simulations of text data with Mapper and reducer on Hadoop Map-Reduce programming paradigm

Searching phase is also separated into two steps of Reducer part. First Reducer stores key-value pairs (keyword, index-value) into Hash table (by applying Hash function on key) with the help of Hashing explained in Algorithm 3.2. Afterwards with the help of key searches the value stored in the Hash table. It is a string matching algorithm so from Reducer's output desired key-value pairs (key, indexes \locations) will be matched. In the proposed technique, iterative function is used to collect all the values for keyword and to search value.

```
Mapper Algorithm


Mapper(LongWritable  key, Text value)
{
//key: starting index value of text content
//value: text content (of size 64MB)
     String s= value.toString();
     char[] words = s.toCharArray();
   For(j =ino; j<words.length; j++)
         {
         If(words[j]!=' ' && words[j]!= '.' &&
   words[j]!=',');
         }
         //calculate the starting index value(address)
   to key
         For each key 'w' the index value
         Output.collect(new Text(temp), new
   LongWritable(ino));
}
```

**Algorithm 3.1** Mapper Algorithm to Assign Index Value to the Words

Time complexity analysis for proposed algorithm, Time is taken to assign a unique integer value for substrings is $O(m*(n - m))$, and Time is taken to assign a unique integer value for pattern $O(m)$. So total pre-processing time is,

$$O(m * (n - m) + m).$$

```
Reducer Algorithm

Reducer(Text  key, Iterator value)
{
//key: word (text)
//value: list of index values for word
    HashMap<String,String>    hm=new      HashMap<String,
String>;
    //Store key value pair into hash table
    // use hashing with chaining technique
    Hm.put(text,s);
          {
          String a= (String) m.getKey();
          //search values for key in the hash table
          String b = (String) m.getValue();
          Output.collect(new Text(a), new Text(b));
          }
}
```

**Algorithm 3.2** Reducer Algorithm to Combine Index Value to the Words and Search a Key

Search time analysis, algorithm searches the location/index value of a string based on the integer value of the string in the Hash table. So that, searching a pattern which is not present in the text gives null value or only one integer value stored in the Hash table.

$$\text{Best case } O(1)$$

If Hash function returns Hash table location that has long linear linked list, then proposed algorithm takes maximum time.

$$\text{Worst case } O((n - m) + 1)$$

## 4 Experiment Specifications and Result Analysis

Experiment setup is done for Hadoop single node cluster (we configured Master node and slave nodes on a single system) and unstructured text data on Hadoop Distributed File System (HDFS). Authors have proposed two different algorithms

for Mapper and Reducer (Algorithms 3.1, 3.2) that runs into Mapper phase and Reducer phase, respectively. Input and output data for the Mapper and Reducer programs is stored in HDFS. The Master node runs master services with Name Node and Job Tracker. The slave node runs slave services with Data node and Task Tracker.

System has Intel Pentium Core, CPU with 3.20 GHz processors and 64-bit operating system with 4 GB (RAM) memory. The operating system is GNU/CentOS-6.5-64, VMware workstation 10.0.1 and Hadoop framework v-1.0.4. We have implemented proposed algorithm using the Map-Reduce [8] programming model and used (JAVA Programming) JDK 1.7 and IDE Eclipse standard kepler-64. Authors have taken number of replicas to 1 (default), and HDFS block size is 64 MB during the tests. File format we were selected for experiment is text (.txt) file format.

Proposed research work is tested on the several unstructured text file of different size. Authors have chosen time constraint for the comparison with three different time context i.e. user time, real time and system time. Proposed algorithm result is compared with Brute Force algorithm. For every file, proposed algorithm is working good and giving better result.

Authors have done experimental and comparative study of proposed algorithm with Brute force algorithm and taken time as a constraint. We have presented the result for Real, User and System time. Real time means actual time from start to finish of the call of processes, CPU busy with outside task (user/unsupervised



**Fig. 2** Experimental result comparison of our algorithm with Brute Force algorithm for user time, real time, system time

mode) but within the process is known as User time and CPU busy with system or kernel task (supervised/kernel mode) known as system time (Fig. 2).

## 5 Conclusion and Future Work

Map Reduce is widely accepted and well-known Programming Model in Hadoop Ecosystem. Model designed with Map Reduce is parallelized and executed on Machines (Commodity hardware). Nowadays, Hadoop framework has been used in wide variety of application such as text retrieval, Text Mining, Search Engines and Machine learning. In the proposed algorithm, author has reduced the searching time of a pattern which increases efficiency of data retrieval. In the research work, author has examined text file 4–10 MB of data that contains more than 30–80 pages of character data and compared with the Brute Force algorithm. Even though Brute Force algorithm does not have searching phase, our algorithm gives better result with less number of comparisons.

For the future work, we can extend our proposed algorithm with changes on hashing function to test it on a real time for Big Data and with different string algorithm Rabin–Karp, Knuth Morris Pratt and Boyer Moore. In this paper, Author has done experiment for unstructured data, but we can do for all type of data. We can use this algorithm for data having key-value pairs and surely will generate very good results as we have used Hashing.

## References

1. Patel, A.B., Birla, M., Nair, U.: Addressing Big Data Problem Using Hadoop and Map Reduce. In: Nirma University International Conference on Engineering, IEEE-Ahemdabad (2012)
2. Pandey, S.K., Dubey, N.K., Sharma, S.: A study on string matching methodologies. Int. J. Comput. Sci. Inf. Technol. **5**(3), 4732–4735 (2014)
3. Dwivedi, K., Dubey S.K.: Analytical review on Hadoop distributed file system. In: 5th International Conference—Confluence The Next Generation Information Technology Summit (Confluence), IEEE (2014)
4. Pandey, S., Tokekar, V.: Prominence of map reduce in big data processing. In: 4th International Conference on Communication Systems and Network Technologies, IEEE (2014)
5. Hinson J.: Simulation of Map Reduce with the Hash of Hashes Technique, pp. 1748–2014. Accenture Life Sciences, Berwyn, PA, USA (2012)
6. Singla, N., Garg, D.: String matching algorithms and their applicability in various applications. Int. J. Soft Comput. Eng. **1**(6), 218–222 (2012)
7. Pandey, S.K., Tiwari, H.K., Tripathi, P.: Hybrid approach to reduce time complexity of string matching algorithm using hashing with chaining. In: Proceedings of International Conference on ICT for Sustainable Development, vol. 1, pp. 185–193 (2015)

8. Yuan, L.: An improved algorithm for Boyer–Moore String matching in Chinese information processing. In: International Conference on Computer Science and Service System (CSSS), IEEE (2011)
9. Katsoulis, S.: Implementation of parallel Hash Join algorithms over Hadoop. Master of Science, School of Informatics, University of Edinburgh (2011)
10. Fuyao, Z.: A string matching algorithm based on efficient hash function. In: International Conference on Information Engineering and Computer Science, IEEE (2009)
11. Charras, C., Lecroq, T.: Exact String Matching Algorithm Animation Java, Laboratoire d'Informatique de Rouen, Faculté des Sciences et des Techniques, Université de Rouen, Mont-Saint-Aignan Cedex, France. http://www-igm.univmlv.fr/~lecroq/string/node1.html
12. Zha, X., Sahni, S.: Multipattern string matching on a GPU. In: IEEE Symposium on Computers and Communications, IEEE (2011)
13. Sunarso, F., Venugopal, S., Lauro, F.: Scalable Protein Sequence Similarity Search using Locality-Sensitive Hashing and Map Reduce. School of Computer Science and Engineering, The University of New South Wales, Australia (2013)

# Automatic Recognition of Bird Species Using Human Factor Cepstral Coefficients

**Arti V. Bang and Priti P. Rege**

**Abstract**  Identification of bird species based on their song is very important task from biodiversity point of view. In order to develop an automatic system of recognition of bird species, a system using signal processing and pattern recognition techniques has gained huge importance. In this paper, we compare the performance of mel frequency cepstral coefficients and human factor cepstral coefficients combined with time- and frequency-based features. Gaussian mixture models have been used for developing feature models, and maximum likelihood estimation is used for classification. Further, selective features have been used in order to increase the performance of the system. With the proposed method, a maximum accuracy of 97.72% has been achieved for a data set of ten bird species.

**Keywords**  Bird species recognition · Mel frequency cepstral coefficients
Human factor cepstral coefficients · Gaussian mixture modeling

## 1   Introduction

One of the important issues that government organizations have to deal with is the conservation of biodiversity [1]. Birds constitute an important component of biodiversity. Protection of endangered species is an important task required for conservation of biodiversity. It is based on monitoring the region in order to ascertain the presence of the species. Therefore, monitoring birds is an important activity, carried out for conservation of biodiversity.

Bird experts, called as ornithologists, identify the bird species from there sounds. Such tasks rely completely on the knowledge of the ornithologists. Moreover, this is time-consuming and tedious. These limitations of manual observations have

A. V. Bang (✉) · P. P. Rege
Department of Electronics and Telecommunication, College of Engineering, Pune, India
e-mail: arti.bang@viit.ac.in

P. P. Rege
e-mail: ppr.extc@coep.ac.in

given rise to automatic recognition of bird species from their audio recordings. This is a typical pattern recognition problem, which requires preprocessing, feature extraction and classification.

Bird songs are varied and divided into a set of hierarchical structures [2]. The basic sound of the bird song is called as 'syllable' that is made up of elements. A sequence of syllables is called a 'phrase.' A typical combination of phrases that occur repeatedly is called as bird 'song.'

## 2    Literature Survey

In the past, fairly good amount of work has been carried out by various researchers to automatically recognize the bird species. Anderson et al. [3] and Kogan and Margoliash [4] are the pioneers of this work. A number of researchers have, since then, used different parametric representations of the audio signals. Mel frequency cepstral coefficients (MFCC) [5–9] are the most widely used features for bird species recognition. Some of the researchers have combined MFCC with time and frequency or descriptive features [10–13]. Other audio features which have been proposed in the literature are wavelets [14, 15], linear predictive cepstral coefficients [8], tonal features [16], and fft derived features [3, 4]. For recognition purpose, various approaches that are widely applied in speech and audio classification have been used. Researchers have used dynamic time warping [3], hidden Markov model [5, 7, 14], Gaussian mixture model [3, 10], support vector machine [8], and decision trees [17].

In most of the previous work, MFCC and time and frequency features have been used to perform the task of automatic recognition of bird species. In this work, we propose a new feature set called as human factor cepstral coefficients (HFCC) [18] and compare its performance with the widely used MFCC. These feature sets have been combined with descriptive features to obtain better recognition accuracy. Further, the features are judiciously selected for specific classes in order to improve the performance.

## 3    Feature Extraction

### 3.1    Mel Frequency Cepstral Coefficients

MFCC are the most commonly used features in speech and audio classification. Human ear does not resolve frequencies linearly across the audio spectrum. The mel frequency scale is derived from the human perceptual system. Human ear acts like a filter bank. In mel scale filter bank, the center frequencies of first 10 filters are linearly spaced up to 1 kHz. Thereafter, the center frequencies are logarithmically

spaced. In this, the edge frequencies are the center frequencies of adjacent filters. A set of 24 filters is used in this work [19].

$$\text{mel}(f) = 2595\log_{10}\left(1 + \frac{f}{700}\right) \tag{1}$$

The evaluation of MFCC starts with the computation of magnitude of Fourier transform of the audio frame. The spectrum is multiplied by mel scale filter bank. Next, the log energy of each filter bank is computed, and finally, discrete cosine transform is computed to yield cepstral coefficients. Thirteen coefficients have been selected. In order to capture temporal change, delta coefficients (Δ) have also been computed, which gives a total of 26 coefficients.

### 3.2 Human Factor Cepstral Coefficients

HFCC is a biologically inspired feature extraction algorithm, proposed by Skowronsky and Harris [18] for speech recognition. In MFCC computation, filter spacing is determined by mel scale, and bandwidth depends on filter bank parameters. It is conveniently set by the center frequencies of adjacent filters. In HFCC computation, the filter spacing is according to mel scale only, but the difference is, that, the bandwidth is not dependent on the filter bank parameters. Rather, the bandwidth is according to the human auditory system. It uses Moore and Glasberg's [20] critical bandwidth defined by the equivalent rectangular bandwidth (ERB) equation:

$$\text{ERB} = 6.23f_c^2 + 93.39f_c + 28.52 \tag{2}$$

where $f_c$ is the center frequency in KHz. Thirteen cepstral and delta coefficients are computed which gives a total of 26 coefficients.

### 3.3 Descriptive Features

Birds have varied sounds, and the signal model is not known. The spectral characteristics differ largely, and therefore, time and frequency features [10, 11] are derived to parametrize the sound. These features are computed on frame basis and mean and variances of these features are computed over the entire syllable. Frame duration is 10 ms, and 50% overlap is used. Frames are multiplied with hamming window, and discrete Fourier transform (DFT) is applied. These features are:

Spectral centroid (SC): Spectral centroid is associated with brightness of the sound. Brighter sound: higher centroid. It is the center point of the spectrum.

$$SC = \frac{\sum_{k=0}^{N/2} k|X(k)|^2}{\sum_{k=0}^{N/2}|X(k)|^2} \tag{3}$$

where $X(k)$ is the DFT of the frame, $N$ is the size of the DFT and $k$ varies from 0, 1, $...N - 1$.

Bandwidth (BW): It is the width of the frequency band of the frame around the center point of the spectrum.

$$BW = \sqrt{\frac{\sum_{k=0}^{N/2}(k - SC)^2|X(k)|}{\sum_{k=0}^{N/2}|X(k)|^2}} \tag{4}$$

Spectral flux (SF): It is the 2 norm of difference between the spectra of adjacent frames. It measures the difference in spectral shape.

$$SF_i = \sum_{k=0}^{N/2}\left|\left\|X_{(i+1)} - X_i\right\|\right| \tag{5}$$

Spectral roll-off frequency (SRF): It is the point below which the 95% of the power of the frame resides. It is related to the skewness of the spectral shape.

$$SRF = \max\left(M\sum_{k=0}^{M}|X(k)|^2 < 0.95\sum_{k=0}^{N/2}|X(k)|^2\right) \tag{6}$$

Spectral flatness (SFM): Spectral flatness is the ratio of geometric mean to arithmetic mean of the spectrum and is expressed in dB scale. It measures the tonality of the sound. It gives high value for tonal sounds and low value for noisy sounds.

$$SFM = 10\log_{10}\frac{G_m}{A_m} \tag{7}$$

where

$$G_m = \sqrt[\frac{1}{N/2}]{\prod_{k=0}^{N/2}|X(k)|} \tag{8}$$

and

$$A_m = \frac{1}{N/2}\sum_{k=0}^{N/2}|X(k)| \tag{9}$$

Energy (E): All the syllables are normalized to amplitude range of [+1, −1], and the energy of the frame is computed as

$$E = \frac{1}{L}\sum_{n=0}^{L-1} x(n) \tag{10}$$

where $L$ is the length of the frame.

Zero crossing rate (ZCR): A zero crossing occurs when adjacent sample has opposite signs. ZCR is the number of zero crossings in the frame.

$$\text{ZCR} = \sum_{n=0}^{L-1} |\text{sgn}(x(n)) - \text{sgn}(x(n-1))| \tag{11}$$

where sgn, signum function is defined as

$$\text{sgn}(x(n)) = \begin{cases} +1 & x(n) > 0 \\ -1 & x(n) < 0 \end{cases} \tag{12}$$

Duration (T): The duration of the syllables is computed since it varies among the species.

Pitch (P): Pitch is the fundamental period of the audio signal. Each individual has a different pitch. Therefore, pitch is considered as an important parameter in analysis of audio signals. Pitch is computed by finding the time lag with the largest autocorrelation coefficient. The mean, maximum, and the minimum pitch values of the syllable are considered for recognition.

These descriptive features constitute a set of 18 features, viz. mSC, varSC, mBW, varBW, mSF, varSF, mSRF, varSRF, mSFM, varSFM, mE, varE, mZCR, varZCR, T, mP, maxP, minP.

# 4 Classification

Feature vectors are modeled using parametric or nonparametric methods. GMM is one of the widely used parametric methods used for modeling the feature vectors because of its ability to form better approximations to arbitrarily shaped probability density functions. GMM uses a mixture of Gaussian densities to model the feature vector. GMM is a weighted sum of Gaussian densities given by

$$p\left(\frac{x}{\lambda}\right) = \sum_{i=1}^{M} w_i b_i(x) \quad i = 1, 2, \ldots, M \tag{13}$$

where $x$ is a $d$-dimensional feature vector, $b_i(x)$ is the density of $i$th mixture component, and $w_i$ is the weight vector of $i$th component, and $M$ is the number of

Gaussian components of the mixture model. $b_i(x)$ is a multivariate Gaussian density function given by

$$b_i(x) = \frac{1}{2\pi^{d/2} \sum_i^{1/2}} \exp\left\{ -\frac{1}{2}(x - \mu_i) \sum_i^{-1}(x - \mu_i) \right\} \qquad (14)$$

where $\mu_i$ is the mean vector, and $\sum_i$ is the covariance matrix. Each bird species, $s$, is represented by a GMM and is denoted by $\lambda_s$. The parameters of GMM are:

$$\lambda_s = \left\{ w_i \mu_i \sum_i \right\} \quad i = 1, 2, ..., M \qquad (15)$$

The mixture weights satisfy the constraints $\sum_{i=1}^{M} w_i = 1$. The parameters of GMM are obtained using expectation maximization (EM) algorithm, iteratively. Cluster tool [21] is used to estimate GMM parameters of all the classes. For recognition, maximum likelihood (ML) estimation is used. The objective is to find the bird model $\lambda_s$, which gives the maximum a-posterior probability for the test sequence $X$.

$$s = \operatorname{argmax} p\left(\frac{\lambda}{X}\right) = \operatorname{argmax} p(\lambda_s) p\left(\frac{X}{\lambda}\right) \qquad (16)$$

where $s$ denotes the index of the bird syllable achieving the maximum a-posterior probability, $p(\lambda_s)$ is the a priori probability.

## 5    Database

The bird sound recordings have been collected from Internet Web site www.zeno-canto.org. All the bird sounds belong to the region of India. The database of ten bird species has been prepared, as shown in Table 1. The size of the database is about 1 GB, and duration is about 2 h. About 62% of the total data is used for training and 32% for testing. The training and testing syllables belong to different individuals. Such a system is called as individual independent recognition system.

All the recordings are converted to standard wave file format, 16 kHz sampling rate and 16 bits per sample. Different files are accompanied with different noise sounds. The background noise is removed using spectral subtraction algorithm [22].

Bird song file is segmented into syllable files as depicted in the following algorithm.

1. Compute the spectrogram of the song with a hamming window of size 160, FFT with size 512 and 50% overlap.

**Table 1** Database used in the study

| Name of the class | Total no. of samples | No. of test samples |
|---|---|---|
| Red-vented bulbul (CB) | 364 | 154 |
| Greater coucal (GC) | 603 | 247 |
| House crow (HC) | 453 | 160 |
| Hume's leaf warbler (HLW) | 222 | 89 |
| Hill partridge (HP) | 407 | 174 |
| Asian koel (AK) | 215 | 83 |
| Common myna (CM) | 251 | 88 |
| Rose-ringed parakeet (RRP) | 325 | 125 |
| House sparrow (HS) | 471 | 156 |

2. Locate the point with the highest amplitude. Mark the points on either side of this point until the magnitude of the points drops below a certain threshold. These points constitute a syllable. Store this as a separate file.
3. In the spectrogram, equate the magnitude of these points to zero.
4. Repeat step 2 and 3 to locate another syllable, and so on.

## 6 Results and Discussion

DF, MFCC, and HFCC features are computed at frame level with a frame duration of 10 ms and 50% overlap. Block diagram of evaluation of MFCC and HFCC features is as shown in Fig. 1. The results are evaluated by combining DF with MFCC and HFCC, separately. Thus, the length of the feature vector is 44. The feature vectors are modeled using GMM. The aim of using GMM is because of its ability to form better approximations to arbitrary-based densities. Optimal number of mixture models has been derived for each species which range from one to four. The parameters have been derived using EM algorithm. For recognition, ML estimation is used. About 62% of the total data is used for training and 32% for testing.



**Fig. 1** Block diagram of computation of MFCC and HFCC parameters

The species recognition accuracy (SRA) is evaluated as follows:

$$SRA = \frac{CS_A}{NS_A} \times 100 \qquad (17)$$

where $CS_A$ is the number of correctly classified syllables of a species, and $NS_A$ is the total number of test syllables of that species. The reliability (Rel) of a species is expressed as follows:

$$Rel = \frac{CS_A}{MS_A} \times 100 \qquad (18)$$

where $CS_A$ is the total number of correctly classified syllables of a particular species, and $MS_A$ is the total number of syllables recognized as that species. Overall accuracy (OA) of the system is given as follows:

$$OA = \frac{C_A}{M_A} \times 100 \qquad (19)$$

where $C_A$ is total number of correctly classified syllables of all the species, and $N_A$ is the total number of test syllables.

The OA obtained for MFCC + DF is 91.97% and that of HFCC + DF is 94.39%. It can be seen that the performance of HFCC is remarkably better than MFCC. This is because of the fact that the bandwidth of the filters of HFCC filter bank is derived from the ERB relation rather than filter parameters.

The confusion matrix with HFCC + DF features is as shown in Table 2. It is observed that the SRA of four species, viz. HLW, IP, AK, and CM is below 90%. The OA can be further improved by judiciously selecting discriminative features for certain species. The algorithm is as follows:

**Table 2**  Recognition result with HFCC + DF (44) features

|      | Number of correctly recognized syllables in the class | | | | | | | | | | RSA |
|------|------|------|------|------|------|------|------|------|------|------|------|
|      | CB   | GC   | HC   | HLW  | HP   | IP   | AK   | CM   | RRP  | HS   |      |
| CB   | 148  | 0    | 2    | 0    | 0    | 0    | 0    | 0    | 2    | 2    | 96.10 |
| GC   | 0    | 246  | 0    | 0    | 0    | 1    | 0    | 0    | 0    | 0    | 99.60 |
| HC   | 0    | 0    | 160  | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 100.00 |
| HLW  | 0    | 0    | 0    | 72   | 0    | 0    | 0    | 0    | 15   | 2    | 80.90 |
| HP   | 0    | 0    | 0    | 0    | 170  | 0    | 4    | 0    | 0    | 0    | 97.70 |
| IP   | 0    | 0    | 4    | 0    | 0    | 37   | 0    | 0    | 3    | 0    | 84.09 |
| AK   | 0    | 0    | 3    | 0    | 21   | 0    | 59   | 0    | 0    | 0    | 71.08 |
| CM   | 13   | 0    | 0    | 0    | 0    | 0    | 0    | 74   | 1    | 0    | 84.09 |
| RRP  | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 125  | 0    | 100.00 |
| HS   | 0    | 0    | 0    | 1    | 0    | 0    | 0    | 0    | 0    | 155  | 99.36 |
| Rel. | 91.93 | 100.00 | 94.67 | 98.63 | 89.01 | 97.37 | 93.65 | 100.00 | 85.62 | 97.48 | |
| OA: 94.39% | | | | | | | | | | | |

*Algorithm:*

1. Derive the GMM for each species.
2. Input the test bird syllable.
3. Perform recognition using three features mean pitch, maximum pitch and minimum pitch. The confusion matrix is as shown in Table 3. It can be seen that the SRA and the reliability of GC and HLW species are maximum. Hence, if the test syllable belongs to these species, the syllable is correctly classified with maximum accuracy and reliability, hence stop. Else, the syllable does not belong to GC and HLW class. Hence, remove these two classes.
4. Perform recognition using all 44 features and eight classes. The confusion matrix is as shown in Table 4. Now the recognition accuracy and reliability of HC, RRP, and HS are maximum. If the syllable is classified to either of these classes, stop. Else, remove these three classes.

**Table 3** Recognition results with three pitch features and 10 classes

| | Number of correctly recognized syllables in the class | | | | | | | | | | RSA (%) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | CB | GC | HC | HLW | HP | IP | AK | CM | RRP | HS | |
| CB | 116 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 23 | 15 | 75.32 |
| GC | 0 | 246 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 99.60 |
| HC | 0 | 0 | 133 | 0 | 2 | 16 | 9 | 0 | 0 | 0 | 83.13 |
| HLW | 0 | 0 | 0 | 87 | 0 | 0 | 0 | 0 | 0 | 2 | 97.75 |
| HP | 0 | 0 | 6 | 0 | 89 | 0 | 76 | 0 | 3 | 0 | 51.15 |
| IP | 0 | 0 | 2 | 0 | 0 | 42 | 0 | 0 | 0 | 0 | 95.45 |
| AK | 0 | 0 | 0 | 0 | 23 | 1 | 57 | 0 | 2 | 0 | 68.67 |
| CM | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 76 | 8 | 0 | 86.36 |
| RRP | 12 | 0 | 3 | 0 | 0 | 0 | 0 | 0 | 106 | 4 | 84.80 |
| HS | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 154 | 98.72 |
| %Rel | 87.88 | 100.00 | 91.72 | 97.75 | 78.07 | 71.19 | 40.14 | 100.00 | 74.65 | 88.00 | |

**Table 4** Recognition results with 44 features and excluding GC and HLW species

| | Number of correctly recognized syllables in the class | | | | | | | | RSA (%) |
|---|---|---|---|---|---|---|---|---|---|
| | CB | HC | HP | IP | AK | CM | RRP | HS | |
| CB | 148 | 2 | 0 | 0 | 0 | 0 | 2 | 2 | 96.10 |
| HC | 0 | 160 | 0 | 0 | 0 | 0 | 0 | 0 | 100.00 |
| HP | 0 | 0 | 170 | 0 | 4 | 0 | 0 | 0 | 97.70 |
| IP | 0 | 4 | 0 | 37 | 0 | 0 | 3 | 0 | 84.09 |
| AK | 0 | 3 | 21 | 0 | 59 | 0 | 0 | 0 | 71.08 |
| CM | 4 | 0 | 0 | 0 | 0 | 74 | 1 | 0 | 87.50 |
| RRP | 13 | 0 | 0 | 0 | 0 | 0 | 125 | 0 | 100.00 |
| HS | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 156 | 100.00 |
| %Rel | 91.93 | 94.67 | 89.01 | 100.00 | 93.65 | 100.00 | 95.42 | 98.73 | |

**Table 5** Recognition results with 41 features, excluding pitch features, and five species

|  | Number of correctly recognized syllables in the class | | | | | RSA (%) |
|--|-----|-----|-----|-----|-----|--------|
|  | CB | HP | IP | AK | CM | |
| CB | 148 | 0 | 3 | 0 | 3 | 96.10 |
| HP | 0 | 173 | 1 | 0 | 0 | 99.43 |
| IP | 0 | 0 | 40 | 0 | 4 | 90.91 |
| AK | 0 | 4 | 5 | 74 | 0 | 89.16 |
| CM | 2 | 0 | 0 | 0 | 86 | 97.73 |
| %Rel | 98.67 | 97.74 | 81.63 | 100.00 | 92.47 | |

5. Perform recognition using all the features except mean pitch, maximum pitch, and minimum pitch, i.e., 41 features and remaining five classes. The confusion matrix is as shown in Table 5.

The pitch of different bird species is quite different, except HP and AK. Hence, their SRA is low when pitch features are included. Therefore, at the end, pitch features are excluded, and the SRA of these species is enhanced. Finally, the OA obtained is 97.72%.

# 7    Conclusion

In this paper, we have addressed the problem of bird species recognition. The work is based on the perceptual audio features such as MFCC and HFCC. These features have been combined with DF, independently. The features have been modeled using GMM. Recognition is achieved using ML estimation algorithm. The algorithms were implemented in MATLAB environment, and performance of 97.72% has been achieved.

There is no standard data set available for bird species recognition. Our aim is to collect more recordings and also increase the number of bird species. We will also go for more efficient and robust feature extraction and classification approaches to improve performance under noisy conditions.

# References

1. Brandes, T.S.: Automated sound recording and analysis techniques for bird surveys and conservation. Bird Conserv. Int. **18**, S163–S173 (2008). doi:10.1017/S0959270908000415
2. Kumar, A.: Acoustic communication in birds, differences in songs and calls, their production and biological significance. Resonance **6**, 44–54 (2003)
3. Anderson, S.E., Dave, A.S., Margoliash, D.: Template-based Automatic recognition of birdsong syllables from continuous recordings. J. Acoust. Soc. Am. **100**(2), 1209–1219 (1996)

4. Kogan, J., Margoliash, D.: Automated recognition of bird song elements from continuous rcordings using dynamic time warping and hidden Markov models: a comparative study. J. Acoust. Soc. Am. **103**(4), 2187–2196 (1998)
5. Kwan, C., Ho, K., Mei, G., Li, Y, Ren, Z.: An automated acoustic system to monitor and classify birds. EURASIP J. Appl. Signal Process. **2006**, 1–19 (2006). Article ID 96706
6. Lee, C.H., Han, C.H.: Automatic classification of bird species from their sounds using two-dimensional cepstral coefficients. IEEE Trans. Speech Audio Process. **16**(8), 1541–1550 (2008)
7. Trifa, V.M., Kirschel, A., Taylor, C.E.: Automated species recognition of antbirds in a Mexican rainforest using hidden Markov models. J. Acoust. Soc. Am. **103**(4), 2424–2431 (2008)
8. Castano, G.V., Rodriguez, G., Castillo, J., Lu, K., Rios, A., Bird, F.: A Framework for bioacoustical species classifications in a versatile service-oriented wireless mesh networks. In: 18th European Signal Processing Conference, Aug 23–27 (2010)
9. Quian, K., Zhiang, Z., Ringeval, F., Schuller, B.: Bird sounds classification by large scale acoustic features and extreme learning machine. In: IEEE Global Conference on Signal and Information Processing (GlobalSIP) (2015)
10. Somervuo, P., Harma, A., Fagurland, S.: Parametric representations of bird sounds for automatic species recognition. IEEE Trans. Audio Speech Lang. Process. **14**(6), 2252–2263 (2006)
11. Fagurland, S.: Bird species recognition using support vector machine. EURASIP J. Adv. Signal Process. Article ID 38637, 8 p (2006). doi:10.1155/2007/38637
12. Briggs, F., Laxminarayanan, B., Neal, L., Fern, X.Z., Raich, R.: Acoustic classification of multiple simultaneous bird species: a multi-instance multi-label approach. J. Acoust. Soc. Am. **131**, 4640–4650 (2012)
13. Mporas, I., Ganchev, T., Kocsis, O., Fakotakis, N., Jahn, O., Riede, K.: Automated acoustic classification of bird species from real-field recordings. In: 24th IEEE International Conference on tools with Artificial Intelligence (2012). doi:10.1109/ICTAI.2012.110
14. Selin, A., Turunen, J., Tanttu, T.: Wavelets in recognition of bird sounds. EURASIP J. Adv. Signal Process. Article ID 51806, 9 p (2007). doi:10.1155/2007/51806
15. Bang, A.V., Rege, P.P.: Classification of bird species based on bioacoustics. Int. J. Adv. Comput. Sci. Appl. **4**(1), 184–188 (2014)
16. Jancovic, P.: Automatic detection and recognition of tonal bird sounds in noisy environments. EURASIP J. Adv. Signal Process. (2011). doi:10.1155/2011/982936
17. Vilches, E., Escobar, I.A., Vallejo,E.E.,Taylor, C. E.: Data mining applied to acoustic bird species recognition. In: 18th International Conference in Pattern Recognition (2006). doi: 10.1109/ICPR.2006.426
18. Skowronski, M., Harris, J.: Improving the filter bank of a classic speech feature extraction algorithm. In: IEEE Internatinal Symposium on Circuits and Systems, Bangkok, Thailand, vol. IV, pp. 281–284, May 25–28 (2003)
19. Ganchev, T., Fakotakis, N., Kokkinakis, G.: Comparative evaluation of various MFCC implementations on the speaker verification task. In: Proceedings of the SPECOM, pp. 191–194 (2005)
20. Moore, B.C., Glasberg, B.R.: Suggested formulae for calculating auditory-filter bandwidths and excitation patterns. J. Acoust. Soc. Am. **74**(3), 750–753 (1983)
21. Bouman, C.A., Shapiro, M., Cook, G.W., Atkins, C.B., Cheng, H.: Cluster: an unsupervised algorithm for modeling gaussian mixtures. http://dynamo.ecn.purdue.edu/~bouman/software/cluster (1998)
22. Boll, S.F.: Suppression of acoustic noise in speech using spectral subtraction. IEEE Trans. Acoust. Speech Signal Process. **27**, 113–120 (1979)

# A Novel Meta Crawling Algorithm for Terrorist Network Knowledge Aggregation from the Internet

**R. D. Gaharwar and D. B. Shah**

**Abstract** Nowadays, World Wide Web is flooded with different types of web browsers. Each of these web browsers differs from one another on the basis of search time, search efficiency and the area of the Internet they cover in their search. The popular search engines like Yahoo Search, Google, Dogpile, Bing, Ask, etc., are facing efficiency challenges due to tremendous growth rate of data on the Internet. Moreover, the current state of the art in this field works on the general search results. Hence, special data mining search results require customized search services. Terrorist network mining is one such field of data mining which requires customized search services. This paper presents the algorithm for customized Terrorist Meta Crawler. The algorithm presented in this paper is optimized to search terrorist-related information on World Wide Web using different web services. The last section of this paper presents the comparative effectiveness of Terrorist Meta Crawler along with other popular search engines, and results show Terrorist Meta Crawler is a better solution for customized search.

## 1 Introduction

### 1.1 Search Crawlers

End-users employ search engines to browse their query on World Wide Web. Search engines act as an interface between end-user and information on the Internet

---

R. D. Gaharwar (✉) · D. B. Shah
G H Patel Post Graduate Department of Computer Science and Technology, Sardar Patel
University, Vallabh Vidyanagar, India
e-mail: raina.2611@gmail.com

D. B. Shah
e-mail: dbshah66@yahoo.com

generating the output for user query. Search engines like Yahoo Search, Google, Bing, Dogpile, Ask, etc., are among the most popular search engines used by people around the world. Each of these popular search engines uses special code scriptlets to search user query and generate output. These special code scriptlets are called search crawlers/web spiders. Search Crawler searches the web pages on the Internet for user query. All the related web pages are presented to end-user as search result. These code scriptlets index the end-user query related all web pages for efficient search [1].

## 1.2 Meta Crawlers

Meta Crawler is a software application that uses different search services to accomplish its task. The end-user can post a query on the central user interface provided by Meta Crawler [2]. The user query is transfer to multiple search engines through parallel running programs. Meta Crawler collects all the links returned and checks for the relevance of the web pages returned. The links so collected can be further used to load web pages to generate other related and relevant references for the user query. Meta Crawler acts like an intelligent agent available on the web which engages web crawlers for different search engines. A large number of Meta Crawlers are available on internet for end-users. These Meta Crawlers are used to aggregate the results and effectiveness of different search services.

## 1.3 Terrorist Network Mining

Networks are the graphs consisting of nodes and edges connecting the nodes. The terrorists/terrorist organizations form network through communication links during their operations. Such networks are called terrorist networks. These networks are special network/graph where nodes are terrorists/terrorist organizations, and the links between the nodes are the association between these terrorists/terrorist organizations [3]. Terrorist networks are studied to understand the communicational associations between these nodes. The communicational association forms the patterns which decide the role of each terrorist/terrorist organization in the organized crime. Understanding the role of each terrorist/terrorist organization in the network is important for planning effective counter-terrorist activities. The study of terrorist networks for understanding and aggregation of information is called Terrorist Network Mining.

## 1.4  SPAM Issue

Most of the search services suffer from spamming problems. Spamming is the technique of illegitimately increasing the rank of any page in search result. There are different types of spamming techniques available like hidden text, keyword spamming, Meta tag stuffing, hidden links, etc. [4]. These spamming techniques increase the importance of irrelevant web page in the search operation. One of the solutions to spamming is Meta Crawler. Meta Crawler fires parallel query to multiple search engines; hence, the chances of spamming infected search result decreases. Meta Crawler ensures that search result is not completely depended on the output of a single search engine. Hence, the diversifications of different search engines are kept under check. This leads to efficient and effective search results.

## 2  Previous Findings

The several researchers are working in this area. The study says that some authors discussed the important problem related to the fusion of output returned by different search engines. They classified the potential fusion cases for Meta search engines in seven categories, i.e. non-conflict equivalent case, conflict equivalent case, non-conflict inclusion case, conflict inclusion case, disjoint case, non-conflict overlap case and conflict overlap case.

Later on, these authors also categorized Meta search engines into categories like partially homogeneous, heterogeneous, etc., based on the type of fusion case they employ. Yang et al. also stated the necessary constraints in Meta search engines for merge methods. Finally, they concluded that due to the unavailability of previous fusion methods, the reliability of existing merge algorithm cannot be checked completely [5].

The authors stated that any Web Crawler differs from search engines in many different ways. Any search engine does three activities: indexing, storage and retrieval. However, Meta Crawler does not store anything. It just uses the different search services as platform for the retrieval of web pages and hence does not maintain any Internet database of its own. The authors also stated if any Meta Crawler simply collates the reference returned by search engines then its probability of giving obsolete and unrelated references also increases. Hence, Meta Crawler should not just collate the resultant references but should also return only quality references which are useful to end-users. Moreover, authors concluded that though the Meta Crawler presented by them is very useful and popular; new web services will follow and will continuously evolve [6].

Some authors came up with the information retrieval model which uses its own indexing and filtering schemes. The output of the search is clustered among different groups depending on its similarities and differences. They used different

types of clustering parameters like depth of indexes minimum similarity value, size of the indexing vocabulary, data slot size and rounds value of similarity. However, authors have discussed their concern about the success of this clustering technique on a small group of document. They stated that it might be possible that clustering technique is not as effective on small group of document as on large group [7].

The authors designed intelligent interface by providing users with the option to search any word, all words from the query or searching for the results that do not have search words at all. Meta Crawler Softbot tries to induce search engines with the features that they did not have before like Lycos search engine do not support phrase web searching capability, but this softbot simulates Lycos to perform phrase searching. Moreover, it uses novel algorithm for the detection of duplication of references. It compares domain in URL or content of the web page. This reduces the time for duplication detection. The authors concluded that this architecture is adaptable, portable and scalable, but there is still scope for improvement [8].

Hence, in this paper, authors present a novel Meta Crawlers algorithm for Terrorist Network Knowledge Aggregation from the Internet.

## 3   Algorithm

A search engine can be denoted with 4-tuple as $SE = \langle U, Q, I \rangle$ where $U$ is a unique identification number for any web page on the server, $Q$ is the set of queries using that a particular search engine can give search results, $I$ is the internal indexing algorithm used by that particular search engine to rank web pages in the search result.

The algorithm for the generation of search engine specific URL for search engines like Yahoo Search, Dogpile, Bing, Ask and MySearch is as follows:

Algorithm 1:
Input: Q where Q is user query
Procedure:
```
Let

SE be an array indicating 5 different search engines
sufix = ∑si where i = 1,2,...,n and s = fixed string
prefix = fixed string
URL = ∑"Ti" where Ti C Q
index = 1 where index indicates the number of web pages fetch
for any particular search engine
link_queue be an queue storing generated search engine
specific URL
do while SE has element
```

```
step-1:
```

```
URL = prefix + URL + index + suffix;
index ++;
endqueue(link_queue,URL);
repeat step-1 3 times
```

```
end do
```

Output: list of search engine specific URLs

Any Meta Crawler can be denoted with 5-tuple as $MC = \langle S,Q,C,I,r \rangle$ where $S$ is the list search engines that particular Meta Crawler may use to get search results, $Q$ is the set of queries that Meta Crawler give output for, $C$ is the collation algorithm for fusion of result from different search engines, $I$ is an indexing algorithm used to give a unique identification number to each page in result set, $r$ is ranking algorithm deployed by Meta Crawler for generate rank for web pages.

In case of Terrorist Meta Crawler, $S$ is the set of search engines like Yahoo Search, Dogpile, Bing, Ask and MySearch. Any query that is given to Terrorist Meta Crawler will be search on above-mentioned search engines. For generation of search engine specific queries, Terrorist Meta Crawler uses Algorithm 1. This algorithm will return the list of 15 URLs stored in a link array; three URL for each search engine will be generated and stored.

$Q = Q1 + Q2 + Q3 + Q4 + Q5$ where $Qi$ is the set of queries of above-mentioned search engines. Hence, query set for Terrorist Meta Crawler is the union of query set of each individual search engine.

$C = R1 + R2 + R3 + R4 + R5$ where $Ri$ is the result set returned by each search engine. Final result set generated for Terrorist Meta Crawler will be the union of result set generated by each individual search engine.

Terrorist Meta Crawler uses minimum rank removal technique for filtration of web pages in the result set. Terrorist Meta Crawler allocates ranks to each web page. Rank basically denotes the significance of each reference for a particular user query. Rank is allocated as follows:

Rank = Number of times that reference appears in the result set.

Algorithm 2:
Algorithm for Terrorist Meta Crawler is as follows:
Input: link_queue having list of 15 search specific URLs.
Procedure:

```
Let
```

```
While(not empty(link_queue))
url = Get(link_queue)
SourceCode = Download(url)
If (Search(SourceCode == true))
```

```
End_queue(link_queue,url)
url_queue = Extract(url)

for each u in url_queue

End_queue(link_queue,u)
```

Output: list of terrorist-related URLs

Remark 1: Here,

*End_queue() is the function that adds element at the end of the queue
*Get() is the function that removes and returns the first element of the queue
*Search() is the function that parses the source code of any web page to find
whether it has relevant terrorist-related data
*Download() is the function that the downloads the source code for any web page
from the Internet
*Extract() is the function that the URLs from the source code of any web page

## 4   Experiment, Result and Evaluation

Experiment is carried out using Algorithm 1 and 2, and result is collected. Table 1
shows results of Algorithm 1, i.e. the search engine specific URLs generated by
Terrorist Meta Crawler.

Table 2 shows the top 10 web page links which are used to aggregate the
terrorist organization related information.

Following Fig. 1 shows the comparative evaluation of the outputs of different
search services with Terrorist Meta Crawler for terrorist-related information.
Terrorist-related queries were fired on Ask, Yahoo Search, MySearch, Bing,
Dogpile and Terrorist Meta Crawler (TMC) which produce 61, 55, 75, 62 and 56%
relevant links, respectively, whereas Terrorist Meta Crawler (TMC) gave 94%
relevant links. This figure graphically shows that the number of the relevant links
searched by Terrorist Meta Crawler (TMC) is considerably high. Terrorist Meta
Crawler searches for terrorist network related information more accurately with
compared to other general web search services.

**Table 1** Showing the results of Algorithm 1

| No. | URLs | Search engine |
|---|---|---|
| 1 | http://www.ask.com/web?q=Terrorist+attacks+in+india&page=1&qid=D923E7B97651A84EC7A97378902968C17&o=0&l=dir&qsrc=998&qo=pagination | Ask |
| 2 | http://www.ask.com/web?q=Terrorist+attacks+in+india&page=2&qid=D923E7B97651A84EC7A97378902968C17&o=0&l=dir&qsrc=998&qo=pagination | ASK |
| 3 | http://www.ask.com/web?q=Terrorist+attacks+in+india&page=3&qid=D923E7B97651A84EC7A97378902968C17&o=0&l=dir&qsrc=998&qo=pagination | ASK |
| 4 | https://in.search.Yahoo Search.com/search;_ylt=A2oKmJBeH0RXu3kAHia7HAx.:_ylu=X3oDMTEzZZNnZms4BGNvbG8Dc2zBHBvcwMxBHZ0aWQDBHNIYwNwYWdpbmF0aW9u?p=Terrorist+attacks+in+india&fr=sfp&fr2=sb-top-in.search&b=1&pz=10&bct=0&xargs=0 | Yahoo! Search |
| 5 | https://in.search.Yahoo Search.com/search;_ylt=A2oKmJBeH0RXu3kAHia7HAx.:_ylu=X3oDMTEzZZNnZms4BGNvbG8Dc2zBHBvcwMxBHZ0aWQDBHNIYwNwYWdpbmF0aW9u?p=Terrorist+attacks+in+india&fr=sfp&fr2=sb-top-in.search&b=11&pz=10&bct=0&xargs=0 | Yahoo! Search |
| 6 | https://in.search.Yahoo Search.com/search;_ylt=A2oKmJBeH0RXu3kAHia7HAx.:_ylu=X3oDMTEzZZNnZms4BGNvbG8Dc2zBHBvcwMxBHZ0aWQDBHNIYwNwYWdpbmF0aW9u?p=Terrorist+attacks+in+india&fr=sfp&fr2=sb-top-in.search&b=21&pz=10&bct=0&xargs=0 | Yahoo! Search |
| 7 | http://search.mysearch.com/web?ts=1464082527144&tpr=10&q=Terrorist+attacks+in+india&page=1&ots=1464082547632 | MySearch |
| 8 | http://search.mysearch.com/web?ts=1464082527144&tpr=10&q=Terrorist+attacks+in+india&page=2&ots=1464082547632 | MySearch |
| 9 | http://search.mysearch.com/web?ts=1464082527144&tpr=10&q=Terrorist+attacks+in+india&page=3&ots=1464082547632 | MySearch |
| 10 | https://www.bing.com/search?q=Terrorist+attacks+in+india&qs=AS&sk=AS2&pq=terrorist+a&sc=8-11&sp=3&cvid=137c7a93c77e494aac07d0eb62d1f23d&first=1&FORM=PERE | Bing |
| 11 | https://www.bing.com/search?q=Terrorist+attacks+in+india&qs=AS&sk=AS2&pq=terrorist+a&sc=8-11&sp=3&cvid=137c7a93c77e494aac07d0eb62d1f23d&first=11&FORM=PERE | Bing |
| 12 | https://www.bing.com/search?q=Terrorist+attacks+in+india&qs=AS&sk=AS2&pq=terrorist+a&sc=8-11&sp=3&cvid=137c7a93c77e494aac07d0eb62d1f23d&first=21&FORM=PERE | Bing |
| 13 | http://www.dogpile.com/info.dogpl.control/search/web?&q=Terrorist+attacks+in+india&qsi=1&fcop=results-bottom&fpid=2 | Dogpile |
| 14 | http://www.dogpile.com/info.dogpl.control/search/web?&q=Terrorist+attacks+in+india&qsi=11&fcop=results-bottom&fpid=2 | Dogpile |
| 15 | http://www.dogpile.com/info.dogpl.control/search/web?&q=Terrorist+attacks+in+india&qsi=21&fcop=results-bottom&fpid=2 | Dogpile |

**Table 2** Showing the results of Algorithm 2

| Sr. No. | SiteUrl | Rank |
|---|---|---|
| 1 | www.satp.org/satporgtp/countries/india/database/index.html | 10 |
| 2 | www.ndtv.com/topic/india-terror-attacks | 9 |
| 3 | www.theguardian.com/world/mumbai-terror-attacks | 7 |
| 4 | en.wikipedia.org/wiki/Terrorism_in_India | 5 |
| 5 | www.timesofindia.indiatimes.com/topic/terror-attack | 5 |
| 6 | en.wikipedia.org/wiki/List_of_terrorist_incidents_in_India | 3 |
| 7 | www.ibtimes.com/major-terrorist-attacks-india-over-last-20-years-timeline-1752731 | 3 |
| 8 | www.chacha.com | 2 |
| 9 | en.wikipedia.org/Terrorist_incidents_in_India | 2 |
| 10 | http://www.indiabix.com/group-discussion/terrorism-in-india/ | 1 |

**Fig. 1** Comparative evaluation of the outputs of different search services with Terrorist Meta Crawler



## 5   Conclusion

This paper presents the need for a customized search services for the special data mining application like Terrorist Network Mining. In this paper, optimized algorithm for such Terrorist Meta Crawler application is presented. First algorithm presents the logic for generating five popular search engine specific URLs. This algorithm generates three URLs for first three search pages of search engines like Yahoo Search, Dogpile, Bing, Ask and MySearch. These URLs will be passed as input for second algorithm which extracts the source code and parses these pages for getting other terrorist-related links. Hence, the output of second algorithm is terrorist-related links. During the experiment, authors found out that Terrorist Meta Crawler gives effective output with 94% accuracy. This paper presents the comparative evaluation of the outputs of different search services with Terrorist Meta Crawler which shows the effectiveness of the newly developed Terrorist Meta Crawler as customized terrorist-related search service.

# References

1. Gaharwar, R.D., Shah, D.B.: Architectural design of Meta Crawler for Terrorist Network Mining. Commun. Appl. Electron. **5**(8), 37–40 (2016)
2. Etzioni, O.: Moving up the information food chain: deploying softbots on the world wide web. AI Mag. **18**(2), 11 (1997)
3. Gaharwar, R.D., Shah, D.B., Gaharwar, G.K.S.: The study of multi-search services for Terrorist Network Mining. Commun. Appl. Electron. **5**(8), 37–40 (2016)
4. Links & Law: Information about legal aspects of search engines, linking and framing. http://www.linksandlaw.com/technicalbackground-search-engine-spamming.htm
5. Yang, X., Zhang, M.: Necessary constraints for fusion algorithms in meta search engine systems. In: Proceedings of the International Conference on Intelligent Technologies, pp. 409–416 (2000)
6. Selberg, E., Etzioni, O.: Multi-service search and comparison using the MetaCrawler. In: Proceedings of the Fourth International WWW Conference, Boston (1995)
7. Laria, V.G., Griffiths, R., Winstanley, G.: Application of a clustering algorithm to recover topic content in an unstructured text-based environment. https://www.researchgate.net/profile/Graham_Winstanley/publication/2504349_Application_of_a_Clustering_Algorithm_to_Recover_Topic_Content_in_an_Unstructured_TextBased_Environment/links/02e7e52f231abb5cac000000.pdf
8. Selberg, E., Etzioni, O.: The MetaCrawler architecture for resource aggregation on the Web. IEEE Expert **12**(1), 11–14 (1997)

# Password Security Using Bcrypt with AES Encryption Algorithm

**Narander Kumar and Priyanka Chaudhary**

**Abstract** With the advancement of technology, the Internet has become a widely used tool of communication. Million numbers of individual all over in the world can get the utilization of technology. Novel issues like cyber stalking have been increasing worldwide global attention. Cyber stalking can be explained as threatening behavior or undesirable advances intended for another using the Internet and other way of online communications, so client authentication in computer systems is an essential feature in the present time for avoiding cyber stalking. In this paper, we have scheduled a technique utilizing Bcrypt hashing technique with AES encryption for securing an online account and reducing cyber criminal activity.

## 1 Introduction

In the evaluation of new technologies and advancement has enhanced our lives incalculable manner. Likewise, expand dependency on IT and communication approach for dynamic and active field arrangements has its own particular reaction. The impact of computerized data advancements upon the world surely postures unlimited advantages for the nationals of the developing worldwide town. The dark side of it, as anyone might expect, is misappropriating of Information Technology for criminal exercises. Cyber stalking is another sort of violation that existed subsequent to the late 1990s that raised a significant universal criminological type of problem in 2004. Generally, cyber stalking depicts the utilization of ICT keeping

N. Kumar (✉) · P. Chaudhary
Department of Computer Science, B.B. Ambedkar University (A Central University),
Lucknow, Uttar Pradesh, India
e-mail: nk_iet@yahoo.co.in

P. Chaudhary
e-mail: cpriyanka22@gmail.com

in mind the end goal to bug one or more martyr [1]. Furthermore, molestation implies any conduct that causes the martyr trouble whether deliberate or not. Cyber stalking frequently discovers their martyr online [2] as they utilize computer and network systems for criminal exercises, as these advances can without much of a stretch be misusage to startle, threaten, force, annoy, and deceive clueless clients. Password-based authentication protocols cannot depend with respect to persevering put away data on the customer side. Rather, they depend on clients' capacity of exact review of a secret key data. It is primarily because of this exact review prerequisite that clients regularly pick basic and low entropy passwords that are anything but difficult to recollect [3]. The vulnerability of passwords turns into the feeble connection of the system, which assailants misuse by dispatching offline or online dictionary attack.

The password hashing strategy is superior than encryption of secret key since hashing has contained single functionality we cannot find plain text content from its hash implies the plain password that builds hash cannot be reconstructed from its containing hash value. The method of hashing is delicate to dictionary assault. Dictionary assault is a strategy for recouping password from known password. So it is conceivable to split hash password by utilizing pre-computed hash value. A hashing algorithm is exceptionally deterministic as they deliver same hash value for same inputted content. Raw hashes have likewise helpless against rainbow tables, a technique for adjusting a requirement for a pre-calculation of hashes and the clearly substantial function is important to keep a whole dictionary containing hashes [4]. For averting from this type of issues, salting acts as a rescue. Salt guarantees that assailants cannot utilize specific assaults like lookup tables and rainbow tables to break huge collection of hashes rapidly, yet it does not keep them from running a dictionary assault or brute force attack on each hash independently. High-end graphic cards and custom equipment could be processed on billions of hashes each second, so these assaults are still extremely powerful. To make these assaults less powerful, they can utilize a method known as key stretching for whatever length of time that an assailant can utilize a hash to check whether a password is correct or wrong, they could be run a dictionary or brute force attack on the hash. The objective of the proposed work is to fortify existing password-based verification conventions against brute force attack for securing client accounts that are the best approach for the user side to memorize and securing password utilizing hashing technique with AES for online account protection. Thus, it is more secure to assault from instance hacking, phishing, fraud.

## 2   Review of Work

P. Sriramya and R. A. Karthika implemented an algorithm which is used in Bcrypt algorithm with salt hash technique provides more security to users for online shopping [5]. Yu-Chi Chen et al. formulize the blind decoding schema, i.e., planned as technique being secured end client protection in online looking for electronic

proof [6]. Salmin Sultana et al. studied on the security of wireless sensor system, wireless sensor network quick created and utilized as a part of numerous divisions. Subsequently, the need for security turns out to be an exceptionally vital component. There are numerous advances being accessible to give security against the assailants, one of the best innovations is cryptography [7]. Halderman et al. [8] studied, an approach, cryptographic hash function have utilized to prepare secure passwords for subjectively various records, and customer requires to hold short secret passwords [9]. These types of procedure perform task absolutely on the user side, server side adjustment is not required. Khiyal et al. [10] planned a system of password hashing, i.e., process for safe passwords. A recent approach taking into the advanced encryption standard for securing a password is designed in [11]. A symmetric cipher technique that depends on the Rijndael technique is designed in [12] yet this technique begins with 200 bits. M.S.H. Khiyal et al. planned an instrument for securing passwords, i.e., using MD5 and SHA1 cryptographic hash function [13]. The function is utilized to infer one or more secret keys from a secret string. It depends on memory-hard function which offers some additional insurance against assaults utilizing custom hardware [14]. The SHA-192 can be utilized as a part of numerous functions such as an open key cryptosystem, advanced sign-cryption, content verification, approach random generator and in security engineering of forthcoming remote gadgets like programming characterized radio, and so forth [15]. SSH protocol is planned as a substitution for the current rsh, rlogin, rcp, rdist, and telnet protocol [16]. A paper is proposed to deal with security against phishing assaults with "BogusBiter" is discussed in [17]. F. Mwagwabi et al. examination approaches are scheduled during this study have additionally appeared to be a valuable approach to clarifying IS security behavioral intentions [18]. S. Farmand et al. proposed exploration are a way to deal with upgrade the current graphical password strategies and oppose against assaults like shoulder surfing [19].

From an existing review of the work, we have found some vulnerabilities salt guarantees that intruders cannot use specific attacks like lookup tables and rainbow tables to break huge collections of hashes rapidly; however, it does not keep them from running dictionary or brute force attacks on each and every hash independently.

## 2.1 Brute Force Attack

At the point, when a brute force attack to each possible set of characters up to a given length. This kind of assault is computationally excessive, the minimum expert with respect to as hashes cracked per processor time; however, they will dependably and consequently perform the password.

## 2.2 Salt Collision

Salt collision happens when two passwords encrypted utilizing with related same salt value. For developing dictionary attacks, an intruder can be grouped with ciphered passwords through the salt and hash each and every applicant password from a dictionary just per salt. The outcome speedup can be determined as follows:

$$\frac{\text{Number of Password}}{\text{Number of Different Salt}} \qquad (1)$$

In the event that salts are produced by a random number generator, the accepted number of various salts for $n$ password entries with $s$ salts is

$$\text{EV}(n, s) = \sum_{i=0}^{n-1} \left(\frac{s-1}{s}\right)^i = s - (s-1)^n s^{1-n} \qquad (2)$$

## 3 Proposed Method

In this proposed method, we use Bcrypt hashing algorithm using the AES encryption technique if an attacker can utilize a hash to check whether a password guessing theory is correct or incorrect, they can run a word reference table or brute force attack on the hash. The hash function that will be decided by this value.

In proposed model, defined in Fig. 1, user opens a login form. Store user information in a list including with user name and password, as well as the necessary information with a password. We apply Bcrypt hashing algorithm. After applying Bcrypt hashing, we use AES encryption algorithm. The list of data would be saved in the format of windows registry. A user could be open his list by entering a valid password and then applying decryption techniques, and also, check the user is valid or not and provides access to the list of data.

Bcrypt algorithm that is discussed in [6] in which initial key is a key stretching as a password then aggressor will be first to try each word in the dictionary or common password list and then try to check all possible character combinations for longer password. In our proposed method, this problem can be reduced using with AES encryption technique.

## 4 Implementation

We designed our proposed algorithm using with NetBean IDE 8.0 software. The below figures define the implementation of the proposed algorithm by using different numbers of text data values and sizes of a wide range. Figure 2 is defined the

**Fig. 1** Flowchart of proposed algorithm

login page of the user. Figure 3 demonstrates that generating salt with a hash of the corresponding user password with the help of Bcrypt hashing algorithm. Figures 4 and 5 demonstrates that application of AES encryption technique to a generating of hash corresponding user password.

## 5 Result and Discussion

From the outcome, demonstrated database password security utilizing hashing and salting pattern give a stronger security with the end goal that the original password is never put away. Regardless of the possibilities that the password store is compromised, just the hashes get to be a public. The password length is not stored and could not be estimated, making password cracking that much harder. There is no requirement for a secret, as none is utilized to hash the password.

**Fig. 2** Login page

**Fig. 3** Generated hash using Bcrypt hashing algorithm



**Fig. 4** Category and mode selection of AES technique



**Fig. 5** Encrypted password using AES and Bcrypt algorithm



The performance matrices are defined encryption time and throughput time. The encryption time is defined as, the time is taken for generating a cipher text from plaintext, and throughput time is defined as (Figs. 6 and 7; Tables 1 and 2).

$$\text{Throughput} = \frac{\text{Size of Encrypted Text in MB}}{\text{Time Required for Encryption in Seconds}} \qquad (3)$$

**Fig. 6** Execution time for proposed algorithm

**Fig. 7** Throughput time of proposed mechanism

**Table 1** Execution time of proposed algorithm

| Data text (character) | Proposed algorithm execution time (microsecond) |
|---|---|
| Data 1 | 2200 |
| Data 2 | 2436 |
| Data 3 | 2661 |
| Data 4 | 3084 |
| Data 5 | 3572 |

**Table 2** Throughput of proposed algorithm

| Plain text (in terms of character) | Data set (size in MB) | Execution time (seconds) | Throughput |
|---|---|---|---|
| Data 1 | 0.000967 | 0.002 | 0.4835 |
| Data 2 | 0.000984 | 0.00244 | 0.4032 |
| Data 3 | 0.000965 | 0.00267 | 0.3711 |
| Data 4 | 0.000990 | 0.00309 | 0.3203 |
| Data 5 | 0.001090 | 0.00350 | 0.3114 |

# 6 Conclusions

Protection and accuracy are the two major fields of overall system customers. This type of problems is illuminated by the study of cryptographic methods. The storage of password security is an essential part of information protection, as most systems nowadays require a validation technique utilizing passwords. The technique of hashing is usually utilized for transforming plain content passwords into set of character probably it cannot be decrypted by intruder due to the way of their one-sided cipher technique. Be that as it may, with time, the attacks got to be possible through the using of word reference tables and rainbow tables. In this paper, we have intended to the utilizing Bcrypt hashing algorithm with AES for providing the user's account protection. Our future work would explore this concept, and a combination of algorithms will be applied either sequentially or parallel, to set up a more secure environment for data storage and retrieval.

# References

1. Bocjj, P.: The Dark Side of the Internet: Protecting Yourself and Your Family from Online Criminals, 2nd edn, pp. 159–161. Greenwood Publishing Group, Westport, CT (2006)
2. Morley, D.: Understanding Computers in a Changing Society, 3rd edn, pp. 196–199. Course Technology Cengage Learning, Boston, MA (2008)
3. Buxton P.: Egg rails at password security, Netimperative, June, 24, (2002)
4. Zombie PCs for Rent Information Security&p%5Bne%wsletterId%5D=609, September 2004
5. Dorrans, B.: ASP.NET Security, Wiley (John Wiley & Sons, Ltd), ISBN:978-0-470-74365-2, 2010
6. Sriramya, P., Karthika, R.A.: Providing password security by salted password hashing using Bcrypt algorithm. J. Eng. Appl. Sci. **10**(13), 5551–5556 (2015)
7. Chen, Y.C., Horng, G., Huang, C.C.: Privacy protection in on-line shopping for electronic documents. In: 5th International Conference on Information Assurance and Security, pp. 105–108 (2009)
8. Sultana, S., Ghinita, G. et. al.: A lightweight secure scheme for detecting provenance forgery and packet drop attacks in wireless sensor networks. IEEE Trans. Dependable. Secure Comput. **12**(3), 256–269 (2015)
9. Halderman, J.A., Waters, B., Felten, E.: A convenient method for securely managing passwords. In: Proceeding of the 14th International World Wide Web Conference, pp. 471–479 (2005)
10. Khiyal, M.S.H., Khan, A., Bibi, N., Ashraf, T.: Analysis of password login phishing based protocols for security improvements. In: Proceeding of IEEE 5th International Conference on Emerging Technologies (ICET 2009), pp. 376–379 (2009)
11. Stallings, W.: Data and computer communications, Pearson Education, Inc., Eighth Edition, ISBN: 0-13-243310-9, (2007)
12. Islam, M.N., Mia, M.M.H., Chowdhury, M.F.I., Matin, M.A.: Effect of security increment to symmetric data encryption through AES methodology. In: Nineth ACIS International Conference on Software Engineering. Artificial Intelligence. Networking and Parallel/Distributed Computing, pp. 291–294 (2008)
13. Zhao, Z., Dong, Z., Wang, Y.: Security analysis of a password-based authentication protocol proposed to IEEE 1363. Theor. Comput. Sci. 352, 280–287 (2006)
14. Khiyal, M.S.H., Khan, A., Bibi, Ashraf, N.T.: Analysis of password login phishing based protocols for security improvement. In: Proceeding of IEEE 5th International Conference on Emerging Technologies (ICET 2009), pp 376–379 (2009)
15. Lakshmanan, T., Muthusamy, M.: A novel secure hash algorithm for public key digital signature schemes. Int. Arab J. Inf. Technol. 262–267 (2012)
16. Ylonen, T.: SSH secure login connections over the internet. In: Proceedings of the USENIX Web Security, Privacy & Commerce, 2nd edn
17. Yue, C., Wang, H.: Anti-phishing in offense and defense. In: Proceedings of the 24th Annual Computer Security Applications Conference (AC-SAC'08), pp. 345–354 (2008)
18. Mwagwabi, F., McGill, T., Dixon, M.: Improving compliance with password guidelines: how user perceptions of passwords and security threats affect compliance with guidelines. In: 47th Hawaii International Conference on System Sciences, pp. 3188–3197 (2014)
19. Farmand, S., Zakaria, O.B.: Improving graphical password resistant to shoulder-surfing using 4-way recognition-based sequence reproduction (RBSR4). In: 2nd IEEE International Conference on Information Management and Engineering (ICIME), pp. 644–650 (2010)

# A Novel Cost Minimization Approach for Data Transaction in the Context of IoT

Ranjit K. Behera, Nicky Kumari, Avinash Mahala
and K. Hemant Reddy

**Abstract** With the advent of pervasive computing and ubiquitous services in today's life, there are numerous data-centric services boosted by IoT. Moreover, mobile devices and their utility have resulted in a huge amount of digital data content which fails to maintain QoS. Standalone IoT cannot handle all these huge data, which resulted in the merging of IoT with cloud computing. Similarly, the retrieval of required data has become more time-consuming with the exponential increase of IoT devices. In order to remove these shortcomings, we propose an IoT-enabled three-layered architecture to process huge data and support high-end processing, which is completely suitable for data transaction supporting IoT services. In this three-layered novel architecture, each layer is collectively responsible for data emission, data retrieval, or data processing, thus completing the data transaction process. A case study on smart healthcare system presents the efficacy of the model.

**Keywords** Data clouds · FOG computing · Datacenter · Contention cost and smart healthcare

R. K. Behera (✉) · N. Kumari · A. Mahala · K. Hemant Reddy
Department of Computer Science and Engineering, National Institute of Science and
Technology, Palur Hills, Berhampur 761008, India
e-mail: ranjit.behera@gmail.com

N. Kumari
e-mail: nist.nicky@gmail.com

A. Mahala
e-mail: avinashmahala6122@gmail.com

K. Hemant Reddy
e-mail: khemant.reddy@gmail.com

# 1 Introduction

Since the existence of Networking, the Internet has proved its boon toward the mankind in several ways [1]. The Internet is inevitable and in demand. And the same is proliferating exponentially each day. Looking into today's scenario, we can assume that the human dependence on the Internet has shifted from human-to-human or human-to-machine communication to machine-to-machine communication [2]. And that is where IoT stands. The term 'Things' in IoT refers to any physical world object whether it is an interactive or non-interactive object. Anything ranging from a smart active device to passive day-to-day entities like 'a can of beverage' or 'a bottle of milk' can be a part of Internet. In order to make such day-to-day entities as a part of the interactive world, we can attach a communication means, mainly Radio frequency identification (RFID) tags. The basic working principle lies in Machine-to-machine (M2M) communication [3], but not restricted to it. Basically, M2M communication refers to the interaction between two different machines or objects without any human involvement. Our main contribution through this paper is to present a novel cost minimization approach for data transaction among IoT devices. Therefore, we have tried to enhance the data transaction services of IoT over the Internet using the basics of cloud computing, FOG computing, and the community-based data-centric services through the proposed three-tier architecture. Also, we have shown the efficacy of the proposed architecture using a smart healthcare system.

The remainder of this paper is organized as follows. Section 2 presents a detailed study of literature. Section 3 presents the system model and its underlying principles and design. A detailed case study of the proposed model is presented in Sect. 4. Finally, the conclusions are presented in Sect. 5.

# 2 Related Work

In this segment, we deliberate the former studies which were done in this realm and details of our study presented in three subclasses indicating the prior research in the corresponding domain.

## 2.1 *Cloud-Only Versus FOG Computing*

A recent research has contrasted between the efficiency of data analysis in cloud-only computing architecture and cloud-cum-FOG computing architecture [1]. The rapid increase in the IoT devices has led to the proliferation of the data generated each unit of time [4]. But in traditional cloud-only computing architecture, data generated from the users are directly pushed to the cloud as a result of

which data traffic gradually rises in the network due to which inconsistency arises resulting in the degradation of QoS. In one of its revolutionary research works, Cisco has proposed the concept of FOG computing [5]. It says that whenever a distributed computing architecture handles billions of Internet-connected devices, that architecture is known as FOG computing. The concept is actually based on the technology of edge computing where the services are hosted with the help of so-called edge devices like gateways, routers, and access points who are not only capable of intermediary network routing but also capable of storing intermediary data and take valuable decisions of whether to forward these data packets to cloud or to drop back and send the decisions to the user after analyzing. To support IoT-enabled applications, Parwekar [6] has discussed the integration of cloud and IoT.

## 2.2 Data Clouds

In the recent past, various data-centric services were incapable of being accommodated in the network due to the present Internet architecture, i.e., host-centered architecture. Therefore, a variety of Information-centric networking (ICN) [7] concepts came into existence and tried to solve this problem. However, IoT demands advanced data-centric services to be held. In this regard, a number of ICN architectures have been established and recommended. In data-oriented network architecture (DONA) [8], flat names were associated with the data items. According to the authors, for the naming of data items, a group of Resolution handlers (RHs) have been introduced into the network. Data item to be advertised is initially registered by the owner with the local RH, which will then be distributed among the hierarchy of other RHs. Based on the demand of data retrieval by the user, the RHs are requested which will find the matching owner and forward the requests. Finally, the data is received by the concerned user from the relevant owner. The Content-centric networking (CCN) [9] has hierarchical naming of data items which are extremely different from DONA. Here, the communication between the users and the owner is made with the help of Content routers (CRs). Likewise, a new architecture called as Data Clouds [2] used the concept of communities which also governs the hierarchical naming of data items where a community stands for similar user interests in data contents which use a network consisting of data-centric services of interest. Due to the presence of name resolution and data routing within each community data dissemination among the users become possible in an efficient manner. This community-oriented communication could also improve management over the Internet which develops mechanisms to design user mobility and in-network caching and preserve user security.

## 2.3   Contention Cost Task Scheduling Algorithm

In the present day, the advancement of pervasive and ubiquitous services tends to increase several data-centric services that are boosted by IoT. Moreover, mobile devices and their utility have resulted in a huge amount of digital data content. Standalone IoT cannot handle all these huge data which resulted in the integration of IoTs with cloud computing. Vast data content makes it difficult to provide quality of service requirements. Therefore, these shortcomings are addressed by the cloud computing paradigm. The main idea is to move a heavy data to the high processing, storage data centers in the cloud. IoT also generates multimedia data which requires more processing power, storage space, and scheduling resources. Therefore, it is very difficult to perform an efficient task scheduling process in the cloud. In [10], the authors propose an approach of contention cost scheduling algorithm which has a potential to solve the above shortcomings. The main idea is to schedule tasks in clouds so as to minimize the contention cost as well as the resource cost. This model consists of three layers (1) cloud provider layer, which contains data centers, (2) cloud customer layer, where 'things' of IoT reside and be able to send a request through sensors, and (3) broker layer. Each individual server in the cloud is known as a Virtual machine (VM). In the broker layer, there is a number of VMs functioning as a centralized management node, called as a task scheduler. This task scheduler receives all computational requests of users and decides whether the received request should be added to the service queue or not enabling efficient task processing in the cloud.

## 3   IoT-Enabled Cloud Model

After going through numerous research works in this paper, we have described a detailed architecture that could support the infrastructure of IoT on multiple levels. The proposed architecture consists of three tiers—Tier 1, Tier 2, and Tier 3. Each layer has its own features and control functions enabling the entire infrastructure of IoT to perform data transaction. The following subsections described the architecture of three layers (Fig. 1).

   **TIER 1 User Layer**: This layer is the initial layer or the layer where the user generates crude data through various technologies like RFIDs, Barcodes. Each moment data generated is exponentially booming. It is really a tedious job to operate on this ever-growing data. In order to ease such task, first of all, we need to organize the data in an efficient manner. Now, this can be done through the concept of Data Clouds [2]. According to the concept of Data Clouds, the authors have extended the idea of 'community-oriented communication' and proposed the concept of Data Clouds. According to this concept, various users having common interests are congregated into diverse groups, and users within the same group perform intercommunication. Hierarchical Naming Convention is used according to

**Fig. 1** 3-Tier architecture of IoT-enabled cloud model



[2] in order to name and address the data content. Different communities can communicate with each other through Designated forwarding and caching node (DFCNs). This device is responsible for holding all the addresses within a particular community. It also has the capability of storing or caching frequently used data/information as and when requested by the clients. Now, a number of such communities connected with each other through DFCNs are again connected with a Community rendezvous point (CRP) that keeps track of all the communities. Such device is responsible for enabling intercommunication between the communities. Data generated in this tier is pushed to the next tier through CRPs, DFCNs, and Forwarding nodes (FNs). These devices are only responsible for forwarding the data to the desired location. Data Clouds provide a feasible platform for

disseminating data in data-centric services and ICN-based architectures of the IoT. Each CRP being referred belongs to individual administrative domain.

**TIER 2 FOG Layer**: This layer is termed as the From cOre to edGe (FOG) layer. This tier consists of the FOG devices which act as an intermediary between the lower layer and the upper layer serving for cloud purposes. The term FOG computing was coined by Cisco in 2012 [7]. It is defined as a distributed computing paradigm that authorizes the network devices at different hierarchical stages with multiple degrees of computational as well as storage ability (Fig. 2).

The devices present over here are given enough intelligence to examine whether or not the incoming data from the Tier 1 is required to send to the cloud layer or Tier 3 or not. Data arrived here are classified into two categories: simple data and complex data. The data those are not so large and needs immediate processing for real-time application, and can be processed easily by the devices present in the FOG layer. Such data is processed here in this layer without any further chaos of pushing to the cloud. This helps in lowering the data traffic and enhancing the performance of the end applications. Useful information generated here is either pushed back to the user layer on-demand of different users and for corresponding applications or is pushed up to the cloud storage for further insight processing. Insight processing is the technique of extracting business intelligence reports from the accumulated data and processing on the same. Hence, we can assume that this layer acts as a filter that



**Fig. 2** **a** User layer architecture, **b** FOG layer architecture, **c** Cloud layer architecture

filters out data to be processed in the cloud layer enhancing the performance in the cloud as well as the experience of the user.

**TIER 3 Cloud Layer**: This is the most vital layer that is responsible for heavy data processing and finds out customer insights and extracting business intelligence reports alike data mining [11]. Data received here is the complex data that requires thorough analysis and processing. The various crude data are accumulated in the cloud storage, commonly known as data centers where high-end servers with high-storage and computing capacity are available to work on data. First of all, whenever data arrives from the lower tier, a scheduling agent is responsible for pushing these data into queues such that they will be in a task queue which would be easier to be pushed. Now a dispatcher is available to take the task queue and is scheduled by a scheduler for cloud processing. Here, a contention-cost algorithm [10] is applied so as to preserve data integrity. Generally, due to a huge rush of data and limited resources available, it is a common problem to have network contention. This leads to data congestion and unavailability of real-time services. In order to eliminate such problems, we use a contention-cost algorithm that can easily schedule the various tasks containing data to be processed on and hence provide reliable and real-time services having low-latency to improve the performance. After that data arrived in this cloud layer is also encrypted for security purposes. Various cryptographic technologies like RSA [12] and others are used for maintaining the data security and integrity.

## 4 Case Study

The detailed architecture that we discussed above can be implemented in several IoT-related applications such as in industries, healthcare, traffic management system, and many more. Here, we have taken a case study of the smart healthcare system and tried to show how the implementation of such architecture could prove to be very economical and reliable. Starting from the very first layer or the first tier, whenever a patient gets admitted to a hospital, first of all, he is kept in a particular room if required. In such case, he or she is surrounded by lots of RFID [13] devices, mostly active RFID devices—maybe in his or her body (inbuilt or wearable), or in the patient's surrounding. Apart from an active RFID device, we can also use passive type RFIDs that requires a constantly monitoring active RFID reader (Fig. 3).

Now, these RFID devices constantly monitor the surrounding of the patient, trying to maintain the atmospheric constraints like pressure, temperature, heart rate, respiration rate, breathing rate, blood pressure, and others. Constant measurement of such parameters is highly essential for preventing any kind of catastrophic incident.

In such way, every room in the hospital generates data. Also, there can be lots of other users like a community of doctors, a community of families, friends and relatives, nurses, ward boys who are also responsible for the maintenance of the

**Fig. 3** Smart healthcare system

hospital scenario. It is the constant effort of all the people involved that makes the healthcare system more effective. Now, these people at regular intervals communicate with each other in order to provide better care to the patients. The crude data generated at this level is moved toward the local data centers those can be termed as the FOG devices. Now, these devices with their limited storage and computing capability analyze the type of medical data received from various communities. And further finds out which data is simple that can be processed by the same local data center and which needs to be forwarded to the cloud layer. This filtering of data saves time as well as decreases the data traffic and congestion between the FOG layer and the cloud layer. Thus the performance of computation gets increased in the cloud layer. In this layer, gateways can generate requests to the local medical authorities in case of any emergency and notifications could also be sent so that special care could be taken off the case. After performing on the local transaction, Tier 2 gateways are responsible for routing the required data to the cloud for further analyzing and processing. The in-depth analyzing of the data and processing it generate several medical histories of a patient that may prove to be useful in future instances. Hence, these data are saved in the remote data centers, and reports are sent back to the local data centers helping in the maintenance of continuous health assistance to the patients. Scheduling of data and tasks in this cloud layer eliminate the possibility of loss of data integrity and data contention. The contention-cost algorithm constantly improvises the data processing performance leading to the improvement in the better services to the patient.

## 5 Conclusion

Contemporary IoT demands expansion of data-centric services over traditional Internet architecture. Through this paper, we have come up with an architecture that integrates various interrelated computing areas supporting IoT infrastructure aiming to achieve data transaction efficiency and transparency. The problem of datamandering has been solved with the involvement of community-based, i.e., data clouds architecture in the user layer. Then comes the advent of FOG computing as it is preferable for computing of various data in the edge devices itself rather than cloud-only computing. Problems like data congestion, service latency, data traffic, and many others may decrease with this technique. Data retrieval is no more a difficulty when this type of architecture is going to be deployed. Another major issue of increasing cost by the cloud providers can be resolved with the help of a novel contention-cost algorithm in which a proper task scheduling is done in the cloud layer. Our further work mainly focuses on justifying the proposed architecture through a detailed simulation so that the practicality of the proposed model can be validated.

# References

1. Sarkar, S., Chatterjee S., Misra, S.: Assessment of the suitability of fog computing in the context of internet of things. IEEE Transactions on Cloud Computing (2015)
2. Yue, H., Guo, L., Li, R., Asaeda, H., Fang, Y.: Data clouds: enabling community-based data-centric services over the internet of things. IEEE Internet Things J. **1**(5), 472–482 (2014)
3. Salam, S.A., Mahmud, S.A., Khan, G.M., Al-Raweshidy, H.S.: M2M communication in smart grids: implementation scenarios and performance analysis. In: Wireless Communications and Networking Conference Workshops (WCNCW), 2012 IEEE, pp. 142–147. IEEE (2012)
4. Cisco delivers vision of fog computing to accelerate value from billions of connected devices. Press Release Cisco. [Online] (2014, Jan). Available: http://newsroom.cisco.com/release/1334100/CiscoDelivers-Vision-of-Fog-Computing-to-Accelerate-Value-from-Billionsof-Connected-Devices-utm-medium-rss
5. Bonomi, F., Milito, R., Natarajan, P., Zhu, J.: (2014). Fog computing: a platform for internet of things and analytics. In: Big Data and Internet of Things: A Roadmap for Smart Environments, (pp. 169–186). Springer International Publishing
6. Parwekar, P.: From internet of things towards cloud of things. In: 2011 2nd International Conference on Computer and Communication Technology (ICCCT), IEEE, pp. 329–333 (2011)
7. Xiaoke, J.: Information-centric networking. China Commun. **3** (2015)
8. Miorandi, D., Sicari, S., De Pellegrini, F., Chlamtac, I.: Internet of things: vision, applications and research challenges. Ad Hoc Netw. **10**(7), 1497–1516 (2012)
9. Lee, U., Rimac, I., Hilt, V.: Greening the internet with content-centric networking. In: Proceedings of the 1st International Conference on Energy-Efficient Computing and Networking, pp. 179–182. ACM (2010)
10. Hung, P.P., Aazam, M., Huh E.-N.: A novel contention-cost scheduling algorithm in internet of things environment
11. Fayyad, U., Piatetsky-Shapiro, G., Smyth, P.: From data mining to knowledge discovery in databases. AI Mag. **17**(3), 37 (1996)
12. Somani, U., Lakhani, K., Mundra, M.: Implementing digital signature with RSA encryption algorithm to enhance the data security of cloud in cloud computing. In: 2010 1st International Conference on Parallel Distributed and Grid Computing (PDGC), IEEE, pp. 211–216 (2010)
13. Want, R.: An introduction to RFID technology. IEEE Pervasive Comput. **5**(1), 25–33 (2006)

# Comparison between Genetic Algorithm and PSO for Wireless Sensor Networks

**Pritee Parwekar, Sireesha Rodda and S. Vani Mounika**

**Abstract** One of the most promising algorithms for network optimization is the particle swarm optimization (PSO) and genetic algorithm (GA). The paper is about comparing these two as applied to wireless sensor networks. If a sink is placed at a longer distance from the sensors then the battery life (energy) drains faster, and it reduces the life of the network. Our analysis shows that optimized clustering technique of sensors can minimize the communication distance and can help to increase the network stability. GA and PSO can optimize the cluster formation of sensors. Simulation results have shown us that PSO performs better than GA for clustering algorithms in wireless sensor networks.

**Keywords** Wireless sensors · Network · Ad hoc networks · Clustering
Genetic algorithm · Particle swarm optimization

## 1 Introduction

Wireless sensor networks (WSN) have evolved as one of the most significant technologies in the twenty-first century and have started becoming the foundation for a few more. In the past decades, it has been the focus of attention for both the academia and the industry worldwide. A WSN usually comprise of sizable no of cost efficient, low-powered wireless nodes, having limited computation capabilities. These sensor nodes typically communicate in short range and collaborate to accomplish the network function, for example, environmental information

P. Parwekar (✉) · S. Vani Mounika
Anil Neerukonda Institute of Technology and Sciences, Visakhapatnam, India
e-mail: pritee.cse@anits.edu.in

S. Vani Mounika
e-mail: sambangimounika3@gmail.com

S. Rodda
GITAM University, Visakhapatnam, India
e-mail: sireesha@gitam.edu

gathering, surveillance, and in manufacturing processes. The WSN rides on the philosophy of collaborative network capability for completing the tasks for the required mission.

The future WSN applications are expected to incorporate a benchmarked combination of hardware and software. But as we stand today, the network designers are still juggling between the tradeoffs that they have to adopt so as to minimize the cost of deployment and maximize performance. Designers of wireless embedded systems, therefore, are required to assess these tradeoffs and select the optimal sensor hardware. The sensors are a compact, small, battery-powered device, and therefore have limited energy resource. Therefore, energy consumption remains the most important issue to be tackled in sensor networks. The data collected by individual sensors in a WSN must be transmitted to another more robust node usually called as the sink which aggregates data, and finally transmits the same to the base station for data processing or use. Scalability is another aspect of such networks where they should be able to accommodate more area of interest with help of additional sensors.

Energy exhaustion of sensors depends on two factors, namely the distance of communication between the sensors and sink node, and also the period for which the sensor is active. The energy reserve of each sensor is directly dependent on the power consumption pattern, and this is the most important factor affecting the network lifetime. Therefore, limiting the long-range data exchanges between the sensors and their respective sink nodes helps to prolong the network life and effectively maximize network performance. Clustering has evolved as the most effective method to efficiently increasing the overall network life.

## 1.1 Related Work

The Sound Surveillance System (SOSUS) was developed by the US Military in the 1950s. This can be considered as the first wireless network resembling a modern WSN. This was used to detect and track Soviet submarines. Later on, Distributed Sensor Network (DSN) program was started by DARPA in 1980 for implementing distributed/wireless sensor networks. The Carnegie Mellon University and the Massachusetts Institute of Technology: Lincoln Laboratory have participated actively in the program. WSN technology started leaving the military domain and entered the scientific research in civilian domain. Universities started exploring the potential of WSNs in different applications.

Low-energy adaptive clustering hierarchy (LEACH) (Heinzelman et al's paper [1]) is one of the first modern routing algorithms designed for wireless sensor networks which enable collection of data and delivery to the sink node through the networks of sensors. Clustering has been the central theme of this protocol, herein a hierarchical approach has been adopted. LEACH algorithm is distributed and autonomous, and it does not require control information from the base station. Simulation results have shown that LEACH protocol outperforms conventional

routing protocols. Most of the research has therefore taken LEACH as the benchmark protocol. A-LEACH, LEACH-A, LEACH-B, LEACH-C, C-LEACH, LEACH-E, LEACH-F, I-LEACH, LEACH-L, LEACH-M, M-LEACH, LEACH-S, TL-LEACH, V-LEACH are improvement of LEACH protocol [2–16]

Heinzelman et al. [17] have established a predetermined number of optimal clusters which is about 5% of the overall number of nodes. They have used the communication energy model. However, several other factors such as density of the sensor, the sink position have been overlooked. Particle swarm optimization (PSO), a population-based heuristic optimization technique is popular as parameter estimation and localization algorithm [18]. This technique addresses the computation complexity issue; however, in the basic form, this technique converges prematurely. Leung et al. [19] has propsoed a model in which linearly decreased inertia weight of each particle can be automatically calculated according to fitness value. [20] proposed a medhod to adjusts the population size automatically to select the candidate location. Tillett et al. [21] have shown an approach called the particle swarm optimization wherein the field of sensor nodes is equally sized. PSO is an evolutionary programming technique, wherein interaction and behaviour of ants and termites in a colony have been the mimicked through a programming technique to find a good solution. This method is not suggested for heterogeneous networks. As shown by Hussain et al. [22] in WSN, GA-based hierarchical clusters prolong the network lifetime. Ostrosky and Rabani [23] have provided the solution through assumption of a large data set having $n$ points and further partitioning this data into $k$ disjoint clusters whilst minimizing the overall distance between cluster heads. Solution is obtained using polynomial-time approximation scheme. Agarwal and Procopiuc [24] have further developed this method, named as k-clustering where the nodes in a workspace are separated in k clusters, and distance metric determines the cluster membership. Both these methods assume a fixed number of clusters. The aim here is to solve problem having initially unknown number of clusters.

## 2   Current Work

Both the GA and PSO optimization techniques are used to find the cluster heads in the current work. A set of regular nodes collectively form into small clusters, and every cluster has a cluster head. On selection of the cluster head, the nearby nodes establish communication only with it. Each cluster head, in turn, aggregates data collected by the sensors within the cluster and transmits this aggregated data to the sink.

## 3   Assumptions

In this paper, the sensors are static, and all the sensors are homogenous. The sensors are spread in a hostile environment, and the sink is out of this environment. All nodes can become cluster heads and adjust transmission power depending on

distance. GPS assisted by terrestrial triangulation methods for finer accuracy (A-GPS) is used for precise localization of sensor. This paper implements the proposed algorithm with GA [25] and compares it with PSO for achieving minimum overall communication distance and also reduction in energy and thereby maximizing the network lifetime.

## 4    Proposed Solution

This paper is using genetic algorithm (GA) and particle swarm optimization (PSO) to optimize the number of clusters in an arbitrary network. When the cluster head is selected, all its nearby sensors will be connected to it. In each cluster, the cluster head aggregates the data that is collected from all the sensors and transmits it to the base station. Figure 1 explains clustering process. Both the techniques use the same objective function described below, and the results are compared.

### *4.1    Problem Representation*

Cluster head selection is important considering minimization of the distance between the sensors and the CH. Binary representation (as tabulated below) has been adopted, wherein each bit is related to one sensor. A "1" identifies the sensor as a cluster head and a "0" as a regular node. Table 1 shows the nodes with their values.



**Fig. 1**  Clustering technique in WSN

**Table 1** Chromosome representation in genetic algorithm

| S0 | S1 | S2 | S3 | S4 | S5 | S6 | S7 | S8 | S9 |
|----|----|----|----|----|----|----|----|----|----|
| 0  | 1  | 0  | 0  | 1  | 0  | 0  | 1  | 0  | 0  |

To begin with, the population comprises of random individuals amongst which cluster heads are selected by GA. The nearest cluster head is found using a deterministic method by the regular nodes. The governing factor is reducing the total transmission distance. Lesser the number of cluster heads in the network the better is the energy efficiency since the cluster heads drain more energy as compared to regular nodes. Hence, the number of cluster heads also has been factored in the governing fitness function equation. This leads us to the following consolidated equation as fitness function [25]:

$$f = W * (D - di) + (1 - W) * (N - Hi) \tag{1}$$

where $D$, total distance of all nodes to the sink; $di$, sum of the distances between regular nodes and cluster heads coupled with cluster heads to the sink; $Hi$, count of cluster head; $N$, node count; $w$, weight.

A given objective function has all parameter fixed except the $di$ and $Hi$. Shorter distances, therefore, improve the fitness value. The GA maximizes the fitness function value for a good optimal solution. The range of w is $0 \leq W \leq 1$ depends on the consideration of distance or the cost induced by cluster heads. If $W = 1$ the network is optimized with respect to the distance of communication, and if $W = 0$ then the number of the cluster heads are optimized.

## 5 Experimental Results

The experimental setup involves 200 nodes in a 2D simulation.

- GA and PSO have been implemented for optimizing the sensor network in MATLAB R2012b for minimizing the distance and the number of cluster heads to ensure more lifetime.
- Experimental results show that PSO algorithm is performing better than GA and is efficient and adaptive to multiple network topologies.
- It has been proven that this approach is quick in finding solutions. For a 100-node problem, below 120 generations, a good solution is achieved. Experiments indicate that the scaling window plays an important role in the quality of the solution found.
- In an event of a single node near a sink, the node directly transmits data to the sink.
- When the sink is closer to the middle of the network, it has been found that more cluster heads are required.

- In a region having higher density of sensors, the middle node is usually elected as a cluster head.
- No two cluster heads are co-located. The GA tends to elect one of them as a CH to avoid duplication of communication

When the count of node doubles, the population size is required to be doubled to maintain comparable performance. About 10% of the nodes assume the role of cluster heads. This percentage varies depending on the distribution of the network. On an average, the overall transmission distance is reduced by 80% compared to direct transmission (Table 2).

Table 1 illustrates number of iterations required with the increase in count of nodes. The experiment is done with 25, 50, 75, 100 and 200 in all cases PSO has taken minimum number of iterations to converge to the fitness function compared to GA (Figs. 2, 3, 4).

**Table 2** Values considered for the experiment and iterations for convergence

| No. of nodes | PSO_Iteration | GA_Iteration |
|---|---|---|
| 25 | 4 | 7 |
| 50 | 8 | 18 |
| 75 | 6 | 7 |
| 100 | 6 | 31 |
| 200 | 7 | 12 |



**Fig. 2** Comparison of GA and PSO

**Fig. 3** GA fitness versus generations



**Fig. 4** PSO fitness versus generations

# 6  Conclusions and Future Work

In this paper, GA and PSO methods are compared to optimize the distance between sensors in a WSN. Fitness function is used to select cluster heads after an initial random selection. It has emerged that PSO is more effective for converging to the solution. Further, a better and more useful solution can be obtained in case of uniform and non-uniform topologies. Scalability of the network has also been indicated. This optimization technique can also be utilized for base station construction to maximize transmission coverage and minimize overheads. Future work would involve evaluating these concepts more objectively.

# References

1. Handy, M.J., Haase, M., Timmermann, D.: Low energy adaptive clustering hierarchy with deterministic cluster-head selection. In: 4th International Workshop on Mobile and Wireless Communications Network, 2002, pp. 368–372. IEEE (2002)
2. Qiao, X., Yan, C.: A control algorithm based on double cluster-head for heterogeneous wireless sensor network. In: 2nd International Conference on Industrial and Information Systems, vol. 1, pp. 541–544 (2010)
3. Singh, S.K., Singh, M.P., Singh, D.K.: A survey of energy efficient hierarchical cluster-based routing in wireless sensor networks. Int. J. Adv. Netw. Appl. **02**(02), 570–580 (2010)
4. Kumar, N., Kaur, J.: Improved leach protocol for wireless sensor networks. In: 2011 7th International Conference on Wireless Communications, Networking and Mobile Computing (WiCOM), pp. 1–5. IEEE (2011)
5. Mehta, R., Pandey, A., Kapadia, P.: Reforming clusters using C-LEACH in wireless sensor networks. In: 2012 International Conference on Computer Communication and Informatics (ICCCI), pp. 1–4. IEEE (2012)
6. Xu, J., Jin, N., Lou, X., Peng, T., Zhou, Q., Chen, Y.: Improvement of LEACH protocol for WSN. In: 2012 9th International Conference on Fuzzy Systems and Knowledge Discovery (FSKD), pp. 2174–2177. IEEE (2012)
7. Kaur, P., Katiyar, M.: The energy-efficient hierarchical routing protocols for WSN: a review. Int. J. Adv. Res. Comput. Sci. Softw. Eng. **2**(11), 194–199 (2012)
8. Rengugadevi, G., Sumithra, M.G.: Hierarchical routing protocols for wireless sensor network —a survey. Int. J. Smart Sens Ad Hoc Netw. (IJSSAN) **2**(1), 71–75 (2012)
9. Norouzi, A., Zaim, A.H.: An integrative comparison of energy efficient routing protocols in wireless sensor network. Sci. Res. Wirel. Sensor Netw. **4**, 65–67 (2012)
10. Gnanambigai, J., Rengarajan, D.N., Anbukkarasi, K.: Leach and its descendant protocols: a survey. Int. J. Commun. Comput. Technol. 1(3), 15–21 (2012)
11. Aslam, M., Javaid, N., Rahim, A., Nazir, U., Bibi, A., Khan, Z.A.: Survey of extended LEACH-based clustering routing protocols for wireless sensor networks. In: 2012 IEEE 14th International Conference on High Performance Computing and Communication and 2012 IEEE 9th International Conference on Embedded Software and Systems (HPCC-ICESS), pp. 1232–1238. IEEE (2012)
12. Neetika, S.K.: Review on hierarchical routing in wireless sensor networks. Int. J. Smart Sens. Ad Hoc Netw. (IJSSAN) **2**(3), 85–90 (2012)
13. Kaur, R., Sharma, D., Kaur, N.: Comparative analysis of leach and its descendant protocols in wireless sensor network. Int. J. P2P Netw. Trends Technol. **3**(1), 51–55 (2013)

14. Jan, M.A., Khan, M.: A survey of cluster-based hierarchical routing protocols. IRACST Int. J. Comput. Netw. Wirel. Commun. (IJCNWC) **3**, 138–143 (2013)
15. Bhattacharjee, A., Bhallamudi, B., Maqbool, Z.: Energy-efficient hierarchical cluster based routing algorithm in WSN: a survey. Int. J. Eng. Res. Technol. (IJERT) **2**(5), 302–311 (2013)
16. Verma, S., Mehta, R., Sharma, D., Sharma, K.: Wireless sensor network and hierarchical routing protocols: a review. Int. J. Comput. Trends Technol. (IJCTT) **4**(8), 2411–2416 (2013)
17. Heinzelman, W.B., Chandrakasan, A.P., Balakrishnan, H.: An application-specific protocol architecture for wireless microsensor networks. IEEE Trans. Wirel. Commun. **1**(4), 660–670 (2002)
18. Gopakumar, A., Jacob, L.: Performance of some metaheuristic algorithms for localization in wireless sensor networks. Int. J. Netw. Manage. **19**, 355–373 (2009)
19. Leung, S.Y.S., Tang, Y., Wong, W.K.: A hybrid particle swarm optimization and its application in neural networks. Expert Syst. Appl. **39**, 395–405 (2012)
20. Chen, D.B., Zhao, C.X.: Particle swarm optimization with adaptive population size and its application. Appl. Soft Comput. **9**(1), 39–48 (2009)
21. Tillett, J., Rao, R., Sahin, F.: Cluster-head identification in ad hoc sensor networks using particle swarm optimization. In: 2002 IEEE International Conference on Personal Wireless Communications, pp. 201–205. IEEE (2002)
22. Hussain, S., Matin, A.W., Islam, O.: Genetic algorithm for energy efficient clusters in wireless sensor networks. In: Fourth International Conference on Information Technology, 2007. ITNG'07, pp. 147–154. IEEE (2007)
23. Ostrosky, R., Rabani, Y.: Polynomial-time approximation schemes for geometric min-sum median clustering. J. ACM **49**(2), 139–156 (2002)
24. Agarwal, P.K., Procopiuc, C.M.: Exact and approximation algorithms for clustering. Algorithmica **33**(2), 201–226 (2002)
25. Jin, S., Zhou, M., Wu, A.S.: Sensor network optimization using a genetic algorithm. In: Proceedings of the 7th World Multiconference on Systemics, Cybernetics and Informatics, pp. 109–116 (2003)

# Secure Text Dissemination in Delay Tolerant Networks

**Afreen Fatimah and Rahul Johari**

**Abstract** Delay Tolerant networking is a field that lies under the wireless networks that is characterized by intermittent connectivity and frequent disruptions. In this paper, an approach is presented where the DTN nodes can be deployed as spy nodes in highly sensitive border areas. These nodes would act as data acquisition nodes, would acquire data from the surroundings, and transfer it to the respective commanding officer or the monitoring authority. Therefore, they work as spy in border regions, capable of sensing even minute movements of the enemies and alarming the soldiers well in time.

**Keywords** DTN · Security · Spy nodes · Biodegradable

## 1 Introduction

Delay Tolerant networks [1–3] are the category of networks that have low network density and are usually sparse in nature. This connectivity pattern among the nodes varies as the physical and geographical terrains vary. The messages are transferred from one node to another based on the type of contact that the architecture uses. There are majorly 5 types of contacts in DTN [4]. They are as follows:

1. Persistent contact: These contacts always remain available just like the "always-on" internet connections. No connection establishment needs to be done here.
2. On-demand Contacts: These contacts are established on demand basis as and when there is a need to transfer the message.

A. Fatimah (✉) · R. Johari
University School of ICT, GGSIP University, Delhi, India
e-mail: fatimahafreen@yahoo.co.in

A. Fatimah · R. Johari (✉)
GGSIP university, Delhi, India
e-mail: rahul.johari.in@ieee.org

3. Intermittent-Scheduled Contacts: These contacts show an agreement to establish the connection at a particular time and disconnect as and when the task gets over. But in DTNs, this particular time may be delay dependent.
4. Intermittent-Opportunistic Contact: These are not scheduled contacts but function based on the occurrence of any opportunity to transfer the message.
5. Intermittent-Predicted Contact: Unlike the scheduled contacts, they are based on the predictions made as per the historic data. The contacts that have been previously used can be used again for sending messages.

In the below approach presented in Sect. 4, an opportunistic contact had been used so that whenever there is an opportunity, the data are transferred from the DTN node to the satellite.

Security of data that resides within every node of the DTN is of paramount importance in such cases. To address this issue, a biodegradable nature of the node is proposed here. These nodes have the ability to sense variable temperatures, movements, velocity, sounds, etc and then transmit this data to the Commanding Officer through a geostationary satellite, deployed for this purpose. The current work gives detail about the deployment of such nodes in border areas where sensing the movements of enemies is very important to take necessary actions in anticipation. The organization of paper is as follows: Sect. 2 describes the nature of problem envisioned, Sect. 3 gives motivation of the present work, and Sects. 4 and 5 show the approach adopted with detailed algorithms.

## 2 Problem Envisioned

Nowadays, terrorism from across the border is turning out to be of great concern to the security of any nation. Every year, millions of dollars are spent to safeguard the border and other sensitive areas of a country from where important data pilferage can take place. This is done by restricting the enemies or unauthorized people to enter one's territory. So, rather than deploying soldiers in such harsh terrains, it is envisioned to replace men with machine that is to make use of DTN nodes that are capable of data acquisition and would get self-decomposed when the defined TTL gets over or whenever they come in human contact. This approach would be of major help to the defense deployments that have to keep on monitoring the border areas to check that there is no enemy intrusion.

## 3 Motivation for the Proposed Work

As discussed, there have been so many problems in the border areas so as a result, there is a need to administer a security mechanism that can prohibit any intrusion in country's territory by deploying sensor nodes that can sense data well in advance

and inform the base stations which are the Commanding Officer's stations about it. The data acquired by the node are sent to the satellite in just one hop so that during transit, the data are absolutely safe and no tampering or pilferage can happen. The DTN nodes used here are made of biodegradable material, such as magnesium etc that start degrading or dismantling once the defined TTL gets over or it comes under the contact of any human or undesired creature.

## 4 Our Approach

DTN spy nodes architecture includes a Commanding officer (CO), in charge of a particular area, a geostationary satellite which is dedicated for army message retrieval purpose and the DTN sensor nodes that are self-disposable in nature. Figure 1 shows the architecture and Fig. 3 describes the flow chart of the proposed architecture. In this, the commanding officer sends the helicopter carrying nodes to the sensitive areas from where data have to be acquired. The helicopter after dropping the nodes, returns back to the CO, and waits for other assignments to be given. After sending DTN nodes at the desired location, CO generates a message retrieval request to the node via satellite. After this, the satellite sends a signal to the DTN nodes to get the data present inside the nodes.

The Satellites usually gather two different types of data: one is the housekeeping data [5] and other is the actual science data. The housekeeping data are the one that tells about its health, position, temperature of the surroundings, and other similar data. On the other hand, the actual science data are the one that contains the required data. The required data in this case may be enemies or bunker location,



**Fig. 1** The architecture

tank and truck movement, deployment of land mines in the area. This data also contains the measurements and images that are actually required to be studied. Therefore, the major data to be studied are evaluated using the science data, but the housekeeping data also helps to know that how the equipments were functioning when the data were recorded. This also lets the researchers to make sure that the data recorded are accurate.

The nodes on receiving such signals transfer whatever they have sensed in between the refresh intervals. Here the total time to live for the nodes is greater than the refresh time of the node, and this measure has been adopted due to security reasons.

The data gathered by the satellites are sent down to the CO through downlinks in batches and as soon as possible. The data are sent to the ground through antennas which may be high gain or low gain [6]. The high-gain antennas are pointed directly to the receiving stations and send a lot of data very quickly, while the low-gain antennas are not directly pointed to the receiving stations and they send data slowly. In this case, a high-gain antenna is used to receive the signals at the CO's end.

### 4.1 Schema of the DTN Spy Node

The DTN spy node that has been proposed is somewhat similar to the HBE-Ubi-Coin: low power and coin size sensor platform based on MSP430 (CPU) and CC2420 (RF) [7]. This node consists of sensors that can easily sense the environmental changes that are even small in nature. The sensors are made up of five main components viz. Sensor unit, CPU, Memory unit, Energy source and Transceivers. Figure 2 shows a structure similar to that of the UBI coin sensor node. It further comprises of other add-ons like energy producer, position changer, and localization unit. The buffer capacity of the node is 10 KB. This wireless node resembles the shape of coin as they are very small in size. The range of such nodes is 0–50 km, based on the power, cost and bandwidth. But in proposed approach, range of the node is kept as 100m only.

The nodes use ZIGBEE technology; therefore, it is capable of data transfer by consuming very less amount of energy.

### 4.2 Message to be Sent

The message that is sent by the node to the satellite consists of the data acquired by it between the refresh intervals. When the satellite instructs the node to retrieve the information for it, the node generates a message header that is of 4 KB in size along with the actual data. Since the node has a maximum capacity of 10 KB so the message cannot be more than 6 KB in size. This header helps to direct the data packet to the correct destination and prevents it from getting lost inside the network. The message header is depicted in the Table 1.

**Fig. 2** Structure of the node

**Table 1** Message header

| Latitude (512B) | Longitude (512B) | Drop_Time (512B) |
|---|---|---|
| Count of total number of nodes (200B) | | Load_Time of the node (512B) |
| Node_ID (512B) | | Size (300B) |
| TTLn (262B) | Refresh_Time (262B) | CO-ID (512B) |
| Message (6 KB) | | |

# 5  Algorithms to Implement the Bio Degradable Spy Nodes

The algorithms mentioned in Sects. 5.1–5.4 give step by step sequence of implementing the spy node architecture. Figure 3 presents a comprehensive flow diagram with all the steps.



**Fig. 3** Flow chart depicting flow of control for the DTN spy nodes

### 5.1     Algorithm 1: Populating the contents of node

**Notation**
$H$  :   header of message to be sent
$PT$: plain text(information acquired by the node)
**Trigger:** When the Storage of data in the node takes place
1.    Input  header fields of the message [input_($H$)]
Suppose          $l^t$: node's Latitude co-ordinates
                 $l^n$: node's Longitude co-ordinates
                 $t$:   Time to live of the node (actual life of the RFID tag)
                 $n$:  Number of nodes in the delay tolerant network
                 $n_i$:  i$^{th}$ node in the delay tolerant network
                 $t_c$: Refresh time of the node, where $t_c$ <$t$
                 $l_i$: Time at which node is loaded at source
                 $d_i$: Time at which node is dropped at destination
                 $c$: Commanding Officer-ID
                 $N$: ID of the node
                 $m$:  size of the actual message that has to be sent
2.    Storing header fields inside the node
$H$← { $l^t$, $l^n$, $t$, $n$, $t_c$, $l_i$, $d_i$, $c$, $N$, $m$}
3.    Message to be entered in the node [enter($PT$)]
$PT$← {sensor data acquired by the node}

### 5.2     Algorithm 2: Concatenation of the header with plain text

**Notation**
$X$: Prepared sensor data that has to be uploaded to the satellite
**Trigger:** Preparing the message
1.    Concatenate the header $H$ with $PT$ [prepare_the_message ($PT, H$)]
2.    $X$← {prepare_the_message ($PT, H$)}

### 5.3    Algorithm 3: Send the prepared message

**Notation**
$Z$:  data that has been prepared is ready to be transferred / uploaded
$sid$: Satellite ID
**Trigger:** Transferring the message
1.    Attached  message is aligned with the assigned ID of the satellite
2.    $Z$←align($sid, X$)
3.    Transfer($Z$)

---

### 5.4  Algorithm 4: Self decomposing/degradation of the node

---

**Trigger:** node is dropped at the assigned location
if ($n_i \leftarrow$ Comes in contact with any undesired creature)
    self_destroy ($n_i$)
elseif ($n_i \leftarrow$ Does not come in contact with undesired creature)
  if ($t_c < t$)
     Memory_refresh
     $t++$
  else
     self_destroy ($n_i$)
else
  self_destroy ($n_i$)

---

## 6  Acknowledgment and Submission

For the above-mentioned work, we are grateful to Guru Gobind Singh Indraprastha University for providing us the necessary resources in a timely and efficient manner. This is an extended version of the paper "BDN: Bio Degradable Node, Novel Approach For Routing In Delay Tolerant Network" [8].

## 7  Conclusion and Future Work

This seems to be the first, of its kind, approach for acquiring data in DTN. This proposal does not lead to any pilferage of the sensitive data during data transmission. It has been envisaged that if the node carrying the message does not come in contact with anyone (human or animal), it is for sure it will not be going into the hands of any unauthorised person, and even if it does, he would not able to decipher the plaintext.

If such a defense system is implemented on the border areas, the sensitive information will be safe from all types of unintended personnel, thereby safely getting and protecting the details of activities taking place at the border areas without actually deploying soldiers there.

As a part of future work of the present work, the algorithms mentioned would be implemented using the ONE simulator [9], designed and developed in Java, open source programming language. Also to ensure the secrecy of sensitive data, it is proposed to encrypt the data using various encryption techniques, which can be mono-alphabetic or poly-alphabetic substitution cipher.

# References

1. Fall, K.: A delay-tolerant network architecture for challenged internets. In: Proceedings of the 2003 Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications. ACM (2003)
2. Johari, R., Dhama, S.: Routing protocols in delay tolerant networks: application-oriented survey. In: Zeng, Q.A. (ed.) Wireless Communications, Networking and Applications, pp. 1255–1267. Springer, New Delhi (2016)
3. Fatimah, A., Johari, R.: Part: performance analysis of routing techniques in delay tolerant network. In: Proceedings of the International Conference on Internet of Things and Cloud Computing, p. 76. ACM (2016)
4. Cerf, V., et al.: Delay-tolerant networking architecture. No. RFC 4838 (2007)
5. Zacchei, A., et al.: HouseKeeping and science telemetry: the case of Planck/LFI, In: 47th, p. 268 (2003)
6. Taylor, J., et al.: DESCANSO design and performance summary series (2002)
7. http://www.hanback.cn/en/pro_web.asp?id=332&anclassid=&anclass= soebzxidalyum&nclassid=&nclass=
8. Fatimah, A., Johari, R.: BDN: bio degradable node, novel approach for routing in delay tolerant network. In: Proceedings of the International Conference on Frontiers of Intelligent Computing: Theory and Applications (FICTA). Springer (2016)
9. Keränen, A., Ott, J., Kärkkäinen, T.: The ONE simulator for DTN protocol evaluation. In: Proceedings of the 2nd International Conference on Simulation Tools and Techniques. ICST (Institute for Computer Sciences, Social-Informatics and Telecommunications Engineering) (2009)

# Automated Diagnosis of Tachycardia Beats

**Usha Desai, C. Gurudas Nayak, G. Seshikala, Roshan J. Martis and Steven L. Fernandes**

**Abstract** Due to tachycardia, heart generates lethal arrhythmia beats namely *atrial flutter* (AFL), *atrial fibrillation* (A-Fib), and *ventricular fibrillation* (V-Fib). These irregular patterns are very effectively and noninvasively reflected using standard electrocardiogram (ECG). In this study, an automated diagnosis support system (DSS) is developed for accurate discrimination and classification of complete classes of tachycardia beats (atrial as well as ventricular) using higher-order spectra (HOS). In this multiclass diagnosis problem, dimensionality of HOS third-order cumulants is reduced using independent component analysis (ICA) and fed for standard hypothesis test ANOVA ($p < 0.05$). Finally, statistical significant components are subjected for ensemble classification using random forest (RAF) and rotation forest (ROF) classifiers and to realize best performance tenfold classification is performed. Further, the consistency of classifiers is assessed using Cohen's kappa matric. Proposed DSS achieved overall classification accuracy of 99.54% using ROF. Our reported results are highest than published in the earlier works.

**Keywords** Multiclass diagnosis · Ensemble classifiers
MIT-BIH atrial fibrillation database

U. Desai
NMAM Institute of Technology, Udupi, India
e-mail: usha.namait@nitte.edu.in; usdesai2003@yahoo.co.in

U. Desai · G. Seshikala
REVA University, Bengaluru, India
e-mail: seshikala.g@reva.edu.in

C. G. Nayak (✉)
MIT, Manipal University, Manipal, India
e-mail: cg.nayak@maniapl.edu

R. J. Martis
VCET, Puttur, India
e-mail: roshniitsmst@gmail.com

S. L. Fernandes
SCEM, Mangalore, India
e-mail: steven.ec@sahyadri.edu.in

# 1 Introduction

Tachycardia is a precarious cardiac health problem particularly in geriatric individuals (age $\geq 65$ years). If the symptoms are not diagnosed and treated at the initial stage, eventually they can cause heart attack, stroke, or sudden cardiac death (SCD) [1]. WHO [2] statistics, estimates around 17.1 million worldwide deaths are owing to the cardiovascular diseases. Among which, 7.4 million are because of the heart attack and 6.7 million due to the stroke. In addition, this number is estimated to increase toward 23.6 million by 2030 [2]. In fact, due to the dawn of aging population, the occurrence of tachycardia is extremely rising worldwide. Also, the magnitude of geriatric population in India is going to increase by 20% of the total population by 2050. Indeed, the incidence of tachycardia is expected to grow more in Asian-Pacific region [3] and the difficulties associated with the aging people will add huge economic burden to the developing countries.

These atrial and ventricular physiological changes in the cardiac cycles are noninvasively captured using standard tool ECG. Researchers studied different automated techniques to diagnose the tachycardia beats. Christov et al. [4] used P-wave absence features and sequential analysis algorithm and classified with 98.8% of accurate result for diagnosis of three classes of tachycardia viz. A-Fib, AFL, and NSR beats. Tsipouras et al. [5] applied time-frequency analysis method for extraction of transformed domain coefficients and artificial neural network (ANN) classifier to diagnose the same three classes of fatal cardiac beats and achieved 94.2% of precision. Recently, Martis et al. [6, 7] conducted study to identify same classes namely A-Fib, AFL, and NSR beats using bispectrum coefficients and K-Nearest Neighbour (K-NN) classification system and attained 97.65% of accuracy. Besides, using discrete cosine transform (DCT) coefficients and similar classifier, Martis et al. [8] obtained performance with 99.45% of accuracy to diagnose the atrial tachycardia classes. Moreover, these studies conducted in literature considered only the atrial tachyarrhythmia classes. In which, the ventricular tachycardiac condition are not considered for development of the DSS.

Fahim and Khalil [9] discriminated four classes of arrhythmia beats viz. V-Fib, *ventricular flutter* (VFL), *premature ventricular contraction* (PVC), and A-Fib and obtained 97% of classification accuracy. In fact, the study does not reflect the typical main classes of tachycardia beats. Recently, Desai et al. [10, 11] discriminated and classified main classes of tachycardia using recurrence quantification analysis (RQA) and achieved 98.36% of classification accuracy.

Thus, it is evident from literature review that all the research studies reported emphasis on automated classification of either atrial or ventricular tachycardia beats. Our proposed work compares the performance with our earlier work [10] and addresses an automated and robust technique with higher accuracy to diagnose four main classes of tachycardia beats. In addition, we have proposed a unique visualization tool to clearly discriminate normal and tachycardia classes. Figure 1 presents the flow of proposed system. In this work, the preprocessed tachycardia beats are subjected to HOS third-order cumulants and dimensionality is reduced using

ICA. Indeed, the clinical significance ($p < 0.05$) is verified using ANOVA test. Further, the performance is compared using RAF and ROF ensemble classifiers, respectively. Continuing part of the paper is organized as follows: Sect. 2 discusses methodology applied in this study. Sections 3 and 4 illustrate results achieved and its outcomes. Finally, Sect. 5 completes this paper with conclusion.

## 2 Methodology

### 2.1 ECG Database

In this study, altogether 3858 cardiac beats, belonging to four classes of tachycardia namely, AFL (855 beats), A-Fib (887 beats), V-Fib (1108 beats), and NSR (1008 beats) are considered from the three standard databases [10].



**Fig. 1** Flow diagram of automated tachycardia diagnosis system

## 2.2 ECG Decomposition

ECG signals obtained using three different databases are sampled at a common sampling rate of 250 Hz and unwanted frequency components corresponding to baseline wander and high-frequency noise are removed using discrete wavelet transform (DWT) multi-resolution analysis system. Further, filtered ECG signals are subjected for Pan-Tompkin's algorithm for heartbeat segmentation [11–16].

## 2.3 Disease Marking

In this work, for marking the abnormalities nonlinear transform domain technique higher-order spectra (HOS) third-order cumulants coefficients are subjected. Since, cumulants exhibit the total amount of higher-order correlation and also measure the extent of non-Gaussianity. However, for Gaussian signals, cumulants' value is zero. The third-order cumulants are functions of dual time lags (*tou*1 and *tou*2) which are computed from respective ECG beat as follows:

If $[x_1, x_2, x_3, \ldots, x_n]$ represent the $k$ samples of ECG beat, its third-order cumulants $C_3^x$ are given by,

$$C_3^x = E\{x(n)x(n+tou1)x(n+tou2)\} \tag{1}$$

where $E\{\cdot\}$ is the statistical expectation operator and *tou*1 and *tou*2 are the time-lag constraints [17, 18].

Further the cumulants are estimated by replacing expectations by sample averages,

$$C_3^x = (tou1, tou2) = \frac{1}{N_R} \sum_{t \in R} x(t)x(t+tou1)x(t+tou2) \tag{2}$$

where $N_R$ is the number of samples in the region $R$.

In this study, each ECG beat of 150 samples are compressed to 126 samples after computing third-order cumulants. In addition, the classes of tachycardia are discriminated using cumulants three-dimensional graphical plots. Further the dimensionality of cumulants coefficients is reduced using ICA and subjected for ANOVA test to validate against the null hypothesis.

## 2.4 Automated Diagnosis System

In this study, RAF and ROF ensemble classifiers (ECs) are considered to develop the classification system. In this technique, a group of base learners are combined to

provide identical (or dissimilar) ensembles. Conjoining output from various classifiers results into greater prediction accuracy and decreases bias. In current work, decision tree (DT) is taken as the root (or the base) classifier to design ECs RAF and ROF, respectively. During the classifier design, TFCV is applied for dividing testing and training data samples. The classifiers performance is measured using a) class-specific accuracy (%), b) overall accuracy (%), and (c) Cohen's kappa matric ($\kappa$) [13].

## 3   Results

ECG signal is decomposed, and overall 3858 heartbeats belonging to four classes of tachycardia NSR, A-Fib, AFL, and V-Fib, respectively, are considered in this study. Wherein, each beat is of 150 samples. These beats are subjected to disease marking using HOS third-order cumulants, and later the data reduction is achieved using ICA method. Figure 2 represents the third-order cumulants plot and equivalent counter plot for four classes, respectively. Figure 3 represents deviation in accuracy in each fold of TFCV process using RAF and ROF, respectively. It can be inferred that ROF ensemble classifier accuracy is reasonably consistent compared with the RAF performance. This performance can be proved using the corresponding *kappa* value. In fact, it is higher for ROF method. The consistency over tenfold is verified using *kappa* coefficient. The summary of results achieved is presented using Table 1 using class-specific accuracy (%), overall accuracy (%), and equivalent *kappa* value. In present work, using ROF classifier, highest results are achieved with class-specific accuracy of 99.61, 99.55, 99.29, and 99.54% for NSR, A-Fib, AFL, and V-Fib, respectively. Also, 99.54% of overall accuracy is achieved with corresponding *kappa* value of 0.9895 for ROF system. Therefore, our study verifies the superiority of ROF method to achieve the best results for diagnosis of fatal classes of ECG.

## 4   Discussion

Our current paper addresses a technique to discriminate and classify the main classes of tachycardia (atrial and ventricular). Our study reveals that three-dimensional plots obtained using HOS third-order cumulants can finely characterize the normal and three different classes life-threatening of ECG beats. Indeed, our study can provide valuable information, which is very much tedious to obtain using conventional manual ECG interpretation system due to the hidden nonlinearities. Table 2 presents the comprehensive summary of studies conducted on automated diagnosis of tachycardia classes. In current paper, we have proposed unique plots for discrimination of tachycardia classes using HOS third-order cumulants. These plots represent the cardiac signatures, which can be used to

**Fig. 2** HOS cumulants plot and its contour for **a** NSR, **b** A-Fib, **c** AFL, and **d** V-Fib beats

**Fig. 3** Measure of consistency over tenfold cross-validation for RAF and ROF classifiers

**Table 1** Classification results of RAF and ROF ensemble methods

| ECs | NSR (%) | A-Fib (%) | AFL (%) | V-Fib (%) | Overall (%) | *kappa* |
|-----|---------|-----------|---------|-----------|-------------|---------|
| RAF | 99.12 | 98.76 | 99.42 | 98.92 | 99.04 | 0.9853 |
| ROF | 99.61 | 99.55 | 99.29 | 99.64 | 99.54 | 0.9895 |

distinguish NSR and tachycardia class (A-Fib, AFL, and V-Fib) using ECG beats. Our proposed DSS displays that NSR cumulants plot is characterized by clear sharper peak whereas tachycardia classes (irregular beats) typically signify with the indistinct and broader view at the bottom. This is mainly due to generation of heart rates with high frequency during the tachycardia conditions.

**Table 2** Summary of studies conducted on automated diagnosis of tachycardia beats

| Authors | Classes | Method | Accuracy (%) |
|---------|---------|--------|--------------|
| [4] | A-Fib, AFL, NSR | P-wave absence—sequential analysis | 98.8 |
| [5] | A-Fib, AFL, NSR | Time–frequency analysis–ANN | 94.2 |
| [6] | A-Fib, AFL, NSR | Bispectrum–ICA–$k$-NN | 97.65 |
| [8] | A-Fib, AFL, NSR | DCT–ICA–$k$-NN | 99.45 |
| [9] | V-Fib, VFL, PVC, A-Fib | Compressed ECG—data mining | 97 |
| [8] | A-Fib, AFL, V-Fib, NSR | RQA–TFCV–rotation forest | 98.37 |
| Proposed study | A-Fib, AFL, V-Fib, NSR | HOS third-order cumulants | 99.54 |

In this study, these fatal arrhythmia classes are appropriately classified with 99.54% of overall accuracy. Our achieved performance is higher than the earlier reported results. Our study considers in total 3858 ECG beats belonging to NSR and complete classes of tachycardia. In future, study can be performed using large database. In the meantime, this technique can be stretched to other cardiac disorders also.

## 5 Conclusion

Around the world geriatric population is rising at an alarming rate. Tachycardia beats A-Fib, AFL, and V-Fib can cause life-threatening disorder or long-term disability generally at the advanced age. In this study, a system is developed to discriminate and classify normal and three classes of tachycardia beats. Our study confirms the potential of nonlinear technique third-order cumulants discrimination level in diagnosis of fatal tachycardia classes. Proposed methodology using data reduction approach on HOS third-order cumulants and ROF ensemble classification system achieved therapeutically important classification accuracy of 99.54%. Current developed system is low cost and can be effectively used as mass screening device and as an aid for the doctors and researchers to cross-validate results. This diagnosis supporting system contributes in preventing patients from having life-threatening disorders and save life significantly.

## References

1. Goldberger, A.L.: Clinical Electrocardiography: A Simplified Approach. Mosby, St. Louis, MO (2012)
2. World Health Organization: Cardiovascular Disease: Global Atlas on Cardiovascular Disease Prevention and Control. WHO, Geneva (2012)
3. Sasayama, S.: Heart disease in Asia. Circulation **118**(25), 2669–2671 (2008)
4. Christov, I., Bortolan, G., Daskalov, I.: Sequential analysis for automatic detection of atrial fibrillation and flutter. In: Computers in Cardiology 2001. IEEE, Piscataway pp. 293–296 (2001)
5. Tsipouras, M.G., et al.: Classification of atrial tachyarrhythmias in electrocardiograms using time frequency analysis. In: Computers in Cardiology, 2004, pp. 245–248. IEEE (2004)
6. Martis, R.J., et al.: Application of higher order spectra for accurate delineation of atrial arrhythmia. In: Proceedings on Annual International Conference on IEEE Engineering in Medicine and Biology Society, pp. 57–60 (2013)
7. Martis, R.J., et al.: Application of higher order statistics for atrial arrhythmia classification. Biomed. Signal Process. Control **8**, 888–900 (2013)

8. Martis, R.J., et al.: Computer aided diagnosis of atrial arrhythmia using dimensionality reduction methods on transform domain representation. Biomed. Signal Process. Control **13**, 295–305 (2014)
9. Fahim, S., Khalil, I.: Diagnosis of cardiovascular abnormalities from compressed ECG: a data mining-based approach. IEEE Trans. Inf. Technol. Biomed. **15**, 33–39 (2011)
10. Desai, U., et al.: Diagnosis of multiclass tachycardia beats using recurrence quantification analysis and ensemble classifiers. J. Mech. Med. Biol. **16**(1–21), 1640005 (2016)
11. Desai, U., Nayak, C.G., Seshikala, G.: An application of EMD technique in detection of tachycardia beats. In: 2016 International Conference on Communication and Signal Processing (ICCSP), pp. 1420–1424. IEEE (2016)
12. Desai, U., et al.: Decision support system for arrhythmia beats using ECG signals with DCT, DWT and EMD methods: a comparative study. J. Mech. Med. Biol. **16**(1), 1640012 (2016)
13. Desai, U., et al.: Machine intelligent diagnosis of ECG for arrhythmia classification using DWT, ICA and SVM techniques. In: 2015 Annual IEEE India Conference (INDICON), pp. 1–4 (2015)
14. Desai, U., et al.: Discrete cosine transform features in automated classification of cardiac arrhythmia beats. In: Shetty, N., Prasad, N., Nalini, N. (eds.) Emerging Research in Computing, Information, Communication and Applications, pp. 153–162. Springer, New Delhi (2015)
15. Desai, U., Nayak, C.G., Seshikala, G.: An efficient technique for automated diagnosis of cardiac rhythms using electrocardiogram. In: IEEE International Conference on Recent Trends in Electronics, Information & Communication Technology. IEEE (2016)
16. Nayak, C.G., et al.: Identification of arrhythmia classes using machine-learning techniques. Int. J. Biol. Biomed. **1**, 48–53 (2016)
17. Martis, R.J., Acharya, U.R., Adeli, H.: Current methods in electrocardiogram characterization. Comput. Biol. Med. **48**, 133–149 (2014)
18. Nikias, C.L., Mendel, J.M.: Signal processing with higher-order spectra. IEEE Sig. Proc. Mag. **10**(3), 10–37 (1993)

# State-of-the-Art Object-Oriented Metrics and Its Reusability: A Decade Review

**Neelamadhab Padhy, Suresh Satapathy and R. P. Singh**

**Abstract** In this article, the importance is given to the object-oriented metrics and its reusability factor. There has been much work already done, but some improvements are needed in the software industry. Object-oriented metrics is one of the most popular and ongoing studies in different engineering branches including mathematics. Software metrics will be helpful to estimate the reusable code. The objective of this paper was to identify the reusability factor and accessibility from the last decades. A comprehensive survey on metrics and its applications has been carried out for more than one decade, and the main aim of this survey is to point out the reusable factor rather than coding. This paper shows the competence and applicability in a mixture of domains.

**Keywords** Object-oriented metrics · Reusability · Metrics model

## 1 Introduction

The term metrics is all about the measurement in the software code, and nowadays, object-oriented metrics (OOM) plays a crucial role in the industry and it is being used as a research tool. Researchers use it as a well-defined positive method to examine the data organization from the database. Several properties of the metrics have been defined

N. Padhy (✉)
Computer Science and Engineering, Satya Sai University of Technical
and Medical Science (SSSUTMS), Sehore, Madhya Pradesh, India
e-mail: neela.mbamtech@gmail.com

S. Satapathy
P.V.P.Siddhartha Institute of Engineering and Technology,
Kanuru, Vijayawada, Andhra Pradesh, India
e-mail: sureshsatapathy@ieee.org

R. P. Singh
Satya Sai University of Technical and Medical Science (SSSUTMS),
Sehore, Madhya Pradesh, India

so far, but still this measurement is not sufficient. The researchers put more emphasis on its software quality assurance. Whatever the metrics are defined, it is practically used in the programming languages. Significant work has already been done in this field, but the problem lies in the quality. To achieve a quality product is too difficult in the software industry; to fill this gap, we must deliver the good product. This can be achieved if the reusable components are provided so that time, effort and staffing are minimized. In fact, it is really a big challenge in this twenty-first century. Quality development of the software is an ongoing process; it cannot be achieved overnight. There are really two challenges for the programmers in the industry: bugs identification and rectification. Even though automated software does all the things, but measuring the quality attributes is still a challenging task from the code repositories. To find the quality attribute from the software programs, Kemerer and Chidamber [1] have developed a metric called the CK metrics suite. In-depth research has been conducted so far about the measurement and quality product [2–5]. There are other sources available, such as books and the Internet, which describe the metrics measurement [6, 7]. A major attempt has been initiated to produce the eminence of the product. At the end of the software development life cycle (SDLC), quality products will be extensively accepted; quality of the product depends on proper testing of the code. Quality is directly proportional to testability which gives satisfaction to the customer. The role of reusability plays a crucial role in this era. It is not a new concept. It is widely used for estimation of the software assessment. If the component is not reusable, then the whole concept of SDLC will fail. The new product will be developed by the existing one. A survey has been conducted on reusability.

The prime objective of this literature survey is to represent the current state-of-the-art software reusability metrics. Different researchers' views about the reusability are different.

## 2    Research Structure

In this paper, we investigated and formulated the followings:

**Research Goal 1**: What are the most reusable assets rather than coding?
**Research Goal 2**: How to estimate reusability? What are the methods or approaches for reusability?
**Research Goal 3**: How to validate the reusable products? What are the steps required to validate the reusable products?

## 3    Workflow Model for Reusability Evaluation

This is the proposed model for designing and developing the object-oriented paradigm; besides, this model provides the quantification of the reusability factor in the source codes (Fig. 1).

**Fig. 1** Workflow model

**Fig. 2** Metrics classification



Reusability occurs through the inheritance. The object-oriented design includes features such as inheritance, encapsulation, coupling, and cohesion. One of the important properties like encapsulation indicates hides the internal structure of the program. The above research framework demonstrates the quantification of the reusability. First of all, the challenging task is to identify the factors that influence reusability and testing and then to identify the properties of object-oriented design (OOD) metrics. OOD metrics help us to describe the quantification process of reusability as well as establish a multivariate linear model for reusability. The reusability factor can be achieved by using the other metrics: inheritance, cohesion, coupling, and encapsulation. This can be pictorially quantified and represented. During this study, we found different types of metrics. These are broadly classified into two types (Fig. 2).

## 4 Reusability Assets

In this paper, we only focused on the object-oriented metrics in terms of reusability and its assets. Software resources are one of the building blocks of the program paradigm; it can be one of them financial, economical. Not only the programs but also other things are reused. These are listed below. The reusable assets may consist of a single asset or several assets in one asset [8]. In this literature survey, we found 12 items that are reused apart from program code. These are as follows:

1. Used in the data
2. Modules in the program
3. Architecture-driven approach

4. Algorithms used in the program
5. Design patterns
6. Documentation for the project
7. Knowledge requirement
8. Models in the project
9. Planning stage
10. Requirement analysis
11. Service contracts

**Table 1** Reusable properties

| Reusability factor name | Meaning of the factor |
| --- | --- |
| 1. Used in the data project (UD) | This indicates that the data can be reused frequently, thus achieving the target. Data mean an experience that is recorded during the previous projects [9] |
| 2. Modules in the program (MIP) | A project is divided into several modules which contain the set of instructions. "Module" implies a single executable file that is only a part of the application, such as a DLL [10] |
| 3. Architecture driven approach (ADP) | It is the approach which represents the overall structure of the project or a component |
| 4. An algorithm used in the program (AP) | It is the reuse of the algorithms if the same type of problem occurs in the picture. Reusable algorithms are used in software designs [11] |
| 5. Design patterns (DP) | The existing design will be reused if the same type of requirement occurs |
| 6. Documentation in project (DIP) | It is one of the assets in the project. Documentation will be done during the SDLC, and a requirement specification analysis (RSA) will be done. New documents are designed which often share features of the old ones. All these are to reduce time and cost [12–14] |
| 7. Knowledge requirement (KR) | During SDLC, knowledge will be generated and treated as one of the most prominent assets for the software component. The knowledge may represent the experience, idea, or reasoning [15–18] |
| 8. Models in the project (MP) | A model can represent the task of the project, and it can consist of meaningful codes. It should also be able to represent the solutions and insights |
| 9. Requirement analysis (RA) | It is the process of gathering the requirements for the projects. A requrement is a condition or capability that must be met or possessed by a system, product, service, result or component to satisfy a contract, standard, specification, or other formally imposed document [19] |
| 10. Service contracts (SC) | It is the two-way communication between the two parties, i.e. the developers and users who are going to reuse the products. Hence, it is termed the reusable interface. This information can be helpful in predicting where and how the system can be tested, what problems might occur, and how to rectify the problem, after the system is evolved [20] |
| 11. Test cases/test design (TCTD) | After the designing stage, the tester will develop the set of test cases, which is called as test case suite, and it can further be reused. They can be reused many times for different versions belonging to the same family [21–23] |

12.  Test cases/test design

From the above literature survey, we found the property used by the author in different aspects in different contexts (Table 1).

From the above studies, the data sets are created to recognize the most valuable assets for reusability. It is concluded that most of the researchers are using the requirement analysis, to which more attention has been paid in this twenty-first century.

## 5   Data Set for Reuse Assets

See Table 2.

## 6   Performance of Reusability Assets

See Fig. 3.

This graph is called the surface graph, and it represents the assets for the reusability during the last decades. Numerous researchers have studied in the different studies papers per year (SMPY). During the years 1991–1999, more researchers have shown their efforts (e.g., UDPO, MIP, DP, DIP, AP, KR, RA, and TCTD). Interestingly, almost all of them put forth their efforts in the RA, DIP, and ADP. Apart from these, during the period 2000–2009 the researchers extensively used other kinds of assets such as MIP and SC.

**Table 2** Data set of reusable assets

| S. No. | Name of the reusable assets | No. of articles |
|--------|------------------------------|-----------------|
| 1      | UP                           | 05              |
| 2      | MIP                          | 04              |
| 3      | ADP                          | 11              |
| 4      | AP                           | 01              |
| 5      | DP                           | 10              |
| 6      | DIP                          | 09              |
| 7      | KR                           | 07              |
| 8      | MP                           | 01              |
| 9      | RA                           | 19              |
| 10     | SC                           | 02              |
| 11     | TCTD                         | 10              |

**Fig. 3** Reusable assets

## 7   Conclusion and Future Scope

The prime objective of this literature review was to investigate the reusable assets. We have fully focused on the assets and found 11 most reusable attributes in this paper. During the investigation, we found that some of the articles are mentioned as nonvalidated, but they have not proved and claimed as reusable. Hence, extensive survey is required. Our future work is all about the extensive survey concerned with the industry-oriented reusable software analysis and maintenance, because our task is to reduce the maintenance cost and staffing save time. In the industry, about 65% of fund is invested for maintenance purposes. Further, this work is to enhance the understanding of how to maintain the reusable code in the industry. Further in-depth investigation is required in the manufacturing department in the industry where requirement and design are to be more emphasized. We found that most of the researchers observed that reusable assets are validated strongly in academics but poorly in the production sector. From this survey, researcher have shown that only 36% are validated and about 60% are nonvalidated, and the rest of the researchers have shown their kin interest to review studies. Again; comparison between academics and industry, then industry-oriented survey is less than academic. Not only reusability plays a vital role in the industry, but also the new kind of challenge is aging.

## Appendix

Tables 3, 4, and 5 indicates that if studies then that article is represented as ($\sqrt{}$) mark otherwise it is indicated as (**x**). In this paper, we have investigated the limited attributes.

**Table 3** Algorithms in the program (AP)

| SL | Year | Category | Validation done or not | | | | | Not valid | Metrics | Module | Discussion | | Review | Description about an algorithm used in this paper | Reference |
|----|------|----------|------------------------|---|---|---|---|-----------|---------|--------|------------|---|--------|-------------|-----------|
| | | | Industrial case study | Academic case study | Academic experience | Industry experience | Survey | | | | Approach | Mention | | | |
| 1 | 97 | Algorithm | X | X | X | X | X | X | X | X | | √ | | Represented an algorithm and stated that an algorithm can also be reused. Focused on the main target of algorithms | [11] |

**Table 4** Data used in the project (UD)

| S. No. | Year | Category | Validation done or not | | | | | | Metrics | Module | Discussion | | Review | Description about data used in this paper | References |
| | | | Industrial case study | Academic case study | Academic experience | Industry experience | Survey | Not valid | | | Approach | Mention | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 93 | Used | X | X | X | X | X | ✓ | X | X | ✓ | – | – | It is based on the data and focuses on how it is reused | [25] |
| 2 | 94 | Data in | X | X | X | X | X | X | X | X | X | – | ✓ | | [9] |
| 3 | 96 | the | X | X | X | X | X | X | X | X | X | – | ✓ | | [26] |
| 4 | 97 | project | X | X | X | X | X | X | X | X | X | – | ✓ | | [17] |
| 05 | 04 | Algorithm | X | X | X | X | X | X | X | X | X | X | ✓ | | [27] |

**Table 5** Models in the project—MP

| S. No. | Year | Category | Validation done or not | | | | | | Metrics/tools | Module | Discussion | | Review | Description about models used in this paper | Reference |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Models in the project | Industrial case studey | Academic case study | Academic experience | Industrial experience | Survey | Not valid | | | Approach | Mention | | | |
| 1 | 06 | | X | X | X | X | X | X | √ | √ | X | X | - | Refer the models and tools for reorganizing the reusable assets, representing and how to reuse the models | [24] |

# References

1. Kemerer, C.F., Chidamber, S.R.: A metrics suite for object oriented design. IEEE Trans. Softw. Eng. **20**(6), 476–493 (1994)
2. Henry, S., Li, W., Kafura, D., Schulman, R.: Measuring object-oriented design. J. Object Oriented Program **8**(4), 48–55 (1995)
3. Lorenz, M., Kidd, J.: Object-Oriented Software Metrics. Prentice Hall Object-Oriented Series. Prentice Hall, Englewood Cliffs, NJ (1994)
4. Henderson-Sellers, B.: Object-Oriented Metrics: Measures of Complexity. Prentice Hall, Englewood Cliffs, NJ (1996)
5. Singh, Y., Kaur, A., Aggarwal, K.K., Malhotra, A.: Empirical study of object-oriented metrics. J Object Technol. **5**(8), 149–173 (2006)
6. Fenton, N., Pfleeger, S.L.: Software Metrics: A Rigorous and Practical Approach, 2nd edn. International Thomson Computer Press, London (1997)
7. Shepperd, M.J., Ince, D.: Derivation and Validation of Software Metrics. Clarendon Press, Oxford (1993)
8. Hussein, K.: Measuring Reuse Characteristics of Software Components in an Extensible IDE, pp. 16–17. VDM Verlag, Saarbrücken (2008)
9. Visaggio, G.: Process improvement through data reuse. Softw. IEEE **11**(4), 76–85 (1994)
10. Capiluppi, A., Boldyreff, C.: Coupling patterns in the effective reuse of open source software. In: Proceedings of the First international Workshop on Emerging Trends in FLOSS Research and Development (20–26 May 2007). FLOSS. IEEE Computer Society, Washington, DC (2007)
11. Karsten, W.: Reuse of algorithms: still a challenge to object-oriented programming. In: Proceedings of the 12th ACM SIGPLAN Conference on Object Oriented Programming, Systems, Languages, and Applications. ACM, Atlanta, GA (1997)
12. Childs, B., Sametinger, J.: Literate programming and documentation reuse. In: Proceedings of the Fourth International Conference on Software Reuse, pp. 205–214 (1996)
13. Levy, D.M.: Document Reuse and Document Systems, vol. 6, no. 4, pp. 339–348. Electronic Publishing (1993)
14. Gil, J., Barta, D.: A system for document reuse. In: Proceedings of the 7th Israeli Conference on Computer systems and Software Engineering, pp. 83–94. IEEE Computer Society Press, Washington, DC (1996)
15. Arango, G., Schoen, E.: Design as evolution and reuse. In: Software Reusability, 1993. Proceedings Advances in Software Reuse (1993)
16. Hall, P.A.V.: Software components and reuse. Comput. Bull. **3**(4), 14–15 (1987)
17. Yglesias, K.P.: Information reuse parallels software reuse. IBM Syst. J. **32**(4), 615–620 (1993)
18. Soundarajan, S.F.N.: Inheritance: from code reuse to reasoning reuse. In: Fifth International Conference on Software Reuse ICSR'98, p. 206 (1998)
19. IEEE standard glossary of software engineering terminology. IEEE STD 610.12–1990, 1 (1990)
20. Lucas, C., Steyaert, P.: Managing software evolution through reuse contracts. In: Software Maintenance and Reengineering, 1997. EUROMICRO 97 (1997)
21. Mark, F., Lamey, T.: Common Test Patterns and Reuse Test Designs. Microsoft, Redmond (2008)
22. Mraz, T., Anneliese, V., Mayrhauser, R.: Domain based testing: increasing test case reuse. In: Proceedings of the IEEE International Conference on Computer Design, pp. 484–491. Cambridge, MA (1994)
23. Lonngren, D.D.: Reducing the cost of test through reuse. In: AUTOTESTCON'98. IEEE Systems Readiness Technology Conference, pp. 48–53. IEEE (1998)
24. Larsen, G.: Model-driven development: assets and reuse. IBM Syst. J. **45**(3), 541–553 (2006)

25. Jones, C.: Software Return on Investment Preliminary Analysis. Software Productivity Research Inc, Hendersonville (1993)
26. Frakes, W., Terry, C.: Software reuse: metrics and models. ACM Comput. Surv. **28**, 415–435 (1996)
27. Issenin, I., Brockmeyer, E., Miranda, M., Dutt, N.: Data reuse analysis technique for software-controlled memory hierarchies. In: Proceedings of the Conference on Design, Automation and Test in Europe, vol. 1, p. 10202. IEEE Computer Society, Washington, DC (2004)

# Performance Analysis and Comparison of Various Channel Decoding Techniques

**Gupta Akanksha, Jain Anjana and Vyavahare Prakash**

**Abstract** Channel encoding and decoding are performed to remove the influence of the channel on the transmitted data that result bits in error. Some redundant bits are added with the original data to mitigate the effect of error. Various channel coding/decoding techniques are used for error correction and detection. This paper presents the comparative analysis of various decoding techniques for the purpose of selecting the optimum technique for a particular application and fulfills the requirement of the system.

## 1 Introduction

In the present scenario of mobile communication in noise-prone wireless communication channel environment, selection of proper channel codes is one of the emerging challenges. Many error correction and detection techniques have been employed to improve the performance of data transmission with minimum bit error rate probability up to $10^{-6}$ [1]. Coding gain is achieved if the encoded bits are transmitted over the channel. Compared to data bits without coding, lesser $E_b/N_o$ is required for the transmission of coded bits to maintain the constant bit error rate probability. Section 2 presents the classification of coding technique in brief. Various decoding techniques are discussed in Sect. 3. Performance analysis of

G. Akanksha (✉) · J. Anjana (✉) · V. Prakash
Shri Govindram Seksaria Institute of Technology & Science, Indore, India
e-mail: akanksha.19921992@gmail.com

J. Anjana
e-mail: jain.anjana@gmail.com

V. Prakash
e-mail: prakash.vyavahare@gmail.com

various codes is demonstrated in Sect. 4. Section 5 shows the simulation and depicts the results. Finally, conclusion is given in Sect. 6.

## 2 Classification of Codes

Various channel encoding and decoding techniques are available in literature having specific features in the terms of error detection and correction capability, with variation in complexity, memory size, etc. The error correcting and detecting code can be classified on the basis of the way in which data stream is transmitted [2]. If the transmitted data stream is fragmented into fixed sized blocks, then it is called block codes. If the stream of data is directly encoded by the encoder and some sort of memory elements are used to save the previous state and decoded output depends not only on present bit stream but also previous state, this notion is followed in convolution codes.

In the cyclic redundancy check codes, check bits are concatenated at the end of block of message. The receiver can ascertain whether or not the check bits concur with data to detect, if the received bits are in error. The classification of channel coder is shown in Fig. 1. Further, channel codes are classified on the basis of parity bit calculation.



**Fig. 1** Classification of codes based on parity bits calculation

## 3 Categorization of Decoding Techniques

Decoding technique ultimately provides the reliable recovery of transmitted bits with optimized consumption of resources like time and memory size. The various decoding techniques can be categorized as linear block codes, cyclic redundancy check codes and convolution codes shown in Fig. 2. The selection of optimized decoding technique for a particular coding depends on soft/hard decoding, algorithm used, time and system complexity, number of iterations and application (Table 1).

## 4 Performance Analysis of Channel Codes

For the comparison of channel codes, block codes and convolution codes are studied and simulated in the MATLAB by creating the script file using simulating parameter shown in Table 2.

The data stream of 365,760 bits is generated and modulated using quadrature amplitude modulation technique (16-QAM) and transmitted over additive white Gaussian noise channel. Received data are demodulated and compared with the transmitted data to estimate the number of bits in error. Simulation is done to obtain the constellation diagram and BER versus $E_b/N_o$ graph of received bits. Constellation diagram is an indication of signal as a two-dimensional $X$–$Y$ plane. It allocates the complex envelope of each possible signal. In the BER versus $E_b/N_o$ graph, BER changes with respect to the variation in any of the value of simulation parameters.



Fig. 2 Categorization of decoding techniques

**Table 1** Comparison of decoding techniques based on various parameters

| S. no. | Parameter | Linear block codes | CRC codes | Convolution codes |
|---|---|---|---|---|
| 1. | Decoding techniques | Soft/hard decision decoding, syndrome-based/syndrome-less decoding, iterative decoding, ML decoding | Polynomial division by generator polynomial is performed and provides the error detection value called as FCS | ML decoding, MAP decoding, graph based decoding, iterative decoding, sequential decoding |
| 2. | Error detection and correction capability | Single bit error and burst error | Single bit and burst error | Mainly designed for removing burst error |
| 3. | Decoding algorithm used | Concept of parity check matrix, bit flipping algorithm, sum-product algorithm [6], order statistics and generalized minimum distance algorithm for RS code [7] | CRC algorithm | Viterbi algorithm [8], BJR algorithm, stack algorithm, soft output algorithm, feedback decoding algorithm |
| 4. | BER performance | Depends on modulation scheme, channel, block size | Depends on selection of the generator polynomial | Depends on constraint length, block size |
| 5. | Complexity of retrieval of data | Single iteration is required except LDPC code | Single iteration is required | The no. of iterations is required depending on algorithm and decoding technique [9] |
| 6. | Implementation | By different folded version of Xilinx vertex 6 FPGA | Hardware based shift register implementation | Implemented by VHDL and MATLAB toolbox |
| 7. | Time complexity | $O(k)$ operations for each $(n-k)$ parity bits in C | $O(n)$ operations | $O(n)$ or $O(n^2)$ operations depending on decoding technique |
| 8. | Applications | In cryptosystem, magnetic/electrical data storage in hard disk and magnetic taps | 32-bit CRC code for internet applications, CRC-8 for XMODEM protocol | Used in digital video, radio mobile communications, satellite communications, deep space |

**Table 1** (continued)

| S. no. | Parameter | Linear block codes | CRC codes | Convolution codes |
|---|---|---|---|---|
|  |  |  |  | communications and wireless communications [10, 11] |
| 9. | Examples of codes | VRC, LRC, even parity, odd parity, hamming code, LDPC code [12, 13], RS code [14] | CRC-12, CRC-16, CRC-CCITT, CRC-32 | Turbo code, serial/parallel concatenated codes, modified turbo code |

**Table 2** Simulation parameters

| S. no. | Parameters | Value |
|---|---|---|
| 1. | Length of data stream ($N$) | 365,760 |
| 2. | Modulation scheme | 16-QAM (quadrature amplitude modulation) |
| 3. | Size of signal constellation ($M$) | 16 |
| 4. | No. of bits per symbol ($K$) | 4 |
| 5. | Transmission channel | AWGN channel |
| 6. | $E_b/N_o$ | 10 dB |

# 5  Simulation and Results

## 5.1  Error Performance Without Coding

The BER probability can be observed, when the overall communication is carried out. Binary data stream of length of $N = 365,760$ is transmitted without using any encoding and decoding technique before the modulation and after the demodulation, respectively. It is oversampled and modulated using 16-QAM modulation technique, where $M = 16$ (size of signal constellation) and $K = \log_2(M)$, number of bits per symbol is transmitted over AWGN channel with $E_b/N_o = 10$ dB and SNR is calculated by following equation.

$$SNR = E_b/N_o + 10 * \log_{10}(k) - 10 * \log_{10}(nSamp) \tag{1}$$

The constellation diagram in Fig. 3 shows the ideal constellation points and received signal constellation points. The number of bits in error due to the divergence of signal points from the centered ideal constellation points and BER versus $E_b/N_o$ graph are shown in Fig. 3. Results demonstrate that for 10 dB value of $E_b/N_o$, SNR = 16.0206, number of errors = 819, and bit error rate = $2.2 \times 10^{-3}$ are obtained. In the BER versus $E_b/N_o$ graph, performance is measured over the $E_b/N_o$ range of 2–10 dB. Simulating result of BER versus $E_b/N_o$ graph is derived same as theoretical result.

**Fig. 3** Constellation diagram and BER versus $E_b/N_o$ graph, when no coding technique is used



**Fig. 4** Constellation diagram and BER versus $E_b/N_o$ graph, when BCH coding is used

## 5.2 Error Performance with Linearly Block Coding

The same binary data stream of $N = 365{,}760$ is BCH-coded with length of the code word $(n) = 127$, length of the message bits $(k) = 120$, corresponding value of error correcting capability $(t) = 1$ and then modulated [3]. The Galois field Gf (2) is used to perform the BCH encoding and decoding operation [4]. The same value of $E_b/N_o$ and SNR (calculated by Eq. 1) is used for the transmission of signal over AWGN channel. In this case, results demonstrate that the number of bits in error is 379 and bit error rate is $0.001 \times 10^{-3}$. The performance of BER versus $E_b/N_o$ graph shown in Fig. 4 is improved compared to when binary data stream was transmitted without using any coding technique.

## 5.3   Error Performance with Convolution Coding

The convolution coding affixes more number of bits to the original, which increases its error detection and correction capability and complexity. Binary data stream of $N = 365{,}760$ is convolution-coded by using the convolution coding trellis "poly2trellis([5 4],[23 35 0; 0 5 13])" and code rate of value 2/3. For the decoding purpose, the Viterbi decoder with hard decision decoding is used [5]. Same value of parameters, i.e., modulation scheme, channel, $E_b/N_o$, is applied, and SNR value over the AWGN channel is calculated by the Eq. (2).

$$\text{SNR} = E_b/N_o + 10 * \log_{10}(k * \text{coderate}) - 10 * \log_{10}(\text{nSamp}) \qquad (2)$$

Results depict that the number of bits in error is 16 and bit error rate is 4.3748e −005. In the BER versus $E_b/N_o$ graph depicted in Fig. 5, the BER is greatly improved, when it is compared with the previous block coded or data without coding simulation results. The curve of simulated and theoretical result are very close to each other and shows the approximate same response. Comparison of all above discussed conditions are depicted in Fig. 6.

The comparison of BER versus $E_b/N_o$ graph of binary data stream without coding and binary data stream with using coding techniques shown in Fig. 6 depicts that up to 8 dB value of $E_b/N_o$, there is no significance effect of any of the coding technique due to the increased number of bits during coding procedure. As the value of $E_b/N_o$ increases, the BER (bit error rate) performance improves rapidly if the data are convolution-coded or BCH-coded. The convolution codes have better performance on the BER as there are small changes in the $E_b/N_o$ value with block size of 365,760.



**Fig. 5**  Constellation diagram and BER versus $E_b/N_o$ graph, when convolution coding is used

**Fig. 6** Comparison of BER versus $E_b/N_o$ graph, when no coding technique is used and when coding techniques are used

## 6 Conclusion

The forward error correction technique can be realized by using the block codes and convolution codes at the receiver. The comparative analysis of coding technique in terms of BER graph is presented. The simulation result shows that BER performance of convolution-coded data is better than the block-coded data and data without coding. As the length of data stream and $E_b/N_o$ value increase, the performance of BER graph is improved. The increment in the length of data stream increases the complexity of the system and requires more time for simulation. The presented work could be extended by analyzing the performance of LDPC and turbo codes in terms of length of data stream, algorithms, number of iteration required for decoding using EXIT chart.

## References

1. Bahl, L., Cocke, J., Jelinek, F., Raviv, J.: Optimal decoding of linear codes for minimizing symbol error rate. IEEE Trans. Inf. Theory **20**(2), 284–287 (1974)
2. Lin, S., Costello, D.J.: Error Control Coding. Prentice Hall, Englewood Cliffs (1982)
3. Grigorescu, E., Kaufman, T.: Explicit low-weight bases for BCH codes. IEEE Trans. Inf. Theory **58**, 78–81 (2012)
4. Rohith, S., Pavithra, S.: FPGA implementation of (15, 7) BCH Encoder and decoder for text message. Int. J. Res. Eng. Technol. **2**, 209–214 (2013)

5. Viraktamath, S.V., Talasadar, D.G., Attimarad, G.V., Radder, G.A.: Performance analysis of Viterbi decoder using different digital modulation techniques in AWGN channel. IOSR J. Electron. Commun. Eng. **9**(1), 01–06 (2014)
6. Yang, K., Wang, X.: Analysis of message-passing decoding of finite-length concatenated codes. IEEE Trans. Commun. **59**(8), 2090–2100 (2011)
7. Fossorier, M.P.C., Lin, S.: Error performance analysis for reliability based decoding algorithms. IEEE Trans. Inf. Theory **48**(1), 287–293 (2002)
8. Seshadri, N., Sundberg, C.E.W.: List Viterbi decoding algorithms with applications. IEEE Trans. Commun. **42**(2/3/4), 313–323 (1994)
9. Wang, R., Zhao, W., Giannakis, G.B.: Error correction in a convolutionally coded system. IEEE Trans. Commun. **56**(11), 1807–1815 (2008)
10. Hagenauer, J., Imai, H., Wicker, S.B.: Application of error control coding. IEEE Trans. Inf. Theory **44**(6), 2531–2560 (1998)
11. Çalhan, A., Çeken, C., Ertürk, İ.: A teaching demo application of convolutional coding techniques for wireless communication. In: International Conference on Application of Information and Communication Technologies 2009, AICT (2009)
12. Smarandache, R., Vontobel, P.O.: Quasi-cyclic LDPC codes: influence of proto- and Tanner-graph structure on minimum hamming distance upper bounds. IEEE Trans. Inf. Theory **58**, 585–607 (2012)
13. Van Nguyen, T., Nosratinia, A., Divsalar, D.: The design of rate-compatible protograph LDPC codes. IEEE Trans. Commun. **60**(10), 2841–2850 (2012)
14. Jin, W., Fossorier, M.P.C.: Towards the maximum likelihood soft decision decoding of the (255,239) Reed-Solomon code. IEEE Trans. Magn. **44**(3), 423–428 (2008)

# A Non-unit Protection Scheme for Series-Compensated Transmission System Using Fuzzy Inference System

**Praveen Kumar Mishra and Anamika Yadav**

**Abstract**  This paper presents a non-unit protection scheme for series-compensated transmission system (SCTS) using fuzzy inference system. The proposed algorithm uses the fuzzy logic technique to analyze the presence of fault and its direction from the relaying point. In this work, only phase angle of positive sequence current is taken as an input to fuzzy-based fault direction detector (FDD). The output of the proposed scheme will be '−1' for a reverse fault and '1' for forward fault. The performance of proposed algorithm is evaluated on a 735 kV, 60 Hz series compensation transmission system simulated in MATLAB Simulink platform for the wide variation in fault type, fault location, fault inception angle, and system frequency. The proposed scheme is robust, irrespective of change in the location of series capacitor and heavy-load interconnection. The main advantage of the proposed scheme is that it accurately and rapidly identifies the direction of forward as well as reverse close-in faults. The obtained results demonstrate that the proposed algorithm is found to be 100% accurate in large number of fault case studies. The proposed algorithm identifies the presence of a fault in either forward (F) or reverse (R) direction in less than half or full cycle, and results confirm its reliability, accuracy, and security.

**Keywords**  Series compensation · Sequence component · Fault detection
Directional relaying · Fuzzy inference system

P. K. Mishra · A. Yadav (✉)
Department of Electrical Engineering, National Institute of Technology,
Raipur 492010, Chhattisgarh, India
e-mail: ayadav.ele@nitrr.ac.in

P. K. Mishra
e-mail: praveenmishraeee@gmail.com

# 1  Introduction

Nowadays, protection of SCTS becomes very much crucial because it offers a cheap alternative to an existing transmission line for increasing power transmission capacity. With an involvement of series capacitor (SC) in the line, high amount of power flows, so if any type of fault occurs in the system, it may result in widespread disruption of power and revenue loss to the transmission utility. Also, changes in regulatory development and increased electricity demand require raising power transfer capability of the existing transmission system up to the thermal limit. This is feasible by the inclusion of SC compensation devices as an effective viable solution over adding another line in parallel with the network [1, 2]. Use of SC compensation is a commonly applied helps in enhancing the power transfer capability by reducing the effective transmission line length with reduced losses, and also, it enhances the transient and steady-state stability of the SCTS [3]. But the inclusion of a series capacitor in the transmission line needs changes in existing protection concepts due to the hasty change in apparent impedance seen by the relay, voltage inversion, current inversion, and also sub-synchronous resonance [4, 5]. The main challenge for fault direction and detection estimation method is a series capacitor and its protection unit. In the case of distance relay, when any fault occurred in protected line, impedance seen by relay may change because of the series capacitor which may result in mall operation of the relay. Fast and reliable fault detection, direction estimation, and fault classification are important requirements of protection systems. In modern scenario due to increasing complexity of the power system network, it is an important operational requirement to design accurate, fast, and reliable protective relaying scheme [6].

A number of research articles have been published in referred journals regarding the protection of series-compensated lines [7–23]. In [7], fast-relaying scheme based on the equal transfer process of the transmission line is discussed. In [8], author considers an average of voltage and current for relaying scheme. In [9, 10], fault detection and classification techniques using the artificial neural network have been reported. While the ANN-based techniques have been quite successful in protection tasks, sometimes it may stick in local minima. In [11], fuzzy neural network-based distance relaying scheme has been proposed. In [12, 13], support vector machines have been used in combination with wavelet transform for fault classification and section identification. In [14, 15] for detection, classification, and location of the fault, fuzzy logic-based algorithm is used. In [16], for fault classification, a combined wavelet-fuzzy approach has been reported. In [17], a probabilistic neural network-based method has been presented for both section identification and fault classification of the fault using S-transform. In [18], authors proposed a differential equation-based approach to find out relaying impedance for series-compensated line. In [19], authors reported the fault detection algorithm based on negative sequence for all types of faults during the power swing in an SCTS. In [20], authors discussed about differential impedance-based technique for detection and classification of faults in series-compensated double-circuit line and

also tested in real-time digital simulator (RTDS) platform. In [21], authors have proposed big bang-big crunch algorithm for the protection of SCTS lines and compared its results with different techniques like time-delay neural network (TDNN), radial basis function neural network (RBFNN), and nonlinear autoregressive network with exogenous inputs. In many reported research papers, authors did not mention about the fault detection time (Tr in ms) which is the most important task for relaying. In [22], authors have been reported a complete survey of different protection schemes reported till 2014. In [23], authors have been presented a scheme based on the theory of equal transfer process of transmission lines for SCTS to identify the relative position of the fault with respect to the SC and determine the fault location.

Among all the reported techniques, fuzzy logic-based technique is more feasible and reliable for protective relaying task since it does not need to solve complex higher-order equations to reach a trip decision. In this paper, fuzzy logic-based approach is presented for protective relaying of the SCTS which employs simple IF-THEN rule for decision making. The proposed fuzzy logic-based scheme emphasized on the aforesaid problems related to directional relaying and recommends an improved technique for the protection of series SCTS. The proposed technique is designed for the detection of the presence of a fault and also to identify its direction whether it is forward or reverse fault. Proposed scheme is evaluated on a 735 kV, 60 Hz power system network with midpoint series capacitor compensation using MATLAB Simulink software [24]. Inputs for fuzzy logic is based on symmetrical sequence components and fundamental components of current and voltage obtained from data acquisition system taking sampling frequency of 1.2 kHz.

## 2 Series-Compensated Power System Network

A single line diagram of the 735 kV, 60 Hz three-phase SCTS under consideration is shown in Fig. 1. It consists of two sources connected by a three-phase transmission line of 400 km with a series capacitor connected at the midpoint of Bus B4 and Bus B2. Source-1 at Bus B1 represents a power plant consisting of six 13.8 kV, 350 MVA generators which are connected to step up transformer of 13.8/735 kV, 350 MVA. The power is transferred through transmission Line-1 between Bus B3 and B4 and another Line-2 between buses B4–B2 to an equivalent grid represented by source-2 of 735 kV, 60 Hz. The length of the Line-1 and Line-2 is 100 and 300 km, respectively, having line modeled with distributed parameters. The series capacitor provides 40% compensation and an MOV and air-gap arrangement is used for overvoltage protection of the capacitor. Positive sequence reactance of Line-2 is 105.6 $\Omega$ (0.352 $\Omega$/km $\times$ 300 km) and to compensate the Line-2 by fixed capacitor capacity of 40% compensation can be calculated as 62.8 $\mu$F. In this research, the voltage and energy rating of the MOV are considered as 298.7 kV and 30 MJ, respectively. In this SCTS, four loads are connected among which Load-1

**Fig. 1** Series-compensated power system network

of 100 MW at Bus B1 and Load-2 and Load-3 of 132 MW and 332 MVAR are connected at Bus B2, respectively. In this work, a directional relay based on the discrete Fourier transform and fuzzy system is proposed, and the relay is connected to Bus B4 where the 3-phase current and voltage signals are collected and Line-1 is reverse line section and Line-2 is forward line section to be protected. The sampling rate is maintained at 1.2 kHz. The power system network is simulated in MATLAB Simulink environment [24]. When Line-2 is series compensated by 40%, the active power flows increased from 1253 to 1469 MW and also current carrying capability increased from 1395 to 1651 A.

## 3 Proposed Fuzzy-based Directional Relay

Three-phase SCTS is simulated, and three-phase current signals at the relaying are obtained. Positive sequence analyzer is used to get the phase angle of positive sequence current, and input to this sequence analyzer [24] is three-phase current at the relaying point. This phase angle of positive sequence current acts as an input to the fuzzy logic system. The basic general layout of a fuzzy system involves fuzzification, fuzzy inference system (FIS), fuzzy rule base, and defuzzification as shown in Fig. 2 for fault detection. In the first stage, i.e., fuzzification stage, crisp input variables are translated into fuzzy variables. After fuzzification, the inference engine combines the fuzzified inputs with IF-THEN rule base which provides the fuzzy output.

Finally, in the defuzzification stage, the fuzzy output is translated into crisp output. In this paper, for fuzzification, triangular membership function has been chosen.

**Fig. 2** General layout of a fuzzy system

**Table 1** Description of conditions for design of fuzzy-based fault direction detector

| FDD output | Trip condition | Fault | Fault direction |
|------------|----------------|-------|-----------------|
| 0 | TN | No | No fault |
| −1 | TR | Yes | Reverse fault |
| 1 | TF | Yes | Forward fault |

**Table 2** Fuzzy rule base of fault direction detector (FDD)

| Phase angle | IpLOW | IpMID | IpHIGH |
|-------------|-------|-------|--------|
| Trip condition | TF | TN | TR |

Here, input variable for fuzzification is phase angle of positive sequence current, and for fuzzification, it is categorized into three ranges: IpLOW, IpMID, and IpHIGH. Based on the presence of fault, fuzzy output also categorizes into three ranges: for forward fault: trip forward TF (1) or reverse direction fault: trip reverse TR (−1) or no fault condition: trip not TN (0).

Based on the presence of fault and its direction (Dr) from the relaying point, Table 1 exemplifies the details of condition of the system for designing of the fuzzy-based fault direction detector. Fuzzy rule base of fault direction detector is shown in Table 2. In Table 2, if phase angle is IpLOW, then the trip condition will be TF; if phase angle is IpMID, then the trip condition will be TN; and if phase angle is IpHIGH, then the trip condition will be TR. For designing of fuzzy-based directional relay, the phase angle of positive sequence current only is considered as an input variable. The input membership function for fault direction detector (FDD) is shown in Fig. 3.

The flow diagram of the proposed fuzzy-based directional relay is shown in Fig. 4. In this scheme, three-phase current signals are acquired at relaying Bus B4, and with the application of sequence analyzer, we get the phase angle of positive sequence current. Then, the phase angle of positive sequence current in the time domain is fed as input to the fuzzy fault direction detector which determines the presence of a fault and discriminates between forward and reverse fault. The algorithm issues the trip signal to circuit breaker only when the fault occurs in the forward direction.

**Fig. 3** Input membership functions for fault direction detection (FDD)

## 4  Results and Discussion

The feasibility of the proposed scheme has been tested on a 765 kV, 60 Hz SCTS. All the fault simulation studies have been carried out by using MATLAB Simulink software. The fault simulation studies have been carried out under different conditions and under wide variation of fault type (FT), fault inception angle (FIA), fault resistance (FR), fault location (FL), system frequency, the location of series capacitor, heavy-load interconnection, and close-in faults etc. The parameter values which have been chosen for this study are as follows: FIA (0–360° in the steps of 45°), fault resistance (0.001, 10, 20, 50 and 100 Ω), fault location (in between Bus B3 to B4-10 locations and in between Bus B4 to B2-30 locations), system frequency (58–62 Hz in the steps of 1 Hz), location of series capacitor (source side, in midpoint, both the side of source and load with 50% of actual capacity), close-in faults (for +1 and −1 km of fault location from relaying point), and for heavy-load interconnection.

For each of the above case, all the 10 types of faults (AG, BG, CG, ABG, BCG, CAG, AB, BC, CA, and ABC/ABCG) have been considered. Thus, a total of 8 * 5 * 40 * 5 * 3 * 2 = 48,000 cases have been generated for each type of faults. Thus, gross total for all 10 types of faults will be 10 * 48000 + 3 = 480,003. Both forward (F) and reverse (R) side fault cases have been considered for validation.

The response time for FDD can be calculated as:

**Fig. 4** Flow diagram of the proposed fuzzy-based directional relay



$$T_r = T_o - T_i \tag{1}$$

where $T_r$, $T_o$, and $T_i$ are response time, relay operating time, and fault inception time, respectively.

## 4.1 Performance in Case of Varying Fault Type

Effect of varying fault type on the performance of the proposed fuzzy-based algorithm is simulated with varying fault location including different FIA and by

**Table 3** Test results for effect of fault type (FR = 50 Ω)

| FT | FL (km) | FIA (°) | Response time of FDD Tr (ms) | Output of FDD | Dr |
|---|---|---|---|---|---|
| BG | 55 | 45 | 10.417 | 1 | F |
| CG | 175 | 90 | 14.133 | 1 | F |
| ABG | −65 | 135 | 7.05 | −1 | R |
| BCG | 175 | 180 | 4.97 | 1 | F |
| CAG | 205 | 225 | 8.684 | 1 | F |
| AB | −35 | 270 | 8.30 | −1 | R |
| BC | −25 | 315 | 5.0 | −1 | R |
| CA | 135 | 0 | 5.80 | 1 | F |
| ABC | 245 | 45 | 7.01 | 1 | F |

keeping FR = 50 Ω constant. Results are listed in Table 3. Table 3 shows that the proposed algorithm can detect the direction of fault for the wide variation in FL and FIA also.

Figure 5 shows the performance of proposed algorithm for LL fault. Figure 5a, b shows the three-phase voltage and current waveform in line to line (BC) fault. Figure 5c shows the output of FDD. When BC fault occurs in SCTS at −25 km (reverse side) from relaying point Bus B4 with FIA = 315°. The output becomes '−1' (Reverse) at 0.5091 s. In this work, we have taken reference time as 0.5 s which means if FIA = 0°, then fault inception time (FIT) = 0.5 s, here FIT = 0.5141 s for FIA = 315°.

So, the response time of FDD can be calculated from Eq. (1) and is 5.0 ms. Fault directions are identified correctly, i.e., reverse within a half power frequency cycle in most of the cases.

## 4.2 Performance in Case of Different Locations of Series Capacitor

The proposed scheme is also tested for change in the location of the series capacitor. Three different locations of SCs are chosen as follows: (i) at relaying point (i.e., at Bus B4), (ii) at midpoint (i.e., between Bus B4-B2), and (iii) at both the points (i.e., at Bus B4 and at Bus B2), half of SC is connected to the protected line. Some of the test results are shown in Table 4 by keeping FR, FIA, and FL are constant as 30 Ω, 0° and 50 km, respectively, and also it confirms the suitability of proposed algorithm for varying location of the series capacitor. Figure 6 shows the inputs and output of proposed algorithm for ABG fault with FR = 30 Ω, FIA = 0°, FIT = 0.5 s at 50 km from the relaying point, and two series capacitors of C/2 capacity (each) are placed at both end of the line, i.e., at Bus B4 and B2.

**(a)**



**(b)**



**(c)**



**Fig. 5** LL fault in phase C and A at −25 km, FIA = 315°, and FR = 50 Ω: **a** three-phase voltage waveform, **b** three-phase current waveform, **c** output of FDD

**Table 4** Test results for different location of series capacitor

| FT | Response time of FDD Tr (ms) | | | Output of FDD | Dr |
|---|---|---|---|---|---|
| | SC at sending end of the line (40% compensation) | SC at middle of the line (40% compensation) | SC at both sending and receiving end of the line (40% compensation) | | |
| AG | 15.8 | 8.3 | 9.1 | 1 | F |
| AB | 6.6 | 5.8 | 5.8 | 1 | F |
| ABG | 6.6 | 5.8 | 5.8 | 1 | F |
| ABC | 6.6 | 5.0 | 5.0 | 1 | F |

**Fig. 6** Double-line (ABG) fault at 50 km, FIA = 0° and FR = 30 Ω in SCTS compensated at both end with C/2: **a** three-phase voltage waveform, **b** three-phase current waveform, **c** output of FDD

## 4.3 Performance in Case of Heavy-load Interconnection

Another parameter that might have an impact on fault direction estimation is heavy-load interconnection. The impact of this parameter on the performance of proposed algorithm was evaluated by considering different situations. When there is a sudden change in load, the relay may consider it as a faulty situation which is not true. In order to analyze the effect of heavy-load interconnection on the accuracy of the proposed method, a variety of simulations has been carried out, considering different system frequencies. Some of the results are presented in Table 5. From Table 5, it is clear that proposed scheme works well in the case heavy-load inter-connection of $P = 1000$ MW and $Q = 800$ MVAR. Fig. 7a, b shows the three-phase voltage and current waveform, and in Fig. 7c, output of FDD in load is switched ON at 0.5 s and at a system frequency of 60 Hz. In this case when a load of $P = 1000$ MW and $Q = 800$, MVAR is switched ON then currents of all the

**Table 5** Simulation results for the heavy-load interconnection

| Load | | Switching instant(s) | Output of FDD |
|---|---|---|---|
| P (MW) | Q (MVAR) | | |
| 500 | 300 | 0.5 | 0 |
| 800 | 600 | | 0 |
| 1000 | 800 | | 0 |



**Fig. 7** Effect of heavy-load (P = 1000 MW and Q = 800 MVAR) interconnection at Bus B2: **a** three-phase voltage waveform, **b** current waveform, **c** output of FDD

three phases increased to more than 100 A. In Fig. 7b, it can easily visualize that peak of the current of phase B increased to 1792 A from 1622 A, then also proposed relaying algorithm is unaffected, and output of FDD remains unchanged, i.e., '0'.

**Table 6** Performance of FDD for close-in fault condition with variation in fault parameters

| FL (km) | FT | FIA (°) | FR (Ω) | Response time of FDD Tr (ms) | Output of FDD | Dr |
|---|---|---|---|---|---|---|
| −1 | 3-PHASE | 0 | 0.001 | 3.3 | −1 | R |
| | | | 50 | 4.1 | −1 | R |
| | | 45 | 0.001 | 3.717 | −1 | R |
| | | | 50 | 4.517 | −1 | R |
| | | 90 | 0.001 | 3.333 | −1 | R |
| | | | 50 | 3.333 | −1 | R |
| +1 | 3-PHASE | 0 | 0.001 | 5.0 | 1 | F |
| | | | 50 | 5.0 | 1 | F |
| | | 45 | 0.001 | 4.517 | 1 | F |
| | | | 50 | 5.417 | 1 | F |
| | | 90 | 0.001 | 4.934 | 1 | F |
| | | | 50 | 5.834 | 1 | F |



**Fig. 8** Performance during close-in fault: input and output of FDD **a** −1 km, three-phase reverse fault, FIA = 0°, FR = 50 Ω; **b** 1 km, three-phase forward fault, FIA = 0°, FR = 50 Ω

## 4.4   Performance in Case of Close-in Fault

Proposed fuzzy-based technique is also tested for close-in fault conditions. Table 6 shows the test results for close-in faults with wide variation in FIA and FR. Figure 8 shows the performance of proposed algorithm during fault at '−1' km and '+1' km from the relaying point (i.e., Bus B4). Behaviors of the phase angle of positive sequence current are different in both the cases. From the Table 6, it can observe that fuzzy-based direction estimation scheme has a response time of less than half power frequency cycle for all the fault cases. Figure 8a depicts the phase angle of positive sequence current and output of FDD during three-phase fault at −1 km from relaying point, whereas Fig. 8b represents for fault at +1 km from relaying point. In both the cases, FIT is taken at 0.5 s.

## 5   Conclusion

This paper presents a fuzzy-based directional relaying scheme for series-compensated transmission line. The fuzzy-based directional relaying scheme works in the time domain. The problems of change in system frequency, series capacitor location, heavy-load interconnection, and close-in fault on directional relaying in a series-compensated line have been addressed.

## References

1. Vyas, B., Maheshwari, R.P., Das, B.: Protection of series compensated transmission line: issues and state of art. Electr. Power Syst. Res. **107**(1), 93–108 (2014)
2. Jancke, G., Fahlen, N., Nerf, O.: Series capacitors in power system. IEEE Trans. PAS **94**(3), 915–925 (1975)
3. Anderson, P.M.: Power System Protection. IEEE Press, New York (1999)
4. IEEE Standard C37.116.: IEEE Guide for Protective Relay Application to Transmission-Line Series Capacitor Banks (2007)
5. Elmore, W.A.: Line and circuit protection. In: Dekker, M. (ed.) Protective Relaying Theory and Applications 2nd edn. pp. 273–275. New York (Chap. 12) (2003)
6. Phadke, A.G.: Computer Relaying for Power Systems. Wiley, New York (1988)
7. Qi, X., Wen, M., Yin, X., Zhang, Z., Tang, J., Cai, F.: A novel fast distance relay for series compensated transmission lines. Int. J. Electr. Power Energy Syst. **64**, 1–8 (2015)
8. Hashemi, S.M., Hagh, M.T., Seyedi, H.: High-speed relaying scheme for protection of transmission lines in the presence of thyristor-controlled series capacitor. IET Gener. Transm. Distrib. **8**(12), 2083–2091 (2014)
9. Kumari, A., Yadav, A.: ANN based fault detection in series-compensated transmission lines. In: International Conference on Magnetics, Machines & Drives (AICERA-2014 iCMMD). Amal Jyothi College of Engineering, 24–26 July 2014, Kottayam, Kerala, India
10. Verma, A., Yadav, A.: ANN based fault detection & direction estimation scheme for series compensated transmission lines. In: 2015 IEEE International Conference on Electrical,

Computer and Communication Technologies (IEEE ICECCT 2015), Paper ID EE 9037. 05–07, Mar 2015. SVS College of Engineering, Coimbatore, Tamil Nadu, India

11. Dash, P.K., Pradhan, A.K., Panda, G.: A novel fuzzy neural network based distance relaying scheme. IEEE Trans. Power Deliv. **15**(3), 902–907 (2000)
12. Dash, P.K., Samantaray, S.R., Panda, G.: Fault classification and section identification of an advanced series-compensated transmission line using support vector machine. IEEE Trans. Power Deliv. **22**, 67–73 (2007)
13. Parikh, U.B., Das, B., Maheswari, R.P.: Combined wavelet-SVM technique for fault zone detection in a series compensated transmission line. IEEE Trans. Power Deliv. **23**, 1789–1794 (2008)
14. Yadav, A., Swetapadma, A.: Enhancing the performance of transmission line directional relaying, fault classification and fault location schemes using fuzzy inference system. IET Gener. Transm. Distrib. **9**(6), 580–591 (2015)
15. Swetapadma, A., Yadav, A.: Fuzzy inference system approach for locating series, shunt, and simultaneous series-shunt faults in double circuit transmission lines. Comput. Intell. Neurosci. (2015). doi:10.1155/2015/620360
16. Pradhan, A.K., Routray, A., Biswal, B.: Higher order statistics-fuzzy integrated scheme for fault classification of a series-compensated transmission line. IEEE Trans. Power Deliv. **19**(2), 891–893 (2004)
17. Samantaray, S.R., Dash, P.K.: Pattern recognition based digital relaying for advanced series compensated line. Int. J. Electr. Power Energy Syst. **30**, 102–112 (2008)
18. Saha, M.M., Rosolowski, E., Izykowski, J., Pierz, P.: Evaluation of relaying impedance algorithms for series-compensated line. Electr. Power Syst. Res. **138**, 106–112 (2016)
19. Nayak, P.K., Pradhan, A.K.: A fault detection technique for the series-compensated line during power swing. IEEE Trans. Power Deliv. **28**(2), 714–722 (2013)
20. Jena, M.K., Samantaray, S.R.: Intelligent relaying scheme for series-compensated double circuit lines using phase angle of differential impedance. Int. J. Electr. Power Energy Syst. **70**, 17–26 (2015)
21. Deihimi, A., Solat, A.: Optimized echo state networks using big bang-big crunch algorithm for distance protection of series-compensated transmission lines. Int. J. Electr. Power Energy Syst. **54**, 408–424 (2014)
22. Vyas, B., Maheshwari, R.P., Das, B.: Protection of series compensated transmission line: issues and state of art. Electr. Power Syst. Res. **107**(1), 93–108 (2014)
23. Qi, X., Wen, M., Yin, X., Zhang, Z., Tang, J., Cai, F.: A novel fast distance relay for series compensated transmission lines. Int. J. Electr. Power Energy Syst. **64**, 1–8 (2015)
24. MATLAB, Version R2013a. The MathWorks Inc., Natick, MA, USA

# Intelligent System for Team Selection and Decision Making in the Game of Cricket

**Narendra Siripurapu, Ayush Mittal, Raghuveer P. Mukku
and Ritu Tiwari**

**Abstract** The traditional way of team selection in the game of cricket requires lot of expertise and consumes a lot of time. To make this process simpler and easy for the selection committee, an Adaptive Neuro-Fuzzy Inference Model is developed that considers various parameters of a player. Using the player parameters, he/she is clustered and rated by the use of Fuzzy rules into different categories as per his/her performance throughout the career. The player data along with their rating are sent into an Android application that does the task of team selection. This would ease out the work of the selection committee.

**Keywords** Team selection · Adaptive neuro-fuzzy inference model
Fuzzy rules

## 1 Introduction

The traditional way of team selection for any sport is based on the expert panel judgments and most importantly includes negotiative decision making by these experts. It is in fact difficult to predict a player's performance based on their past records. But, the expert panel sitting for the team selection has to consider this fact. This plays a vital role in the selection of the team for any sport for the international competitions. Here, we make use of the advancement in the technology to help the

N. Siripurapu (✉) · A. Mittal · R. P. Mukku · R. Tiwari
Robotics and Intelligent System Design Lab, Indian Institute of Information Technology and Management, Gwalior, India
e-mail: narendra.s1502@gmail.com

A. Mittal
e-mail: ayush2709@gmail.com

R. P. Mukku
e-mail: raghuveerprasadmukku@gmail.com

R. Tiwari
e-mail: tiwariritu2@gmail.com

expert panel to assess the players based on their past performances. We use some nonlinear modeling techniques like neural networks [1] to provide the experts with an analytical aid for such decision making.

Traditional method [2] of selecting a team includes a panel of experts called selection committee, who rate the player based on their performance and then voting is conducted so as to know whether or not to include a player in the squad. After arriving at a final squad of players, negotiative decisions are made so as to recommend which player plays for the team in the respective competition. In this paper, we will discuss how to rate the players using the Neuro-Fuzzy Models and then use them to make the best team out of them.

## 2   Related Work

Kahn [3] implemented multilayer perceptron model to forecast the outcome of the National Football League winners. Flitman [4] used Neural Networks and Linear Program Model to predict the winners in the Australian Football League, and a Neural Network application to model bowler action. Subramanian and Ramesh [5] had designed a system which uses neural networks to predict player's performance from the past record of the players. Kaluarachchi and Varde [6] developed a software tool 'CricAI,' which predicts the probability of victory in a ODI cricket match. Bayesian classifiers were used on input factors such as batting first, winning the toss, day/night effect, and home game advantage. Zualkernan et al. [7] designed a system that uses machine learning techniques and Dynamic Time Warping providing a feedback to the coaches and the players which helps them in improving their technique and also their performance in the game of cricket.

Singh et al. [8] had developed a model to predict the first innings score and outcome of second innings not only on the basis of current run rate but also considering factors such as venue of the match, batting team, and number of wickets. Pathak and Hardik [9] had applied modern classification techniques such as Random Forest, Support Vector Machines, and Naive Bayesian to determine the outcome of a cricket match. Prabhu and Shaila [10] had developed a Cricket Strategy Game which focuses on the captaining, managing, strategizing, and aspects of cricket.

## 3   Methodology

This Adaptive Neuro-Fuzzy Inference System (ANFIS) combines the principles of both fuzzy logic and neural networks and has the capability to have their benefits under a single umbrella. The backbone of this system is the set of IF-THEN rules using fuzzy logics, which have the ability to learn using some nonlinear functions.

## 3.1 Process Flow

The process flow of this intelligent system is described as follows, and the corresponding block diagram is shown in the Fig. 1.

- Created Input Variables for a batsman which include number of matches played, innings played, runs scored, batting average, team strength, and opponent strength.
- Created input variables for a bowler which include number of matches played, innings played, bowling strike rate, economy rate, wickets taken, team strength, and opponent strength.
- Rules have been defined for Fuzzy Inference System (FIS), which are used to rate players for different series and for different opponents.
- Newer players are given moderate rating initially.



**Fig. 1** Flow diagram of the proposed work

- Dataset [11] for various series against different teams for team India has been collected and normalization is done.
- The normalized dataset is then trained, and an FIS is generated using the rules and dataset.
- Finally, once the network is trained, the output of the FIS is the rating of the players.
- The Android application based on player ratings and user choices recommends the team squad.

### 3.2   Technical Details

**Rules for FIS**  The batting average and bowling strike rates of players are calculated using exponential averaging on the previous 5 performances of the player against the respective team. The decay factor used was 4% for each performance considering the latest performance of the player as the most important one.

Rules have been defined to rate players under various circumstances against different opponent teams. The players are rated as Very High, High, Moderate, Low, and Very Low. The players having high range of values for the parameters, matches played, innings played, batting average, runs scored, team strength, and opponent team strength, are rated Very High. The players having next range of values are rated High and similarly Moderate, Low, and Very Low. For example, if matches played are high, innings played are high, batting average is high, runs scored is high, team strength is high, and the opponent team strength is high, then the player is rated Very High.

**Team Selection Using the Android Application**  The Android application developed for the proposed work is named 'MyCricket Auto-Squad.' It was designed in such a way that it works on all the Android versions. The application is user-friendly and provides various options. The application automatically chooses the players from the squad of available players using the ratings given to players.

In the team selection process, the user is given choices to select the series for which he wants to form the team. Here, in the application, we have deployed the team selection process for three different series India versus Australia, India versus South Africa, and England versus Australia. The series can be in any of the three important format of the game, namely Test, One Day Internationals (ODI), and T20.

The teams will be chosen according to the performance of the player against that particular opposition. This is because the players were rated as per their strength and performance against specific oppositions. The application also facilitates the user to choose his alternative among the given popular team formats. The team format selection option facilities the user to strengthen bowling and batting departments according to the game situation. The team format may be as follows:

- 5 Batsmen; 1 Wicketkeeper; 4 Bowlers; and 1 All-rounders
- 5 Batsmen; 1 Wicketkeeper; 3 Bowlers; and 2 All-rounders
- 4 Batsmen; 1 Wicketkeeper; 4 Bowlers; and 2 All-rounders

# 4 Experimental Results

The output of the ANFIS is the rating of the players. Based on these player ratings and user choices, the Android application recommends the best team.

## 4.1 Results of ANFIS

Figure 2 shows the training of the ANIFS. The Training Data Plot contains player's dataset index on its X-axis and their corresponding expected output ratings on Y-axis.

Figure 3 shows the testing of the ANIFS. The Testing Data Plot contains player's data index on its X-axis and their corresponding expected output ratings and actual output ratings on Y-axis. The blue dots indicate expected output ratings, and red dots indicate actual output ratings.

**ANFIS information for** Figs. 2 and 3.



**Fig. 2** Training ANFIS

**Fig. 3** Testing ANFIS

- Number of nodes: 116
- Number of linear parameters: 54
- Number of nonlinear parameters: 90
- Total number of parameters: 144
- Number of training data pairs: 22
- Number of fuzzy rules: 9
- Epochs: 8

## 4.2 Results of Team Selection Using Android Application

Figure 4 shows the teams selected using the Android application. In this case, the user selects Test format for the India vs South Africa series. The team formats 5–1–4–1 and 5–1–3–2 for the teams India and South Africa, respectively, were chosen by the user. Eleven member teams are selected for both the teams as the final result.

## 5 Conclusion and Future Scope

The proposed Neuro-Fuzzy Model rates the players into five categories based on their performances. The categories are as follows: Very High, High, Medium, Low, and Very Low, and using these ratings, the players are selected into the team. The

**Fig. 4** Teams selected



proposed rating method produced similar ratings as that of the Reliance ICC Ratings System [12]. The proposed method takes care of the players' recent performance. The player is judged on all the parameters of his career and then rated. The Android application developed will aid the selectors and recommends them with a team of best eleven players from all the available players for the team. The new players are rated in the same way using their domestic level statistics. They are also considered better performers and are included in the list of the available players for the team.

The Adaptive Neuro-Fuzzy Model used here considers only the players of a team and not the whole cricketing players of all the countries for rating. So there is scope that the rating may be improved by the use of all the players at a single stretch. The parameters used to rate players can be improved, and many more fuzzy rules can be added to this model in that respect to get more clarity about some players' performance.

# References

1. Shukla, A., Tiwari, R., Kala, R.: Real Life Applications of Soft Computing, pp. 41–101. CRC Press, Boca Raton (2010)
2. Team selection. http://en.wikipedia.org/wiki/Indianationalcricketteamselectors

3.  Kahn, J.: Neural network prediction of NFL football games. In: World Wide Web Electronic Publication, pp. 9–15 (2003)
4.  Flitman, A.M.: A hybrid approach utilizing genetically defined neural networks and linear programming. Comput. Oper. Res. **33**(7), 2003–2022 (2006)
5.  Subramanian, I., Ramesh, S.: Prediction of athletes performance using neural networks: an application in cricket team selection. Expert Syst. Appl. **36**, 5510–5522 (2009)
6.  Kaluarachchi, A., Varde, A.S.: CricAI: a classification based tool to predict the outcome in ODI cricket, information and automation for sustainability (ICIAFs). In: 5th International Conference, pp. 250–255. IEEE (2010)
7.  Zualkernan, I.A., Assaleh, K., Dabrai, S.G., Hoque, M.H., Pedhiwala, H.Y.: A wireless electronic training system for cricket. In: IEEE 13th International Conference Advanced Learning Technologies (ICALT), pp. 55–57 (2013)
8.  Singh, T., Singla, V., Bhatia, P.: Score and winning prediction in cricket through data mining. In: International Conference on Soft Computing Techniques and Implementations (ICSCTI), pp. 60–66. IEEE (2015)
9.  Pathak, N., Hardik, W.: Applications of modern classification techniques to predict the outcome of ODI cricket. Proced. Comput. Sci. **87**, 55–60 (2016)
10. Prabhu, S., Shaila, P.: Cricket strategy game: a management game in cricket using C++. In: 2015 International Conference on Nascent Technologies in the Engineering Field (ICNTE), pp. 1–6. IEEE (2015)
11. Dataset for player statistics is collected from. http://www.espncricinfo.com, http://www.cricbuzz.com
12. Reliance ICC Player Rankings website. http://www.relianceiccrankings.com

# A Study of Clustering Techniques for Wireless Sensor Networks

**Neeharika Kalla and Pritee Parwekar**

**Abstract** Wireless sensor network (WSN), in recent times, has been exhibiting potential in variety of applications like military surveillance, disaster management, and wildlife monitoring. In these applications, humans cannot access the place either to deploy the sensors or to monitor the sensors as the environment is inhospitable. Therefore, the sensors are expected to be remotely deployed and to operate in an autonomous mode. To support scalability, sensor nodes are formed or organized into clusters and the clusters formed are without overlap and disjoint. In this paper, classification of clustering attributes is presented from the various published clustering algorithms of WSN. The clustering schemes are compared based on the metrics such as sensor mobility, overlap of clusters, position awareness, efficient energy-based, uniform clustering, and stability of cluster.

**Keywords** Wireless sensor network · Clustering · Base station · Cluster head

## 1 Introduction

Recent advances in wireless sensor network (WSN) have led to the increase in the limited-power, low-cost, tiny-sized, multi-functional sensors that are capable of communicating over short distance range. Sensors are embedded with microprocessor, radio receiver, and power components to enable sensing, data processing, and communication [1]. Sensors are battery-operated devices; hence, decrease in the energy consumption to prolong the lifetime of the network is a major challenge in the field of wireless sensor network (WSN).

In the sensing field, the sensors are deployed randomly along with the base station (BS). The function of the sensor is to sense the physical conditions of the

N. Kalla (✉) · P. Parwekar
Anil Neerukonda Institute of Technology and Sciences, Bheemunipatnam, India
e-mail: neeharika.kalla@gmail.com

P. Parwekar
e-mail: pritee.cse@anits.edu.in

environment like temperature, humidity, pressure, light, and disaster and send these data to the BS. The BS may be placed far away from the sensor nodes. Each and every sensor must report its sensed information to the BS either directly in single step or through multi-hop fashion. The sensors that are near to the BS consume less energy when compared with the sensors that are placed far away from the BS. So, more the energy consumption of the sensors less will be its battery life as the sensors have limited power.

Wireless sensor network have variety of applications like military surveillance, disaster management, health monitoring, air pollution monitoring, and area monitoring. In some applications like disaster management, humans cannot directly interact to deploy the sensors or to monitor the sensors as the environment is so harsh. These applications require thousands of nodes to be deployed. So, the sensors need to be deployed in a remote manner and these are operated autonomously [2]. The WSN must have the feature of scalability to manage the thousands of sensors. To support scalability and to reduce the energy consumption of the sensor nodes, clustering technique is used. A detailed description of clustering is being presented in the next section of the paper.

## 2   Clustering

Clustering is a technique in which the sensors are organized into different number of clusters. In each cluster, one node acts as a cluster head (CH) and the remaining nodes form the members of the cluster. The member nodes of the cluster sense their physical surroundings and transmit the data to the CH. The CH is responsible for aggregating the data that are received and reporting the data to the sink or BS. The cluster members send the data to the CH usually in a single-hop fashion. And it is termed as intra-cluster topology [3].

The main advantage of clustering is that it supports network scalability [2]. The CH can be randomly picked from the set of deployed sensors, or it can also be decided by the network. If the CH is pre-assigned, then it can be provided in advance with more resources [2]. In this, the clusters formed are stable throughout the network as the cluster head is fixed. If the CH is selected randomly, then the clusters are not stable and they change dynamically. When compared with the cluster members, CHs must be provided with more energy and they must always be in active state (Fig. 1).

## 3   Critical Analysis of Clustering Protocols

A brief description of the published clustering protocols is provided, and the protocols are compared against the metrics sensor mobility, overlap of clusters, position awareness, efficient energy-based, uniform clustering, and stability of cluster.

**Fig. 1** Clustering technique in WSN



**Fig. 2** Timeline chart of protocols

The timeline chart in Fig. 2 shows the important protocols used for the clustering from year 2002 to 2016.

### 3.1  Improved EADUC Protocol

In [4], the nodes are randomly deployed and are energy heterogeneous. Base station (BS) is located far from the sensor field. This protocol operates in rounds. After the nodes are deployed, BS broadcasts a signal to all the nodes. The nodes then approximate its distance to the BS based on the received signal strength. Every round comprises of two stages: cluster setup stage and steady-state stage. Setup phase consist of three sub-phases. The first sub-phase is the neighbor node information collection. In the second sub-stage, selection of cluster heads takes place. In the last sub-stage, the cluster formation starts. Once the network is divided into clusters, the steady-state stage begins in which the transmission of data takes place. The steady-state phase consists of number of major slots. After one major slot, the cluster head rotation takes place. An extension of EADUC protocol is proposed in order to reduce the cluster overhead and also to balance the energy.

### 3.2  Cluster Head Rotation Mechanism

In [5], for a random period of time the sensor node goes into idle mode. Node self-elects and declares itself as a cluster head and broadcasts an advertisement message to all the nodes within the range. After receiving the response from the neighboring nodes, it decides whether to continue as a cluster head or to join as a member node in any of the existing clusters. For cluster head rotation, it uses a scheme in which the initial cluster head nominates the successor cluster head based on the information collected by the initial cluster head. This information contains the identification of member nodes sorted as per received signal strength indicator (RSSI) reading and sorted in a packet called NODE packet. This packet contains net location field with the ID of the successor cluster head. The initial cluster head communicates this packet to the successor cluster head. In this way, the cluster heads are rotated.

### 3.3  DUCF Protocol

In [6], initially all the sensors are randomly deployed. RSSI is used to calculate the distance between the sensor nodes. Three metrics are considered for selecting the node as a cluster head namely residual energy of the node, degree of the node, and distance from the node to BS. Similarly two output variables are also used namely

chance and size. This protocol operates in two stages: cluster formation stage and data collection stage. In the cluster formation stage, cluster heads are selected using fuzzy logic. In the data collection phase, a TDMA schedule is generated by each cluster head and the nodes should send the information to the cluster head within that schedule.

### 3.4 Hierarchal Clustering Algorithm with Cluster Head Selection

Generally, a hierarchical protocol works in two layers. The first layer is used for selection of cluster heads, and the second layer is used for routing of the data. The clustering algorithm mechanism mentioned in [7] is used for cluster head selection. In this, cluster head selection is the key point which aims for decreased energy consumption and decrease in delay. Here, three parameters are used namely, remaining energy of the node, distance from the BS, and degree of the node (surrounding nodes). In every round, the weight of each node is calculated using the above-mentioned three parameters. The node with the maximum weight will be elected as cluster head for that cluster.

### 3.5 Mobile Sink-Based Approach

In [8], sensors are randomly deployed and are static. But the sink is mobile. In order to reduce the overhead on the single mobile sink, two mobile sinks are used and they are moved in the counterclockwise direction for every half of the round. LEACH protocol is used for the cluster formation and cluster head selection. The sensing field is equally divided into regions and each region contains clusters which are of unequal size. Initially, the node which is at the center of the region is selected as a cluster head. At the end of each round, the node that is closer to the center point and the node that is having higher residual energy are selected as a cluster head. As the sink is mobile, nodes must set up the route to deliver the data corresponding to the new location of the sink. Only one cluster head is responsible for maintaining the new route information of the mobile sinks.

### 3.6 SECC Protocol

In [9], the protocol consists of two phases. In the first phase, the energy of an individual node and the distance between its adjacent sensor nodes is calculated. Sensor nodes with nearly same average distance range are grouped into same

cluster. In the next phase, energy-aware clusters are formed based on the threshold value. The nodes whose energy value is less than threshold is disabled and do not participate in the sensor operations.

## 3.7 P-LEACH Protocol

In [10], P-LEACH algorithm uses the combination of both PEGASIS and LEACH protocols. From PEGASIS, chain formation implementation is considered and from LEACH, data forwarding through the cluster head is considered. The nodes are divided into clusters based on the sensing range of the REQ message sent by the BS. Then, all the nodes of each cluster send their location and the energy information to the BS. Now for cluster head selection, the BS selects a node with higher remaining energy as a cluster head for each cluster. After the election of the cluster heads, the cluster head with minimum distance to the BS and with maximum energy is elected as a leader.

## 3.8 H-LEACH Protocol

H-LEACH proposed in [11] is the combination of HEED and LEACH protocols. When the setup phase is started, initially the average energy of all the nodes is calculated by using the required parameters. Then, the threshold value is calculated. Now, a random number is picked by the node in the range of 0 and 1. If the random number is less than the threshold and the sensor is having the energy greater than the average energy, then the corresponding node is elected as a cluster head. The energy required for data transmission is deducted from the energy of the node.

## 3.9 Clustering Using Fuzzy Logic

In [12], the algorithm operates in two phases. Residual energy and the required energy are the two linguistic input variables used in the fuzzy logic system. In the first phase, each sensor node calculates the energy required to transmit '$k$' number of bits to the BS based on the received signal strength of the message sent by the BS. Nodes set its timer value inversely proportional to the output variable of the fuzzy logic system called as 'chance'. The higher the value of chance, the more will be the probability of the node to become a cluster head. Each sensor node sets the countdown, and if it reaches zero, the node advertises itself as a CH and clusters are formed based on the distance between the sensors. In the second phase, TDMA schedule is used by the cluster head and sends it to the cluster members.

### 3.10  Area Coverage-Aware Clustering Protocol (ACACP)

In [13], the setup stage consists of information update, sensor activation, CH election, and relay node selection phases for the formation of clusters. Steady-state stage consists of data communication phase in which the active sensor node collects the data from time to time and sends this information to the CH nodes. Then, CHs aggregate the data and transmit these data to the BS through the relay node.

### 3.11  Virtual Grid Margin Optimization and Energy Balancing Scheme (VGMEB)

In [14], the network is organized in the form of grids in order to achieve uniform clustering. Here, each grid of the network uses cluster evaluation model to select the cluster head. The cluster evaluation model is the probability of selecting the cluster head which is located at the center of the grid so that the consumption of energy can be reduced. The weight factor of the node is calculated to find the probability of selection of a cluster head. The selected CH is responsible for forwarding the data to the sink by performing data fusion. As the sink is mobile, it is responsible for informing its location to all the cluster heads when moving from one grid to another.

### 3.12  Hybrid Backtracking Search Optimization Algorithm (BSA)

The protocol in [15] is a hybrid of genetic algorithm and k-means. Selection, crossover, and mutation operations are performed. This algorithm has two selection stages in order to update population 'P' with final trial population 'T' individuals, if they have better fitness value than 'P'. Next, k-means algorithm is run to find the centroids for the clusters. These centroids are mapped to the nearest sensor nodes to obtain new cluster heads and again new clusters are formed.

Table 1 shows the summary of all the clustering protocols with their attributes.

## 4  Conclusion

In this paper, some published clustering techniques have been studied and these algorithms are categorized based on the classification of clustering attributes. Further, the clustering algorithms are compared based on the metrics such as sensor mobility, overlap of clusters, and position awareness. It can be concluded that

**Table 1** Comparison of different clustering protocols based on the metrics

| Technique used | Sensor mobility | Overlap of clusters | Position awareness | Efficient energy-based | Uniform clustering | Stability of cluster | Scalability |
|---|---|---|---|---|---|---|---|
| Improved EADUC protocol [4] | No | No | Not required | Yes | No | Yes | Yes |
| Cluster head rotation mechanism [5] | Stationary | No | Not required | Yes | Either balanced or unbalanced as the cluster formation takes place based on the response of the neighboring nodes | Yes | Yes |
| DUCF protocol [6] | Stationary | No | Not required | Yes | Yes | Yes | Yes |
| Mobile sink approach [8] | Stationary with multiple mobile sink | No | Not required | Yes | No | Yes. In this only the sink is mobile but the sensors are stationary | Yes |
| Hierarchal protocol with CH selection [7] | No | No | Required | Yes and delay is also considered | Clusters can be formed initially by using any approach | Yes | Yes |
| SECC protocol [9] | Yes | Not present as clusters are formed based on the same distance value | Required | Yes | Clusters are formed based on the distance between the sensor nodes. Cluster may or may not be balanced | No. Because new clusters are formed in every round | Yes |
| | No | No | Required | Yes | Clusters are formed based on the sensing range | Yes | Yes |

(continued)

**Table 1** (continued)

| Technique used | Sensor mobility | Overlap of clusters | Position awareness | Efficient energy-based | Uniform clustering | Stability of cluster | Scalability |
|---|---|---|---|---|---|---|---|
| P-LEACH protocol [10] | | | | | | | |
| H-LEACH protocol [11] | No | No | Not required | Yes | Initially clusters can be formed by using any protocol | Yes | Yes |
| Clustering using fuzzy logic [12] | No | No | Not required | Yes | Yes | Only cluster formation is described in this paper | Yes |
| ACACP [13] | No | No | Not required | It focuses on coverage area problem | Clustering is done based on the cover sets | No. In each round, cluster membership evolves | Yes |
| VGMEB [14] | Sensors are stationary, but sink is mobile | No | Not required | Yes | Yes | Yes | Yes |
| BSA [15] | Stationary | No | Not required | Yes | Clustering is done based on the distance between the cluster member and the cluster head | No | Yes |

clustering technique greatly reduces the energy consumed by the nodes and thus extends the stability of the network. Each of the illustrated technique can be treated as a milestone in the domain, a unique solution or a significant improvement over an existing method. By this fastidious review, we try to lend a beneficial resource for researchers and practitioners in the field of wireless sensor networks.

# References

1. Singh, S.K., Singh, M.P., Singh, D.K.: A survey of energy-efficient hierarchical cluster-based routing in wireless sensor networks. Int. J. Adv. Netw. Appl. **2**(2), 570–580 (2010)
2. Abbasi, A.A., Younis, M.: A survey on clustering algorithms for wireless sensor networks. Comput. Commun. **30**(14–15), 2826–2841 (2007)
3. Singla, S., Kaur, K.: Comparative analysis of homogeneous N heterogeneous protocols in WSN. Int. J. Sci. Res. **5**(6), 1300–1305 (2016)
4. Vrinda, G., Rajoo, P.: An improved energy aware distributed unequal clustering protocol for heterogeneous wireless sensor networks. Eng. Sci. Technol. Int. J. **19**(2), 1050–1058 (2016)
5. Pradhan, S., Sharma, K.: Cluster head rotation in wireless sensor network: a simplified approach. Int. J. Sens. Appl. Control Syst. **4**(1), 1–10 (2016)
6. Baranidharan, B., Santhi, B.: DUCF: distributed load balancing unequal clustering in wireless sensor networks using fuzzy approach. Appl. Soft Comput. **40**, 495–506 (2016)
7. Chaubey, N.K., Patel, D.H.: Energy efficient clustering algorithm for decreasing energy consumption and delay in wireless sensor networks (WSN). Int. J. Innov. Res. Comput. Commun. Eng. **4**(5), 8652–8656 (2016)
8. Nagamalar, T., Rangaswamy, T.R.: Energy efficient cluster based approach for data collection in wireless sensor networks with multiple mobile sink. In: 2015 International Conference on Industrial Instrumentation and Control (ICIC) 348–353 (Pune, 2015)
9. Bala Krishna, M., Doja, M.N.: Self-organized energy conscious clustering protocol for wireless sensor networks. In: 2012 14th International Conference on Advanced Communication Technology (ICACT) 521–526 (PyeongChang, 2012)
10. Razaque, A., Abdulgader, M., Joshi, C., Amsaad, F., Chauhan, M.: P-LEACH: energy efficient routing protocol for wireless sensor networks. In: 2016 IEEE Long Island Systems, Applications and Technology Conference (LISAT) 1–5 (Farmingdale, NY, 2016)
11. Razaque, A., Mudigulam, S., Gavini, K., Amsaad, F., Abdulgader, M., Krishna, G.S.: H-LEACH: hybrid-low energy adaptive clustering hierarchy for wireless sensor networks. In: 2016 IEEE Long Island Systems, Applications and Technology Conference (LISAT) 1–4 (Farmingdale, NY, 2016)
12. Singh, M., Soni, G.S., Kumar, V.: Clustering using fuzzy logic in wireless sensor networks. In: 2016 3rd International Conference on Computing for Sustainable Global Development (INDIACom) 1669–1674 (New Delhi, 2016)
13. Nguyen, T.G., So-In, C., Nguyen, N., Phoemphon, S.: A novel energy-efficient clustering protocol with area coverage awareness for wireless sensor networks. Peer-to-Peer Netw. Appl. **10**, 1–18 (2016)
14. Tang, C., Yang, N.: Virtual grid margin optimization and energy balancing scheme for mobile sinks in wireless sensor networks. Multimed. Tools Appl. **76**, 1–20 (2016)
15. Latiff, N.M.A., Malik, N.N.N.A., Idoumghar, L.: Hybrid backtracking search optimization algorithm and K-means for clustering in wireless sensor networks. In: 2016 IEEE 14th International Conference on Dependable, Autonomic and Secure Computing, 14th International Conference on Pervasive Intelligence and Computing, 2nd International Conference on Big Data Intelligence and Computing and Cyber Science and Technology Congress (DASC/PiCom/DataCom/CyberSciTech) 558–564 (2016)

16. Heinzelman, W.B., Chandrakasan, A.P., Balakrishnan, H.: An application-specific protocol architecture for wireless microsensor networks. IEEE Trans. Wirel. Commun. **1**(4), 660–670 (2002)
17. Coyle, E.J., Bandyopadhyay, S.: An energy efficient hierarchical clustering algorithm for wireless sensor networks. In: IEEE INFOCOM 2003. Twenty-second Annual Joint Conference of the IEEE Computer and Communications Societies (IEEE Cat. No. 03CH37428), vol. 3, pp. 1713–1723 (2003)
18. Fahmy, S., Younis, O.: HEED: a hybrid, energy-efficient, distributed clustering approach for ad hoc sensor networks. IEEE Trans. Mobile Comput. **3**(4), 366–379 (2004)
19. Yulin, S., Xiangning, F.: Improvement on LEACH protocol of wireless sensor network. In: 2007 International Conference on Sensor Technologies and Applications (SENSORCOMM 2007) 260–264 (Valencia, 2007)
20. Kim, J.M., Park, S.H., Han, Y.J., Chung, T.M.: CHEF: cluster head election mechanism using fuzzy logic in wireless sensor networks. In: 2008 10th International Conference on Advanced Communication Technology 654–659 (Gangwon-Do, 2008)
21. Chen, G., Li, C., Ye, M., Jie, W.: An unequal cluster-based routing protocol in wireless sensor networks. Wirel. Netw. **15**(2), 193–207 (2009)

# Intensifying the Security of Information by the Fusion of Random Substitution Technique and Enhanced DES

**Rayi Sailaja, Ch. Rupa and A. S. N. Chakravarthy**

**Abstract** Information security is the process of defending information from unauthorized access and loss. Nowadays in most of the organizations, information is accessed through online. So there is a need to safeguard the information from attackers and hackers. Cryptography is the method of storing and transmitting data in an unreadable form so that only indented user can read and process it. There are two categories of cryptography: symmetric key and asymmetric key encryption algorithms. Symmetric key algorithms use single key for both encryption and decryption. RC2, RC4, DES, Triple DES, IDEA etc. are some of eminent existing symmetric key algorithms. In this paper, we tried to enhance the security of DES, which is now considered to be insecure for many applications. We tried to provide double security by adding a random substitution technique before applying EDES algorithm where we modified the round function and increased the number of rounds from 16 to 32. The performance of the proposed technique is compared with traditional DES algorithm by means of avalanche effect and plaintext sensitivity.

**Keywords** Random substitution technique · Symmetric encryption · DES Round function · Avalanche effect

R. Sailaja (✉)
Aditya College of Engineering & Technology, Surampalem, India
e-mail: sailajarai@gmail.com

Ch. Rupa
VR Siddhartha Engineering College, Vijayawada, India
e-mail: rupamtech@gmail.com

A. S. N. Chakravarthy
JNTUK, Kakinada, India
e-mail: asnchakravarthy@yahoo.com

# 1   Introduction

Network security consists of the provisions and policies which are adopted to prevent and monitor misuse, modification, denial of a computer network, and unauthorized access [1].

Cryptography is the art of achieving security by encoding messages to make them non-readable [2]. Cryptography is the practice and study of hiding information [2]. There are two basic types of cryptography: symmetric key and asymmetric key. Symmetric key algorithms are the fastest and most generally used type of encryption [3, 4]. Symmetric cryptography is a type of cryptosystem in which encryption and decryption are performed using the same key. It is well suited for lengthy data streams. It transforms plaintext into cipher text using a secret key and an encryption algorithm. Using the decryption algorithm and the same key, the plaintext is recovered from the cipher text [5, 6].

In this paper, we have discussed about the drawbacks of DES and proposed a random substitution method and an enhanced round function of DES to improve its performance. This paper is organized as follows. Section 2 introduces the Existing System, Sect. 3 introduces proposed method, and Sect. 4 talks about experimental results, concluding remarks are given in Sect. 5.

# 2   Existing System

Symmetric encryption is a variety of cryptosystem in which encryption and decryption are done using the same secret key [6].

The Data Encryption Standard is one of the well-known symmetric key algorithms for encryption of digital data [3]. Data Encryption Standard is adopted in 1977 by the National Bureau of Standards, which is currently the National Institute of Standards & Technology (NIST) as federal information processing standard 46 (FIPS PUB 46) [4, 7]. DES is the block cipher which takes a fixed-length string of plaintext bits, processes it through a sequence of complex operations, and transforms it into another cipher text of same length [8]. In this, the block size used is 64 bits and it also uses a key to customize the transformation so that decryption can only be performed by those who know the secret key which is used to encrypt the data [9]. It uses a key size of 56 bits. Initially, the key is consisting of 64 bits. The bit positions 8, 16, 24, 32, 40, 48, 56, 64 are discarded from the key length [10].

## 2.1 DES Algorithm

As per [11], the main points of the algorithm are as follows:

1. Dividing the plaintext into 64-bit (8 octet) blocks.
2. Perform initial permutation of blocks.
3. Division of blocks into two parts: left and right, named L and R.
4. The steps permutation and substitution are repeated 16 times. So we can say 16 rounds of operation.
5. Left and right parts are re-joined, and then inverse initial permutation is applied to get the cipher text.

## 2.2 DES Advantages

According to [11, 13], DES has the following advantages:

1. DES is been used since 1977, but even today Brute force attack is the most popular attack of DES.
2. It is an ISO and ANSI standard, so everybody can get the knowledge of it and can also implement it.
3. DES is relatively fast in its hardware implementation when compared to software implementation.

## 2.3 DES Disadvantages

According to [4, 13], DES has the following drawbacks:

1. The biggest defect of DES is its key size which is 56 bit.
2. DES software implementation is slow.
3. Since DES is a symmetric encryption, it uses the same secret key for encryption and decryption. If somebody loses that key, then the receiver cannot get the understandable data [11].
4. In traditional DES, operations are always done on only right-hand side of data, which make the work of cryptanalyst very easy.

**Fig. 1** Proposed system (encryption)



## 3  Proposed System

The proposed system provides two-level security. Firstly, we apply the random substitution technique on the given input; then, that cipher text is given as input to the second-level security in which we enhance DES by modifying the round function (Small function), increasing the number of rounds to 32 and also keeping the key size to 64 bit as shown in Fig. 1. We can evaluate its performance and also compare it with traditional DES performance with respective to Avalanche effect and plaintext sensitivity. In cryptography, Avalanche effect is one of the desirable properties. According to it, if the input is changed slightly, the output changes drastically [14].

In the case of extreme quality block ciphers, if we make a small change in either key or the plaintext, it should cause a significant change in the cipher text [14].

### 3.1  Random Substitution Method

1. In this method, we will consider a key string, input string which contains all printable characters and is assumed to be available at both ends. For example,

   **KeyString** = {'ABCDEFGHIJKLMNOPQRS
   TUVWXYZabcdefghijklmnopqrstuvwxyz{#@:$^}*) = !| </> ~ &%)` +\.,'`?
   {0123456789[]-_'}
   **InputString** = {'ABCDEFGHIJKLMNO
   PQRSTUVWXYZabcdefghijklmnopqrstuvwxyz{#@:$^}*) = !| </> ~ &%)
   ` +\.,'`?{0123456789[]-_'}

Plaintext: Welcome

**Input String Index** : 23    31    38    29    41    39    31

Random Cipher Text :  {nc| > Q_

**Key index** : 53    40    29    64    67    17    95

2. To obtain cipher text, we substitute every character of plaintext with a random character from the key string. Key index is the index of symbol taken from key string., i.e., 53 for {.
3. So 'W' whose index is 23 in the input string is substituted with '{' from key string.
4. As we are using random substitution, we will get a different cipher text for each execution.
5. The actual key for this substitution technique is the input string index.
6. For decryption, we will use the input string index which will be the index of the plaintext character in input string, using which cipher text character will be replaced with symbols in input string at decryption. So sending this input string index to the receiver also plays a vital role for providing security.
7. Key Transmission:

   (i) Consider a cover image. Select the pixel values of that image using key string index {53, 40, 29, 64, 67, 17, 95}, i.e., {(5, 3), (4, 0), (2, 9), (6, 4), (6, 7), (1, 7), (9, 5)} and place input string index values at these pixel values, i.e., 23 is placed at (5, 3) pixel value of the cover image.
   (ii) At the receiver, one can get the key string index by comparing cipher text and the key string. Once we get the key string index, we can get the input string index from the same pixel values using which we can get the plaintext of random substitution method.

## 3.2   Enhanced DES Encryption (EDES)

The output from random substitution method will be taken as input or plaintext to the EDES. The overall structure of EDES will remain same as DES. But its round function is modified, and numbers of rounds have been increased to 32.

**Fig. 2** DES round function

**EDES Round Function**

1. Plaintext is subjected to an initial permutation round whose output is divided into two halves say $L_{i-1}$ and $R_{i-1}$.
2. Then, parity check is made on both the halves; if the number of ones in the plaintext is even, then we will make a circular rotate left otherwise circular rotate right on both the halves.
3. Then, 64-bit session key is divided into two halves say $C_{i-1}$ and $D_{i-1}$. Then, both the halves were circularly rotated left.
4. Then, the left-hand side output from step 2 is made XOR-ed with right-hand side of output from step 3 (key), and the right-hand side output from step 2 is made XOR-ed with left-hand side of output from step 3. So key is used at two halves of data where as in traditional DES key is used at only one side of data.
5. Then, the output from the step 5 is given as input to the original round function of DES.

The enhanced round function is shown in Fig. 3. And traditional DES round function is shown in Fig. 2.

**Fig. 3** EDES round function



**EDES Advantages**

1. The key size is maintained as 64 bit in round function.
2. In this algorithm, key is used on both the halves of plaintext in each round where as in traditional DES key is used on only one half which will improve its performance.
3. We are increasing the number of rounds to 32 so Brute force attack on the cipher text will be much more complex.

Finally, we can check the quality of the block cipher with respective to avalanche effect and plaintext sensitivity by observing the change in cipher text when a small change is made in plaintext or key bit sequence.

## 4 Experimental Results

When the button plaintext is clicked, the input is given to random substitution method. It produces random cipher which is converted into binary bits as shown in random output. Then, this random output and key are given as input to EDES which produces second level of cipher in binary as shown in Fig. 4.

**Fig. 4** Encryption



**Fig. 5** Decryption

When the same key is given as input to the cipher of EDES algorithm, encryption process is done in reverse and then substitution method is applied in reverse to get the plaintext successfully as shown in Fig. 5.

## 4.1 Performance Evaluation

We can check the performance of proposed technique with avalanche effect and plaintext sensitivity. If we change one bit in plaintext, the proposed system

**Fig. 6** Avalanche effect for proposed technique and DES



**Fig. 7** Performance evaluation of proposed technique and DES

**Table 1** Performance metrics for proposed technique

| S. No. | Plaintext | Proposed technique | | Traditional DES | |
|---|---|---|---|---|---|
| | | Avalanche effect | Plaintext sensitivity | Avalanche effect | Plaintext sensitivity |
| 1. | IEEE XPLORE | 73 | 57.0313 | 36 | 28.1250 |
| 2. | Hello! how are you? | 95 | 49.4792 | 29 | 15.1042 |
| 3. | What is your name? | 105 | 54.6875 | 32 | 16.6667 |
| 4. | My college name is a.c.e.t. | 128 | 48.4375 | 35 | 13.6719 |
| 5. | Life can only be understood backwards | 161 | 50.3125 | 27 | 8.4375 |

produces 99 bits change as shown in Fig. 6. The change in cipher text for traditional DES is only 34 bits. We can also measure algorithm's performance by plaintext sensitivity. Table 1 shows the effect of avalanche effect and plaintext sensitivity for the proposed technique and the traditional DES. As the number of bits changed due to avalanche effect is more for proposed technique than traditional DES, we can say that the proposed technique performs better than traditional DES. Figure 7 shows the performance comparison of proposed technique and DES in bar chart.

## 5    Conclusion

This paper provides a two-level security by collaborating a random substitution technique and enhanced DES. At the first level, a random substitution technique is applied on the secret message to get a random cipher. In the second level of security, the random cipher is given as input to the enhanced DES algorithm. In the enhanced DES algorithm, the number of rounds is increased to 32, 64-bit key size is considered, and its round function is modified. This will improve the traditional DES to a large extent as the key is operated with both the halves of the input data. It provides better security than traditional DES. So it may resistant to different types of attacks like linear attacks, statistical attacks, Brute force attack, and Meet-In-Middle attack. It provides high security, especially for the applications such as financial banking services as well as in defense services. Its performance is evaluated by using the Avalanche effect and plaintext sensitivity.

## References

 1. https://en.wikipedia.org/wiki/Network_security
 2. https://simple.wikipedia.org/wiki/Cryptography
 3. http://searchsoftwarequality.techtarget.com/definition/cryptography
 4. Ayushi : A symmetric key cryptographic algorithm. Int. J. Comput. Appl. 1(15), (0975–8887) (2010)
 5. Rakeshkumar, SK.: Performance analysis of data encryption standard algorithm & proposed data encryption standard algorithm. Int. J. Eng. Res. Dev. 7(10), (July 2013)
 6. https://en.wikipedia.org/wiki/Symmetric-key_algorithm
 7. National Bureau of Standards—Data Encryption Standard. Fips Publication 46, (1977)
 8. https://en.wikipedia.org/wiki/Data_Encryption_Standard
 9. https://www.utdallas.edu/~edsha/OS2000/des-algorithm-details.txt
10. Ren, W.: A hybrid encryption algorithm based on DES and RSA in bluetooth communication.: In: Second International Conference on Modeling Simulation and Visualization Methods (WMSVM) (2010)
11. http://engineeringfourum.blogspot.in/2014/04/data-encryption-standard-des.html
12. Stallings, W.: Cryptography and Network Security: Principles and Practices, 5th edn. Prentice Hall, New York (1999)
13. Kapoor, P., Mohan, P., Kumar, M.: DES (Data Encryption Standard). http://priyaprakharm rigank.blogspot.in
14. http://en.wikipedia.org/wiki/Avalanche_effect

# Voter Authentication Using Modified Elliptic Curve Cryptography

**K. Sujatha, A. Arjuna Rao, Prathyusha Yejarla and K. J. Sruthi**

> Please note that the LNCS Editorial assumes that all authors have used the western naming convention, with given names preceding surnames. This determines the structure of the names in the running heads and the author index.

**Abstract** Voter authentication in general election is a very expensive and time-consuming issue. At present, identifying the voter is done through their voter ID cards in the presence of local people and this process leads to malpractice and rigging. Some security algorithms are available in similar areas where identification is required. However, they are found to be failed in many situations. When compared to the existing system, proposed system will give better results by reducing security issues. This introduces concept of voter authentication using modified elliptic curve cryptography where the voter can be authentication over Internet by using public communication channel. Then he can cast his vote through Internet. Elliptic curve cryptography is a public key cryptographic algorithm where it uses couple of keys as private and public. Private key is used by the voter and public key is used to authenticate the voter. The private key in ECC is chosen by using cuckoo search optimization technique instead of randomly choosing the values. The proposed methodology is enhanced using live sample database and is found to be secured by all means.

**Keywords** Voter authentication · General election · Security algorithms
Modified elliptic curve cryptography (MECC) · Public key cryptographic
algorithms · Private key · Cuckoo search algorithm

K. Sujatha (✉)
Nadimpalli Satyanarayana Raju Institute of Technology, Visakhapatnam, India
e-mail: sujathakota29@gmail.com

A. Arjuna Rao · P. Yejarla · K. J. Sruthi
Miracle Educational Society Group of Institutions, Vizianagram, India

497

# 1   Introduction

In voting system, there are many different effective changes from beginning to right now like traditional ballot voting to electronic voting and now to the Internet voting. In manual system, voter stamp on the symbol of ones own choice but in voting voter choice should be reliable and verifiable. To resist those attacks, there requires a security algorithm in order to authenticate voter in online. Figure 1 represents process of traditional voting system where the votes are cast by means of ballot paper. In this procedure, first voter is identified through voter ID card and a ballot paper is allocated [1]. Then voter moves to the polling officer to mark the left forefinger with the indelible ink. Voter enters into the screened compartment to select the candidate of their choice by marking the symbol on the ballot paper with rubber stamp. In polling station, there will be a common ballot box inside a screened compartment where voter inserts the ballot paper by folding it.

Problems identified in Identifying Voter in general election.

(a) Identifying the voter is becoming very expensive in elections.
(b) Misuse of official machinery.
(c) Use of caste and religion in election data collection.
(d) Mis-identification of voter in cases of rigging of election and booth capturing.
(e) Misuse of mass media.
(f) Low polling percentage due to problems arising due to mis-identification.
(g) Delay in the disposal of election petitions [1].

Voter authentication is the way of identifying the voter. In general, the voter claims will be represented by specific cards called as voter ID cards, but the actual way of representing a voter is done by a secured mechanism developed in the



**Fig. 1** Architecture of existing voting system

system and then permits the voter to vote through this system. In this process of granting rights to the voter, an identity proof must be submitted to the system in order to verify.

In the process of voter authentication, there are three components involved:

1. Supplicant: In authentication process, this party helps in providing its identity and evidence, which refers this party as a authenticated user through voter identity.
2. Authenticator: The resources required to the client or supplicant are provided by this party which is referred as server. Those requirements are used to identify the authorized user.
3. Security database: The credentials which are used by the user will be checked in this database and provides authentication for a user.

## 2 Modified Elliptic Curve Cryptography

### 2.1 Cuckoo Search Algorithm

In 2009, Yang and Deb developed an algorithm as cuckoo search (CS) which represents new metaheuristic algorithm which is used to solve optimization problems [2, 3]. The thing that is unique is the obligatory brood parasitism behavior of cuckoo birds including the Levy flight behavior of some birds and fruit flies [2]. Every egg in a nest is represented by key solution and each cuckoo egg will give a new solution. The main aim of this algorithm is to exploit a quality or superior solution of cuckoo to substitute a solution with the best solution to rank the nests. In complicated cases like multiple eggs in a nest, a set of solution may be incorporated in this algorithm. This will be the simplest approach as one egg in one nest has been applied.

**Concept of Levy Flight**

Levy flight was introduced by Manelbrot for specific definition of the distribution of step sizes. It can be thought as a random walk where in step dimension has a Levy tailed for probability distribution. To refer discrete space rather than continuous space, Levy flight can be eventually used [2, 3].

**Cuckoo Search Using Levy Flights Process**

The algorithm for cuckoo search is as given below. Concept of Levy flights is inbuilt.

1. Select objective function f(x) where

$$x = (x1, x2, . . . . . . . . . xd)r \qquad (1)$$

2. Initialize with random population n of host nests.
3. While stopping criteria reached
4. By using Levy flights get a cuckoo randomly
5. Quality/fitness should be evaluated as Fi
6. Nest should be chosen randomly among n (say, j)
7. if (Fi > F j),
   New solution replaces j;
   End
8. New nests are built and using function (pa) worse nests are abandoned;
9. Best solutions are traced
10. Find the best to rank the solutions
    [end of step 3 while]
11. Publish postprocess results and visualization
12. End

**Flowchart**

Procedure of cuckoo search is drawn though flowchart shown in Fig. 2. To generate a random number which is used to authenticate, the voter cuckoo search algorithm is used.

## 2.2 Modified Elliptic Curve Cryptography (MECC)

Miller and Koblitz [4] discovered elliptic curve cryptography as an alternative technique to implement public key cryptography. For an elliptic curve, the equation is given as y2 = x3 + ax + b. Here a random value is selected as private key [5, 6]. However, if the parameters that are chosen randomly are not selected properly, this leads in wrong calculations and the cipher text generated will not be decoded to plain text correctly. Hence, a suitable optimization algorithm is required that should properly select the random values. Cuckoo search is used for this purpose which is termed as modified elliptic curve cryptography (MECC). Authentication procedure first decides domain parameters and then computation basing on MECC and Diffie Hellman key exchange protocol [7, 8].

**Fig. 2** Cuckoo search flowchart

## Modified Elliptic Curve Cryptography for Voter Authentication

In this process, MECC generates two keys public and private keys. Voter signs using private key with help of cuckoo search algorithm to select the random number for the private key. Then voter must be authenticated by the admin using the public key generated by MECC. Figure 3 shows how authentication of voter is possible.

**Fig. 3** Voter authentication procedure



Authentication using MECC procedure is shown below

1. Admin takes a point *P* and generates random number Ka using cuckoo search.
2. Then admin computes point by using random number and admin key and sends to voter based on below Eq. (2)

$$Q = \mathrm{Ka}\, P \tag{2}$$

3. Voter generates random number Kv and computes point *R* and sends to admin based on Eq. (3)

$$R = \mathrm{Kv}\, P \tag{3}$$

4. Final point Pa is computed by the admin based on Eq. (4)

$$\mathrm{Pa} = \mathrm{Ka}\, R \tag{4}$$

5. Final point Pv is computed by the voter based on Eq. (5)

$$\mathrm{Pv} = \mathrm{Kv}\, Q \tag{5}$$

6. The concept of shared secret key will be implemented as shown in Eq. (6)

$$\mathrm{Pa} = \mathrm{Ka}\, R = \mathrm{Ka}\, \mathrm{Kv}\, P = \mathrm{Kv}\, \mathrm{Ka}\, P = \mathrm{Kv}\, Q = \mathrm{Pv} \tag{6}$$

Authentication is handled by MECC as voter signs by using voter's private key, and at server end, this is verified using public key of voter.

## 3 Results

In this paper, "modified elliptic curve cryptography" (MECC) gave optimized results when compared to other algorithms such as RSA and ECC. The following graph shows better results for authenticating a person in less time using MECC.

Table 1 shows how MECC is better when compared to other authentication algorithms. MECC takes less time for processing and results are obtained in a short span. Finally, the graph which is in Fig. 4 shows that MECC is better for implementing voter authentication.

**Table 1** Comparison table on RSA, ECC and MECC

| S. No. of trails | Authentication process in ms | | |
|---|---|---|---|
| | RSA | ECC | MECC |
| 1 | 1000 | 300 | 180 |
| 2 | 800 | 250 | 220 |
| 3 | 1100 | 330 | 200 |
| 4 | 1900 | 390 | 240 |
| 5 | 2400 | 450 | 320 |
| 6 | 400 | 100 | 50 |
| 7 | 2200 | 420 | 310 |

**Fig. 4** Graph showing performance of MECC comparable to RSA, ECC

## 4 Conclusion

The developed secured voter authentication system has been tested and implemented on different applications. Security issues were implemented by enforcing both integrity measure and authentication measure of the authenticated voting system. All the issues in manual system are solved by using this automated voter authentication system. This can be used in any system that requires to identify the user securely and with reliability. The application developed stands as a barrier to next applications and those applications are used whenever authentication is required.

## References

1. Sujatha, K., Nageswara Rao, P.V., Arjuna Rao, A., Rajesh, L.V., Vivek Raja, V.: Secured internet voting system based on combined DSA and multiple DES algorithms. In: ICT and Critical Infrastructure: Proceedings of the 48th Annual Convention of Computer Society of India-Vol II Advances in Intelligent Systems and Computing, vol. 249, pp. 643–650 (2014)
2. Barthelemy, P., Bertolotti, J., Wiersma, D.S.: A L'evy flight for light. Nature **453**, 495–498 (2008)
3. Yang, X.-S.: Cuckoo Search via Lévy Flights. Department of Engineering, University of Cambridge, Trumpinton Street, Cambridge CB2 1PZ, UK
4. Amara, M., Siad, A.: Elliptic curve cryptography and its applications. In: 2011 7th International Workshop on Systems, Signal Processing and their Applications (WOSSPA), pp. 247–250 (2011)
5. Qing, Z., Zhihua, H.: The large prime numbers generation of RSA algorithm based on genetic algorithm. In: 2011 International Conference on Intelligence Science and Information Engineering (ISIE), pp. 434–437 (2011)
6. Wang, M., Dai, G., Hu, H., Pen, L.: Selection of security elliptic curve based on evolution algorithm. In: CINC '09. International Conference on Computational Intelligence and Natural Computing, vol. 1, pp. 55–57 (2009)
7. Koblitz, N.: Elliptic curve cryptosystems. Math. Comput. **48**, 203–209 (1987). American Mathematical Society, ISSN: 1088-6842 (online), ISSN: 0025-5718 (print)
8. Miller, V.S.: Use of elliptic curves in cryptography. In: Williams, H.C. (ed.) Advances in cryptology-CRYPT0 '85, LNCS 218, pp. 417–426 (1986). IEEE International Conference, vol. 1, pp. 681–685. IEEE (2011)

# Sensor-based Dam Gate Control System and Alert Using Zigbee

S. Shivani, R. Dharani, J. Annush J. Sharon and S. Mrudhula

**Abstract** Floods are natural phenomenon that brings havoc to human life. It is necessary to bring dam controlling system to overcome the threats posed by floods. This system focuses on using Zigbee for providing availability and the provisioning of an alarm system for intimation of the nearby area in case of any emergency. It has an added advantage of being transparent to the public and also controls the speed of the motor for outflow of water. The system comprises of the dam controlled by a control room and uploading the data to the website such that it can be accessed by the public.

**Keywords** Zigbee · ARM · Raspberry pi · GPU · UART

## 1 Introduction

There is an increasing demand for mechanisms that provide safety for the public. Embedded system is the combination of hardware and software co-design. Embedded systems are now-a-days playing a vital role in Engineering design process for efficient analysis and effective operation. Our work describes the design of an embedded system for the "dam controlling system". At the dam, this system gives the concept of interfacing high-voltage electrical devices such as DC motor, level sensors and flow sensors to Zigbee. This system facilitates us to control

S. Shivani · R. Dharani (✉) · J. A. J. Sharon · S. Mrudhula
Information Technology, Panimalar Institute of Technology, Poonamallee,
Chennai 600123, India
e-mail: dharani.rit@gmail.com

S. Shivani
e-mail: devisukushiv@gmail.com

J. A. J. Sharon
e-mail: annushsharon@gmail.com

S. Mrudhula
e-mail: mrudhulasatheeshnair@gmail.com

the gates of the dam depending on the water level. It consists of a set of sensors connected to a DC motor through an ARM-7 processor. The water level is detected and transmitted to the control room where the motor speed can be controlled by the authorities. The values are updated on the website periodically. The public can view the status of the dam by accessing the website.

## 1.1  Need

With this mechanism, the water wastage can be reduced thus leading to efficient utilization of it. Another need is to keep record of water status of the dam even during emergency situations. The official government website www.chennaimetrowater.tn had failed to provide the dam status from 1 December 2015 to 9 December 2015.

This system provides proper intimation to the authorities even in the absence of mobile network and hence being available during disasters.

Another issue prevailing today in India is the improper intimation to the public, and this system overcomes it by facilitating an alarm system.

## 1.2  Intent

The main objectives of the system are

- To consider more than one parameter for judging the dam level such as water level, flow rate and landscape.
- Providing an alarm system for alerting the nearby locals in case of emergency.
- Ensuring availability by no depending on mobile networks.
- Facilitating transparency to the public on the dam level.

## 2  Literature Survey

The system is mainly concerned with monitoring and controlling water with increased availability and transparency. In 1986, Davidson [1] proposed a control system for efficient working of hydroelectric power plant in North of Scotland with the help of visual display units (VDUs). This proposed VDU system guided us in making a GUI control panel for operator. Another proposal was by Litric [2] that is related to water management in dam using SIMO systems. This SIMO system deals with the real-time calculation of the upstream and downstream flow of the water in dams. Mohani et al. [3] proposed a PC-based dam control model where they introduced the concept of division of reservoirs.

Other supporting literature surveys are "An intelligent SMS-based remote water metering system" [4] proposed by Islam and Wasi-ur-Rahman that deals with a technique having adequate security support, for prepaid billing of water using short message service (SMS). "GSM-based wireless home appliances monitoring & control system" [5] proposed by Al-Ali et al. that monitor and control home appliances locally.

Our system mainly comprises of using graphical processing unit (GPU) for display of the statistics of the dam. It uses Zigbee in order to transmit messages even when there is no mobile network available. An LCD screen is also available to view the readings [6].

## 3 Proposed System

In this section, we discuss the modules in the proposed system, its block diagram and the related functions.

### 3.1 Modules

There are three major modules in the proposed system,

- Dam controlling system
- Control room
- Web server

The dam controlling system is used to monitor the water level of the dam using the sensors. The control room is responsible for the authorities to take action depending on the received water status. The web server is used to host the webpage and update the readings periodically.

### 3.2 Overall Architecture

Figure 1 depicts the overall architectural block diagram for the proposed system. The dam controlling system uses three different types of sensors to monitor the dam. The different sensors are used to detect the level of water, flow rate and landscape measurement. The statistics are then transmitted by Zigbee. The control room receives the statistics by using the Zigbee receiver, and the authority can take necessary action by controlling the motor speed. The authorities can also invoke the alarm system. By the alarm, the general public is aware of the emergency situation.

**Fig. 1** The overall architectural diagram of the proposed system

## 3.3 Working Hypothesis

Figure 2 represents the block diagram for the dam controlling system. It comprises of three sensors the level sensor [7], flow sensor and the landscape sensor. The readings from the various sensors are given to the ARM processor, which then determines the threshold depending upon the nature of the dam. These values are then transmitted to the control room via the Zigbee transmitter. This system also comprises of an alarm system that can be invoked from the control room.

Figure 3 represents the block diagram for the control room. It comprises of Zigbee that receives readings from the dam controlling and also comprises of a controller to graduate the motor speed. These readings from the dam controlling system are given to the raspberry pi processor by which the administrating authority

**Fig. 2** The block diagram of the dam controlling system



**Fig. 3** The block diagram of the control room system

determines the action to be taken. The action may be to invoke an alarm system or to control the motor's speed. These values are then transmitted to the web server from which it is hosted to the cloud. Additionally, this system also comprises of a LCD screen to display the values obtained from the dam.

Figure 4 represents the block diagram for the web server. It comprises of Wi-fi and a personal computer. The public can access the dam status by using their personal computer, mobile devices, etc. via the website.

**Fig. 4** The block diagram of the web server system

## 4  Related Systems

The related systems involved in the proposed system are described in detail as follows.

### 4.1  Zigbee

It is an IEEE 802.15.4 [8] specification for a set of high-level protocols used to create a personal area network (PAN) by using low-power digital radios. The advantages of employing Zigbee are:

- It is less expensive than bluetooth and Wi-fi.
- The transmission distance is in 10–100 m, but long data transmission can be facilitated by employing mesh networks.
- Its network is secured by 128-bit symmetric encryption key.

### 4.2  ARM Processor

It stands for Acorn RISC Machine or Advanced RISC Machine. It is a RISC architecture for computer processors that is configured for various environments.

### 4.3  Raspberry Pi

It is a series of small single-board computers [9]. The main reason for employing raspberry pi as it provides in-built Wi-fi and bluetooth. An additional advantage it also has a random access memory (RAM).

### 4.4  GPU

It stands for graphic processing unit. It is an electronic circuit to manipulate memory for producing output for display.

## *4.5 UART*

It stands for universal asynchronous receiver/transmitter. It is a computer hardware device for asynchronous serial communication in which data format and the speeds can be variable.

## 5 Hardware Specifications

- ARM 7 controller
- Analogue to Digital Converter (ADC)
- Universal Asynchronous Receiver/Transmitter (UART)
- Zigbee
- Alarm System
- Sensors (Level, Flow, Landscape (MEMS))
- Raspberry Pi Processor
- Universal Asynchronous Receiver/Transmitter (UART)
- Zigbee
- LCD
- PC/Desktop
- Web Server

## 6 Conclusion

The proposed system reduces water wastage from dams and ensures the efficient use of the water resources. It also reduces the human resources required to operate the dam. This system mainly focuses on providing availability of dam status to the public and transparency. Also it is possible to save lives in case of heavy rains and floods by gradual outflow of water. The operational time is also reduced. This project is an initiative to merge technology for safety to provide a better living environment.

## 7 Future Work

There is scope to modify our proposed system by better refinement. Since it is very difficult to design a single system for all the dams in India, it is possible to use artificial intelligence [10] which modifies itself according to its surroundings by using feedback from the environment.

On a higher scale, it is possible to have a single integrated control room for all the 5200 (approx.) [11] dams in India using wireless technology that is under control by the central government.

# References

1. Davidson, E.G.: Control of hydro-electric plant. In: IEE Proceedings C (Generation, Transmission and Distribution), vol. 133, no. 3, pp. 145–147
2. Litric, X.: Robust IMC flow control of SIMO DamRiver open-channel systems. IEEE Trans. Control Syst. Technol. **10**(3), 432–437 (2002)
3. Mohani, S.S., Talha, S.M.U., Ahmed, S.H., Ebrahim, M.: Design for irrigation and monitoring system of an automated dam. In: Proceedings of the MultiConference of Engineers and Computer Scientists, vol. II, IMECS 2012, Hong Kong, 14–16 Mar 2012
4. Islam, N.S., Wasi-ur-Rahman, M.: An intelligent SMS-based remote water metering system. In: 12th International Conference on Computers and Information Technology, Dhaka, Bangladesh, 21–23 Dec 2009
5. Al-Ali, A.R., Rousan, M.A., Mohandes, M.: GSM-based wireless home appliances monitoring & control system. In: Proceedings of International Conference on Information and Communication Technologies: From Theory to Applications, pp. 237–238 (2004)
6. Iyer, M., Pai, S., Badri, S., Kharche, S.: Embedded dam gate control system using 'C' and Visual Basic. Int. J. Comput. Appl. **69**(2), 32–37 (2013)
7. Divya, A., Gajalakshmi, M., Kanchana, V., Dharani, R.: Sensor based dam gate controlling with high level protection. In: International Conference on Science and Innovative Engineering, April 2016
8. Kaushal, K., Kaur, T., Kaur, J.: ZigBee based wireless sensor networks. Int. J. Comput. Sci. Inf. Technol. **5**(6), 7752–7755 (2014)
9. Mazidi, M.A., Mazidi, J.G. : Microcontrollers, pp. 284–290, 300–310, 428–430, 441–450. Prentice Hall Publications, Upper Saddle River
10. Deshpande, A., Pitale, P., Sanap, S.: Industrial automation using Internet of Things (IOT). Int. J. Adv. Res. Comput. Eng. Technol. (IJARCET) **5**(2), 266–269 (2016)
11. http://www.cwc.nic.in/main/downloads/new%20nrld.pdf

# Prolonged Network Lifetime to Reduce Energy Consumption Using Cluster-Based Wireless Sensor Network

**Sagargouda S. Patil and Anand Gudnavar**

**Abstract** Efficient utilization of energy for each sensor node is the key to prolong the network lifetime. The balance of energy consumption among the nodes located at central and edge areas also plays another important role in the situation of lifetime extension. In earlier work, LEACH algorithm has been investigated in the existing cluster head selection and the edge sub-clustering scheme is further developed to ease the effect of non-uniform node distribution on the edge, but it has not covered entire coverage area of the network due to failure of some nodes. The proposed work uses partition around medoid algorithm to improve optimization of number and size of sub-clusters along with the communication range for better improvement in cluster head selection and overcoming failure of edge nodes. Thus, it not only improves the energy conservation but also drastically improves the balance energy consumption, which can contribute to longer network lifetime and outage probability.

**Keywords** Clustering · Energy consumption · Node connectivity
Network lifetime · Sensor node

## 1 Introduction

Wireless sensor system comprises of sensor hubs which are set above a large geographical area and have wireless communication. The sensor node is a node which is small, low cost and lows down energy communication [1]. Each sensor node has 3 components, namely (i) sensing data from environment (ii) processing or analysing sensed data (iii) storing data until it sends data to base station. Each

S. S. Patil · A. Gudnavar (✉)
Department of Computer Science and Engineering, Shaikh College
of Engineering and Technology, Belagavi, India
e-mail: anangud@sgibgm.org

S. S. Patil
e-mail: sagapat@sgibgm.org

sensor node is battery-oriented; the capacity of battery is very limited to use, and it is also impossible to replace battery because of harsh environment in wireless sensor network. So, we aim at improving and maximizing the network lifetime. It is very important to inspect and approximate how protracted network is working properly or lifetime of a system. Remote sensor system comprises of sensor hubs and base station. Each sensor node is distributed in a larger geographical environment [2]. It is also used in various fields like war, health care, polyhouse monitoring systems, transport, etc. The key dispute of WSN is to build the lifetime of the system by minimizing the energy consumption [3]. Energy consumption is the amount of consumption of energy or power while transmitting data from one location to another location. Studying from previous years, a variety of changes have been made to limit the energy requirement in WSN, as mainly energy consumption is more for wireless transmission and response. Main approaches proposed so far focused on making the changes at MAC layer and network layer to minimize the energy consumption. Two more major challenges are how to place the cluster heads over the network and how many clusters would be there in a network. If the cluster heads are properly placed over the network and sufficient clusters are formed, it will help to minimize the utilization of vitality and would expand the lifetime of the network. To deal with all the above-mentioned challenges, clustering is always been referred as an effective technique to improve the lifetime of remote sensor system. Clustering is a method which consists of group of sensor nodes and cluster head, where some sensor nodes are selected as cluster head at the time of cluster creation. Sensor node in each cluster transmits their information to separate bunch head, and group head sends such total information to base station (Fig. 1).

Few key terms related to clustering are discussed below:



**Fig. 1** Clustering in sensor network

1. **Sensor node**: A sensor centre point is the principal part of a WSN. Sensor centres perform capacities, for example, detecting, information storing, routing and information preparing.
2. **Clusters**: Clusters are the assorted levelled units for WSN. Significant sensor systems should be separated into bunches to abbreviate assignments, for example, correspondence between the base station and the group head [4].
3. **Cluster head**: Cluster head is the expert of a bunch. Bunch head is habitually required to sort out exercises in the group. These assignments incorporate information accumulation, pressure and correspondence to base station [5].
4. **Base station**: The base station (BS) gives the correspondence join between the sensor framework and the end-customer. It is basically the sink in a WSN.

## 2  Problem Statement

To minimize the vitality utilization of sensor hubs and expanding the system lifetime are the two vital considerations in the WSN. Since sensor nodes are outfitted with battery-powered devices, it is impracticable to change the battery for all nodes if the network is implemented for a huge area with numerous nodes. In the existing method, inefficient selection of size of cluster and cluster head is used and edge connectivity to cover all nodes which affect energy consumption as well as network lifetime to improve network in wireless sensor network.

## 3  Literature Survey

**(1) Low-Energy Adaptive Clustering Hierarchy (LEACH):**

LEACH is one of the primary progressive directing conventions for sensor systems. It is a self-arranging, versatile bunching convention. It decreases the vitality radically. The LEACH is a group-based convention and arbitrarily chooses the cluster-sets out towards every bunch. It is discarded in proposing works due to the following drawbacks.

1. Cluster heads are not consistently dispersed inside group that implies a cluster head, situated at the edges of a group.
2. Many applications need large coverage area, so it does not work well with large area network in WSN.

**(2) Weighted residual energy and distance (WRED) and Edge sub clustering (ESC):**

Clustering is only based on geography effect of poison point process, this leads to consume more energy because PPP is suitable for only finding location and distance of nodes but not the communication range. More complex scheme for

cluster head and cluster selection is that the weighted lingering vitality and separation (WRED) calculation is wanted to draw out and determine the bunch head. Utilizing WRED, hubs with more energy are placed near cluster centre [6]. The edge sub-grouping (ESC) plan is created to advance and enhance the impact of non-uniform circulation of hubs on the edge.

We have two existing algorithms for reducing energy consumption: WRED and ESC, although both have disadvantages like (1) No optimal number of clusters formed leads to increase in energy consumption and leads to decrease in network lifetime. (2) Increased number of hops.

**(3) EECS (energy efficient clustering scheme in wireless sensor networks)**

The hubs with more remaining vitality have more likelihood to be chosen as group heads. On the off chance that hub does not discover with more lingering vitality, it turns into a bunch hub. In group development stage, LEACH utilizes the base separation of hubs to equal their bunch head. In this manner, EECS determines the issue that groups at a more prominent separation from the sink and requires more imperativeness for transmission than those that are ever closer to low message overheads and to the uniform assignment of gathering head stood out from LEACH.

*Advantage*—The calculation builds multilevel groups, and the hubs in every bunch achieve the bunch head by hand-off through different hubs.

**(4) K-means**

Creators have depicted that k-mean grouping scientific system is emphatically indistinguishable to the system bunching issue; some energizing k-implies bunching changes were recently connected to the specially appointed systems. Up till now what makes the view we ponder not the same as others is the truth that the measure of going before information is little than it is in average *k*-implies applications.

Hubs are just ready to quantify separations to their one jump neighbour's positions which are obscure and the system design does not offer the unified learning required for fundamental k-implies calculation applications.

**(5) Survey of various clustering methods**

Ameer Ahmed Abbasi provides a review of a variety of clustering algorithms that are exclusively planned for WSNs. Here, the discussion is about a choice of convergence time algorithms. Convergence time is the time required before all the cluster heads attain conformity about the topology of the WSN. Classified clustering algorithms are in two categories—variable convergence time algorithms and constant [7] convergence time algorithms. Variable convergence time algorithms are helpful when number of nodes in WSN is low, while constant convergence time algorithms are useful when number of nodes in WSNs is high.

# 4  System Architecture

## 4.1  Overview of System Architecture

- **Network initialization**: is an input phase with which the network parameters namely x-dim, y-dim, initial energy of node, number of nodes, MAC type, antenna model, transmission speed, area of simulation are initialized (Fig. 2).
- **Communication establishment**: Once Initializations are completed, the next step is to compute the average, individual energy and location of each and every node along with the distance of all nodes and hence it establishes communication in the network.
- **Cluster formation**: It is done by using input parameter such as energy and distance gathering data from slave nodes to master nodes and transmitting data from master node to base station with minimum distance to reduce energy consumption in wireless sensor network.
- **Result analysis: from our "novel approach" that is PAM algorithm which results in**

  1. **Energy consumption**
     The power utilization level of a hub can be controlled by finding the contrast between the present energy $E_C$ and during imitation energy $E_i$. On the off chance that an energy level of a hub achieves zero, it can't get or transmit any longer parcels. The measure of vitality utilization in a hub can be imprinted in the following record. The vitality level of a system can be controlled by summing the whole hub's energy level in the system [8, 9]. The average power consumption of the application traffic $n$, which is denoted by PC, is obtained as

$$PC = \frac{1}{n} \sum_{i=1}^{n} E_i - E_C.$$



**Fig. 2** System design architecture

2. **Probability of node coverage**
   This mainly tells about how much area is covered by sensors. It is very important to measure in terms of connectivity. Connectivity means that each data collected in the sensor nodes is to stretch connection within the cluster that it should cover the range to all respected nodes in the cluster [10].
3. **Network lifetime**
   Lifetime of wireless sensor node depends on battery life so that those batteries cannot be rechargeable or removable so that the optimal clustering solution to overcome these types of problems to reduce energy consumption in increase network lifetime. Lifetime of network is determined using formulae,

$$\text{Lifetime} = (\text{residual energy} - \text{distance covered})/100$$

## 5   Proposed Research Model

Let us discuss proposed model algorithm to overcome existing problems to efficiently reduce energy consumption and increased in lifetime of network.

**Partition around medoid** is the most common insight of *K*-medoid clustering algorithm. This algorithm is used for partitioning around medoid. Partition around medoid (PAM) algorithm is a clustering algorithm that attempts to minimize the distance between leaf nodes and points from centre [11, 12]. It also makes use of prospect maximization (EM) strategy to unite to a minimum error condition. It starts from initial set of medoids and iteratively replaces a medoid by a non-medoid as and when it improves the total distance of resulting clustering. Integration of range model-based mechanism and distance-based mechanism helps in node distribution technique. Nodes are distributed dynamically that is it is based on non-uniform network.

**Algorithm Steps:**

**Step 1**. Initialization randomly selecting K-**mediod from "n" datapoints (nodes).**
**Step 2**. Associate each data point to closest medoid or nearest medoid (cluster head).
**Step 3**. Network improvements can be done using following steps: **for each medoid "m" and for each nomedoid "o".**

(a) Swapping "m" and "o", re-computing cost.

   Cost is nothing but sum of distances and energy of datapoints to closest medoid.

(b) If any node is out of range, then we communicate edge node and add to respected cluster based on Euclidean distance.

(c) If the total cost of above process is greater than original value, then we can undo using swap process.

**Step 4**. Repeat step until cost decreases, in each cluster make point that minimize sum of distance within cluster and finally reassign each point to respective cluster.
**Step 5**. Selecting appropriate cluster head with respect to medoids.

## 5.1 Advantages of PAM algorithm over LEACH protocol and ESC scheme

1. Enhancement of number and sizes of sub-groups.
2. Reduced energy consumption and increased network lifetime.
3. Solves problem of number of hops and communication range.
4. Better maintainability of system.
5. It finds more accurate centres of clusters.

## 6 Experimental Results

Simulations are conducted using MATLAB (R2013a) tool, and we have used true-time simulator to implement partition around medoid algorithm to get better results. We have assumed following network parameters for our simulation (Table 1):

Figure 3 shows that energy consumption with respect to number of nodes in the simulation. As the number of nodes increases, consumption of energy decreases. Suppose if the number of nodes is 40, consumption of energy is around 15 J in the existing system likewise if nodes are 100, consumption of energy in existing system is 26 J, but using our proposed algorithm we are consuming few joules of energy compared to the existing system about 13%.

Figure 4 shows that energy consumption with respect to the number of nodes is for optimal cluster formation. Here, choosing $k$-value is key thing to efficiently convey optimization scenario in wireless sensor network. With the lack of optimization in the existing system to overcome this, PAM algorithm can be used by

**Table 1** Network parameters

| | |
|---|---|
| Number of nodes | 100 |
| Area of simulation | $100 \times 100$ m |
| Initial energy of node | 100 J |
| $k$-value | Trial and error |
| Number of sink | 1 |
| Deployment model | Random |

**Fig. 3** Energy consumption in WSN



**Fig. 4** Optimal cluster formation in WSN

changing $k$-value using trial and error method until we get optimized $k$-value. In Fig. 4, we are taking up to $k = 9$ iterations with respect to energy consumption, we get reduction in energy consumption and efficient optimization at $k = 5$ that increase in optimization of number and size of sub-cluster.

Figure 5 shows that lifetime of network in percentage (%) with respect to number of nodes in the simulation. Lifetime of network increases when energy of node decreases, at the node, 100 lifetimes is about 0.78% in existing system but using our proposed system with 100 nodes lifetime of network increased about 0.9 that means it increased by 0.12% drastically in network.

Figure 6 shows that probability of node coverage in percentage with respect to number of nodes in the simulation. It has been concluded that percentage of

Fig. 5  Network lifetime in WSN



Fig. 6  Probability of node coverage in WSN

probability of node coverage in proposed system should be increased when compared with an existing system. The main logic involved in this graph is increasing connectivity of all nodes which are there in network that is covering all the nodes which are on edges and non-communication range. When we are considering 60 nodes in existing system, it covers around 0.72% of coverage, but using our coverage algorithm, we are covering around 0.2% of improvement for better communication in network.

# 7    Conclusion

Energy factor in wireless sensor network is a major challenge; most of sensor nodes are equipped with very limited power. The lifetime of wireless sensor nodes depends upon battery lifetime and these batteries are neither easily removable nor rechargeable. The paper proposes a protocol, partition around medoid, which helps in formation of efficient optimization number and size of sub-clustering that confirms minimum energy consumption. Connectivity among nodes on edges or out of range is very important in consideration for better communication establishment among all nodes, which is evident through the simulation that exhibits better results compared to existing LEACH and ESC algorithms. Hence improvement in reduction of energy consumption and increased network lifetime increases overall performance of wireless sensor networks.

# References

1. Takaishi, D., Nishiyama, H., Kato, N., Miura, R.: Towards energy efficient big data gathering in densely distributed sensor networks. IEEE Trans. Emerg. Top. Comput. **2**(3), 388–397 (2014)
2. Shah, S.H., Khan, F.K., Ali, W., Khan, J.: A new framework to integrate wireless sensor networks with cloud computing. In: IEEE Aerospace Conference, March 2013, pp. 1–6 (2013)
3. Akyildiz, I.F., Su, W., Sankarasubramaniam, Y., Cayirci, E.: A survey on sensor networks. IEEE Commun. Mag. **40**(8), 102–114 (2002)
4. Tao, Y., Zheng, Y.: The combination of the optimal number of cluster-heads and energy adaptive cluster-head selection algorithm in wireless sensor networks. In: International Conference on Wireless Communications, Networking and Mobile Computing, Wuhan, China, pp. 1–4 (2006)
5. Chan, T.J., Chan, M.U., Huang, Y.F., Lin, J.Y., Chen, T.R.: Optimal cluster number selection in ad-hoc wireless sensor networks. WTOC **7**(8), 837–846 (2008)
6. Dorsey, D.J., Kam, M.: Non-uniform deployment of nodes in clustered wireless sensor networks. In: 43rd Annual Conference on Information Sciences and Systems, (CISS), Baltimore, MD, USA, pp. 823–828 (2009)
7. Heinzelman, W.R., Chandrakasan, A., Balakrishnan, H.: An application specific protocol architecture for wireless microsensor networks. IEEE Trans. Wirel. Commun. **1**(4), 660–670 (2002)
8. Tripathi, R.K., Singh, Y.N., Verma, N.K.: Two-tiered wireless sensor networks-base station optimal positioning case study. IET Wirel. Sens. Syst. **2**(4), 351–360 (2012)
9. Mathapati, B.S., Patil, S.R., Mytri, V.D.: Energy efficient cluster based mobility prediction for wireless sensor networks. In: IEEE International Conference on Circuits, Power and Computing Technologies (ICCPCT), (2013)
10. Elbhiri, B., El Fkihi, S., Saadane, R.: A new spectral classification for robust clustering in wireless sensor networks. In: IFIP WMNC (2013)
11. Jawad Ali, S., Roy, P.: Energy saving methods in wireless sensor networks. In: IDE0814, (May 2008)
12. Saini, M., Saini, R.K.: Solution of energy-efficiency of sensor nodes in wireless sensor networks. Int. J. Adv. Res. Comput. Sci. Softw. Eng. **3**(5) (2013)

# Energy-Efficient Hybrid Protocol for Routing Based on Mobile Data Collectors in Wireless Sensor Network

B. A. Mohan and H. Sarojadevi

**Abstract** Wireless Sensor Network is increasingly becoming popular solution for wide range of real-life applications. Energy is a primary factor in sensor networks as battery is used for powering sensor nodes, which is non-replaceable. Choosing energy-efficient technique for routing can increase the network lifetime. In this paper, a new energy-efficient clustering routing protocol based on multiple mobile data collectors is proposed. The network life is based on rounds, and in each round, clustering is applied. This protocol is hybrid because it uses both centralized and distributed clustering approach for multiple rounds after which left out nodes directly forward the data to base station (BS). Clusters are grouped into sector, and a mobile node is assigned for data collection. The mobile nodes in the network simultaneously collect the data, which in turn reduces latency. Simulation reveals that the proposed protocol with multiple data collectors improves energy consumption compared to node density-based clustering and mobile collection (NDCM).

**Keywords** WSN · Mobile data collector · Energy efficiency · NDCM

## 1 Introduction

Wireless Sensor Networks (WSNs) are progressing as state-of-the-art solution to wide range of applications of monitoring, event detection, and target tracking in domains such as precision agriculture, pollution prevention, buildings health monitoring, intrusions, fire/flood emergencies, and surveillance [1, 2].

WSNs are made up of sensors characterized by small in size, lightweight, low energy consumption, and low computation capacity and deployed in inaccessible

B. A. Mohan (✉) · H. Sarojadevi
Nitte Meenakshi Institute of Technology, Bangalore, India
e-mail: mohan.ba@nmit.ac.in

H. Sarojadevi
e-mail: sarojadevi.n@nmit.ac.in

places where humans cannot reach or impossible to replace battery, so conserving energy of the nodes to extend overall lifetime of the network is very important.

Clustering is an approach for energy conservation, where nodes in a network are grouped into clusters and each cluster will have a designated node called cluster head (CH) which is responsible for collecting data from group of member nodes and removing the redundant data and forwarding it to the base station or sink, instead of all nodes sending data directly to the sink which incurs more communication cost for the energy constraints [3–5].

The conventional sensor network architecture comprises of densely deployed static nodes over a network area, leading to the problem of network isolation as the connecting node dies in the network. To address this problem WSN architectures are exploiting mobility [6–10].

The structure of the paper is as follows: Sect. 2 analyzes the working principle of node density-based clustering and mobile collection protocol. Section 3 discusses the network model used and working principle of proposed protocol. Section 4 compares the energy consumption, throughput, and network lifetime of the proposed and existing protocol like NDCM and, Sect. 5 concludes the proposed protocol.

## 2  Analysis of Node Density-based Clustering and Mobile Collection

NDCM combines cluster routing and mobile element-based data collection in WSNs. The CHs in the network first gather data from the cluster members and wait for mobile element (ME) to visit for collecting the data. A new cluster head selection method based on density of the nodes is proposed. Thus, a node at the center of densely populated nodes is more likely to be elected as a cluster head, which in turn improve the efficiency of intra-cluster communication.

This protocol solves most of the problems of traditional WSNs like energy consumption, funnelling effect, and network isolation but fails to address issues like latency, energy consumption, and throughput.

As this protocol employs single mobile node for collecting the gathered data from cluster heads, there is a latency introduced for delivering the data to the BS [11]. This problem can be solved by introducing multiple mobile nodes in the network and defining the mobility path for the mobile elements [12].

## 3  Energy-Efficient Hybrid Protocol with Mobile Data Collectors (EEHPMDC)

The working of the protocol is broadly classified as cluster formation, transmission of data from sensor nodes to cluster head, and data forwarding from cluster head to base station via a mobile node which is clearly depicted in Fig. 2 using a flowchart.

Cluster Formation—The sensor nodes are randomly deployed and expected to send their location and energy level to the base station, which is used to divide network into clusters and group clusters into sectors as shown in the Fig. 1. Each sector in the network is assigned a mobile node for collecting data. Base station selects and broadcasts 2 cluster heads (primary and secondary) to all the nodes. This protocol uses hybrid approach for cluster head selection. First, it uses centralized-CH selection method for initial two rounds and distributed-CH selection method for the rest of the rounds. Modified LEACH protocol is used for cluster head selection in which every node checks its random values against threshold value, if its random values are less than threshold, then that will be elected as cluster head. Cluster head will broadcast the advertisement message (ADV message) to neighbor nodes using MAC protocol to form the cluster. The normal nodes which receive ADV message from multiple cluster heads will decide to join one of the clusters [1, 3, 5, 13–17].

Data forwarding from CM to CH—sensed data from cluster members—are transmitted to cluster head using time division multiple access schedule, which divided the timeframe into slots/frames and assigned to member nodes for transmission of data. In a cluster, only one sensor node will be active at a time and use the channel for data transmission. Once the data are received at cluster head from all of its members, it forwards to the mobile node after removing redundant data.

Data transmission from cluster heads to sink—mobile node moves around the sector in a given pattern for collecting data from the cluster heads and forwarding it



**Fig. 1** Network model with multiple mobile nodes

to base station. The base station uses the location information of the cluster heads to derive the path for the mobile nodes. Base station uses trajectory and integral linear programming for deriving path for mobile node. The steps are summaries as follows.

   I. Base station calculates path for mobile nodes.
  II. The mobile nodes follow the given path to reach every cluster head in sector, collect data, and preserve energy of the cluster heads.
 III. The mobile node sends the data to BS for further processing.

Trajectory of mobile nodes is found in subsequent rounds as follows. Mobile nodes will periodically update the location information of cluster heads in sector to the base station for calculating new path using linear programming. The calculated path information is updated in the mobile node (Fig. 2).

## 4    EEHPMDC Versus NDCM

This section discusses simulation setup, assumptions made, and performance analysis of proposed EEHPMDC protocol. A comparison with the existing NDCM protocol with single mobile node is presented.

The proposed model is simulated using NS-2.35 as shown in Fig. 3 on Ubuntu 12.04 platform, and the parameters assumed are shown in Table 1.

Simulation is carried out for 50, 60 and 70 no. of nodes for studying energy consumption which is inversely proportional to network lifetime and the throughput.

Results in the Fig. 4 indicate the energy consumed for proposed protocol (EEHPMDC) is constant as number of nodes in the network is increased from 50 to 70, but for the NDCM, protocol energy consumed increases as nodes in the network is increased. Similarly, Fig. 5 shows the network lifetime remains constant for proposed protocol, whereas for NDCM network, lifetime decreases as the nodes are increased from 50 to 70.

Figure 6 shows the throughput of proposed protocol verses NDCM protocol, in which throughput of the proposed protocol is better and remains constant when the nodes increased from 50 to 60 and throughput increases when the nodes is increased to 70.

## 4.1    Analysis of Results

The proposed protocol is a new protocol which uses multiple mobile nodes in network, due to which data collection effort is share among the mobile nodes. Cluster heads are not participating in the transmission of data to base station. This

**Fig. 2** Forming sectors and clusters

**Fig. 3** Network deployment, sectors, and clusters formation

**Table 1** Parameters assumption for network simulation using NS-2

| Parameters | Values |
| --- | --- |
| No. of nodes | 70 |
| Simulation area | 600 m × 600 m |
| Energy model | <50 J |
| Wireless interface | WirelessPhy |
| MAC type | Mac/Sensor |
| Queue type | Droptail/Priorityqueue |
| Queue length | 50 packets |
| Antenna | Omni-Antenna |
| Propagation type | TwoRayGround |
| Routing protocol | EEHPMDC, NDCM |
| Data packet size | 1000 bytes |
| Simulation time | 300 s |

results in overall reduction in energy consumption and increase in the network lifetime. The proposed protocol is topology independent which makes overheads not depending on the number of nodes. The graphs show a trend with a minimum variation as the numbers of nodes is increased. Throughput performance increases with the increase in the number of nodes, as in Fig. 6.

Fig. 4 Energy consumption
verses no. of nodes



Fig. 5 Network lifetime
verses no. of nodes



Fig. 6 Throughput verses
number of nodes

## 5 Conclusions

In this paper, a hybrid routing protocol is proposed which employs multiple mobile nodes for data collection unlike node density mobile collector (NDCM). Deploying multiple mobile nodes is a specific feature of this proposed protocol. The simulation results establish that the usage of multiple mobile nodes can be improved in terms of throughput and latency. Network which is divided into sectors has dedicated mobile nodes for collecting data from the clusters within that sector. As a result, mobile nodes in the network will simultaneously collect the data contrary to single mobile node moving around the network for collecting the data which reduces latency in data collection. The results show that the energy consumption decreases and the throughput and network lifetime increase with an increase of nodes in the network for the proposed protocol.

## References

1. Kaur, P.: Wireless sensor networks: a survey. Int. J. Adv. Res. Comput. Sci. Softw. Eng. **5**, 988–994 (2015)
2. Xinhua, W., Jianbing, C.: Protocol for WSN. In: International Conference on Internet Technology and Applications (iTAP), IEEE (2011). doi:10.1109/ITAP.2011.6006320
3. Zhang, M., Babaei, A., Agrawal, P.: A cluster-based hybrid access protocol for wireless sensor networks. In: Sarnoff Symposium (SARNOFF), 2012 35th IEEE, pp. 1–5 (2012)
4. Gnanambigai, J., Rengarajan, N., Navaladi, N.: A clustering based hybrid routing protocol for enhancing network lifetime of wireless sensor network. In: 2nd International Conference on Devices, Circuits Systems, pp. 2–5 (2014)
5. Shilpa P. Kamble, N.M. Thakare, S.S. Patil: A review on cluster-based energy efficient routing with hybrid protocol in wireless sensor network. In: 2014 International Conference on Computer Communication and Informatics (ICCCI-2014), pp. 5–8 (2014)
6. Mohan, B.A., Sarojadevi, H.: A review on mobile data collector. Int. J. Latest Trends Eng. Technol. **7**(4), 240–247 (2016)
7. Safdar, V., Bashir, F., Hamid, Z., Afzal, H., Pyun, J.Y.: A hybrid routing protocol for wireless sensor networks with mobile sinks. In: ISWPC 2012 Proceedings, pp. 1–5 (2012). doi:10.1109/ISWPC.2012.6263665
8. Anastasi, G., Conti, M., Di Francesco, M.: Data collection in sensor networks with data mules: an integrated simulation analysis. In: 2008 IEEE Symposium on Computers and Communications, pp. 1096–1102. IEEE (2008). doi:10.1109/ISCC.2008.4625629
9. Basagni, S., Carosi, A., Petrioli, C.: Controlled vs. uncontrolled mobility in wireless sensor networks: some performance insights. IEEE Veh. Technol. Conf. (2007). doi:10.1109/VETECF.2007.70
10. Abdulaziz, M., Simon, R.: Mobile data collection using multi-channel network coding in wireless sensor networks. In: 2015 IEEE 40th Conference on Local Computer Networks (LCN), pp. 205–208. IEEE (2015). doi:10.1109/LCN.2015.7366307
11. Zhang, R., Pan, J., Liu, J. Xie, D.: A hybrid approach using mobile element and hierarchical clustering for data collection in WSNs. In: 2015 IEEE Wireless Communications and Networking Conference (WCNC), pp. 1566–1571. IEEE (2015). doi:10.1109/WCNC.2015.7127701

12. Mohan, B.A., Sarojadevi, H.: A hybrid approach for data collection using multiple mobile nodes in WSN (HADMMN). IEEE Int. Conf. Recent Trends Electron. Inf. Commun. Technol. **5**, 736–739 (2016). doi:10.1109/RTEICT.2016.7807916
13. Liu, X.: Atypical hierarchical routing protocols for wireless sensor networks: a review. IEEE Sens. J. **15**, 5372–5383 (2015). doi:10.1109/JSEN.2015.2445796
14. Mohan, B.A; Sarojadevi, H.: The study of distributed and centralized cluster formation protocol in WSNs, vol. 3, pp. 8–12 (2014)
15. Mohan, B.A, Sarojadevi, H.: Energy efficient clustering scheme with secure data aggregation for mobile wireless sensor networks (EECSSDA). In: 2016 International Conference on Innovations in Information, Embedded and Communication Systems, pp. 790–794 (2016)
16. Prabowo, S., Abdurohman, M., Erfianto, B.: (EDsHEED) Enhanced simplified hybrid, energy-efficient, distributed clustering for wireless sensor network. In: 3rd International Conference on Information and Communication Technology, pp. 97–101 (2015)
17. Singh, J., Mishra, A.K.: Clustering algorithms for wireless sensor networks: a review. In: 2015 2nd International Conference on Computing for Sustainable Global Development (INDIACom), vol. 2, pp. 637–642 (2015)

# Enrichment of UML by Introjection of Functional Model

Shivanand M. Handigund, B. N. Arunakumari
and Ajeet Chikkamannur

**Abstract** The de facto status of unified modeling language (UML) remains unaltered even after 20 years of advances in other technologies. The non-ascension in the status indicates some lacuna in UML, though it has undergone up-gradation in subsequent versions. Initially, UML is developed by three amigos to widen the scope and enhance the richness of diagramming languages. The clairvoyant study based on the number of involved diagrams indicates that, instead of unifying the semiotics of all three amigos diagramming languages, the diagrams of three languages are unified, discarding some of the diagrams with apprehension of descension in UML richness. The functional model which is quintessence of any open system has been discarded costing utilities of UML. This paper formulates the richness through logical positivism, unifies the semiotics of participating diagramming languages, and enhances through the reintroduction of modified functional model in the form of work process flow diagram.

**Keywords** Functional model · Unified modeling language · Semiotics Syntactics · Semantics · Pragmatics · Business process · Work process Open system · Class diagram · Data flow diagram

S. M. Handigund (✉)
Department of Information Science & Engineering, Vemana Institute of Technology, Bengaluru 560034, India
e-mail: smhandigund@iitbombay.org

B. N. Arunakumari
Department of Computer Science & Engineering, Research Resource Center, Visvesvaraya Technological University, Belagavi 590018, India
e-mail: arunakbn@gmail.com

A. Chikkamannur
Department of Computer Science & Engineering, R. L. Jalappa Institute of Technology, Bengaluru 561203, India
e-mail: ajeetac@rediffmail.com

# 1  Introduction

The goal of this paper is represented in state-of-the-art form of vision mission and objectives as follows:

*Vision*: To enhance the richness and ramification of unified modeling language (UML) from de facto to global standard.
*Mission*: To resuscitate the detrimental functional model into salubrious model to enhance the UML status.
*Objectives*:

- To formulate richness through logical positivism and provide ramification to UML.
- To transform the detrimental factors into incremental factors with respect to richness of UML.
- To expand the scope (utility) of UML to explicitly cover the open systems.

## 1.1  Motivation

Originally, the UML developed by three amigos is an attempt of mere unification of the diagramming pragmatics of their modeling languages, but there is no attempt for unification of trinity of diagramming semiotics viz., syntactics, semantics & pragmatics. So in their attempted unification to retain the some of the richness, the diagrams including functional model are deleted. The reason may be the difficulty in mapping the process of functional model to syntactics, semantics, and pragmatics of other languages participated in unification. The deletion of pragmatics results in non-availability of the language for specific pragmatics. Unfortunately, the information of the functional model is incorporated only in one of the diagramming languages. The functional model is the only diagram that identifies the recursion, iteration, concurrencies of the activities, the flows to and fro the environment to the system. Therefore, there is quintessence of desiderate UML with introjection of functional model. This paper attempts to resuscitate the loss by reintroducing it (functional model) by replacing the 'process' of functional model with work process activities. The quintessence of the process is the retainment of modularity. This may be achieved by representing the embedded behavioral requirements between (i) single read and single write, (ii) single read and many writes, (iii) many reads and single write, and (iv) many reads and many writes. The increasing complexity of business information system outdates (i). The use of data store embedding many reads and many writes transforms (iv) to (ii) or (iii) and (iii) highly inefficient as number of reads used for single purpose and cannot be converted into other perspective views. Therefore, (ii) is the ideal replacement of modular process to be incorporated as pragmatics in one of the other two diagramming languages. In addition, this module work process is converted into part of usecase and/or part of work in the business process.

Though, UML is considered as de facto standard multilevel diagramming language, it does not possess ramification of diagramming pragmatics. Each diagram is to be designed from scratch. This is herculean task of higgledy-piggledy designing diagrams for the requirements. Moreover, the input for all these diagrams is to be repeatedly abstracted from the analyzed dossier. This knocks the analysis as onerous task. Furthermore, the purpose of the unified modeling language is to enhance the utility i.e. utility $\propto$ richness of the language but the achievement of au-courant UML is (utility $\otimes$ richness) = constant, which has jeopardized the aim of UML.

## 1.2  Literature Survey

UML was developed by unifying the modeling techniques of three amigos viz. object modeling technique (OMT) [1], object-oriented analysis and design with applications (OOADA) [2], and object-oriented software engineering [3]. In this unification [1–6], the developers instead of unifying the semiotics of the languages have summed up their diagrams i.e. enriching the number of diagrams. This enhances number of synonymies and heteronyms which dampens the richness of the language. A diagramming language used for information engineering should ease the process of understanding and communication between stakeholders. The virtues of the language are determined by its capability of coding the information in different abstraction levels which needs more number of syntactics, less number of semantics. The abstraction levels determine the richness of the language by [7]. To retain the richness of the diagramming language, the functional model (detrimental to richness) of OMT [6] is discarded in the au-courant UML. As a result, the information about the interacting actors is lost which undermine applicability of UML to open systems. In our research, we have attempted to enhance the richness of UML without losing the interaction with various actors through the resuscitation of functional model by transforming synonymous syntactics into an existing semantics. The fact that the retainment of de facto standard status of UML for long time inspite of its perennial updation reveals that the scope of UML pragmatics is limited in its present form. Moreover, the UML pragmatics is loomed. Any language semiotics should be organizable in stratiform ramification levels that converts the extant lower level semantics into higher level syntactics so as to concurrently enhance the syntactics and diminishes the semantics. The hierarchical levels in any diagramming language enhance the richness of the language, as lower level semantics are converted to higher level syntactics. In UML [8] work process diagram which is a pragmatics at the lower level (semantics in the proposed higher level diagram) has been transformed into syntactics at our proposed work process flow diagram (WPFD). To audit the correctness and completeness of software development, the proposed WPFD serves as milestone as system and environment are embedded in a single diagram. All work processes should satisfy flow constraints [9] with provider actors as sources and consumer actors as sinks. If there is

no provider and consumer of information then $f(V) - f(s) - f(t) = 0$. Here $V$ is set of work processes, $s$ is set of sources, and $t$ is the set of sinks. The capacity constraints are roles and responsibilities of each work processes. Moreover, they enhance the richness of the language as the syntactic of WPFD (higher level) is gemmated into pragmatics of work process diagram at the lower level.

## 1.3 Taxonomy

*Actor*: An actor is a person, file, or environment system which interacts with system.

*Attribute*: An attribute is syntactic representation of a class/attribute characteristic.

*Actors' interface attributes*: Subset of the actor's attributes participating in the system which is either provided/consumed by the system in the actor's language.

*Business process*: Process that takes zero/more input and produces one/more output that is of some value to the society. In this business process, the information flows from an individual/department to another individual/department.

*Class*: Comprises cohesive set of attributes of an entity (that satisfies good database design principles) along with methods that define high cohesive subset of attributes of the class which are very loosely coupled (that satisfy good software engineering principles).

*Object method*: A behavioral thread that defines high cohesive subset of attributes of the class.

*Pragmatics*: An architectonic way of representing the perspective view of the information system.

*Syntactics*: The primitive symbol used in UML.

*Semantics*: The valid combination of syntactics that serves as constraint.

*Semiotics*: It is a group name for trinity of syntactics, semantics, and pragmatics.

*Work process*: A set of activities performed by an individual with tools and techniques taking input from other work processes or actors and transforming the result to other work processes or actors.

## 2 Methodology

Our proposed methodology is to strengthen the UML by reintrojection of functional model in the form of au-courant work process flow diagram (WPFD). The input to this diagram is semiotics of functional model, class, and data flow diagrams, and the output is resuscitated WPFD. The work process flow diagram represents a perspective view of business process through the reticulation of activities of each work process. Normally, work process activities means the object methods of other classes in which the work process attributes are referenced. Since, all these work

processes are components of single business process, the ensuing WPFD should represent all interrelated work processess in a reticular directional graph analogical to network flow graph in which each vertex represents either 'work process' or an 'inanimate class', each directed edge super scribing on it the 'attribute names' represents data flow, the 'environment' symbol represents actor, and the rectangle with three compartments represent 'class'.

## 2.1 Procedure for Abstraction of Tagmemes of Work Process Flow Diagram

*Input*: class and data flow diagrams, semiotics of functional model [10]
   // The following steps from 1 to 8 identifies work process

1. Get class diagram and store the information in five column table as shown. For each method of the class, enter class name, definitional attributes, object method, signature attributes, and initially fifth and sixth columns as blank and as shown in Table 1.

Let i represents a row explaining parameters of each method and k represent the column number containing parameter values. Let (i, j) represent a pair of consecutive rows.

2. Sort the records of table in the lexicographic order of class name.
3. Organize $i_2$ and $i_4$ columns attributes in the lexicographic order of attribute names (it eases the process of intersection)
4. Initialize $C_{i5} \leftarrow i_2$
5. While ($i_1 = j_1$)
   $C_{i5} \leftarrow C_{i5} \cup j_2$
   $l \leftarrow i$
   $i \leftarrow j$
   $j \leftarrow j + 1$
   end while
   for (u = l to i)
   $u_5 \leftarrow C_{i5}$
   $i \leftarrow j$
   $j \leftarrow j + 1$
   Repeat step 4 until $j \leftarrow \varphi$
6. Initialize $i \leftarrow 1$; $j \leftarrow i + 1$

**Table 1** Class table

| Class name $\langle i_1 \rangle$ | Defined attributes $\langle i_2 \rangle$ | Object method $\langle i_3 \rangle$ | Signature attributes $\langle i_4 \rangle$ | Class attributes $(C_{i5} = \cup i)$ $\langle i_5 \rangle$ | Array of class attributes pair $\langle i_6 \rangle$ |
|---|---|---|---|---|---|

7. While $(i_1 \neq j_1)$
   For $(i = 1$ to $n)$
   For $(j = 1$ to $n$ & $j \neq i)$
   Take $i_5 \cap j_4$
   Store $i_5$ $(\cap j_4)$ in matrix
   $I6(i, j)$
   end j end i
8. Consider matrix $I6(i, j)$ and $[I6(i, j)]^{-T}$
   $N(i, j) \leftarrow I6(i, j) \parallel [I6(i, j)]^{-T}$
   $\forall (i, j)$ and $(i, j + n) \neq \varphi$
   If $((i, j) \neq \varphi$ & $(i, j + n)) \neq \varphi$
   then take $i = i_1$ and $i_1$ is intervivo (work process)
   // The steps 9–10 identifies data flows and actor interface attributes
9. If $i$ is intervivo $(i, j) \neq \varphi$ then $\exists$ flow $N(i, j)$ from $j$ to $i$
   else if $(i, j + n) \neq \varphi$ then $\exists$ flow $N(i, j)$ from $i$ to $j$
10. $N(i, j) \notin \cup C_{i5}$ then $N(i, j)$ are interface attributes of actor $j$
    //Step 11 identifies class (data store)
    if $i \in \{I\}$, $j \in \{J\}$ and
10. $[\{i\}] \geq 2$ & $[\{j\}] \geq 2$
    $i \in [\{I\}] \exists$ flows from $i$ to $j$
    $\forall j \in J$ and vice versa
    Then introject additional vertex $k$ (data store) and transform the flow
    $(i, j)$ to $(i, k)$ & $(k, j)$

## 2.2 Design of New Work Process Flow Diagram

The *syntactics* of the work process flow diagram:

*Work process*: Process should have been denoted as oval shape but is representation syntactic for usecase in UML. Since it gives rise to heteronym, which dampens the richness of the language, and process symbol is represented as 'predefined process' from well-known project management body of knowledge (PMBOK) [11] it has used as symbol for work process. This is because PMBOK is de facto standard project development process that equally valid for software development. Therefore, the section of people who are engaged in software development activities may well understand the same syntactics without harming UML.



```
|| <<Work
   process name>> ||
```

*Data flow*: Data flow is input flow and output flow, respectively for the process at the head of arrow and at the tail of the arrow. In UML, message starts with verb

represented by an arrow, transition is represented by arrow without any list of attributes. If adopted in WPFD, it dampens the richness through the enhancement of synonyms. Hence, syntactics of data flow denoted as directed arrow with super scribed list of attributes.

List of attributes

*Environment (Actor)*: Actor symbol is represented by 'environment' symbol adopted from PMBOK.

<<Actor>>

*Class (data store)*: Data store is analogous to class except the behavioral aspects. Therefore, data store is represented by class symbol. Data store is represented as sequential, indexed, hash file. The semantics and pragmatics of functional model are analogous to data flow diagram, and hence is adopted in WPFD.

<Inanimate
class name>

## 2.3 Formulation of Richness of UML

English is the ideal language that helps to support our formulation. The English language has millions of syntactics with lot of flexibility for their use i.e. the semantics is less. Though, we feel there exist more synonyms and heteronyms, the clairvoyant study reveals that each synonymous word in strict sense gives different meanings. This has served as the base for our formulation.

The richness of the language depends on the

$$
\text{Richness of the language} \propto \left\{ \frac{\text{syntactics} - (\text{synonyms \& heteronyms})}{\frac{1}{\text{semantics}}} \right\} \quad (1)
$$

$$
\text{Richness of the language} = k \left\{ \frac{\text{syntactics} - (\text{synonyms \& heteronyms})}{\frac{1}{\text{semantics}}} \right\} \quad (2)
$$

Let $R_0$ be the richness of the UML after combining all the trinity of diagramming languages minus functional model and is given by the following equation

$$R_0 = k\left[\text{Synt}_{\text{OMT}} + \text{Synt}_{\text{OOADA}} + \text{Synt}_{\text{Usecase}}\right] + \left[\frac{1}{\text{Sem}_{\text{OMT}} + \text{Sem}_{\text{OOADA}} + \text{Sem}_{\text{Usecase}}}\right]$$

$$- \left[\text{Syno}_{\text{UML}} + \text{Heto}_{\text{UML}}\right] - \left[\left[\text{Synt}_{\text{FM}} + \frac{1}{\text{Sem}_{\text{FM}}}\right] - \left[\text{Syno}_{\text{FM}} + \text{Heto}_{\text{FM}}\right]\right]$$

$$\text{Synt}_{\text{UML}} = \left[\text{Synt}_{\text{OMT}} + \text{Synt}_{\text{OOADA}} + \text{Synt}_{\text{Usecase}}\right], \text{Sem}_{\text{UML}} = \left[\text{Sem}_{\text{OMT}} + \text{Sem}_{\text{OOADA}} + \text{Sem}_{\text{Usecase}}\right]$$

$$(3)$$

$$(3) \Rightarrow$$

$$R_0 = k\left[\text{Synt}_{\text{UML}} + \frac{1}{\text{Sem}_{\text{UML}}} - \left[\text{Syno}_{\text{UML}} + \text{Heto}_{\text{UML}}\right]\right] - \left[\text{Synt}_{\text{FM}} + \frac{1}{\text{Sem}_{\text{FM}}} - [0+1]\right] \qquad (4)$$

After introduction of resuscitated functional model (work process data flow diagram) in the UML, the richness of the language $(R_0^+)$ is increased. The richness of language flows from the additional syntactics namely work process, class, environment, and data flow. Rectangle is used to represent class and data store. In the object-oriented technology (OOT), the data store is represented as class. Data flow and message are represented by directed arrow. Data flow contains superposed attributes list. Message starts with transitive verb; therefore, its synonym is eliminated. Actor is represented by rectangle inscribed with two rectangular frameworks. The 'process' symbol which represent a work process (a pragmatics of lower level work process diagram) is used in WPFD as a syntactics in tune with hierarchical ramification. The number of semantics representing this has been transformed to syntactics. Thus, syntactic is increased by one each of the above symbol and is given by following equation

$$R_0^+ = k\left[\left[\left[\text{Synt}_{\text{UML}} + 4\right] + \frac{1}{\text{Sem}_{\text{UML}} - 1}\right] - \left[\text{Syno}_{\text{UML}} + \text{Heto}_{\text{UML}}\right]\right] + \left[\left[\text{Synt}_{\text{FM}} + \frac{1}{\text{Sem}_{\text{FM}}}\right] - [0+1]\right]$$

$$(5)$$

From Eqs. (4) and (5) $R_0^+ > R_0$. This shows that with true unification of three amigos' diagramming languages, the functional model which was hitherto dampening the richness is actually enhances the richness of UML.

## 2.4   Illustrative WPFD

The following is the illustrative WPFD for banking environment. In the diagram, both consumer and provider actors have been shown with duplicates to feel the analogy with WPFD as network flow diagram so as to satisfy the flow constraints. Since, the nature and purpose of two diagrams are different with interchange of vertices and edges, the flow constraints are applied to vertices (process) as shown in Fig. 1.

**Fig. 1** Work process flow diagram for banking system

## 3    Conclusion

UML has not retained its true name as it does not contain unification of semiotics of participant diagramming languages which enhances its richness. Attempt has been made to formulate richness on the pedestal of mathematical rigor and facilitate the ensuing WPFD to be ramified extant work process diagram. The functional model was detrimental to richness which we transformed into incremental to richness by enhancing syntactics with 'work process' as process, dataflow as an 'arrow', data store as 'transient class', and actor as 'environment'. The design of WPFD requires least effort as it is a common perspective view abstractable from class diagram in the presence of environment (actor). To ease the understanding of WPFD, the symbols are chosen from de facto standard project management symbol as depicted in PMBOK a sequel of standard managerial activities. UML is made more realistic as with our innovation, it enhances the applicability of UML to open systems. The WPFD can be ramified to work process diagrams at the lower level.

# References

1. Rumbaugh, J., Blaha, M., Premerlani, W., Eddy, F., Lorensen, W.E.: Object-Oriented Modeling and Design. Prentice Hall, Englewood Cliffs (1990). ISBN 0-13-629841-9
2. Booch, G.: Object-Oriented Analysis and Design with Applications, 3rd edn. Pearson Education India, Noida (2009). ISBN 8131722872
3. Jacobson, I., Christerson, M., Jnsson, P., Overgaard, G.: Object Oriented Software Engineering. Addision-Wesely, New York (1992)
4. Booch, G., Jacobson, I., Rumbaugh, J.: The Unified Modeling Language User Guide, 2nd edn. Pearson Publisher (2007). ISBN 9788131715826, 8131715825
5. Rumbaugh, J., Jacobson, I., Booch, G.: Unified Modeling Language Reference Manual, 2nd edn. Addison Wesley, Pearson Education, Boston (2005)
6. Eriksson, H.-E., et al.: UML 2 toolkit, vol. 26. Wiley, Hoboken (2003). ISBN 978-0-471-46361-0
7. Stephenson Smith S.: New International Webster Comprehensive Dictionary of English Language (Deluxe Encyclopedic Edition), Trident Press International (2003). ISBN 9781582795577
8. Handigund, S.M., Kshama, S.B., Ranjana, N.: An ameliorated methodology for design and development of a work process diagram to be incorporated in UML. In: Proceedings of IEEE International Conference on Advances in Computing, Communications and Informatics (ICACCI) 2013, pp. 184–190. doi:10.1109/ICACCI.2013.6637168 (2013)
9. Coremen, T.H., et al.: Introduction to Algorithms, 2nd edn. PHI Learning Private Limited, New Delhi (2009). ISBN 978-81-203-2141-0
10. Pressman, R.S.: Software Engineering: A Practitioner's Approach, 7th edn. McGraw Hill Higher Education, New York (2010)
11. Project Management Institute: A Guide to the Project Management Body of Knowledge (PMBOK® Guide)—Fifth Edition 2012, Kindle Edition, Publisher: Project Management Institute (2010). ISBN-13: 978-1935589815

# Fault Detection and Classification on Distribution Line with Penetration of DFIG-Driven Wind Farm Using Fuzzy System

**Anamika Yadav and Chitrarth Rangari**

**Abstract** This paper describes fault detection and classification scheme for protection of doubly feeded distribution line using fuzzy logic system with penetration of the wind farm driven by doubly fed induction generator. The simulation study of doubly feeded distribution line system consists of 120-kV, 60-Hz source and 9-MW wind farm connected to distribution line of 30 km length, which is modelled using pi block in Simulink toolbox of MATLAB 2013a. The proposed method has been exhaustively examined with large variety of fault situations with different fault parameters like all ten types of fault, fault inception angle and fault location. The proposed scheme also identifies the evolving fault and classifies the faulty phase(s) as well. Simulation study shows that the developed scheme works accurately for huge number of fault cases with half cycle detection time.

**Keywords** Fault detection · Fault classification · Wind farm · Fuzzy system Distribution line

## 1 Introduction

As we know, the quantity of nonrenewable energy sources are present in fixed amount below the earth and it is decreasing day by day. Power generation mainly depends upon the nonrenewable sources, which are depleting day by day. In view of this, use of the renewable energy sources like wind energy or solar energy is increasing. With penetration of renewable energy resources at distribution level, the protection of distribution system is greatly affected. Thus, the protection of distribution lines with penetration of wind farm has raised the interest of researchers. Monitoring strategy of short-circuit fault between turns of the stator windings and open stator phases in doubly fed induction generator by fuzzy logic technique has

A. Yadav (✉) · C. Rangari
Department of Electrical Engineering, National Institute of Technology,
Raipur, Chhattisgarh 492010, India
e-mail: ayadav.ele@nitrr.ac.in

been proposed in [1]. Online monitoring of induction motors is becoming increasingly important. Knowledge-based fuzzy logic approach helps in diagnosing the induction motor faults [2]. Another method for induction motor fault detection using electrical signature analysis techniques is discussed in [3]. The main difficulty in this task is the lack of an accurate analytical model to describe a faulty motor having different types of faults like winding faults [4]. Implementation of broken rotor bar fault detection in an inverter-fed induction motor using motor current signal analysis (MCSA) and prognosis with fuzzy logic is reported in [5]. Method to diagnose static and dynamic air gap eccentricities in doubly fed induction generators operated for closed-loop stator power control by using a modified control technique to enable detection and isolation of this fault from electrical unbalances in the stator and rotor and load torque oscillations is discussed in [6].

Above methods do not deal with fault detection and classification in distribution line. Different methods are used for detection and classification of fault in distribution line [7–17]. Artificial neural network and wavelet transform has been combined in order to detect and classify fault in [7] [8]. The demerit of using ANN is that it requires training or learning, which is very difficult and time taking process [9]. Techniques based on fuzzy system have been proposed by many researchers in [10–13]. Another method is proposed for fault detection, classification and location in transmission lines, which is reported in [13]. Artificial neural-network-based fault location scheme has been proposed for power distribution lines using the frequency spectra of fault data [14]. Combination of ANN and fuzzy system can be used for fault diagnosis approach based on combined wavelet transform and adaptive neuro-fuzzy inference system for fault section identification, classification and location in a series compensated transmission line [15]. Another method is proposed for detection of islanding, and fault disturbance in power system is reported in [16].

This paper deals with fuzzy system implementation for fault detection and classification in the doubly feeded distribution line. The method has been tested for different cases by varying the length of the distribution line and fault inception angle. The proposed scheme also identifies the evolving fault and classifies the faulty phase(s). This method is giving efficient result for all the ten types of faults. The different sections are present in this paper as Sect. 2 describes about the simulated system. Section 3 describes the fuzzy inference system, and FIS is divided into two parts: first FIS is for detection of fault and second FIS is for ground fault detection. Section 4 discusses the test results for different fault scenarios followed by conclusions (Table 1).

## 2 Doubly Feeded Distribution Line System

The system to be simulated is a 120-kV, 60-Hz, 9-MW wind farm driven by doubly fed induction generator. It consists of 30 km distribution line. Simulation is carried out in MATLAB/SIMULINK 2013a environment. A 9-MW wind farm driven from

**Table 1** Parameters of DFIG

| S. No. | Parameters | Values |
|---|---|---|
| 1. | No. of wind turbine | 06 |
| 2. | Nom. power, L-L volt. and freq. [Pn (VA), Vsnom (Vrms), Vrnom (Vrms), fn (Hz)] | 1.5e6/.9VA, 575 V, 1975 V, 60 Hz |
| 3. | Stator [Rs, Lls] (p.u.) | 0.023 Ω, 0.18 H |
| 4. | Rotor [Rr′, Llr′] (p.u.) | 0.016 Ω, 0.16 H |
| 5. | Magnetizing inductance Lm (p.u.) | 2.9 H |



**Fig. 1** Doubly fed distribution line model simulated in MATLAB

doubly fed induction generator is connected at the sending end, and 120 kV three phase programmable voltage source is connected at the receiving end as shown in Fig. 1. A 9-MW wind farm consisting of six 1.5-MW wind turbines connected to a 575-V bus exports power to 120 kV grid through 30 km line. The pi-section line is used in distribution line system. The parameters of doubly fed induction generator are given in Table 2, and line parameters of distribution line are given in Table 3. Three buses are placed in between sending end and receiving end having voltages 120 kV, 25 kV and 575 V, respectively, and the measurements are being taken at 25-kV bus.

The DFIG stator is connected directly to the grid. Back-to-back VSCs are included in the rotor circuit. A braking resistor is provided in the dc-link bus as a form of protection to dissipate excess energy during a grid fault. The resistor is connected to the dc-link bus in series with an IGBT and is referred to as dc crowbar. The DFIG consists of stator windings connected directly to the grid and wound rotor windings connected to a power converter. Torque is created by the interaction of the rotor magnetic field with the stator magnetic field. The magnitude of the generated torque is dependent on both the strength of the two magnetic fields [17]. Two types of subsystems are shown in Fig. 1, one is for fault phase classification, and the other is for ground detection. In the fault phase classifier subsystem, the

**Table 2** Parameters of distribution line

| S. No. | Parameters | Units | Values |
|---|---|---|---|
| 1. | Line length | km | 30 |
| 2. | Positive Seq. R1 | Ω/km | 0.1153 |
| 3. | Positive Seq. L1 | Henry/km | 1.05e-3 |
| 4. | Positive Seq. C1 | Farad/km | 11.33e-009 |
| 5. | Zero Seq. R0 | Ω/km | 0.413 |

**Table 3** Fuzzy rules for fault classification

| Parameters | LI | MI | HI |
|---|---|---|---|
| LV | Lowtrip | Lowtrip | Hightrip |
| MV | Lowtrip | Lowtrip | Hightrip |
| HV | Lowtrip | Lowtrip | Hightrip |

values of each phase voltage and current at the 25-kV bus are recorded. To find the fundamental values, these values are passed from a low pass filter of order 2, having a cut-off frequency of 400 Hz and then by a zero-order hold block and a discrete Fourier transform (DFT) block with sampling frequency of 1.2 kHz. In the ground detection, subsystem, the zero-sequence components of voltage and current are taken by using sequence analyzer. Now the values of each phase voltage and current are sent to the fuzzy system controller with rule viewer block, which examines faulty phases and differentiates between different types of fault. The output is seen from the respective scopes.

## 3 Fuzzy Inference System

The flowchart of the proposed scheme is illustrated in Fig. 2. Two different fuzzy inference systems are created. One is for detection of faulty phase(s), and the other is for detection of ground. The rules and procedures to obtain results are given in tables below.

### 3.1 FIS for Faulty Phase Detection

The values of voltage and current of each phase are taken as the inputs. FIS file used is Mamdani type, and triangular membership function is used. The ranges of membership function are taken as low, medium and high. The output is named as 'Trip', and it is having triangular membership functions. Ranges are defined as lowtrip, vergetrip and hightrip. Table 3 is showing the rules used.

**Fig. 2** Flowchart of proposed scheme

## 3.2 FIS for Ground Detection and Fault Classification

The zero-sequence components of voltage and current are taken by sequence analyzer, and then, they are used as inputs in the FIS. Mamdani-type FIS file is used with three ranges, taken as low, medium and high, with triangular membership functions. The output 'GND' is defined with triangular membership functions with ranges defined as NG, VG and G. Table 3 shows the rules used.

# 4    Results and Discussion

## 4.1    Variation in Fault Location

The proposed scheme is applied for all types of fault with fixed value of fault resistance of 0.001 Ω at different location in distribution line (Table 4).

The fault inception angle is taken as 0° at 1.5 s for all type of fault. Now fault is created near the 25-kV bus, and length of the distribution line is taken 30 km as shown in Table 5. Figure 3 shows the case of AG fault out of all the cases discussed in Table 5. Figure 3 shows the occurrence of AG fault at 3 km from the 25-kV bus on distribution line with fault resistance of 0.001 Ω. It requires 0.011 and 0.005 s to detect the phase A and ground G, respectively.

**Table 4**  Fuzzy rules for ground fault

| Parameters | LI0 | MI0 | HI0 |
|---|---|---|---|
| LV0 | NG | G | G |
| MV0 | NG | G | G |
| HV0 | NG | G | G |

**Table 5**  Simulation results of fault on distribution line near 25-kV bus

| Type | Fault location (km) | Fault resistance (Ω) | Phase detection output (time required to detect the faulty phase/ground in s) | | | |
|---|---|---|---|---|---|---|
| | | | A | B | C | G |
| AG | 3 | 0.001 | 1 (0.011 s) | 0 | 0 | 1 (0.005 s) |
| BG | 6 | 0.001 | 0 | 1 (0.009 s) | 0 | 1 (0.003 s) |
| CG | 9 | 0.001 | 0 | 0 | 1 (0.014) | 1 (0.003) |
| AB | 12 | 0.001 | 1 (0.008 s) | 1 (0.009 s) | 0 | 0 |
| ABG | 15 | 0.001 | 1 (0.011 s) | 1 (0.010 s) | 0 | 1 (0.005 s) |
| BC | 18 | 0.001 | 0 | 1 (0.009 s) | 1 (0.008 s) | 0 |
| BCG | 21 | 0.001 | 0 | 1 (0.011 s) | 1 (0.010 s) | 1 (0.007 s) |
| CAG | 27 | 0.001 | 1 (0.013 s) | 0 | 1 (0.015 s) | 1 (0.005 s) |
| ABC | 29 | 0.001 | 1 (0.013 s) | 1 (0.011 s) | 1 (0.015 s) | 0 |

**Fig. 3** AG fault at 3 km with fault resistance of 0.001 Ω

## 4.2 Variation in Fault Inception Angle

The fault inception angle is varied from 0° to 360°, the fault resistance is fixed at 0.001 Ω, and results for different types of fault on distribution line near the 25-kV bus are shown in Table 6. All the types of fault are tested at different inception angles, and the results are shown in Table 6. Out of all the cases discussed above, Fig. 4 shows the occurrence of ABC fault at 29 km from the 25-kV bus on distribution line with fault inception angle 135° and fault resistance 0.001 Ω. It requires 0.013, 0.010 and 0.0010 s to detect the phases A, B and C, respectively. From the output result, it can be noticed that the implementation of fuzzy logic scheme quickly detects the faulty phase and classifies it. It gives accurate result from the quarter cycle time to half cycle time, which is very beneficial thing for distance relay protection.

## 4.3 Evolving Fault on Distribution Line

Evolving fault can be defined as a fault, which originates in one phase of a distribution line and after a few cycles spreads into another phase, e.g. a single line to

**Table 6** Simulation results of faults at different fault inception angles

| Type | Fault resistance (Ω) | Fault inception angle (°) | Phase detection output (time required to detect the faulty phase/ground in s) | | | |
|------|------|------|------|------|------|------|
| | | | A | B | C | G |
| AG | 0.001 | 0 | 1 (0.011 s) | 0 | 0 | 1 (0.005 s) |
| AB | 0.001 | 45 | 1 (0.008 s) | 1 (0.008 s) | 0 | 0 |
| ABG | 0.001 | 90 | 1 (0.015 s) | 1 (0.013 s) | 0 | 1 (0.005 s) |
| ABC | 0.001 | 135 | 1 (0.013 s) | 1 (0.010 s) | 1 (0.010 s) | 0 |
| BG | 0.001 | 180 | 0 | 1 (0.009 s) | 0 | 1 (0.003 s) |
| BC | 0.001 | 225 | 0 | 1 (0.010 s) | 1 (0.011 s) | 0 |
| BCG | 0.001 | 270 | 0 | 1 (0.011 s) | 1 (0.010 s) | 1 (0.008 s) |
| CG | 0.001 | 315 | 0 | 0 | 1 (0.014 s) | 1 (0.003 s) |
| CA | 0.001 | 360 | 1 (0.010 s) | 0 | 1 (0.009 s) | 0 |



**Fig. 4** ABC fault at 29 km from 25-kV bus with fault inception angle 135° and fault resistance 0.001 Ω

ground fault gets converted into double line to ground fault. From Table 7, we can notice that after some cycles, the AG fault in the distribution line gets converted into ABG fault.

Some of the cases are tested in which first fault is created on 3 km from the 25-kV bus at 1.5 s. The second fault is created on 3 km from the 25-kV bus at

**Table 7** Simulation results of evolving fault on distribution line

| S. No. | Type of first fault | Phase detection output (time required to detect the faulty phase/ground in s) | | | | Type of second fault | Phase detection output (time required to detect the faulty phase/ground in s) | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | A | B | C | G | | A | B | C | G |
| 1 | AG | 1 (0.011) | 0 | 0 | 1 (0.005) | ABG | 1 (0.011) | 1 (0.027) | 0 | 1 (0.005) |
| 2 | BG | 0 | 1 (0.008) | 0 | 1 (0.002) | BCG | 0 | 1 (0.008) | 1 (0.031) | 1 (0.002) |
| 3 | AG | 1 (0.011) | 0 | 0 | 1 (0.005) | ACG | 1 (0.011) | 0 | 1 (0.027) | 1 (0.005) |
| 4 | AG | 1 (0.011) | 0 | 0 | 1 (0.005) | ABCG | 1 (0.011) | 1 (0.027) | 1 (0.031) | 1 (0.005) |

**Fig. 5** Evolving fault ABG at 3 km from 25-kV bus and fault resistance 0.001 Ω

1.52 s. Figure 5 shows the occurrence of single line to ground fault, i.e. AG fault near the 25 kV bus on distribution line at fault resistance 0.001 Ω. Here, phase A and ground G get high after 0.011 and 0.005 s, respectively. In case of an evolving fault; the inital single line to ground fault is converted into double line to ground fault after some time which is clearly seen in the Fig. 5 where phase B also gets high after 0.027 s.

## 5  Conclusions

The proposed scheme detects and classifies the type of fault at different position on doubly fed distribution line with penetration of the wind farm driven from doubly fed induction generator (DFIG) with implementation of fuzzy inference system. The FIS type used here is Mamdani. Triangular types of membership functions are used. For phase classification, we are using fundamental values of voltage and current. For ground fault detection, zero-sequence components of voltage and current are taken as inputs. It has been noticed that the system will work for fault resistance at 0.001 Ω, i.e. on low impedance faults. This scheme is applied with various fault inception angles. The evolving fault has also been tested with the system.

From the obtained results, it is observed that the system works efficiently in various cases to classify and detect the type of fault. Various simulation results

confirm the suitability of the proposed scheme for the detection and classification of fault in a doubly fed distribution system.

# References

1. Merabet, T., Bahi, N.: Condition monitoring and fault detection in wind turbine based on DFIG by the fuzzy logic. Energy Procedia **74**(2015), 518–528 (2015)
2. Bhardwaj, K., Agarawal, A.: Fault diagnosis of three phase induction motor using fuzzy logic controller and sequence analyzer. Int. J. Electr. Instrum. Eng. **2**(2), 112–118 (2012)
3. Somaya, A., Shehata, M., Hamdy, S., El-Goharey, M., Marei, I., Ibrahim, A.: Detection of induction motors rotor/stator faults using electrical signatures analysis. RE&PQJ **1**(11), 382–387 (2013)
4. Mini, V.P., Sivakotaiah, S., Ushakumari, S.: Fault Detection and Diagnosis of an Induction Motor Using Fuzzy Logic, SIBIRCON-2010, pp. 459–464, Russia, 11–15 July 2010
5. Akar, M., Ankay, I.C.: Broken rotor bar fault detection in inverter-fed squirrel cage induction motors using stator current analysis and fuzzy logic. Turk. J. Electr. Eng. Comput. Sci. **20** (Supl 1), 1077–1089 (2012)
6. Toliyat, H.A., Vivek, Sundaram, M.: Air gap eccentricity fault detection in DFIG-based wind. In: International Conference on Energy Conversion Systems, Monaco (2011)
7. Yadav, A., Swetpadma, A.: A novel transmission line relaying scheme for fault detection and classification using wavelet transform and linear discriminant analysis. Ain Shams Eng. J. **6** (1), 199–209 (2015)
8. Silva, K.M., Souza, B.A., Brito, N.S.D.: Fault detection and classification in transmission lines based on wavelet transform and ANN. IEEE Trans. Power Deliv. **21**(4), 3700–3705 (2006)
9. Yadav, A., Dash, Y.: An overview of transmission line protection by artificial neural network: fault detection, fault classification, fault location, and fault direction discrimination. Adv. Artif. Neural Syst. **2014**, 1–20 (2014)
10. Das, B., Reddy, J.V.: Fuzzy-logic-based fault classification scheme for digital distance protection. IEEE Trans. Power Deliv. **20**(2), 609–616 (2005)
11. Adhikari, S., Sinha, N., Dorendrajit, T.: Fuzzy logic based on-line fault detection and classification in transmission line. SpringerPlus **5**, 1–14 (2016)
12. Swetapadma, A., Yadav, A.: Fuzzy inference system approach for locating series, shunt, and simultaneous series, shunt faults in double circuit transmission lines. Comput. Intell. Neurosci. **2015**, 1–12 (2015)
13. Yadav, A., Swetapadma, A.: Enhancing the performance of transmission line directional relaying, fault classification and fault location schemes using fuzzy inference system, IET Gener. Transm. Distrib. **9**(6), 580–591 (2015)
14. Aslan, Y., Yağan, Y.E.: Artificial neural-network-based fault location for power distribution lines using the frequency spectra of fault data. Electr. Eng. **99**, 301–311 (2017)
15. Swetapadma, A., Yadav, A.: High speed directional relaying using adaptive neuro-fuzzy inference system and fundamental component of currents. IEEJ Trans. Electr. Electron. Eng. **10**, 653–663 (2015)
16. Ray, P.K., Panigrahi, B.K., Rout, P.K., Mohanty, A., Dubey, H.: Fault detection in IEEE 14-bus power system with dg penetration using wavelet transform. In: First International Conference on Advancement of Computer Communication and Electrical Technology, pp. 1–5, Murshidabad, India, Oct 2016
17. Snyder, M.A.: Development of simplified models of doubly-fed induction generators (DFIG). Department of Energy and Environment Division of Electric Power Engineering Chalmers University of Technology Goteborg, pp. 13–75, Sweden (2012)

# Text Summarization with Automatic *Keyword* Extraction in Telugu e-Newspapers

**Reddy Naidu, Santosh Kumar Bharti, Korra Sathya Babu and Ramesh Kumar Mohapatra**

**Abstract** Summarization is the process of shortening a text document to make a summary that keeps the main points of the actual document. Extractive summarizers work on the given text to extract sentences that best express the message hidden in the text. Most extractive summarization techniques revolve around the concept of finding keywords and extracting sentences that have more keywords than the rest. Keyword extraction usually is done by extracting relevant words having a higher frequency than others, with stress on important one's. Manual extraction or annotation of keywords is a tedious process brimming with errors involving lots of manual effort and time. In this work, we proposed an algorithm that automatically extracts keyword for text summarization in Telugu e-newspaper datasets. The proposed method compares with the experimental result of articles having the similar title in five different Telugu e-newspapers to check the similarity and consistency in summarized results.

**Keywords** Automatic keyword extraction · e-newspapers · NLP
Summarization · Telugu

R. Naidu (✉) · S. K. Bharti · K. S. Babu · R. K. Mohapatra
National Institute of Technology, Rourkela, Odisha 769008, India
e-mail: naidureddy47@gmail.com

S. K. Bharti
e-mail: sbharti1984@gmail.com

K. S. Babu
e-mail: prof.ksb@gmail.com

R. K. Mohapatra
e-mail: mohapatrark@nitrkl.ac.in

# 1   Introduction

Many popular Telugu e-newspapers are freely available on the internet, such as Eenadu, Sakshi, Andhrajyothy, Vaartha, and Andhrabhoomi. The extraction of all the relevant information from these newspapers is a tedious job for people. So, there is a need for a tool that extracts only relevant information from these data sources. To get the required information, we need to mine the text from newspapers. Natural Language Processing (NLP) is a powerful tool for text mining. Text mining deploys some of the techniques of NLP and includes tasks like Automatic Keyword Extraction and Text Summarization. Summarization is a process where the most prominent features of a text are extracted and compiled into a short abstract of the original wording [1]. According to Mani and Maybury [2], text summarization is the process of "distilling the most important information from a text to produce an abridged version for a particular task and user". Summaries are usually around 17% [3] of the original text and yet contain everything that could have been learned from reading the original article.

The Telugu language is the second most popular language in India just after Hindi, and it has got importance over other Indian languages as there are about 75 million native Telugu speakers. Telugu ranks fifteenth in the Ethnologue list of most-spoken languages worldwide [4]. Telugu has rich agglutinative characteristics which motivated us to consider Telugu as the topic language over other Indian languages.

The rest of this paper is organized as follows: Related work is mentioned in Sect. 2. Proposed Scheme and implementation details of the paper are presented in Sect. 3. Experimental results are shown in Sect. 4. Finally, the conclusion of the paper presented in Sect. 5.

# 2   Related Work

In recent times, for the English language, many authors suggested the procedure for automatic keyword extraction in their state-of-the-art work [1, 5, 6]. Based on previous work done toward automatic keyword extraction from the text for its summarization, extraction techniques can be categorized into four approaches, namely simple statistical approach, linguistics approach, machine learning approach, and hybrid approaches as discussed in subsequent sections. These are the approaches used for the Text Summarization in English language. In this paper, the proposed work follows a mixed approach of machine learning and statistical methods.

### 2.1 Simple Statistical Approach

These statistical methods are unprocessed, simplistic which do not require training data. They mainly focus on statistics derived from non-linguistic features of the document text such as the position of a word within the document, the term frequency, and inverse document frequency. The methods of this approach include word frequency, term frequency (TF) or term frequency-inverse document frequency (TF-IDF), word co-occurrences, and PAT-tree [5].

### 2.2 Linguistics Approach

This approach uses the linguistic features of the words in the sentences and articles. It includes the lexical analysis, syntactic analysis, discourse analysis. TreeTagger, WordNet, Electronic dictionary, N-grams, POS pattern, etc., are the primary resources of lexical analysis while Noun Phrase (NP) chunks (parsing) belong to syntactic analysis.

### 2.3 Machine Learning Approach

These approaches consider supervised or unsupervised learning from the examples, but previous work on keyword extraction prefers supervised method. The article is first converted into a graph. Each word is treated as a node, and whenever two words appear in the same sentence, the nodes are connected with an edge for each time they appear together. Then, the number of edges connecting the vertices is converted into scores and is clustered accordingly. The cluster heads are treated as keywords. Bayesian algorithms use the Bayes classifier to classify the word into two categories: keyword or not a keyword depending on how it is trained.

### 2.4 Hybrid Approach

These approaches combine any of the above two mentioned methods or use heuristics, such as position, length, layout feature of the words, HTML tags [7] around the words. These algorithms are designed to take the best features from above mentioned approaches.

# 3 Proposed Scheme

This section deals with explicit details on the approach that was used for automatic keyword extraction and summarization. The proposed algorithm follows a mixed approach of machine learning and statistical method. A Telugu POS tagger [8] was used to identify appropriate POS information in Telugu text. Then, a statistical method is used to extract keywords. In this scheme, we are not considering the stop words to summarize the news. First, it extracts the keywords by the automatic keywords extraction algorithm and then to summarize the article, the keywords are determined to summarize the article. Finally, the summarization algorithm accordingly chooses sentences to form the necessitated summary. This algorithm applies to a single article at a time.

## 3.1 Automatic Keyword Extraction

The main aim of automatic keyword extraction is to point out a set of words or phrases that best represent the document. To achieve this, a hybrid extraction technique has been proposed. The model is shown in Fig. 1.

**Keyword Annotation**: For this model, there is a need of human intervention for an annotation to train the proposed algorithm. The human annotators analyze documents and select probable keywords. These keywords are supplied to the POS tagger on the documents, and the output is provided to the next section of the model.

**Telugu POS Tagging**: In this paper, we have used a Telugu POS tagger [8] to analyze accurate POS information of any given text based on its context (relationship with adjacent and related words in a phrase, sentence, or paragraph). The Telugu POS tagger followed the Indian language standard tagset [9], which comprise 21 tags.

**Learning Probability Distribution**: Due to lack of reliable human annotators and it is a tedious job, e-newspapers clippings are used as our training dataset. The articles were considered as the target document and the headlines as the keywords,



**Fig. 1** Model for learning of probability distribution

thus eliminating the need of human annotators. The training dataset was analyzed, and the number of nouns, verbs, adverbs, adjectives, etc. that appeared as a keyword was found in the headlines. The Algorithm 1 is deployed for finding probability distribution values of keywords.

---

**Algorithm 1:** Find Probability Distribution

**Data**: $dataset$ := dataset of Telugu e-news articles
  $keyword$ := Human annotated set of keywords for each article.
**Result**: P($tag$) : Probability Distribution of Tags
$count=0$
**while** $tag\ in\ taglist$ **do**
  | Tag_count($tag$)=0
**end**
**while** $newsarticle\ in\ dataset$ **do**
  | **while** $keyword\ in\ newsarticle$ **do**
  |   | $tag$=POS_tag($keyword$)
  |   | Tag_count($tag$)=Tag_count($tag$)+1
  |   | $count=count+1$
  | **end**
**end**
**while** $tag\ in\ taglist$ **do**
  | P($tag$)=Tag_count($tag$)/$count$
**end**

---

Algorithm 1 infers P(tag) from the Telugu e-newspaper dataset. The value of the count variable is initialized to 0 which stores the number of keywords that has been scanned by the algorithm. It finds POS tag of the keyword for each keyword in every news article of the Telugu dataset and the count for that POS tag is increased by 1. Once this terminates, probability distribution, P(tag), is determined by dividing the tag count by the total number of keywords. This is used as a probabilistic measure to detect keywords.

## 3.2 Extraction

Extraction (Testing) model is shown in Fig. 2. The articles are supplied to the POS tagger on the documents. The score is calculated for each text, and few top scored texts are selected as a keyword.

**Keywords Extraction**: The output file from the POS Tagger is now forwarded to the model for extraction. Unlike TF-IDF (keeping the count of the number of times a particular word has appeared), we keep count of the word-tag pair. i.e. [Book, Noun] and [Book, Verb] are treated differently. When a count of the entire document is taken, the keywords are ranked by the Eq. 1.

**Fig. 2** Model for keywords extraction

$$Score = P(tag) * Count(word, \ tag) \tag{1}$$

where *P(tag)* is the probability of a tag being a keyword and *count(word, tag)* is the number of times the word has appeared in the current document.

---

**Algorithm 2:** Extract Keywords

---

**Data**: *doc* := Input Article
          P(Tag) := List of Trained Probabilities
          *Num_Keywords* := Required Number of Keywords
**Result**: *Keywords*[]
*pos_doc*:=pos_tagger(*doc*)
*top*:=0
**while** *word in pos_doc* **do**
    *flag*:=0
    **for** $i \leftarrow 0$ **to** *top* **do**
        **if** *word.text=wordset*[*i*].*text and word.tag=wordset*[*i*].*tag* **then**
          | *wordset*[*i*].*count*:=*wordset*[*i*].*count*+1 *flag*:=1
        **end**
    **end**
    **if** *flag=0* **then**
        *wordset*[*top* + 1].*word*:=*word.word wordset*[*top* + 1].*tag*:=*word.tag*
        *wordset*[*top* + 1].*count*:=1 *wordset*[*top* + 1].*score*:=0 *top*:=*top*+1
    **end**
**end**
**for** $i \leftarrow 0$ **to** *size* **do**
  | *wordset*[*i*].*score*:=*wordset*[*i*].*count**P(*wordset*[*i*].*tag*)
**end**
sort_desc(*wordset.score*)
**for** $i \leftarrow 0$ **to** *Num_Keywords* **do**
  | *Keywords*[*i*]:=*wordset*[*i*]
**end**

---

Algorithm 2 takes single document article, the number of keywords to be extracted and a probability distribution table trained during the training as an input for extracting keywords. The output of the algorithm will be saved in an array Keywords[]. Wordset[] is another array of structures that keep the record of the words that have already been scanned, and a number of times that word-tag pair has been scanned. The input file is fed to the POS tagger to get the POS tag values. The algorithm then courses through the file, updating existing records in the way and creating new ones when needed. When the algorithm is done with parsing the file, the scores are updated. Once the scores are set, the array is sorted according to the scores of each word-tag pair. The top score value of few texts is then extracted as keywords.

## 3.3 Summarization

With the help of algorithms explained so far, a set of word-tag pair keywords is attained and their respective scores. For summarization, the proposed algorithm suggests that one derives from many sentences for a keyword from the article as it is proportional to the score it received. It can derive these sentences by any means, be it through clustering means or crude scoring. The added advantage of this algorithm is the simple statement that "Not all keywords are equal". So it helps while selecting the keywords by differentiating them.

To see the working procedure of the proposed scheme, let us assume a single document news article on Tamil Nadu ex-CM Jayalalitha's (Amma) death. Possible keywords would be listed in Fig. 3. Assume that one need to summarize it in twenty sentences. Given the individual scores of the keywords, as shown in Fig. 3, we shall extract sentences for each of the keywords using Eq. 2. Finally, the extracted number of sentences is shown in Fig. 3 to get the desired summary of the document.

$$NS = \left\lceil \frac{(Keyword\ score * No.\ of\ sentences\ required)}{(Total\ score\ of\ all\ the\ keywords)} \right\rceil \qquad (2)$$

where NS = number of Sentences needed in summary using each keyword.

| Keyword | అమ్మ అస్తమయం | నింగికేగిన అమ్మ | మరణించిన అమ్మ | అమ్మ మహారాజ్ఞి ఇకలేరు |
|---|---|---|---|---|
| Keyword Score | 3.5 | 2 | 3 | 1.5 |
| NS | 7 | 4 | 6 | 3 |

Fig. 3 Score of each keyword and required number of sentences on each keyword

# 4    Results and Discussion

This section evaluates the quality of the keywords produced with the proposed algorithms. The section starts with article collection from different e-newspapers followed by results and discussion.

## 4.1    Article Collection

After analyzing the several e-newspapers of Telugu, we collected data from five different e-newspapers namely, Eenadu, Sakshi, Andhrajyothy, Vaartha, and Andhrabhoomi. Our dataset included almost 450 articles from each e-newspapers ranging from the 1st of October 2016 to 6th of December 2016. We have collected 150 articles with a total of 1223 keywords, 140 articles with a total of 1148 keywords, 100 articles with a total of 915 keywords, 70 articles with a total of 785 keywords, 50 articles with a total of 568 keywords from Eenadu, Sakshi, Andhrajyothy, Vaartha, and Andhrabhoomi e-newspapers, respectively (Table 1).

## 4.2    Experimental Results

In this paper, work has experimented in a machine with the following configuration:

– *System Configuration*: Intel(R) core(TM) i7-4770 CPU @ 3.40 GHz with 4 GB RAM and minimum 20 GB memory space.

**Table 1**  Number of times a tag has been found in different e-Newspapers headlines and their probability measures

| Newspaper | VM | NN | PRP | RB | INJ | JJ |
|---|---|---|---|---|---|---|
| *Tag count* | | | | | | |
| Eenadu | 350 | 245 | 255 | 152 | 154 | 76 |
| Sakshi | 317 | 185 | 250 | 195 | 78 | 85 |
| Andhrajyothy | 278 | 195 | 142 | 105 | 47 | 73 |
| Vaartha | 207 | 143 | 103 | 168 | 65 | 95 |
| Andhrabhoomi | 177 | 125 | 93 | 75 | 44 | 55 |
| *Probability measures* | | | | | | |
| Eenadu | 0.286 | 0.200 | 0.208 | 0.124 | 0.126 | 0.062 |
| Sakshi | 0.275 | 0.161 | 0.217 | 0.169 | 0.067 | 0.074 |
| Andhrajyothy | 0.303 | 0.213 | 0.155 | 0.114 | 0.051 | 0.079 |
| Vaartha | 0.263 | 0.182 | 0.131 | 0.214 | 0.082 | 0.121 |
| Andhrabhoomi | 0.311 | 0.220 | 0.163 | 0.132 | 0.077 | 0.096 |

**Table 2** Results in terms of confusion matrix, *accuracy*, *precision*, *recall*, *f-score*

| Newspaper | $T_p$ | $F_p$ | $T_n$ | $F_n$ | Accuracy | Precision | Recall | f-score |
|---|---|---|---|---|---|---|---|---|
| Eenadu | 198 | 51 | 912 | 62 | 0.907 | 0.795 | 0.761 | 0.777 |
| Sakshi | 165 | 41 | 887 | 54 | 0.917 | 0.800 | 0.753 | 0.776 |
| Andhrajyothy | 183 | 35 | 653 | 44 | 0.913 | 0.839 | 0.806 | 0.822 |
| Vaartha | 174 | 27 | 545 | 39 | 0.915 | 0.865 | 0.816 | 0.840 |
| Andhrabhoomi | 153 | 27 | 352 | 36 | 0.889 | 0.85 | 0.809 | 0.829 |

- *Operating System used*: Windows 7 Professional X 32.
- *Software Package used*: In Windows, Python Interpreter.

To evaluate the performance of proposed algorithm and compare the results with existing work, three parameters are considered, namely *precision, recall,* and *f-score* (Table 2).

*Precision* is a measure of result relevancy, while *recall* is a measure of how many truly relevant results are returned. *Precision(P)* is defined as the number of true positives over the number of true positives $(T_p)$ plus the number of false positives $(F_p)$. *Recall(R)* is defined as the number of true positives $(T_p)$ over the number of true positives $(T_p)$ plus the number of false negatives $(F_n)$.

For testing, we ran our algorithm on a set of Telugu e-newspaper articles from the above mentioned five e-newspapers. However, this time, the headlines were not provided to the algorithm. We collected the content of the article with similar article title in all five e-newspapers (same news in all five e-newspapers on the same date) to check the accuracy of proposed algorithm. The input was the same newspaper clippings, and the end target was to extract words that were present in the headlines. The result of Table 2 shows how effectively our proposed algorithm work on all five e-newspapers, and all of them attains almost similar *accuracy*, *precision*, *recall,* and *f-score* values. As of my knowledge, there is no reported work for Text Summarization in Telugu Language. So, the results are not compared to any of the work.

## 5    Conclusion

In this paper, the proposed work dealt with interdependent algorithms in keyword extraction and text summarization. The keyword extraction algorithm found the top scored keywords as efficiently as human do. With the help of keyword extraction algorithm, a summarization algorithm was proposed that introduced a concept that primary objective is "not all keywords are equal".

# References

1. Litvak, M., Last, M.: Graph-based keyword extraction for single-document summarization. In: Proceedings of the Workshop on Multi-source Multilingual Information Extraction and Summarization, pp. 17–24. ACL (2008)
2. Mani, I., Maybury, M.T.: Advances in Automatic Text Summarization, vol. 293. MIT Press, Cambridge (1999)
3. Thomas, J.R., Bharti, S.K., Babu, K.S.: Automatic keyword extraction for text summarization in e-newspapers. In: Proceedings of the International Conference on Informatics and Analytics, pp. 86–93. ACM (2016)
4. http://www.ethnologue.com/statistics/size
5. Chien, L.F.: Pat-tree-based keyword extraction for chinese information retrieval. In: ACM SIGIR Forum, vol. 31, pp. 50–58. ACM (1997)
6. Giarlo, M.J.: A comparative analysis of keyword extraction techniques (2005)
7. Humphreys, J.K.: An HTML keyphrase extractor. Department of Computer Science, University of California, Riverside, CA, USA, Technical Report (2002)
8. Reddy, S., Sharo, S.: Cross Language POS taggers (and other tools) for Indian languages an experiment with Kannada using Telugu resources. In: Proceedings of IJCNLP Workshop on Cross Lingual Information Access: Computational Linguistics and the Information Need of Multilingual Societies. Chiang Mai, Thailand (2011)
9. Bharati, A., Sangal, R., Sharma, D.M., Bai, L.: Anncorra: annotating corpora guidelines for pos and chunk annotation for indian languages. Technical Report. Technical Report (TRLTRC-31), LTRC, IIIT-Hyderabad (2006)

# Analysis and Implementation of Reliable Spectrum Sensing in OFDM Based Cognitive Radio

P. Vijayakumar, Jeswin George, S. Malarvizhi and A. Sriram

**Abstract** Cognitive radio (CR) provides better utilization of the spectrum by sensing the presence of free spectrum and assigning it to the unlicensed user temporarily when licensed users are not using the spectrum. Cyclostationary sensing method will be used here because it can differentiate between noise and the primary user (PU) signal. This technique is independent of the type of modulation used so at the receiver side no prior information on PU signal is required. Cyclostationary spectrum sensing is the preferred method for the sensing of multicarrier modulated signals, so also used for the sensing of OFDM signals and its standards. The aim of this project is to find the presence of a PU signal in a given spectrum. For this purpose probability of detection of the PU is calculated for various values of signal-to-noise ratio (SNR). Different OFDM modulated signals can also be differentiated using this method.

**Keywords** Cognitive radio · Cyclostationary · Cyclic autocorrelation function Spectrum sensing · OFDM

## 1 Introduction

The wireless communication systems need to be more flexible and more robust nowadays to meet the ever increasing demands of customers. The spectrum requirement of different technologies is different, and demand for spectrum varies

P. Vijayakumar (✉) · J. George · S. Malarvizhi · A. Sriram
SRM University, Chennai, India
e-mail: vijayakumar.p@ktr.srmuniv.ac.in

J. George
e-mail: georgejeswin22@gmail.com

S. Malarvizhi
e-mail: malarvizhi.g@ktr.srmuniv.ac.in

A. Sriram
e-mail: sriram.a@ktr.srmuniv.ac.in

from user to user which leads to a lack of spectrum in the market for different users for best possible use. The radio frequency (RF) spectrum can be divided into the unlicensed and the licensed spectrum. Licensed spectrum is the one which has been bought by a user for his use alone and hence this spectrum must not have any interference with other user's transmissions. Unlicensed spectrum is the one which has not been put up for sale by the government, and anyone can use it. So the unlicensed spectrum has the limitation of the high amount of interference from different users. But the whole of the licensed spectrum is not used most of the time [1]. The unlicensed users can exploit the unused idle portions of the spectrum known as spectrum holes provided that they have a minimum or negligible interference to the licensed user of that spectrum. The licensed user of the spectrum is known as the primary user (PU), and the unlicensed user who is trying to access the unlicensed spectrum is called the secondary user (SU).

CR technology provides the solution to the problem of spectrum scarcity. CR can be used efficiently to utilize the considered licensed frequency spectrum. CR technology hence helps in sharing the licensed spectrum with an unlicensed user in an opportunistic manner through dynamic spectrum access techniques. Spectrum management by CR comprises the following functions: Spectrum sensing, spectrum decision, spectrum sharing and spectrum mobility [2]. Spectrum Sensing is the foremost part of the CR, and it is the process of finding a spectrum hole or finding an unused spectrum. Spectrum decision is the process of finding the best frequency for use among the available frequencies obtained from spectrum sensing. Spectrum sharing is the process of coordination among transmitters to share the free spectrum. Spectrum mobility is the process of vacating the channel by the unlicensed user when the licensed user's presence is detected.

## 2 Spectrum Sensing

There are many schemes used for spectrum sensing by CR such as energy detection, matched filter detection, waveform based sensing and cyclostationary spectrum sensing [2]. The energy detection method is used mainly for high signal to noise ratio (SNR) conditions. The matched filter method is very complex and perfect information on the PU's signal is also required [3]. The waveform based sensing technique also requires a fixed pattern for the PU signal which must be known at the receiver side. Because of these constraints, cyclostationary spectrum sensing method is chosen which not only performs well under low SNR conditions but also does not need any prior information on the PU signal [4].

## 2.1 Two-State Hypotheses

Spectrum sensing can be simplified as an identification problem by modelling it as a hypothesis test [5]. So the sensing decision is made based on the two hypotheses:-

$$H0 : y(n) = w(n). \tag{1}$$

$$H1 : y(n) = x(n) + w(n). \tag{2}$$

where, '$x(n)$'is the PU signal transmitted by the PU transmitter. '$y(n)$'is the signal received by the secondary user(SU) receiver. '$w(n)$'is the additive white Gaussian noise(AWGN). $n$ is the time instant. 'H0' indicates the absence of PU signal in the spectrum and 'H1' indicates the presence of PU signal in the spectrum.

There are two important cases in this model [6]: The probability of detection ($P_d$) and the probability of false alarm ($P_f$). The probability of detection ($P_d$) is the case where PU is present in the spectrum, i.e., $P(H1/H1)$. The probability of false alarm ($P_f$) is the case where a false detection of PU has been made in the spectrum, i.e., $P(H1/H0)$.

Here $P_d$ is mainly considered because it gives the probability of correctly detecting the presence of PU in the concerned spectrum. The sensing schemes have the aim to maximize the detection probability in a spectrum for various values of SNR.

## 2.2 Cyclostationary Spectrum Sensing

Cyclostationarity spectrum sensing is the method for detecting the presence of PU in the spectrum by making use of the Cyclostationarity features of the received PU signal. Cyclostationary features occur due to the periodic properties of the signal or also in its statistics like mean and autocorrelation. These features are also deliberately added to aid in spectrum sensing. In the concerned spectrum the presence of PU is detected using the cyclic autocorrelation function (CAF). This method can also be used to differentiate between noise and PU's signal as noise is a wide-sense stationary (WSS) process with no correlation or periodicity while modulated signals are Cyclostationary with cyclic autocorrelation due to the redundancy of signal periodicities. As a result, cyclic detectors can operate even in low SNR conditions. This resulting periodic nature of signals can also be used to differentiate between different modulation schemes [7].

For a process to be cyclostationary these conditions must be satisfied [8]:

$$E\{x(n + T)\} = E\{x(n)\}. \tag{3}$$

$$\mathrm{Rx} = E\{x(n+T)x(n)\}. \tag{4}$$

From the above two equations we can infer that for a cyclostationary process both the mean and auto-correlation function Rx must be periodic with a period T.

CAF is used to detect the Cyclostationary signal present in a given spectrum. And CAF is represented in terms of Fourier coefficients as [8]:

$$R_X^{\alpha}(\tau) = \frac{1}{M} \sum_{n=0}^{M-1} x(n) \cdot x^*(\mathrm{n} - \tau)e^{-\frac{j2\pi \alpha n}{M}}. \tag{5}$$

$$\text{Cyclic frequency } \alpha = 1/T. \tag{6}$$

where, $\alpha$ is the cyclic frequency of the signal; $\tau$ is the lag parameter; $M$ is total number of symbols; $T$ is symbol duration; $x(n)$ is the symbol at time instant $n$.

The CAF depicts the periodicity in the signal using the lag parameter $\tau$ and cyclic frequency $\alpha$. The cyclic frequency $\alpha$ is the reciprocal of the OFDM symbol duration. This method is suitable for sensing of OFDM based signals as it can take advantage of the cyclic prefix (CP). The CP is the last portion of the data block copied to the front of the block. In addition to combating the problem of inter-symbol interference (ISI), CP can also aid in cyclostationary spectrum sensing. This method takes advantage of the high autocorrelation between the data part and the CP of the OFDM symbol [9].

## 3   Implementation Description

To implement cyclostationary spectrum sensing the simulation environment was created using the following parameters as shown in Table 1.

To implement sensing PU transmitter is modeled which transmits 16-QAM OFDM modulated signal as PU signal through the channel. And a corresponding SU receiver is also modeled at the receiver side, which detects whether the received signal is that of PU or not.

**Table 1** Parameters of 802.11a used in simulation

| Alpha $\alpha$ | 0.25 MHz |
|---|---|
| CP length | 12 |
| FFT | 64 |
| Number of (data + pilot) subcarriers | 52 |
| $M$ (for QAM modulation) | 16 |
| Number of iterations | 1000 |
| Lag parameter $\tau$ | 52 |
| OFDM symbol duration | 4 μs |
| Channel used | AWGN channel |

Figure 1 shows the steps taken in this work to implement cyclostationary spectrum sensing. To model the PU signal, at the transmitter side the information bits are 16-QAM modulated. The modulated serial data is converted to parallel data, and inverse fast Fourier transform (IFFT) operation is done on the modulated complex data. The transformed result is arranged serially. Then cyclic prefix insertion is carried out on the serial data by prepending data from end to the front. The threshold for this simulation set-up is set by finding the average of maximum and minimum values of the CAF corresponding to the OFDM signal hence obtained and which is further used as PU signal. The OFDM signal obtained is passed through an AWGN channel and received at the SU receiver.

At the SU receiver, the received signal's CAF is calculated using the parameters shown in Table 1. Then maximum value of cyclic autocorrelation function (CAF) is compared with the previously calculated threshold value. If the maximum value of cyclic autocorrelation function is greater than that of the threshold value, then it means that the signal belongs to the primary user (PU) otherwise it is not the PU signal. This method is effective against noise signals because noise signals do not have any periodicity.

## 4   Simulation Results

Figure 2 shows a comparison between CAF of 16-QAM OFDM signal and BPSK OFDM signal. This simulation can be used to prove that cyclostationary method can be used to differentiate between two different modulated OFDM signals.

Figure 3 shows the variation of the maximum value of CAF with respect to SNR (in dB) of 16-QAM OFDM signal and BPSK OFDM signal. Maximum CAF value of BPSK OFDM signal is lower as compared to that of 16-QAM OFDM signal



Fig. 1  Implementation diagram for cyclostationary spectrum sensing

**Fig. 2** Comparison of CAF of 16-QAM OFDM signal and BPSK OFDM signal



**Fig. 3** Variation of maximum value of CAF with respect to SNR (in dB)

because BPSK signals are antipodal signals whose values are either +1 or −1 so its signal strength is less than that of 16 QAM signals which are complex signals.

Figures 4 and 5 show the variation of the probability of detection of 16-QAM OFDM signal and the variation of the probability of detection $P_d$ of BPSK OFDM

**Fig. 4** Probability of detection of 16 QAM OFDM signal with respect to SNR (in dB)



**Fig. 5** Probability of detection of BPSK OFDM signal with respect to SNR (in dB)

signal in the spectrum with respect to SNR. It is inferred that the probability of detection increases with the increase in SNR.

## 5 Conclusion

This paper realizes the spectrum sensing of OFDM based PU signals. The simulation results show that the Cyclostationary spectrum sensing method implemented here can be used for spectrum sensing without any prior information on the modulated PU signal and also used for detection in low SNR conditions. These results can be used to show that why in this project Cyclostationary spectrum sensing technique is preferred over energy detection method, matched filter method, and waveform based sensing. The simulation for the probability of detection for 16-QAM OFDM modulated signal and BPSK OFDM modulated with respect to SNR(in dB) using Cyclostationary spectrum sensing technique has been shown. This method was also used to differentiate between different modulated OFDM signals. This work will be further extended to the hardware implementation of cyclostationary spectrum sensing and also for sensing in VANET.

## References

1. Akyildiz, I.F., Lee, W.-Y., Vuran, M.C., Mohanty, S.: NeXt generation/dynamic spectrum access/cognitive radio wireless networks: a survey. Comput. Netw. **50**, 2127–2159 (2006)
2. Yucek, T., Arslan, Hu: A survey of spectrum sensing algorithms for cognitive radio applications. IEEE Commun. Surv. Tutor. **11**, 116–130 (2009)
3. Anupama, Er., Vipin G.: Spectrum sensing: first step towards nurturing cognitive radio networks. IJSRD Int. J. Sci. Res. Dev. **4**(02), 539–542 (2016)
4. Jain, S., Raghunadh, M.V.: Spectrum sensing based on cyclostationary detector using USRP. Int. J. Eng. Trends Technol. (IJETT) **11**(4), 188–191 (2014)
5. Liang, Y.-C., Zeng, Y., Peh, E.C.Y., Hoang, A.T.: Sensing-throughput tradeoff for cognitive radio networks. IEEE Trans. Wirel. Commun. **7**(4),1326–1337 (2008)
6. Turunen, V., Kosunen, M., Huttunen, A., Kallioinen, S., Ikonen, P., Pärssinen, A., Ryynänen, J.: Implementation of cyclostationary feature detector for cognitive radios. In: Proceedings of the 4th International Conference on CROWNCOM (2009)
7. Mohapatra, S.G., Mohapatra, A.G., Lenka, S.K.: Performance evaluation of cyclostationary based spectrum sensing in cognitive radio network. In: International Multi-conference on Automation, Computing, Control and Compressed sensing (2013)
8. Syrjala, V., Valkama, M., Allen, M., Yamamoto, K.: Simultaneous transmission and spectrum sensing in OFDM systems using full-duplex radios. In: IEEE 82nd Vehicular Technology Conference (VTC Fall), pp. 6–9 (2015)
9. Turunen, V., Kosunen, M., Vaarakangas, M., Ryynänen, J.: Correlation-based detection of OFDM signals in the angular domain. IEEE Trans. Veh. Technol. **61**(3) 951–958 (2012)

# List Colouring of Graphs Using a Genetic Algorithm

**Aditi Khandelwal, Pallavi Jain and Gur Saran**

**Abstract**  The List Colouring Problem (LCP) is an NP-Hard optimization problem in which the goal is to find a proper list colouring of a graph $G = (V, E)$ such that the number of colours used is minimized. This paper deals with the development and implementation of a Genetic Algorithm for the list colouring of graphs. Initial solutions are generated using two heuristics in equal proportion. Experiments have been carried out on randomly generated planar graphs with a list of colours for each vertex. Results for 32 graphs with size varying from 10 to 750 show that the proposed Genetic Algorithm implementation tends towards lower values of the number of colours used.

**Keywords**  List colouring · Genetic Algorithm

## 1    Introduction

Let $C$ be a set of colours, and for each $v \epsilon V(G)$, let $L: V(G) \rightarrow 2^c$ be a function assigning to each vertex $v \epsilon V(G)$ a list of colours $L(v) \subseteq C$. If $\exists$ a function f: $V(G) \rightarrow C$ such that $f(v) \epsilon L(v)$ for all $v \epsilon V(G)$ and $f(u) \neq f(v)$ for $uv \epsilon E(G)$, then $G$ is said to be *L-colourable* [1]. In LCP, a vertex colouring of an *L-colourable* graph $G$ is an assignment of colour to each vertex $v_i \epsilon V(G)$ from its list assignment $L_i$ so as to obtain a proper colouring of $G$ with the least number of colours [1]. LCP is an NP-Hard problem [1]. In this paper, a Genetic Algorithm is developed to optimize the number of colours being used. A large number of practical life problems can be formulated as LCP [1, 2].

A. Khandelwal (✉) · P. Jain · G. Saran
Department of Mathematics, Dayalbagh Educational Institute, Agra, India
e-mail: aditikhandelwal8193@gmail.com

P. Jain
e-mail: pallavijain.t.cms@gmail.com

G. Saran
e-mail: gursaran@dei.ac.in

In 1997, Dimitris and Michael [3] analysed a greedy list colouring algorithm that successfully finds a list colouring of a graph $G$ if no list $L_v$ assigned to vertex $v \in (V(G))$ is empty. They have also proved that almost all graphs with $n$ vertices and 1.923 $n$ edges are 3-colourable. Maya and Constantin [4], in 2005, have designed a heuristic for list colouring of graphs. The colour maximally applicable is chosen and is assigned to vertices without violating the requirements of list colouring. The coloured vertices and the colours used are removed, and the same procedure continues till all vertices are coloured.

## 1.1  Organization

The rest of the paper is organized as follows. Section 1.2 describes LCP. In Sect. 2, heuristics for generating initial population are described. Implementation details of LCP are presented in Sect. 3. Section 4 is devoted to computational experiments, and the results followed by conclusions in Sect. 5.

## 1.2  Genetic Algorithm for List Colouring Problem

Genetic Algorithm (GA) mimics the process of natural evolution to generate solutions to optimization problems. The basic GA is described in [5], its adaptation to the List Colouring Problem is described in Fig. 1 and its implementation details are presented in Sect. 3. In the context of LCP, two construction heuristics are proposed to generate initial population. These are described in Sect. 2.

## 2  Heuristics for Generating Initial Population

In this section, we describe two construction heuristics that generate a feasible initial population. The first heuristic *H1* generates solutions by the hybridization of Welsh Powell Algorithm and random assignment from their respective lists while *H2* by the hybridization of Welsh Powell Algorithm and Greedy Colouring Algorithm. In the implementation of the GA for LCP, solutions in the initial population are generated using both the heuristics in equal proportion.

A. **Heuristic 1 (*H1*)** *H1* is a heuristic which begins with arranging the vertices in descending order of their degrees. Then vertices with the same degrees are selected in an array, and their random permutation is taken. From the permuted list of vertices, each vertex is coloured randomly from its respective list. Figure 2 presents the implementation of *H1*.

---

**Pseudocode of List Colouring Problem (LCP)**

---

**Step 1:** Initialize population size *ps*, crossover rate *pc,* mutation probability *pm* termination criteria *tc*.

**Step 2:** Generate initial population *pop* of size *ps*.

**Step 3:** *bestcost* =least *number of colours* used in solutions generated in Step 2.

**Step 4:** **while** termination criteria *tc* not met

**Step 5:** Apply binary tournament operator on *pop* to obtain Intermediate population *interPop*.

**Step 6:** Among remaining individuals in *interPop,* crossover operator is applied and the solution obtained goes to *Check* and *Repair* to maintain proper colouring , giving *childPop*.

**Step 7:** Apply mutation operator on *childPop* with probability *pm* on each individual, further sending each to *Check* and *Repair* to maintain proper colouring , giving *NewPop*.

**Step 8:** *pop = NewPop*

**Step 9:** Update *bestcost* if the least cost solution in *pop* is smaller than current *bestcost*.

**Step 10:** **end while**

---

**Fig. 1** Pseudocode of LCP

---

**Procedure *H1:* colouring (vertex)**

---

**Step 1:** *degree= {d_G(u)* such that *u ε V(G)}*

**Step 2:** sort *degree* in descending order

**Step 3:** **if** $v_c = \phi$ for every *v ε V(G)*

**Step 4:** *colour=* random colour to the *v*

**Step 5:** update *done*(matrix storing the coloured vertices with their corresponding colour*)*

**Step 6:** **else**

**Step 7:** **for** *j*= 1 to *|N(v)|*

**Step 8:** *c=v_c* for every *v ε N(v)*

**Step 9:** **end for**

**Step 10:** *l= l_v - c*

**Step 11:** **if** *l =*$\phi$

**Step 12:** *l= l_v*

**Step 13:** *colour=*random colour from *l*

**Step 14:** update *done*

**Step 15:** *Repair (vertex, colour)*

**Step 16:** **else**

**Step 17:** *color=*random colour from *l*

**Step 18:** **end if else**

**Step 17:** **end if else**

---

**Fig. 2** Heuristic *H1*

B. **Heuristic 2 (*H2*)** *H2* is a heuristic which begins with arranging the vertices in descending order of their degrees. Then vertices with the same degrees are selected in an array, and their random permutation is taken. From the permuted list of vertices, each is coloured greedily. Figure 3 presents the implementation of heuristic *H2*.

## 3 Implementation Details of LCP

In this section, we describe the implementation of the algorithm presented in Fig. 1.

### 3.1 Solution Representation and Initial Population

Each individual (solution) in the initial population represents the colours assigned to vertices of the graph *G* satisfying proper colouring. The initial population with *ps* individuals is generated using each of the heuristics *H1* and *H2* in equal proportion.

| | |
|---|---|
| **Procedure *H2*** *greedycolouring(v)* | |
| **Step 1:** | *degree* = $\{d_G(u)$ such that $u \in V(G)\}$ |
| **Step 2:** | sort *degree* in descending order |
| **Step 3:** | **if** $v_c = \phi$ for every $v \in V(G)$ |
| **Step 4:** | *colour* = random colour to the *v* |
| **Step 5:** | update *done* |
| **Step 6:** | **else** |
| **Step 7:** | **for** *j* = 1 to $|N(v)|$ |
| **Step 8:** | $c = v_c$ for every $v \in N(v)$ |
| **Step 9:** | **end for** |
| **Step 10:** | $l = l_v - c$ |
| **Step 11:** | **if** $l = \phi$ |
| **Step 12:** | $l = l_v$ |
| **Step 13:** | *colour* = random colour from *l* |
| **Step 14:** | update *done* |
| **Step 15:** | *Repair (vertex, colour)* |
| **Step 16:** | **else** |
| **Step 17:** | *greedyc* = colours used ∩ *l* |
| **Step 18:** | **if** *greedc* = $\phi$ |
| **Step 19:** | *greedyc* = *l* |
| **Step 20:** | **end if** |
| **Step 21:** | *colour* = random colour form *greedyc* |
| **Step 22:** | update *done* |
| **Step 23:** | sort *done* in ascending order of $v \in V(G)$ |
| **Step 24:** | **end if** |
| **Step 25:** | **end if** |

**Fig. 3** Heuristic *H2*

## 3.2  *Evaluation*

The function evaluates the fitness of each individual which is the number of colours used.

## 3.3  *Selection*

Binary tournament operation is applied on the population for selection. In this way, any individual of initial population can have at most two copies in the intermediate population.

## 3.4  *Crossover*

Two-point crossover is used with probability $pc$ taken as 0.25. A *Check* function is applied after the crossover, and *Repair* function is applied if necessary.

## 3.5  **Check** *Function*

This function checks

 (i)  That each individual satisfies proper colouring.
(ii)  Colour assigned to each vertex is available in its respective list.

## 3.6  *Repair of Solutions*

The *Repair* function repairs colouring of neighbouring vertices if proper colouring is not achieved.

## 3.7  *Mutation*

Mutation operator, with mutation rate $pm$ taken as 0.05, recolours a vertex randomly. A *Check* function is applied after the mutation, and *Repair* function is applied if necessary.

## 3.8   Termination Criterion

The above GA terminates if there is no improvement in solutions for 50 generations.

## 4   Experiments and Results

### 4.1   Test Set

The test set of planar random graphs was generated using Mathematica 10.3 with number of vertices ranging from 10 to 750. Experiments for population size and convergence have been carried out on these graphs. The graphs were stored with their list assignment in Matlab files using the following nomenclature: for planar graph with 10 vertices, the notation is p10 (Table 4).

### 4.2   Pilot Experiment for Determining Population Size

Preliminary set of graphs p25, p50, p100, p150, p200, p250, p300, p325, p400, p450, p679, p700 and p750 was chosen randomly to test for population size. Thirty trials were carried out on each graph for each population size 30, 60 and 90 (mean values listed in Table 1). Two-way ANOVA with replication was applied. The $p$ value of columns in Table 2 shows that there was significant difference between the means.

TUKEY's HSD test was then used to compare the means between two population sizes for these instances. Table 3 shows that there is significant pairwise difference between the population size 60 and 90.

A population size 30 is chosen in the interest of performance.

**Table 1**   Mean values for different population sizes

| Population size | 30 | 60 | 90 |
|---|---|---|---|
| Mean | 5.935897 | 5.792308 | 5.74359 |

**Table 2**   Results of TWO-WAY ANOVA

| Source of variation | SS | $df$ | MS | $F$ | $p$ value | $F$ critical |
|---|---|---|---|---|---|---|
| Sample | 1621.93 | 12 | 135.1608 | 1329.277 | 0 | 1.760 |
| Columns | 7.796581 | 2 | 3.898291 | 38.33884 | 1.11E-16 | 3.003 |
| Interaction | 3.003419 | 24 | 0.125142 | 1.230749 | 0.203964 | 1.526 |
| Within | 115 | 1131 | 0.10168 | | | |

**Table 3** Results of TUKEY's HSD

| Test between populations | 30 and 60 | 30 and 90 | 60 and 90 |
|---|---|---|---|
| TUKEY's HSD value | 11.90999 | 8.892791 | 3.017197 |

## *4.3 Final Experiment*

The final experiments were carried out with *ps* = 30 for 30 trials with input list instances generated as described in Sect. 4.3.1.

### Generation of List of Colours on Each Vertex

This procedure is used to generate list of colours available on vertices of each graph. For an input graph *G*, upper bound and lower bound to the number of available colours in its list are given. Then a colour is chosen randomly from a pre-assigned list of available colours. That colour is assigned to random percentage of vertex lists. This procedure is followed until all colours are assigned or upper bound is reached.

Table 4 details the least number of colours used and the average time taken to generate the result for the various test graphs.

**Table 4** Final results for population size 30

| Graph instances | Least number of colours | Average time (s) |
|---|---|---|
| p10 | 3 | 8.535033 |
| p25 | 3 | 19.7271 |
| p50 | 3 | 39.78613 |
| p100 | 5 | 73.9514 |
| p125 | 5 | 99.34663 |
| p150 | 4 | 109.6963 |
| p175 | 5 | 147.92 |
| p200 | 5 | 153.1717 |
| p225 | 5 | 186.06 |
| p250 | 5 | 241.79173 |
| p275 | 6 | 227.3277 |
| p300 | 6 | 307.065 |
| p325 | 5 | 317.20783 |
| p350 | 5 | 313.2547 |
| p375 | 6 | 78.41407 |
| p394 | 6 | 5261.765 |

(continued)

**Table 4** (continued)

| Graph instances | Least number of colours | Average time (s) |
|---|---|---|
| p425 | 6 | 5254.094 |
| p450 | 6 | 102.522 |
| p475 | 6 | 89.1745 |
| p500 | 6 | 99.606 |
| p525 | 6 | 125.75 |
| p550 | 7 | 105.486 |
| p575 | 7 | 5287.6 |
| p591 | 6 | 143.9307 |
| p600 | 6 | 119.9597 |
| p625 | 6 | 157.167 |
| p650 | 6 | 203.363 |
| p675 | 7 | 249.71 |
| p679 | 7 | 172.825 |
| p700 | 6 | 166.3597 |
| p725 | 6 | 251.437 |
| p750 | 7 | 157.775 |

## Convergence of Solutions

Figure 4 shows that for larger and smaller graphs, the population mean tends to converge after 10–15 generations. Figure 5 shows that the least number of colours used converges after approximately ten generations.



**Fig. 4** Mean of the number of colours assigned in each generation with population size 30

**Fig. 5** Minimum number of vertices in each generation with population size 30



## 5 Conclusion

In this paper, we have implemented a GA for the List Colouring Problem. Two heuristics were used to generate the initial population. Preliminary experiments were performed to select the population size based on which the population size of 30 was selected for the final experiments.

Results show that the proposed GA implementation trends towards lower values of the number of colours used. The optimality of these results has not been tested in this paper as the optimal values for the test graphs could not be obtained. As a future extension, local search operators can be incorporated in each generation step to improve the solutions. Further, for the initial solutions, other heuristics can be designed to obtain possibly better results.

## References

1. Baber, C.L.: An Introduction to List Colorings of Graphs, M.Sc. thesis, Faculty of the Virginia Polytechnic Institute and State University (2009)
2. Laurenta, B., Hao, J.K.: List graph colouring for multiple depot vehicle scheduling. Int. J. Math. Oper. Res. **1**(1/2) (2009)
3. Achlioptas, D., Mollo, M.: The Analysis of a List-Coloring Algorithm on a Random Graph (1997)
4. Satratzemi, M., Tsours, C.: A heuristic algorithm for the list coloring of a random graph. In: The 7th Balkan Conference on Operational Research. Constanta, Romania (2005)
5. Deb, K.: Multiple-Objective Optimization Using Evolutionary Algorithms. Wiley, New York (2008)

# Waveform Generation and Reception of IEEE 802.11p Standard for CR-VANET Application

**Ponnusamy Vijayakumar, Ravi Ranjan and S. Malarvizhi**

**Abstract** Vehicular ad hoc network (VANET) is a technique used to improve road safety and reduce traffic jams. IEEE standardized VANET as IEEE 802.11p and vehicular environment having capacity and capability to support Intelligent transport system (ITS). The purpose of work was to simulate IEEE 802.11p waveforms using MATLAB and AWR Virtual system simulator (VSS). The performance of the generated waveform is analyzed by using cancellation plot and Bit error rate (BER) graph. The main contribution of this work is to provide co-simulation environment of 802.11p waveform generation and reception for VANET application on MATLAB and AWR tool. This work also incorporated the multiband transmission by spectrum sensing to enable cognitive radio-enabled 802.11p waveform generation. Since the allocated band for the VANET application is not sufficient for the high dense environment like in urban city especially for safety application, multiband CR-enabled 802.11p waveform generation is required.

**Keywords** OFDM · Intelligent transport system (ITS) · Vehicular adhoc network (VANET) · Vehicle-to-vehicle (V2V) · Vehicle-to-infrastructure (V2I) Bit error rate (BER)

## 1 Introduction

The vehicular communication system is a type of vehicular communication network where communication exists between vehicles and various roadside units. This type of communication is called Dedicated Short-Range Communication (DSRC) which

P. Vijayakumar (✉) · R. Ranjan (✉) · S. Malarvizhi
SRM University, Chennai, India
e-mail: vijayakumar.p@ktr.srmuniv.ac.in

R. Ranjan
e-mail: raviranjan549@gmail.com

S. Malarvizhi
e-mail: malarvizhi.g@ktr.srmuniv.ac.in

exchange the instant warning messages and the traffic information like vehicles position. [1]. VANET-based IEEE 802.11p was standardized by IEEE for the vehicular communication network. Intelligent transport systems (ITS) used VANET for their application [2]. IEEE 802.11p provides operability at high user mobility as compared to IEEE 802.11a to support vehicular communication. In this work, an IEEE 802.11a PHY layer modified its parameter to achieve an 802.11p PHY layer model [3]. The bandwidth of the signal decreased to 10 MHz in IEEE 802.11p compared to IEEE 802.11a signal bandwidth of 20 MHz which makes the communication more effective and efficient for high user mobility and fading channel such as reducing intersymbol interference due to the multipath channel with improved guard interval by double [4, 5]. OFDM technique is used for IEEE 802.11p standard PHY layer for advanced vehicular communication devices. Study the specification of an IEEE 802.11p PHY layer model operated in high mobility, and more multipath environments have potentially improved [6]. Simulations result of Bit Error Rate (BER) performance at various Signals-to-Noise Ratio (SNR) on the different mobility level is analyzed in the literature. They demonstrate that rearranging pilot pattern can offer better results which provide the best performances by adding pilot's symbols. Channel estimation technique such as Least square (LS) and minimum mean square error (MMSE) [7]. Doppler shift is caused by high velocity in vehicular communication. Doppler shift results in variation in carrier frequency offset that induces inter-carrier interference in OFDM systems [8, 9]. Coexistence issues between digital video broadcasting terrestrial (DVB-T2) and IEEE 802.11p transmission in the TV white spaces (TVWS) in the UHF band are discussed. The main outcomes of that work are the transmitter power maximization and signal bandwidth of an 802.11p signal waveform in the neighboring channel of a present active DVB-T2 system at the same time protecting the error rate graph for a different mode of IEEE 802.11p in the white space TV band [10, 11]. CR-VANET improves the overall performance of the new emerging future vehicular applications such as infotainment and public safety communication [12].

## 2 IEEE 802.11p Waveform Design

Orthogonal frequency division multiplexing (OFDM)-based multiplexing technique is implanted for VANET application.

Figure 1 provides the block diagram. Random bits are taken as information bits for a single user, representing multiple concatenated Physical layer convergence protocol (PLCP) service data units (PSDUs), specified as a binary vector stream. When transmitted, the periodic signal or non-random data components of the signal cause spectrum peaks that are larger than average power. These strong discrete frequency components can cause interference, and long sequences of certain values cause a problem in clock recovery. So, scrambler is used to randomly rearrange the values of information bits in a data block to avoid the above issue. Scrambled data feed into forward error correction (FEC) block to perform convolutional encoding,

and it is carried out for coding the scrambled data block. Initial scrambler state of the data scrambler for each packet generated is specified as an integer from 1 to 127, or as an $N_P*N_{users}$) matrix of integers with values from 1 to 127. $N_P$ is the number of packets, and $N_{users}$ is the number of users. Next operation can be interleaving employed to mitigate the effect of burst errors. In this way, serial to parallel block is used to generate N serial data symbols into corresponding N different parallel symbols. And then constellation mapper is used to map bits to the symbol. The adaptive modulation is used with various modulation scheme and code rate (1/2 or 3/4). In this, each OFDM symbol has allotted four pilot subcarriers. Pilot signals are used for frequency offset tracking and calculation of phase noise. The position of pilots set as −21, −7, 7, and 21 in overall OFDM subcarrier. At each time, duration user data are modulated with different orthogonal subcarriers and IFFT operation is performed to create OFDM symbol. The operation of IFFT block shows the allocation of the subcarrier each one orthogonal to each other and no interference exists between adjacent subcarriers and subchannel. Another benefit of IFFT is that it does not require N oscillators for N subcarrier transmission. The beginning of each OFDM symbol prepended with a copy of end called as a Cyclic prefix (CP). The CP is added after the IFFT operation at the transmitter, and then CP is removed in order to get the original signals back at the receiver successfully. This guard interval is used to provide protection against the effect of ISI. Parallel to serial block is used to convert N parallel data symbols into corresponding N serial data symbols and passed through the fading channel. At receiver side, serial data stream is converted into N parallel data symbol to perform all reverse operation and recover information data stream. Parallel data are passed out from serial to parallel block and the CP are removed. Channel estimation is used to find channel condition and corresponding parameter to implantation of adaptive modulation. Equalization done in frequency domain is used to transmit signal recovery; zero-forcing equalizer performs calculation of the inverse of the channel response in frequency domain. The signal is passed through fading channel. The output of zero-forcing equalizer produces equalizes the linear phase distortions to enable to recover the transmitted signal. It is followed by constellation de-mapper, which produces output as the sequence of a bit stream of corresponding symbol stream and same scheme is used for constellation de-mapper and mapper. Again parallel to serial block is used to produce serial output for decoding and descrambling. Its working principle is same as the maximum-likelihood decoding. At the final stage, data stream block provides information bit sequence.

Figure 2 describes PPDU frame structure. In a frame structure, it consists of two-part preamble and data. Preamble is subdivided into training sequence and a Signal field (SIG). Training sequence is further divided into the long and short type called as LTS and STS. The Legacy long training field (LLTF) is the second field in the 802.11P. Channel estimation, frequency offset estimation, and time synchronization rely on the LLTF. The Legacy signal (LSIG) field is the third field of the 802.11P OFDM symbol legacy preamble. Service field contains 16 zeros to initialize the data scrambler. PSDU consists of variable-length field containing the PLCP service data unit (PSDU). Tail bits are required to terminate a convolutional

Fig. 1  Transceiver system block diagram



Fig. 2  PPDU frame structure for IEEE 802.11p

code. The field uses six zeros for the single encoding stream. Pad bits have variable-length field to ensure that the data field contains an integer number of symbols. The long OFDM training symbol $L_k$ consists of 52 subcarriers. $\Delta f$ represents carrier frequency spacing and $T_{GI2}$ represents second guard interval time duration. $N_P$ is the number of packets. The LLTF is two OFDM symbol long and follows the LSTF of the preamble in the packet structure for the entire frame.

$$\text{LLTF} = \sum_{k=-Np/2}^{Np/2} L_k e^{-j2\pi k\Delta f(t-T_{GI2})}. \tag{1}$$

The sequence $S_k$ uses 12 subcarriers as pilot out of available 52 subcarriers per 10-MHz channel bandwidth segment. The Legacy short training field (LSTF) frame structure having subframe shown in Fig. 3 is the first field of the 802.11P OFDM symbol. It is for the start of packet detection, for coarse frequency correction, and for setting the AGC.

$$\text{LSTF} = \sum_{k=-Np/2}^{Np/2} S_k e^{-j2\pi k\Delta ft}. \tag{2}$$

**Fig. 3** Block diagram for IEEE 802.11p waveform design

In a communication system, fading is generated due to the propagation through multipath propagation in vehicular communication systems. Delay profile model is used in this work, where parameters considered are as follows: breakpoint distance is 10 m, RMS delay spread ($T_{RMS}(ns)$) is 50, maximum delay ($T_S(ns)$) is 450 and corresponding frequency is denoted as $f_s$, Rician K factor is equal to 6, number of clusters is 6, number of taps is 36, and carrier frequency (GHz) is taken as 5.9. The first tap is considered as Line of sight (LOS) or dominant link between the transmitter and receiver, whereas all another taps are Nonline of sight (NLOS). As a result, the first tap exhibits Rician behavior, while the others exhibit Rayleigh behavior. Receiver sensitivity is defined as minimum input signal required to reproduce output signal at a specified SNR level. Greater Rx sensitivity shows that device is able to perform weak signal reception, which results in more transmission distances that can be supported.

# 3   AWR Simulation of IEEE 802.11p

The IEEE 802.11p PHY layer has similar parameter specifications as specified in IEEE 802.11a with some improvisation. Increase in guard interval gives the better performance against ISI.

Figure 3 provides the block diagram of IEEE 802.11p waveform design and analysis as described in Sect. 1. In this work, all parameters setup are described in Table 1 according to IEEE 802.11p specification.

## 4   Results and Discussion

The purpose of work was to simulate IEEE 802.11p waveforms using MATLAB and AWR Virtual system simulator (VSS). The performance of WAVE depends on the Physical layer (PHY). An IEEE 802.11a PHY is modified to achieve 802.11p PHY model and better communication performance.

Figure 4 shows the IQ plot of transmitted and received signal for MCS type 7 which uses of 64 QAM modulation scheme.

Figure 5 provides the transmitted power spectrum of OFDM-based IEEE 802.11p. Result shows that power level ranges from −1.859 to −6.759 dB around bandwidth of 8 MHz in place of total transmitted bandwidth of 10 MHz. Spectral leakage per 5 MHz band is −11.7 dB.

Figure 6 provides the plot of BER w.r.t receiver sensitivity. The OFDM modulator attenuates before going to power amplifier. 5-dB fading margin and 10-dB fading margin are assumed for two different levels of fading occurrence in the channel with attenuation level of 13.5 and 20 dBm, respectively. The reference curves are shown in BER graph with transmitter used without power amplifier and BER results with power amplifier loading at approximately 5 and 10 dB fading margin. Evaluation of the receiver sensitivity depends on the signal strength required to achieve the BER for vehicular communication. The receiver sensitivity is 76.5 and −70.6 dB for 10 and 5 dB fading margin, respectively, for BER of 1e-005.

**Table 1** Specification difference between IEEE 802.11p and IEEE 802.11a standard

| Parameter | 802.11p | 802.11a |
|---|---|---|
| $N$: Total no. of subcarrier | 64 | 64 |
| $N_{SD}$: no. of subcarrier | 48 | 48 |
| $N_{SP}$: no. of pilot subcarrier | 4 | 4 |
| $N_{ST}$: no. of total subcarrier | 52 | 52 |
| $N_{STF}$ : no. of subcarrier for STF | 12 | 12 |
| $\Delta F$: Frq. spacing between subcarrier (MHz) | 0.15625 | 0.3125 |
| $T_{SIGNAL}$: OFDM symbol period (μs) | 8.0 | 4.0 |
| $T_{SYM}$: symbol interval (μs) | 8 | 4 |
| Total time slot duration (μs) | 13 | 9 |

Fig. 4 IQ plot of IEEE
802.11p waveform



Fig. 5 Transmitter power spectrums

**Fig. 6** BER versus receiver sensitivity

## 5 Conclusion

Co-simulation of the 802.11p waveform on the cross-platform of MATLAB and AWR tool is done for providing a co-simulation-based VANET research and application development. The power spectrum analysis of the simulated waveform shows that the waveform is generated within the allocated bandwidth without much spectral leakage.

From the BER analysis, it is concluded that even at low transmit power, we can achieve good performance of around 10–5. Here, 70% of work is simulated in AWR and remaining 30% is developed in MATLAB. As a future work, we are working on hardware in loop simulation of the 802.11p and analysis of the performance of waveforms in a mobile environment.

## References

1. Standard Specification for Telecommunications and Information Exchange Between Roadside and Vehicle Systems—5 GHz Band Dedicated Short Range Communications (DSRC)

Medium Access Control (MAC) and Physical Layer (PHY) Specifications, ASTM DSRC STD E2313–02 (2002)

2. Molisch, A., Tufvesson, F., Karedal, J., Mecklenbrauker, C.: A survey on vehicle-to-vehicle propagation channels. IEEE Wirel. Commun. **16**(6), 12–22 (2009)
3. IEEE Std. 802.11-2007, Part 11: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications, IEEE Std. 802.11 (2007)
4. Hartenstein, H., Laberteaux, K.P.: VANET-Vehicular Applications and Inter-Networking Technologies. Wiley, New York (2010)
5. Intelligent Transport Systems (ITS); Access layer specification for Intelligent Transport Systems Operating in the 5 GHz Frequency band, ETSI EN 302 663V1.2.1, Aug 2016
6. Dietzel, S., Petit, J., Kargl, F., Scheuermann, B.: In-network aggregation for vehicular ad hoc networks. IEEE Commun. Surv. Tutor. **16**(4), 1909–1932 (2014)
7. Sassi, A., Elhillali, Y., Rivenq, A.: A symbol-based estimation technique for inter-vehicular communication performance optimization. IJCSI Int. J. Comput. Sci. Issues **10**(2), 3 (2014)
8. Adeyemo, Z.K., Akande, D.O., Ojo, F.K., Raji, H.O.: Comparative evaluation of fading channel model selection for mobile wireless transmission system. Int. J. Wirel. Mob. Netw. (IJWMN) **4**(6), 127–138 (2012)
9. Feldman, J., Abou-Faycal, I., Frigo, M.: A fast maximum-likelihood decoder for convolutional codes. IEEE Veh. Technol. Conf. **1**, 371–375 (2002)
10. Fadda, M., Murroni, M., Popescu, V.: Interference issues for VANET communications in the TVWS in urban environments. IEEE Trans. Veh. Technol. **65**(7), 4952–4958 (2016)
11. Chen, J., Liu, B., Zhou, H., Wu, Y., Gui, L.: When vehicles meet TV white space: a QoS guaranteed dynamic spectrum access approach for VANET. In: IEEE International Symposium on Broadband Multimedia Systems and Broadcasting (BMSB), pp. 1–6. Beijing, China, 25–27 June 2014
12. Singh, K.D., Rawat, P., Bonnin, J.-M.: Cognitive radio for vehicular adhoc networks (CR-VANETs): approaches and challenges. EURASIP J. Wirel. Commun. Netw. **2014**, 49 (2014)

# Mutual Correlation-based Optimal Slicing for Preserving Privacy in Data Publishing

K. Ashoka and B. Poornima

**Abstract** Privacy preservation is a substantial concern for the organizations that publish/share personal data for informal analysis. Several anonymization algorithms such as generalization and Bucketization are developed as a solution to this Privacy Preserving Data Publishing (PPDP). Latest research has shown that generalization loses significant amount of information, particularly for high dimensional data. However, Bucketization does not prevent membership disclosure. In this paper, we propose a novel approach that makes use of Information Gain of the attributes with respect to sensitive attributes, which gives the effectiveness of an attribute in classifying the data, which is two-way association among attributes. We show that our approach preserves better data utility and has lesser complexity than earlier techniques. Our proposed technique is theoretically analyzed, and mathematical analysis outstrips past works with sufficient experiments.

**Keywords** Privacy preserving data · Data anonymization · Data perturbation Data utility

## 1 Introduction

Nowadays with the rapid growth in computing, networking and database technologies result into collection and integration of massive amount of digital data. For the purpose of analysis or due to government policies, this data need to be shared/published among various parties. When these data consist of personal sensitive information, privacy of an individual will become a critical concern. To resolve these privacy concerns, many PPDP algorithms are developed by researchers. Some of the well-received techniques are $k$-anonymity [1–4], $\ell$-diversity [5, 6],

K. Ashoka (✉) · B. Poornima
Bapuji Institute of Engineering and Technology, Davangere, Karnataka 577004, India
e-mail: ashoka_kkd1@yahoo.com

B. Poornima
e-mail: poornimateju@gmail.com

$t$-closeness [7], $m$-invariance [8], Personalized Privacy [9], Slicing [10], etc. It has been shown in [11, 12] that $k$-anonymity losses significant information loss and suffers from homogeneity and background knowledge attack. $\ell$-diversity has better data utility than $k$-anonymity, but it suffers from skewness attack [7]. As the Quasi-identifiers are published as it is, this technique cannot avoid membership disclosure [13]. $t$-Closeness [7] requires the distribution of sensitive attribute in each bucket to be close to the distribution of the attribute in the original table. This condition significantly damages data utility as well as the correlation between QID and SAs. $m$-Invariance [8] is a dynamic data re-publishing model as it allows both record insertions and deletions with multiple release of data. But $m$-invariance does not guarantee privacy if the life span is broken. That is, it happens if a record reappears after its first life span.

These limitations are eliminated in personalized privacy [9], but it suffers from 'play safe problem' where the record owner may play safe by fixing his/her guarding node as 'Any disease' which will keep him/her in safer privacy zone. But this play safe will result in inaccurate results in many data mining tasks. Slicing [10] is more powerful technique when privacy protection and information loss is concerned, it partitions the data table both vertically and horizontally. Vertical partitioning is done by grouping the attributes based on correlation among the attributes. They use mean square contingency coefficient to measure the correlation among the attributes, and based on these correlation values, attribute clustering will be done using $k$-medoid method, to partition attributes into columns, which will be computationally very expensive. Horizontal partitioning is done by grouping tuples into buckets, for which they use 'Mondrian' algorithm [4], which is not an optimal algorithm for tuple partitioning.

An improvement in the $k$-anonymity/$\ell$-diversity was suggested in [14, 15], where the researchers presented a systematic clustering technique for $k$-anonymization. In this method, grouping of the tuples will be done systematically to satisfy $k$-anonymity/$\ell$-diversity, which has better data utility and execution time. An independent $\ell$-diversity principle to avoid corruption attack [16] was proposed in [17] that integrates perturbation and generalization to conserve more data utility. Preserving privacy for unstructured textual medical data was discussed in [18, 19] where sensitive association rules are sanitized by altering the support and confidence of related items.

Big data analytics is a trending area that makes remarkable revolution in traffic control, disease outbreak perception, smart grids, product recommendation, etc. Privacy preserving for big data analytics to safeguard differential privacy for individual data contributors was done in [20]. The authors propose a generic framework to generate analysis results on a sampled databases. A differential privacy approach for enhanced classification accuracy was proposed in [21], in which the authors present a non-interactive algorithm to satisfy $\varepsilon$-differential privacy. Preserving privacy for outsourced multimedia material is proposed in [22]. The authors use privacy preserving framework based on robust hashing and partial encryption techniques.

All the approaches for PPDP assume that the attributes of the input microdata are classified into three categories. (1) Identifying Attributes: that distinctively identify an individual, for example, Social Security Number; they are untied from the published table. (2) Quasi-identifiers (QID): which can be used by the opponent to link these values to a publicly accessible external database (voter list) to identify an individual, for example, Gender, Zipcode, and Birthdate. (3) Sensitive attributes (SA): which are unknown to the adversaries and are used in data mining and statistical analysis, for example, Disease, Income. An illustrative original microdata table is shown in Table 1.

## 1.1 Motivation and Paper Outline

To overcome the limitations deliberated above, we are motivated to propose a novel approach for PPDP. We used the concept of entropy of the data set to find the effectiveness of an attribute to classify the data, which in turn gives us the correlation among QIDs and SAs. By using this correlation, the input microdata table is sliced vertically. For horizontal partitioning, we have used optimum partitioning algorithm for $\ell$-diversity.

At first, we begin with the care full analysis of entropy and formalize the concepts that motivated the new technique that gives the efficacy of the attributes to classify data. We concentrate on the accuracy of classification application that provides a simple framework that partitions the data table, vertically by correlation in terms of Information Gain from the entropy and horizontally by optimal partitioning for $\ell$-diversity.

**Table 1** Input microdata table

| Name | Gender | Education | Race | Occupation | Salary (k) |
|------|--------|-----------|------|------------|------------|
| Anil | M | Bachelor | White | Tech-support | $\leq 50$ |
| Basu | F | Bachelor | White | Tech-support | $\leq 50$ |
| Chandru | M | Bachelor | White | Sales | $>50$ |
| David | M | Bachelor | Asian | Professional | $>50$ |
| Edwin | M | Masters | Black | Professional | $>50$ |
| Fatima | F | Masters | Black | Professional | $\leq 50$ |
| Geeta | F | Masters | Black | Sales | $>50$ |
| Henry | M | Bachelor | Asian | Tech-support | $\leq 50$ |
| Imthiyaz | M | Masters | Black | Tech-support | $>50$ |
| John | M | Masters | Asian | Professional | $>50$ |
| Kate | F | Masters | Asian | Tech-support | $>50$ |
| Lilly | F | Bachelor | Asian | Sales | $>50$ |
| Manu | M | Masters | White | Sales | $>50$ |
| Nancy | F | Bachelor | Asian | Professional | $\leq 50$ |

We developed an entropy-based partitioning algorithm for slicing the table that preserves considerable amount of information in the microdata. Our preliminary result specifies that the proposed method can improve data utility, which in turn increases the precision of several data mining tasks. To the best of our knowledge, this is the first of its kind which takes account of entropy-based Information Gain for correlation in slicing the table.

The rest of the paper is organized as follows. In Sect. 2, we formalize the concept of Information Gain. Section 3 gives the entropy-based slicing algorithm. Experimental evaluation of our approach is given in Sects. 4 and 5 concludes the paper with future research directions.

## 2   Information Gain

Entropy is a measure of (im)purity of an arbitrary collection of examples. If the target attribute in the input microdata can take on c different values, then entropy of data S relative to this c-wise classification is given by

$$\text{Entropy}(S) = \sum_{i=1}^{c} -p_i \log_2(p_i) \tag{1}$$

where $p_i$ is the proportion of $S$ belonging to class $i$.

The effectiveness of an attribute in classifying the data is measured by 'Information Gain' (IG). For a data set S and attribute A, the Information Gain is given by

$$\text{IG}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{values}(A)} \frac{|Sv|}{|S|} \text{Entropy}(Sv) \tag{2}$$

where values $(A)$ is the set of all possible values for attribute $A$, and $Sv \subseteq S$ for which attribute A has value v. Therefore from Eq. (2), IG$(S, A)$ is the information provided about the target function value, given the value of some attribute $A$.

For illustration, consider the microdata in Table 1, for the attribute Gender which have the values male and female, there are 14 examples (tuples). For salary value greater than 50 k, out of 14 examples, there are 9 positive (tuples 3, 4, 5, 7, 9, 10, 11, 12, and 13) examples and 5 negative (tuples 1, 2, 6, 8, and 14) examples. Information Gain is given by

$$\text{IG}(S, \text{Gender}) = \text{Entropy}(S) - \sum\nolimits_{v \in \{\text{male,female}\}} \frac{|Sv|}{|S|} \text{Entropy}(Sv)$$

$$\text{Entropy}(s) = -(9/14)\log_2(9/14) - (5/14)\log_2(5/14) = 0.940$$

and

$$S_{\text{male}} \leftarrow \{6\,\text{positive},\ 2\,\text{negative}\}\ (\text{out of 8 male examples})$$

$$S_{\text{female}} \leftarrow \{3\,\text{positive},\ 3\,\text{negative}\}\ (\text{out of 3 female examples})$$

Therefore

$$\text{IG}(S, \text{Gender}) = 0.940 - (8/14)\text{Entropy}(S_{\text{male}}) - (6/14)\text{Entropy}(S_{\text{female}})$$

If the entropy of male and female examples is calculated as above, we get

$$\text{IG}(S,\ \text{Gender}) = 0.940 - (8/14)0.811 - (6/14)1.00 = 0.048$$

Similarly for the other attributes, Education, Race, and Occupation the IG values are calculated, as given below.

$$\text{IG}(S, \text{Education}) = 0.151$$

$$\text{IG}(S, \text{Race}) = 0.029$$

$$\text{IG}(S, \text{Occupation}) = 0.246$$

According to Information Gain measure, the attribute occupation has highest value; hence, it has highest prediction accuracy (so the correlation) of target attribute. Therefore, while dividing the table vertically, the attributes Occupation and Salary are kept in same column and others in other columns. We may also add another attribute for the above column by calculating second-level Information Gain of the remaining attributes for various values of occupation, that is, for Tech-support, Sales, and Professional. The horizontal partitioning was done using

**Table 2** Sliced table

| Gender, education, race | Occupation, salary |
|---|---|
| F, Masters, Black | Tech-support, $\leq 50$ k |
| M, Bachelor, Asian | Tech-support, $\leq 50$ k |
| M, Bachelor, White | Sales, >50 k |
| M, Bachelor, White | Professional, >50 k |
| M, Masters, Black | Professional, >50 k |
| F, Masters, Black | Professional, $\leq 50$ k |
| F, Bachelor, White | Sales, >50 k |
| F, Bachelor, Asian | Sales, >50 k |
| M, Masters, Black | Professional, >50 k |
| F, Bachelor, Asian | Tech-support, $\leq 50$ k |
| F, Masters, Asian | Professional, $\leq 50$ k |
| M, Masters, Asian | Sales, >50 k |
| M, Bachelor, Asian | Tech-support, >50 k |
| M, Masters, White | Tech-support, >50 k |

optimal partitioning method to satisfy $\ell$-diversity. For example, when Table 1 is sliced by this technique, the output table is shown in Table 2.

## 3 Entropy-slicing Algorithm

Our proposed algorithm has two parts. In the first part, the efficacy of an attribute in the form of Information Gain based on entropy to classify the data is calculated for each QIDs. That will be used to determine the attribute partitioning (column partitioning) criteria. Next in the second part, we did the optimal tuple partitioning for $\ell$-diversity slicing. The complexity of the step-2 (calculation of IG values for each QIDs) is $O(m)$, where m is the number of QID attributes. The tuple partition (steps-5–8) takes $O(n \log n)$. The generalization in step-9 takes $O(n \log n)$. Therefore, the overall complexity of the algorithm is $O(n \log n)$ which is efficient than slicing [10].

Algorithm Entropy-Slicing
Input: Private Microdata table T, Number of attributes per column k, l .
Output: The publishable table T*.
1.  If (n < THmin) then return with warning message // T should contain minimum THmin records
2.  Calculate IGi (Information Gain) for all QIDs. 1≤i ≤ m.
3.  Vertically partition T, with k no. of attributes based on IG.
4.  Initialize T*=$\phi$
5.  For each tuple ti ∈ T  (1 ≤i ≤n)
6.      Search ti.QID for matching bucket.
7.      If found  Put ti in the Bucket
8.       Else initialize new bucket with ti in T*
9.  Perform generalization, if tuple insertion violates l-diversity.
10. Randomly permute tuples within each bucket.

## 4 Experiments

In this section, we experimentally estimate the efficacy of our approach as related to $k$-anonymity, $\ell$-diversity, and Slicing. We have taken adult data set from UCI machine learning repository [23]. This adult data set contains fifteen attributes. After removing tuples with missing and null values, there are around 45 k effective tuples in total. In our experiments, we have used Age, Education, Marital status, Race, Sex, and Country as Quasi-identifiers. The attribute Occupation is treated as sensitive attribute. The value of $\ell$ is set to 4. We did the experiments for

classification accuracy and computational efficiency. We used weka tool for evaluation of classification accuracy for Decision Tree and Naïve Bayes method. Learning the classifier was done with tenfold cross-validation.

Figure 1 gives the accuracy of classification for Naïve Bayes and Decision Tree classifiers. We observe from the results that entropy-based slicing is at par with normal slicing, and its performance is better compared to $k$-anonymity and $\ell$-diversity.

We also investigated the computational cost of our technique with $k$-anonymity, $\ell$-diversity, and Slicing. Figure 2 shows the experimental results for execution time in seconds with the work load of 10, 20, and 40 k records. The figure unveils that our entropy-based slicing algorithm shows better performance than normal slicing with respect to computational cost.

## 5 Conclusion and Future Scope

In this paper, we presented a new technique 'Information Gain based slicing' for preserving privacy in microdata publishing, which can preserve the privacy of multiple heterogeneous sensitive attributes. The general method proposed by our work is that by considering the effectiveness of an attribute in terms of Information Gain from entropy for partitioning the attributes, slicing was done. Our experiments shows that Information Gain-based slicing has better classification accuracy and computational cost when compared to normal slicing.

Our work motivates several directions for upcoming research. We did the work for single release data, whereas sequential release and multiple release data privacy are still in its infant stage. Research on hybrid technique that combines the plus points of anonymization, slicing, perturbation, etc. can be considered.

Finally, while there exist abundant anonymization algorithms, providing privacy for unstructured data is still in its suckling stage. Therefore, more stress needs to be given in designing effective algorithms for privacy preserving in an unstructured data.
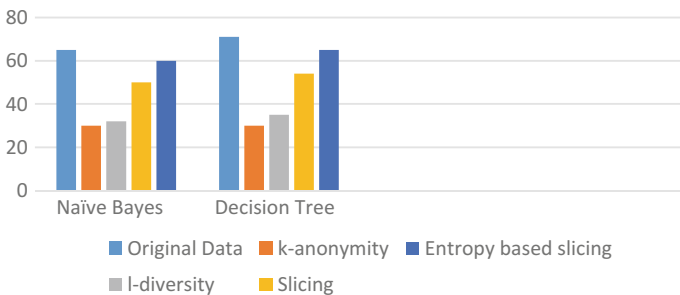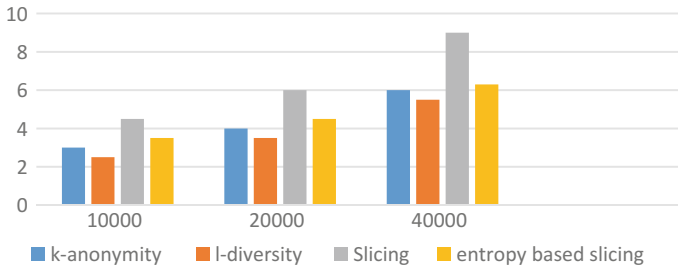


**Fig. 1** Classification accuracy

**Fig. 2** Computational cost—no. of records versus execution time (s)

# References

1. Sweeney, L.: k-anonymity: a model for protecting privacy. Int. J. Uncertain. Fuzziness Knowl. Based Syst. **10**(5), 557–570 (2002)
2. Samarati, P.: Protecting respondents' identities in microdata release. IEEE Trans. Knowl. Data Eng. **13**(6), 1010–1027 (2001). doi:10.1109/69.971193
3. LeFevre, K., DeWitt, D.J., Ramakrishnan, R.: Incognito: efficient full-domain K-anonymity. In: Proceedings of the 2005 ACM SIGMOD International Conference on Management of data (SIGMOD '05). ACM, New York, NY, USA, pp. 49–60 (2005)
4. LeFevre, K., DeWitt, D.J., Ramakrishnan, R.: Mondrian multidimensional k-anonymity. In: 22nd International Conference on Data Engineering (ICDE'06), pp. 25–25. doi:10.1109/ICDE.2006.101 (2006)
5. Machanavajjhala, A., Gehrke, J., Kifer, D., Venkitasubramaniam, M.: $\ell$-diversity: privacy beyond k-anonymity. In: Proceedings of International Conference Data Engineering (ICDE), p. 24 (2006)
6. Domingo-Ferrer, J., Torra, V.: A critique of k-anonymity and some of its enhancements. In: Proceedings of the 3rd International Conference on Availability, Reliability and Security (ARES), pp. 990–993 (2008)
7. Ninghui, L., Tiancheng, L., Venkatasubramanian, S.: t-closeness: privacy beyond k-anonymity and $\ell$-diversity. In: Proceedings—International Conference on Data Engineering, pp. 106–115 (2007)
8. Xiao, X., Tao, Y.: m-invariance: towards privacy preserving re-publication of dynamic datasets. In: ACM SIGMOD International Conference on Management of Data, pp. 689–700 (2007)
9. Xiao, X., Tao, Y.: Personalized privacy preservation. In: Proceedings of ACM International Conference on Management of Data (SIGMOD), Chicago, IL (2006)
10. Li, T., Li, N., Zhang, J., Molloy, I.: Slicing: a new approach for privacy preserving data publishing. IEEE Trans. Knowl. Data Eng. **24**(3), 561–574 (2012)
11. Aggarwal, C.: On k-anonymity and the curse of dimensionality. In: Proceedings of International Conference on Very Large Data Bases (VLDB), pp. 901–909 (2005)
12. Kifer, D., Gehrke, J.: Injecting utility into anonymized data sets. In: Proceedings of ACM SIGMOD International Conference on Management of Data (SIGMOD), pp. 217–228 (2006)
13. Nergiz, M.E., Atzori, M., Clifton, C.: Hiding the presence of individuals from shared databases. In: Proceedings of ACM SIGMOD International Conference on Management of Data, pp. 665–676 (2007)
14. Kabir, M.E., Wang, H., Bertino, E.: Efficient systematic clustering method for k-anonymization. Acta Inf. **48**(1), 51–66 (2011). doi:10.1007/s00236-010-0131-6

15. Pavan, R., Bhaladhare, A.N.D., Devesh, C.: Jinwala: novel approaches for privacy preserving data mining in k-anonymity model. J. Inf. Sci. Eng. **32**(1), 63–78 (2016)
16. Tao, Y., Xiao, X., Li, J., Zhang, D.: On anti-corruption privacy-preserving publication. In: Proceedings of ICDE 08, Cancun, April 7–12, pp. 725–734. Washington, DC, USA (2008)
17. Zhu, H., Tian, S., Lü, K.: Privacy-preserving data publication with features of independent $\ell$-diversity. Comput. J. **58**(4), 549–571 (2015)
18. Fengli, Z., Yijing, B.: ARM-based privacy preserving for medical data publishing. In: Cloud Computing and Security: First International Conference, ICCCS 2015, Nanjing, China, August 13–15. doi:10.1007/978-3-319-27051-7_6 (2015)
19. Sánchez, D., Batet, M., Viejo, A.: Utility-preserving privacy protection of textual healthcare documents. J. Biomed. Inf. **52**, 189–198 (2014). doi:10.1016/j.jbi.2014.06.008
20. Fan, L., Jin, H.: A practical framework for privacy-preserving data analytics. In: Proceedings of the 24th International Conference on World Wide Web (WWW '15), pp. 311–321. ACM, New York (2015)
21. Zaman, N.K., Obimbo, C., Dara, R.A.: A novel differential privacy approach that enhances classification accuracy. In: Desai, E. (ed.) Proceedings of the Ninth International C* Conference on Computer Science and Software Engineering (C3S2E '16), pp. 79–84. ACM, New York. doi:http://dx.doi.org/10.1145/2948992.2949027 (2016)
22. Weng, L., Amsaleg, L., Furon, T.: Privacy-preserving outsourced media search. IEEE Trans. Knowl. Data Eng. **28**(10), 2738–2751 (2016). doi:10.1109/TKDE.2016.2587258
23. Lichman, M.: UCI Machine Learning Repository. (http://archive.ics.uci.edu/ml). Irvine, CA: University of California, School of Information and Computer Science (2013)

# Scale Space-based Capillary Detection to Determine Capillary Density and Avascular in Nailfold Capillary Images Using USB Digital Microscope

**H. S. Ajaya, H. R. Shreyas, Vikram Manikandan, K. V. Suma and Bheemsain Rao**

**Abstract** Nailfold capillaroscopy (NC) is a non-invasive technique used for detecting multiple medical disorders such as connective tissue, rheumatic and systemic diseases. Current NC methods involve acquiring high-resolution images from expensive video capillaroscope for analysis. NC analysis on low-resolution images obtained from a low-cost hardware is a challenging task. Scale space capillary detectors (SCD) is proposed for detection of capillaries under such conditions, which is a unique combination of anisotropic diffusion and Harris corner detector. SCD is followed by Ordinate clustering algorithm (OCA) which is used to eliminate outliers in the detected capillaries. Nailfold capillary images are obtained from a low-cost digital microscope under varying lighting conditions to form a custom database. Experimental results show the promising performance of the proposed algorithm with a high true-positive detection rate and a low false-negative detection rate.

**Keywords** Nailfold capillary · Anisotropic diffusion · Ordinate clustering algorithm · Bi-cubic interpolation · Harris corner detector

H. S. Ajaya
International Institute of Information Technology, Bangalore, India
e-mail: ajayahs@gmail.com

H. R. Shreyas · V. Manikandan
University of Southern California, Los Angeles, CA, USA
e-mail: shreyas.ravindranath23@gmail.com

V. Manikandan
e-mail: manickav@usc.edu

K. V. Suma (✉)
Ramaiah Institute of Technology, Bangalore, India
e-mail: sumakv@msrit.edu

B. Rao
PES University, Bangalore, India
e-mail: drbheemsainrao@gmail.com

# 1 Introduction

Nailfold capillaroscopy (NC) is a non-invasive clinical aid which is used in the diagnosis and prognosis of different diseases. Nailfold capillary image is obtained at the nailbed of distal fingers of non-dominant hand, that is, typically fourth finger of left hand. Morphological aspects such as capillary shape, orientation and architecture are found to give information on the health of a person [1]. Capillary density is a morphological feature which refers to the number of capillaries present in the given area. This parameter is crucial in diagnosis of diseases such as hypertension. It is observed that hypertensive subjects have a rarefaction of capillaries, that is, the capillary density is low. In addition, distance between capillaries can be used to reinforce this aspect. The custom database of nailfold capillary images was acquired using USB digital microscope [2]. Sample images are shown in Fig. 1.

In this paper, we are proposing a fully automated method that does not involve image cropping. Instead, ordinate clustering helps in isolating the two layers of capillaries. Only the upper layer is selected, and capillary tips are detected using scale space detection. It is necessary to process the image so that only corner features corresponding to capillary tips are present [3].

Corners can be important features in capillary images which can be used to identify the capillary peaks and further to determine the capillary density. Many corner detection algorithms are available. Trajkovic et al. [4] proposed a fast corner detector which is suitable for images with high level of light intensity. The nailfold capillary images have a low brightness and hence cannot use fast corner detector. Wang and Brady have come up with real-time corner detector which takes the difference of the second tangential derivative with the edge strength [5].

Anisotropic diffusion introduced by Perona and Malik formed a new definition of scale space [6]. A new class of algorithms was introduced, in which the diffusion coefficient is chosen to vary spatially in such a manner as to encourage intra-region smoothing rather than inter-region smoothing. This technique aims at reducing the



**Fig. 1** Sample capillary images in different health conditions

image noise with the additional benefit of edge preservation which is crucial for image interpretation [7].

This paper is divided into following sections—Sect. 1 gives an introduction to the use of nailfold capillary image as a clinical aid and an overview of the technique used. Section 2 elaborates the automated capillary analysis system proposed, while Sect. 3 describes the computation of the capillary density and avascular followed by results and discussion in Sect. 4. Conclusions are drawn in Sect. 5.

## 2 Automated Capillary Analysis System

A custom database of 60 images of nailfold capillary images belonging to healthy, diabetic and hypertensive subjects were collected. The subjects were in the age group of 18–60 years and consisted of 28 females and 32 males. Ethics committee clearance was obtained from M S Ramaiah Medical College and Teaching hospitals. Informed consent of the subjects was taken. The room temperature was maintained around 25 °C, and the subjects were seated with their hands placed on the table at their heart level. The digital microscope has a ×200 magnification, and a resolution of 640 × 480 was selected.

The fundamental techniques used for image pre-processing such as optimum color channel selection, bi-cubic interpolation, Harris corner detection, and scale space construction by anisotropic diffusion so on play an important role in obtaining significant results as shown in Fig. 2. The images that are taken using the microscope usually contain red, green, and blue channels. The RGB color space is represented normally using a 3D matrix. These images contain large amount of data that require large computational time for processing as well. In order to improve



**Fig. 2** General block diagram of the proposed system

computational time and efficiency, we select only the optimum color channel for processing which reduces the 3D matrix to a 2D matrix.

It is found that the green channel has the highest degree of contrast within the capillaries and the background from [6]. The red and blue channels tend to be noisy and as a result are removed. This technique results in an image which contains only the green channel for further processing.

Resizing an image reduces the sharpness and quality of the image, and since we are using detected corners of capillary tips, we do not want to lose this feature from the images by just resizing it. Bi-cubic Interpolation is used to modify the size of an image. The process of determining estimate values at unknown points using known data is called as interpolation. The available interpolation algorithms are of two kinds, adaptive and non-adaptive. Since all pixels in the image are treated equally, we have used the non-adaptive algorithm viz., bi-cubic interpolation algorithm. This algorithm uses adjacent pixel parameters to compute interpolation; more the number of adjacent pixels, better the accuracy [1]. Here, the ones that are closer are weighted higher and farther the pixel, lower the weight.

$$\text{Model}: f(x, y) \sum_{i=0}^{3} \sum_{j=0}^{3} x^i a_{ij} \tag{1}$$

where $f(x, y)$ is a bilinear surface

Approximation:

$$\partial_x f(x, y) = [f(x+1, y) - f(x-1, y)]/2 \tag{2}$$

$$\partial_y f(x, y) = [f(x, y+1) - f(x, y-1)]/2 \tag{3}$$

$$\partial_{xy} f(x, y) = [f(x+1, y+1) - f(x-1, y) - f(x, y-1) + f(x, y)]/4 \tag{4}$$

Features in an image appear at different scales or resolution. A scale space is a family of images derived from the same original image but filtered from fine to coarser resolutions. The filtering process is generally referred to as diffusion as the image features, particularly edges, are thought to be diffusing through the scale space from higher concentration in the finer resolutions to lower concentrations at coarser resolutions. Filtering is said to be isotropic if the diffusion coefficient remains constant throughout the scale space. Isotropic diffusion does not respect the object boundaries and at coarser scales, the edges become spatially distorted from their original position. In order to avoid this effect, Perona and Malik introduced anisotropic diffusion (AD) wherein the diffusion coefficient varies through the scale space. In this paper, we construct a scale space for the NC images using AD in order to analyze the statistics of nailfold capillaries. The AD is given by Eq. (5)

$$I_t = \frac{\partial I}{\partial t} = c(x, y, t)\Delta I + \nabla c \cdot \nabla I \tag{5}$$

where $I_t$ is the filtered image at scale t, div is the divergence operator; also used are Gradient and Laplacian operators while $c(x, y, t)$ is the diffusion coefficient.

Figure 3 shows the visualization of smoothing achieved by AD through the scale space at $t = 1$, $t = 5$ and $t = 10$. Comparing with the original image, the filtered images are less noisy while the corners corresponding to capillary tips are preserved.

Harris corner detector (HCD) is applied on the image filtered using AD. Only the corners that exceed a certain threshold on the Harris corner metric are taken as capillary tips. We can see in Fig. 4 the output of HCD with and without AD. It can be seen in Fig. 4a that HCD results in a large number of false-positive detections, leading to very low accuracy. The combination HCD + AD, as shown in Fig. 5b, results in a very low false-positive rate. Thus, the detection accuracy of capillaries is greatly increased. The diffusion coefficient function used for AD is given in Eq. (6).

$$c(x; y; t) = e^{(t)} \qquad (6)$$

where the threshold value that controls the diffusion, =25 in all our experiments.

We use Ordinate clustering algorithm (OCA) to categorize the capillaries into upper and lower layers. In OCA, a capillary called the Probe capillary (PC) is picked randomly from the set of detected capillaries. The Y-coordinate distance of the PC with all other capillaries is computed. A threshold is kept on the ordinate distances which divide the set of capillaries into two layers, namely layer 1 and layer 2. Following this, the PC has to be assigned to one of these layers. To do this, the mean ordinate distance of the PC with each layer is computed. The PC is assigned to the layer with the least mean. The complete set of capillaries is divided into layer 1 and layer 2 subsets. The subset with the lowest ordinate belongs to the lower layer and can be discarded. The remaining capillaries belong to the upper layer and are used for further analysis. This accounts for increase in detection accuracy.



**Fig. 3** Anisotropic filtered scale space at $t = 1$, 5, and 10 in comparison with the original image. **a** 2D image and **b** corresponding 3D plots. Observe that corners corresponding to capillary tips are preserved through the scale space while others get smoothed out

**Fig. 4** The detected corners **a** without anisotropic filtering and **b** with anisotropic filtering. Only the required capillary tips are detected after anisotropic filtering. A human nailfold bed will have multiple layers of capillaries present, and the capillaries in the upper layer are of interest to be detected for medical purposes. We can see in Fig. 5 the nail bed which consists of capillaries in two layers. After the detection of the first layer, it is necessary to eliminate the capillaries that lie in the lower layers

**Fig. 5** Capillaries in *upper* and *lower layers*, as seen from the digital microscope



By maintaining uniformity in applying the immersion gel, the region of interest which is cropped out can be maintained constant as well. Hence, the whole setting need only be programmed once.

In order to determine the optimum threshold value for the Harris corner metric, we plot the average TPR, FPR, and FNR obtained by taking two random images against different threshold values, as shown in Fig. 8; two random images are taken to ensure that threshold so obtained does not become specific to the database. We observe from the plot that highest TPR is obtained between =0:8 and =0:9 but =0:8 gives lowest FNR. Hence, the threshold value is taken as 0:8 (mean of all points).

In another experiment, we compare the detection rates obtained with and without using OCA. Since OCA is mainly used to exclude the outliers, FPR is the most relevant metric to grade its performance. It is clearly seen from Table 1 that there is

**Table 1** Comparison of detection rates with and without OCA

| Attributes | Without OCA | With OCA |
|---|---|---|
| TPR | 68.79 | 81.86 |
| FPR | 31.2 | 8.14 |
| FNR | 2.08 | 2.08 |
| TNR | 97.98 | 97.98 |
| Accuracy | 83.38 | 94.92 |

a significant decrease in FPR when OCA is used. This shows that OCA effectively excludes outliers in the initial detection of capillaries.

## 3   Capillary Density and Detection of Avascular Regions

In this experiment, we determine the capillary density and distance between capillaries which can be used to detect diseases. Capillary density is calculated once the number of capillaries in an image is known. In order to accommodate for the faulty detections from the algorithm, we specify the capillary density with a tolerance value. We define this tolerance value as the difference between FPR and FNR. Hence, capillary density is given in Eq. (7).

$$\text{Capillary density} = (\text{number of capillaries in the first layer})/\text{bounded area} \qquad (7)$$

Distance between capillaries is used to find the probable avascular regions. First, the X-axis separations between adjacent capillaries are calculated using the location of detected capillaries. If the distance crosses a threshold value, here taken as 40 pixels, then the region between the corresponding capillaries is said to be avascular. A set of samples with probable avascular regions is shown in Fig. 9. An apparent avascular region can arise if a capillary is not detected (false negative). Hence, the probability of a detected avascular region being truly avascular is inversely proportional to the FNR of the algorithm. The proposed algorithm has a very low FNR, hence a high probability of true avascular region.

## 4   Results and Discussion

Figure 2 shows the block diagram of the automated capillary analysis system. Optimum color selection chooses green channel which is less noisy. Bi-cubic interpolation helps in estimating values based on surrounding $4 \times 4$ pixel values. This is used to resize the image without reducing the sharpness of the features. Further, Harris corner detector finds the intersection of two edges. These points, called as corners, are the capillary tips which can be used to calculate capillary density and inter-capillary distance. Figure 3 shows anisotropic filtered scale space and illustrates that corners corresponding to capillary tips are preserved while smoothing out others. OCA is used to separate the upper and lower layer of capillaries. Table 1 shows how the capillary detection can be drastically improved by using OCA. Sensitivity increased from 68.79 to 81.86% while high specificity is maintained at 97.98%. From Fig. 4, it can be seen that the lower layer of capillaries is not detected, when OCA is employed. Figures 6 and 7 illustrates the erroneous detection when OCA is not used. Figure 8 demonstrates the effect of different threshold values on the capillary detection rate. A threshold value of 0.8 is observed

to give lowest FNR; hence, this threshold is employed. Figure 9 shows the avascular regions detected.

## 5 Conclusion

A novel approach to detect the nailfold capillaries is proposed which uses scale Space capillary detection (SCD) to detect capillaries in a human nail bed, OCA to isolate and procure the cardinal row of capillaries, capillary density is used to detect drop in capillary count in a circumscribed area, and distance between the capillaries is used to detect probable avascular regions. Capillary density and distance between the capillaries play a significant role in identifying the disease of a subject using the detected capillaries. The proposed combination of system performs extremely well in detecting capillaries and isolating the fundamental row of capillaries; the average detection rate of 91.86% was achieved when OCA was followed by scale space detection and 68.79% when only scale space detection was performed. We can also see a drastic fall in the false-positive rate from 31.2 to 8.14% when OCA was used. We can deduce that the proposed SCD detects capillaries from an image far more efficiently than the other corner detectors such as FAST detector, SUSAN detector and Minimum Eigen detector from the experiments conducted. Therefore, a future research scope would be to use a higher-resolution digital microscope to obtain better quality images. The images from these microscopes could be used to determine the arteriolar, venular, and apical limbs of the capillary, the dimensions of the capillary such as the width and height of the capillary, and capillary architecture which is the shape of the capillary which could vary from elongated, enlarged to bushy capillary.



**Fig. 6** Shows the detected capillaries **a** before applying OCA and **b** after applying OCA. The *lower layer* capillaries are eliminated using OCA in (**b**)

O Erroneous detection

**Fig. 7** Shows erroneous detection in **a** before OCA, which are eliminated in **b** after OCA

**Fig. 8** Graph to determine the optimum threshold value





Possible region of Avasculature

**Fig. 9** Shows the probable avascular regions detected in some sample images

# References

1. Cutolo, M., Sulli, A., Secchi, M.E., Paolino, S., Pizzorni, C.: Nailfold capillaroscopy is useful for the diagnosis and follow-up of autoimmune rheumatic diseases. A future tool for the analysis of microvascular heart involvement? Oxf. J. Rheumatol. **45**(4), iv43–iv46 (2006)
2. Vasdev, V., Bhakuni, D.S., Bhayana, A., Kamboj, P.: Nailfold capillaroscopy: a cost effective practical technique using digital microscope. Indian J. Rheumatol. **6**(4), 185191 (2011)
3. Harris, C., Stephens, M.: A combined corner and edge detector. In: Proceedings of Alvey Vision Conference, vol. 15, pp. 147–151 (1988)
4. Trajkovic, M., Hedley, M.: Fast corner detection. Image Vis. Comput. **16**(2), 7587 (1998)
5. Wang, H., Brady, M.: Real-time corner detection algorithm for motion estimation. Image Vis. Comput. **13**(9), 695703 (1995)
6. Perona, P., Malik, J.: Scale-spaceandedge detection using anisotropic diffusion. IEEE Trans. Pattern. Anal. Mach. Intell. **12**(7), 629–639 (1990)
7. Srivastava, A., Bhateja, V., Tiwari, H.: Modified anisotropic diffusion filtering algorithm for MRI. In: Proceedings of (IEEE) 2nd International Conference on Computing for Sustainable Global Development 1885–1890 (2015)

# Augmenting Non-verbal Communication Using a Tangible User Interface

Suraksha Devi and Suman Deb

**Abstract** Since the dawn of the human civilization, non-verbal communication has been a pervasive phenomenon. With the advancement of technology, this natural interaction blended with computing facility turned into a scientific tangible interaction. Computer-aided non-verbal communication can be used as a vital interface between computers and specific areas of modern life like human psychology such as emotions, social attitudes, and even for educational games. Investigative, design-focused studies have shown that TUIs provide learning aids due to the added haptic dimension and the better approachability to the shared space that can be used in supportive circumstances. In this paper, a tangible user interface has been proposed comprising of several tablets/smartphones connected to a single server. This system would help out the teachers to attain feedback (whether the students are able to understand what the teacher is trying to convey. If yes up to what level) from the students in an elementary classroom. The experiments conducted with the prototype have shown significant interaction improvements among the children and the teachers in the classroom sessions. The experiments on modeling and improving the prototype are going on to improvise the interaction specially for mobile devices and tablets.

S. Devi (✉) · S. Deb
Department of Computer Science and Engineering, National Institute of Technology
Agartala, Agartala, India
e-mail: suraksha.malik5@gmail.com

S. Deb
e-mail: sumandebcs@gmail.com

# 1 Introduction

Communication is multifaceted. It is a process by which the information can be shared among the individuals verbally or non-verbally. The relative contribution of non-verbal to verbal interactions could not be quantified, but a lot more meaning is provided by the non-verbal mode of interaction that what is thought by the people. Non-verbal communication is the silent form of communication among the parties without using any form of speech and is often used to make an expression of thoughts and make your conveying message to be more appealing and interesting to whom you are speaking. It can regulate relationships and can support or even replace verbal communications in many situations.

Non-verbal communication could turn out to be a beneficial for classroom session, but for the unexperienced mentors, it could turn out a barrier as the teacher would not be able to recognize whether the students are understanding what is been conveyed. A tangible user interface could be used in the classroom session to remove such barriers. As per Ullmer and Ishii [1], "TUIs challenge to eliminate this discrepancy between the input and output and introduce innovative options for collaboration which brings the physical and digital worlds together." Corporeal and intellect features in input–output are accentuated by the TUIs. The corporeal depiction of real objects is often united in a close manner to the TUIs.

Controller and description are put together in a solitary device by the tangible user interface [2]. The notions of "societal collaboration" and "association" are underwired by these interfaces [3, 4], and it has been proved that such structures are idyllic for tentative and inventive actions and acquiesce students to spawn creations that could be versioned by any sort of means [5]. Inherent enthusiasm is animated, and superlative experience of learning is provided by them both in verbal and non-verbal communication. Hypothetical works on practicing the TUIs in non-verbal communication have been emerging slowly. There have been many studies, and research conducted on the use of tangible interface for non-verbal communication.

This paper highlights the importance of non-verbal communication in classroom and how a tangible user interface could be used in classroom sessions to decrease the barriers between the teacher and the students when it comes to non-verbal communication. It was envisioned to know whether TUI could be used to enhance non-verbal communication; in other words, does this interface help out the teachers to teach the students in a better way, when it is difficult for the teacher to understand the student's cues, thereby making the classes to be interactive rather than uninteresting. Contemplating that the interface could be efficaciously implemented into the classroom, a system has been proposed comprising of tablets/smartphones with an installed android application provided to every student in the classroom and all tablets and smartphones connected to a single server (computer) that could only be accessed by the teacher. This system enables the tutor to attain the feedback (of the level of understanding of what is being taught and what actually are the feeling of the students in that particular classroom session). This feedback would help out the

teachers to make the classroom sessions to be more interesting and interactive for the students so that they could understand the lectures in a better way.

As found in the study, nowadays 85% of the people have their own smartphones/ tablets. So, it was decided to use an android smartphone/tablet for the system. The experiment was conducted over 15 specially abled students of class six of the "FERRANDO rehabilitation center Agartala," and it was found that the system generated output was 82% accurate. The affordance and acceptability of the system was taken into consideration. The results obtained by the performance evaluation of the system confirmed our claim that the system was accurate and efficient enough to be used by the students and the teachers in elementary classrooms. The methodologies suggested in this paper should act as an important precursor for further work and development of better educational software and contents for similar portable systems that could be effectively implemented in elementary classroom sessions to make the learning process to be interesting and interactive for the children.

## 2   Literature Survey

In a classroom, the verbal communication that occurs has a relatively small percentage and it primarily stimulates the intellective domain for the pupil, while on the other hand, the non-verbal collaboration (comprising of approximately 93% of the entire communiqué) inspires the sensuality and insolence of the students (sentimental area) about the objects [6]. Thus, the mentors should wary about their non-verbal prognosis and should observe the non-verbal prompts of the students, in order to study the impression that the non-verbal modes could have on the learning process among the students.

Communication is a twofold interaction method [7]. This declaration turns out to be valid in the teaching space, where as a tutor, we attempt to interrelate in a clear and effective way. Radford [8] states that active interaction during the classroom sessions is quite perilous, which may not always turn out to be beneficial for the students. According to Miller [9], "transmission of knowledge takes place through the communication methods which is operative and fostered by competently transferring the messages…" Nevertheless, a greater part of verbal and non-verbal cues is created in a lively teaching–learning environment by the mentor and envisioned to be processed by the students. Through this construal of the non-verbal prompts, this apparently solitary interaction can be turned into a more interactive, twofold process. The capability and well-being of the children with dispensation of what the mentor has taught can be construed from their non-verbal prompts, which could thereby enable the tutor to carry out the conversation depending on what type of cue was observed.

Angelo and Cross [10] quantified that "instructors can study a lot about how the students are learning and responding to the teaching tactics by observing closely the learning procedure the students are following…" The feedback from the non-verbal

prompts of the children could be used by the teachers for influencing subsequent interactions [8] and altering the method followed to deliver lectures whenever required [11]. Webb et al. [12] specified that "by observing and interpreting body and facial expressions of the students, the decision could be made by the discerning teacher whether the comprehension should be checked, what type of different instructions must be provided, or whether there is a requirement to assign more practices to the students…" Thus, by observing the non-verbal cues that are used within a classroom by the tutors, the teaching methods could be reinforced to make the learning procedures more resourceful and active [10].

Angelo and Cross [10] even discussed various annotations that should be performed by the tutor during the teaching–learning sessions into the classrooms "teachers monitor and answers to the student's queries, comments and expressions while teaching in an automatic manner. This involuntary data collection is an intuitive and couched process. Teachers hinge on their vision of student learning and …" Neill and Caswell [11] have discussed about the inexperienced instructors. The inexperienced mentors are unaware of the extent of the behavior of the non-verbal individuals. This comes with experience. As demonstrated by a study by Webb et al. [12], an evaluation was carried out among the "expert" and "non-expert" teachers to judge students based on visual, non-verbal behavior and it was observed that the experienced teachers were much more precise as compared to the "non-expert" teachers.

In a classroom comprising of active learning, the circumstances (verbal and non-verbal communication) are utilized, making the classroom environment truly busy. The experienced tutor has established the capability to include psychological multiple tasking during the teaching–learning sessions; talking, perceiving various conditions all around, intermingling, and enabling learning in an active manner. These skills could not be developed at one go; instead, it takes experience to develop such skill to handle the class of students.

Radford [8] records that the teaching skills such as observing the students during the class are to be learned and the mentor could be easily comfortable with these skills with some experience and training. Consequently, the teachers who have attained expertise could easily handle the innumerable and multifaceted information, which could be interpreted and used to understand the behavior of the students, communal data, and events going on in the class in an effective manner as compared to the teachers who are less experienced. When a teacher fails to record, understand, and identify that a pupil is unable to understand the concept, then it could lead to frustrating and boring classes for the students. Hence, it is very imperative for a teacher to be aware of what the student is learning, how is he learning, and whether he is able to understand what the teacher is trying to convey.

When we summarize about the study that we demonstrated above, we could see that the teacher may sometimes due to lack of experience may not be able to understand whether the students are able to understand what is being taught. In such situations, a tangible user interface can serve as a boon. The tangible user interface such as I/O Brush [13] can be used in the classroom sessions to make the class to be more interactive and beneficial for the students.

Although a number of publications can be found that describe different kinds of strengths of non-verbal communication for learning environments, there is still a lack of knowledge on using tangible interface as a tool non-verbal communication for classroom teaching. Based on the study of the research work that has already been accomplished, few research objectives have been drawn that will be focused. The underlying primary research objectives have been formulated to drive this research work from the studies carried out so far. These objectives are as follows:

- Is the proposed system affordable, such that it could be used by every student individually?
- Are there any noteworthy improvements in the performance of the students upon using the system proposed in classroom teaching–learning?

## 3 Proposed System

In the initial design phase of the proposed system, we have modeled a hand in Blender 2.78 capable of performing different hand movements. Further, an android application was created using Android Studio consisting of various non-verbal cues such as happy, angry, and confused. The android app also consisted of an understanding meter using which the students can specify how much they were able to understand a particular topic that was taught. In the final phase of the design, the hand modeled in blender was transferred to the android application. This created a gaming-like environment for the students. Each and every student was provided with a tablet or a smartphone with the respective android application installed, and all of the tablets and the smartphones were connected to a single server (computer) that could only be accessed by the teacher. This system usually targets the special abled people (deaf and dumb) but could be even implemented in the regular classroom sessions making the learning process to be more interactive in a playful manner. The students with the help of the hand model can click on the non-verbal cues available on the screen of the tablet/smartphone. The feedback can be collected by the teacher. This would help the teachers to understand whether the students are able to understand the lecture delivered and what are the actual feelings of the students regarding the lecture delivered, and then, accordingly, the tutor can improve the classroom teaching–learning process.

## 4 Experimental Setup

The experiment for this study involved sessions of two groups (Group A and Group B) of partakers. Group A was subjected to the traditional method of teaching–learning followed in the classrooms over the years, and Group B was made to use the application developed to provide the students feedback. The partakers selected

for this study were 15 specially abled students of class six of the "FERRANDO rehabilitation center Agartala." The sessions planned were on the "solar system." In every session, a small portion of the topic was covered and based on that a test of 10 marks containing ten questions each carrying one mark was scheduled. The partakers were inquired for their consensus at the commencement of the session. The allotment of the partakers was done by the technique of random matching [14, 15]. Through this process, the partakers of the groups were harmonized based on their prior knowledge on the topic that was recognized through the fallouts of the pre-tests. This was accomplished by picking the first two uppermost scoring students and allotting them groups A and B, respectively. In the next round of allotment, the order was reversed to B and A. This process was repeated unless no more partakers were left out and each time the order was altered.

## 5  Data Analysis

For each session, the students were imperiled through the standard model of two tests of ten marks, one earlier and one later to the session each consisting of ten questions carrying one mark each. The pre-test specified the understanding students had before the session was started. Post-test results specified their understanding echelons after the session was completed. The performance is then measured by calculating the difference between the pre- and post-test fallouts using the formula specified below:

$$\text{Performance Evaluation} = \text{Post - test result} - \text{Pre - test result} \tag{1}$$

This was continued for about one month, and 14 different sessions were organized. The average results are listed below (Table 1)

## 6  Result Analysis

Students under treatment Group B were initially at the same level as the Group A, and then, in some sessions their performance degraded. But as the sessions progressed, they quickly caught up. In the end, it could be seen that the average performance of Group B was significantly better as compared to Group A. This can be attributed to the fact that students were initially not comfortable with using the tangible user interface system setup, but they got used to it as more sessions were conducted and also the teachers used the feedback that was provided by the students by using the TUI application and modified the lectures which made the performance of Group B comparatively better. The final results show a situation that the Group B had a better overall performance than the Group A (shown in Fig. 1).

**Table 1** Performance evaluation of the students of Group A and Group B

| Sessions | Performance evaluation | |
|---|---|---|
| | Group A | Group B |
| 1 | 6.25 | 6.5 |
| 2 | 5.75 | 5.25 |
| 3 | 7.0 | 7.50 |
| 4 | 6.75 | 6.5 |
| 5 | 3.5 | 5.5 |
| 6 | 6.0 | 7.5 |
| 7 | 5.5 | 5.0 |
| 8 | 7.75 | 7.0 |
| 9 | 4.25 | 6.50 |
| 10 | 5.75 | 7.50 |
| 11 | 6.25 | 7.25 |
| 12 | 8.75 | 9.0 |
| 13 | 6.75 | 6.0 |
| 14 | 4.50 | 6.75 |

**Fig. 1** Line chart showing the difference in performance among the Group A and the Group B



## 7  Conclusion

This experiment along with our system setup shows how tangible user interface environments can be used in educational field for an effective classroom teaching–learning. This is a step in the course of conveying immersive learning environment to the day-to-day teaching–learning process in elementary classrooms especially for the specially abled. This paper should act as a significant forerunner for researchers observing into the field of applied real-life practice of tangible user interface in classrooms. For further work, there is scope in complex dimension of students' activities and attention points via analytic arrangements applied in the tangible user

interface itself. The experiments showed above can be achieved on a larger group of students for longer period to find out the long-term effects.

Declaration: We hereby declare that the sample survey conducted in our work is with the due permission from respondent and their parents. Permission from both major and minor respondents and their parents are taken. We shall be solely responsible in case any dispute arises due to this.

# References

1. Ullmer, B., Ishii, H.: Emerging frameworks for tangible user interfaces. IBM Syst. J. **39**(3–4), 915–931 (2000)
2. Ullmer, B., Ishii, H.: Emerging frameworks for tangible user interfaces. In: Carroll, J.M. (ed.) Human Computer Interaction in the New Millenium, pp. 579–601. Addison-Wesley, Boston (2001)
3. Hornecker, E., Buur, J.: Getting a grip on tangible interaction: a framework on physical space and social interaction. Proc. CHI **2006**, 437–446 (2006)
4. Marshall, P., Rogers, Y., Hornecker, E.: Are tangible interfaces really any better than other kinds of interfaces? In: Workshop on Tangible User Interfaces in Context and Theory, ACM CHI (2007)
5. Carroll, J.M.: Becoming social: expanding scenario-based approaches in HCI. Behav. Inf. Technol. **15**(4), 266–275 (1996)
6. McCroskey, J.C., Richmond, V.P., McCroskey, L.L.: Nonverbal communication in instructional contexts. In: Manusov, V.L., Patterson, M.L. (eds.) The SAGE Handbook of Nonverbal Communication, pp. 421–436. Sage Publications, Thousand Oaks, CA (2006)
7. Suinn, R.M.: Teaching culturally diverse students. In: McKeachie, W.J., Svinicki, M.D., Hofer, B.K. (eds.) McKeachie's Teaching Tips: Strategies, Research, and Theory for College and University Teachers, 12th edn, p. xxii. Houghton Mifflin, Boston (2006). (407 p.)
8. Radford, K.W.: Observing the class. Educ. Canada **30**, 36–39 (1990)
9. Miller, P.W.: Nonverbal Communication, 3rd edn. NEA Professional Library, West Haven (1988)
10. Angelo, T.A., Cross, K.P.: Classroom Assessment Techniques: A Handbook for College Teachers, 2nd edn. Jossey-Bass Publishers, San Francisco (1993)
11. Neill, S.R., Caswell, C.: Body Language for Competent Teachers. Routledge, London; New York (1993)
12. Webb, J.M., Diana, E.M., Luft, P., Brooks, E.W., Brennan, E.L.: Influence of pedagogical expertise and feedback on assessing student comprehension from nonverbal behavior. J. Educ. Res. **91**(2), 89–97 (1997)
13. I/O Brush: Drawing with Everyday Objects as Ink Kimiko Ryokai, Stefan Marti and Hiroshi Ishii MIT Media Laboratory 20 Ames Street Cambridge, MA 02139 USA
14. Ronald Aylmer Fisher. The design of experiments (1935)
15. Kirk, R.E.: Experimental Design. Wiley Online Library, New York (1982)

# An Approach to En Route Environmentally Sustainable Future Through Green Computing

**Bhubaneswari Bisoyi and Biswajit Das**

**Abstract**  Green computing aims in attaining sustainable future by implementation of practices that uses computing resources economically and eco-friendly manner. The prime aspiration is to minimize the practice of hazardous materials and capitalize on energy efficiency throughout the product's lifespan. Green computing encourages recyclability of obsolete products and wastages released from factory. With due course of time, companies where information technology is implemented have realized that going green is beneficiary both in terms of maintaining public relations and reduced overheads. This paper focuses on various factors that motivate companies to implement green computing and put into practice e-waste recycling process. The main motive behind this study is to promote e-waste management as a factor of green computing. The study also focuses on attainting sustainability using this approach.

## 1  Introduction

With the commencement of industrial revolution, exploitation of fossil fuels and other non-renewable resources had began which unknowingly polluted the environment leading to rise in carbon footprint globally. It took somewhat a long period for realizing that damage has been done towards our atmosphere. Today, it is our prime responsibility to develop measures that are environment friendly and will lead to attend sustainable future. It is high time to adapt green technology to meet the need of the hour. With rise in innovation, invention and other technological

B. Bisoyi (✉) · B. Das
KIIT University, Bhubaneswar, India
e-mail: bhubaneswari.bisoyi@gmail.com

B. Das
e-mail: biswajit@ksom.ac.in

development such as computer and several other devices such as mobile phone, tablets, camera that are closely associated with our daily routine consume relative amount of energy. The period of innovation cycle has reduced with advancement in technology; therefore, the rate of obsolesce of electronic products has increased which leads to electronic waste (e-waste) [1, 2].

The technological progression provides ample scope for exploring energy management and e-waste management in order to facilitate a green computing approach in the computing sector. Recycling of e-waste products and eco-friendly approach is considered to a viable option for implementing green computing and for sustainable future. The area of focus in this paper is on realizing green computing through e-waste management. The various methods used for making this approach viable have been discussed in this paper.

We can define green computing as the practice that leads to use computing resources in an efficient manner and thereby minimizing adverse impact on environment. The goal of green computing is not limited to reduction in carbon footprint but also to reduce usage of hazardous materials throughout the supply chin and also in the manufacturing process [3].

## 2 Concept of Green Computing

The term green computing can also be called as green IT that aims in attaining economic feasibility and enhances effective use of computing devices. Green computing practices comprise of improvisation in production practices for sustainable environment, energy efficient processors and better technique for disposal of electronic wastes (e-wastes) and recycling of these waste products. To promote green computing practices globally, we can implement four approaches as mentioned below [4–6]:

Green use: Reducing power spending of computers and secondary devices attached with it and utilizing these devices in an eco-friendly approach.
Green discarding: Re-using existing electronic devices or correctly discarding of or recycling of product.
Green design: Planning to use devices such as energy efficient computers and other devices such as printer, server, projector.
Green development: Reducing waste during the process of production of computers and other devices attached with it to minimize environmental impact.

## 3 Factors Motivating to Adopt Green Computing

The following development that influences data centres and also impacts electronic devices such as desktop computers for adopting green computing practices [7, 8].

### 3.1 Rise in Usage of Internet

With increase in dependability on electronic data, the number and size of data centre have increased. This growth is due to quick espousal of communication through Internet and media, automation of business processes and function carried out such as retention of all accounts, disaster recuperation and many more. The percentage growth in Internet usage annually is more than 10% leading to increase in building up data centre at a rate of 20% CAGR.

### 3.2 Mounting Power Density of Equipments

Even though improvisation in server CPU in several cases has enabled superior performances with low energy consumption per CPU, but taken as a whole, server energy consumption has increased with demand in installing CPU with bigger memory capacity to minimize the use of floor space, the form factor of servers has been reduced. This leads to rise in power density of data centre. The density has increased more than ten times as it was in 1996–2007.

### 3.3 Escalation in Requirement of Energy for Maintenance of Data Centre

The rise in power density of server has led to an associated boost in heat density of data centre. In order to cool the server requirement of power, 1–1.5 W of power is required for cooling per unit of heat. When we calculated the cost of maintenance of data centre, it exceeds the cost of equipment.

### 3.4 Precincts of Power Supply and Access

Renowned companies such as Google, Microsoft and Yahoo with requirement of big data centre need energy supply for sustainable operation of data centre, but they are unable to get adequate power requirement in American cities. Therefore, they have set up new data centres near Columbia River to draw energy for hydroelectric source.

## 3.5  Low Rate of Server Utilization

The major problem in data centre is issue related to efficiency of energy use. The average server utilization rate is about 5–10% of bigger data centre. Server utilization is low which means increasing the cost of energy supply, maintenance and operation.

## 3.6  Emerging Consciousness About Impact of IT on Atmosphere

There exists a proportional relationship between carbon emissions with usage of energy. In the year 2007, about 44 million servers have been set up worldwide, and the power consumed for these data centre is about 0.5% of total electricity produced. Their combined carbon footprint has reached 80 metric megatons of carbon dioxide in the Netherland and Argentina. This figure is expected to reach about 340 metric megatons by year 2020. This pollution in our environment would be only due to electronic products [9–11].

## 4  Electronic Waste (e-Waste)

With revolution in the information technology, many new technologies have been developed and range of product is available at an affordable price to public. The life cycle of these innovate product is very less, and they become obsolete with due course of time.But on the other hand, it lends to uncontrolled resource consumption and alarming e-waste management. The problem of e-waste is faced both by developed countries and developing countries. The whole range of electronic and electrical products such as computers, mobile, refrigerators, television and many more contains toxic materials that are hazardous to environment. Many of the developments used in our daily work are unsustainable and pose severe confrontation to environment and towards human health. Therefore, our focus should be towards optimal and efficient use of renewable resources, reducing waste products and built an environmentally sustainable recycling process for disposal of various categories of waste products. This initiative would help in achieving sustainable economic growth with enhancing living style.

The issue of e-waste has been handled by European Union (EU) and other developed countries by defining policies and also by adopting various methods for recycling and disposing of waste products. This type of waste product was by EU as "Waste Electrical and Electronic Equipment" (WEEE). In our country, WEEE is called as e-waste [12, 13].

### 4.1 Duration of Electronic Equipments

The useful life period of all electronic products has relatively decreased due to change in features of equipment and competency. The life period of central processing unit has reduced from 4 to 6 years in 1997 and to 2 years in 2005. The average lifetime for various electronic devices is mentioned below for the device which will operate without getting obsolete [13, 14] (Table 1).

### 4.2 Estimation of Total e-Waste Produced Yearly

The total e-waste excepted during the period 2014–2016 for the different electronic devices in metric tonne is mentioned below in Table 2. Figure 1 shows the e-wastage of various devices from 2014 to 2016.

## 5 Recycling

At present, scenario with technological innovation and modern marketing strategy to mitigate customer demand has led to rapid yield of electrical and electronic equipments (EEE). With reduction in cost of electronic product and increase in providing value in terms of various features, old electronic products are getting obsolete and leading to generate e-waste. Recycling of these electronic wastes can minimize the pollution caused by these wastes [15, 16] (Fig. 2).

**Table 1** Devices and their average lifespan

| Device | Avg. lifetime in years |
|---|---|
| Desktop | 5 |
| Laptop | 4 |
| Television | 10 |
| Mobile phones | 6 |
| Printer | 4 |

**Table 2** Estimation of total e-waste produced yearly

| Device | 2014 | 2015 | 2016 |
|---|---|---|---|
| Desktop PCs | 59,558 | 66,429 | 67,102 |
| Laptops | 12,640 | 15,248 | 20,673 |
| Mobile phones | 22,919 | 23,101 | 28,827.5 |
| Televisions | 130,200 | 145,800 | 168,000 |
| Total e-waste generated annually | 225,317 | 250,578 | 284,602.5 |

**Fig. 1** e-wastage of various devices from 2014 to 2016



**Fig. 2** Recycling of e-waste

## 5.1 Formal Recyclers

The registered organizations having proper licensing for execution of e-waste recycling are called as formal recyclers. Environmental Protection Bureaus are the authorizing organization for issuing licences. The State Environment Protection Administration (SEPA) provides licences for recycling of materials and that classified under hazardous waste are also sometimes necessary [17, 18].

## 5.2 Informal Recyclers

The participants are considered outside the official regulatory bodies. The activities carried out by these parties cannot be tracked as no official record is being maintained; therefore, they are called illegal [17, 19, 20].

The prime focus of this informal recycler is to remove material from waste product and refurbish. The reason behind increase in informal recycling is mentioned below:

i. Customer unwillingness to pay for disposing of their electronic products.
ii. Unregulated importance of second hand products.
iii. Due to lack of awareness among customers, collectors, recyclers about vulnerable exposure to wastes.
iv. Lack of proper tactics, infrastructure and organization to execute extraction and disposal processes.
v. Lack of interest in e-waste.

The practical implementation for recycling of e-waste through devices such as e-dustbin using sensor system can segregate wastes that can be recycled and reused. This will contribute towards reduction in pollution to our environment and shall motivate for societal development.

## 6 Conclusion

There are many factors that are associated with the end-of-life disposal of e-waste management. The complexity of the materials used, occurrence of hazardous substances, lack of awareness, legislative requirements, availability of technologies, supply chain uncertainty are some of the major issues pertaining with e-waste management. Hence, it is a challenge to us to establish proper line-up to pave a sustainable path of future to ensure green computing. Progress in research and development is increasing, and formal recycling is gaining pace every day. This, on the one hand, ensures recyclability and reusability, thereby enhancing the life cycle of the electronic product. So informal sector is beneficial as well as it is harmful. Integration of informal sector with formal sectors and training them properly could be a viable solution. Awareness among the consumers is increasing. Earlier, the consumers used to care about only speed and price while buying computers. But as Moore's Law marches on and computers commodities, consumers will be choosy enough about being green. There are problems, and there are solutions to this e-waste problem. There are economical and social issues that drive the skilled informal sector which is the real deal to the environmental problem. Hence, the sustainability of the e-waste management system will depend on how well this informal sector is tackled. In other words, the future lies in their hands, and

arguably, it can be commented that optimization and negotiation between the government and this thriving sector will certainly lead us to a sustainable future.

# References

1. Anam, A., Syed, A.: Green computing: e-waste management through recycling. Int. J. Sci. Eng. Res. **4**(5), 1103–1106 (2013)
2. Debnath, B., Baidya, R., Biswas, N.T., Kundu, R., Ghosh, S.K.: E-waste recycling as criteria for green computing approach: analysis by QFD tool. In: Maharatna K., Dalapati G., Banerjee P., Mallick A., Mukherjee M. (eds) Computational Advancement in Communication Circuits and Systems. Lecture Notes in Electrical Engineering, vol. 335. Springer, New Delhi (2015)
3. E-waste, The escalation of global crisis by TCO (2015). http://tcodevelopment.com/news/global-e-waste-reaches-record-high-says-new-un-report/. Accessed 5 Oct 2015
4. Manomaivibool, P.: Extended producer responsibility in a non-OECD context: the management of waste electrical and electronic equipment in India. Resour. Conserv. Recycl. **53**(3), 136–144 (2009)
5. Baud, I., Grafakos, S., Hordjik, M., Post, J.: Quality of life and alliances in solid waste management. Cities **18**, 3–12 (2001)
6. Empa, E.: E-waste pilot study Delhi: knowledge partnerships with developing and transition countries. St. Gallen: Empa; 2004. Fagerberg Jan, Mowery David C, Nelson Richard R., The Oxford Handbook of Innovation. Oxford University Press (2006)
7. Hazardous Wastes (Management and Handling) Amendment Rules, 2003. www.cpcb.nic.in. Accessed Aug 2010
8. Sunil, Heart: Environmental impacts and use of brominated flame retardants in electrical and electronic equipment. Environmentalist **28**, 348–357 (2008)
9. Xia, Huo, et al.: Elevated blood lead levels of children in Guiyu, an electronic waste recycling town in China. Environ. Health Perspect. **115**(7), 1113–1117 (2007)
10. Amit, Jain, Rajneeth, Sareen: E-waste assessment methodology and validation in India. J. Mater. Cycles Waste Manag. **8**, 40–45 (2006)
11. Sepúlveda, A., Schluep, M., Renaud, F.G., Streicher, M., Kuehr, R., Hagelüken, C., Gerecke, A.C.: A review of the environmental fate and effects of hazardous substances released from electrical and electronic equipments during recycling: examples from China and India. Environ. Impact Assess. Rev. **30**, 28–41 (2010)
12. Srivastava, P.K., et al.: Stakeholder-based SWOT analysis for successful municipal solid waste management in Lucknow, India. Waste Manag. **25**, 531–537 (2004)
13. Kiruthiga, P., Kumar, T.V.: Green computing—an ecofriendly approach for energy efficiency and minimizing E-waste. Int. J. Adv. Res. Comput. Commun. Eng. **3**(4), 6318–6321 (2014)
14. Ladou, J., Lovegrove, S.: Export of electronics equipment waste. Int. J. Occup. Environ. Health **14**(1), 1–10 (2008)
15. Lout, Q., Wong, M., Cai, Z.: Determination of polybrominated diphenyl ethers in freshwater fishes from a river polluted by e-waste. Talanta **72**(5), 1644–1649 (2007)
16. Agarwal, S., Nath, A.: Green computing—a new horizon of energy efficiency and electronic waste minimization: A global per-spective. In: Proceedings of IEEE in International Conference on Communication Systems and Network Technologies, pp. 688–693 (2011)
17. Chen, A., Dietrich, K.N., Huo, X., Ho, S.: Developmental neurotoxicants in e-waste: an emerging health concern. Environ. Health Perspect. **119**, 431e8 (2010)
18. United Nations University (UNU), Solve the E-waste Problem (StEP), Massachusetts Institute of Technology (MIT), National Center for Electronics Recycling (NCER). World e-waste

map reveals national volumes, international flows. 2013. https://www.vie.unu.edu/file/get/11505.pdf. Accessed 14 Oct 2014

19. Balde, C., Kuehr, R., Blumenthal, K., Fondeur Gill, S., Kern, M., Micheli, P., et al.:. E-waste statistics Guidelines on classification, reporting and Indicators. United Nations University (2015)

20. Barba-Gutiérrez, Y., Adenso-Diaz, B., Hopp, M.: An analysis of some environmental consequences of European electrical and electronic waste regulation. Resour. Conserv. Recycl. **52**(3), 481–495 (2008)

# Energy Aware Task Scheduling Algorithms in Cloud Environment: A Survey

Debojyoti Hazra, Asmita Roy, Sadip Midya and Koushik Majumder

**Abstract** Cloud computing is a developing area in distributed computing and parallel processing domain. Popularity of cloud computing is increasing exponentially due to its unique features like on-demand service, elasticity, scalability, and security. Cloud service providers provide software, platform, high-end infrastructure, storage, and network services to its customers. To provide such services to its customers, all cloud resources need to be utilized in the best possible way. This utilization is efficiently handled by task scheduling algorithms. Task schedulers aim to map customer service requests with various connected resources in a cost-efficient manner. In this paper, an extensive study of some scheduling algorithm that aims to reduce the energy consumption, while allocating various tasks in cloud environment is done. The advantages and disadvantages of these existing algorithms are further identified. Future research areas and further improvements on the existing methodologies are also suggested.

D. Hazra · A. Roy · S. Midya (✉) · K. Majumder
Department of Computer Science and Engineering, West Bengal University of Technology,
Kolkata, India
e-mail: sadip20@gmail.com

D. Hazra
e-mail: debojyoti.hazra22@gmail.com

A. Roy
e-mail: asmitaroy2002@gmail.com

K. Majumder
e-mail: koushik@ieee.org

# 1 Introduction

Cloud computing is paradigm shift from traditional distributed computing that is capable of delivering both application as service over the Internet as well proving hardware and software services [1–5]. Users send some cloud service requests to a cloud server via Internet. Cloud service providers alias CSPs provide several cloud services like Software, Platform, and Infrastructure as a service to its customers on basis of user demand and in pay-as-you-use basis. To provide these services, a CSP has a large pool of computing resources. When any user request is generated for a CSP, it assigns the request to any of its resources and sends the result back to the user. The process of assigning the resource for a request is performed dynamically. Resources are managed systematically to provide real-time service and improve system performance. Enhancement of system performance is dependent on various factors like scheduling of incoming task requests, workflow, and jobs.

Users' requests are referred to as tasks [4, 5]. To enhance system performance, CSPs need to automate task and resource mapping process. This automation is done by task scheduling algorithms. Task scheduling procedures play a major role in deciding the quality of experience (QoE) for users [6, 7]. Task schedulers reside on a CSP's main system where task requests appear. On receiving a task, the scheduler selects a resource according to the task's need and scheduling policy. Then, that task is mapped with the resource for execution. When execution is over and outcome is sent back to user, then user-resource mapping is unmapped. In case, there is no task scheduler resource mapping and utilization will not be satisfactory.

This paper mainly focuses on reduction in energy consumption by resources. According to [8], energy consumed by various cloud data centers is estimated to be 26 GW. This is 1.4% of worldwide energy consumption, and this value is increasing day by day with an increase rate of 12% per year. From this large amount of energy consumption by cloud servers, 34% energy is consumed by processors. Processors always consume some amount of static energy. The energy required to complete the scheduled task varies depending on which frequency a processor is operating. When frequency is high, then energy consumption is high and tasks are executed faster. In low frequency, energy consumption is low and tasks need more time to successfully execute.

The remaining paper is arranged into following sections. In Sect. 2, literature survey on different existing energy aware task scheduling is done. In Sect. 3, a comparative study is done along with it some future research areas are identified and Sect. 4 concludes the paper.

## 2 Literature Review of Existing Task Scheduling Algorithms

### 2.1 Energy Aware Genetic Algorithm

In [9] authors proposed energy aware genetic algorithm for task scheduling, which aims in reducing the consumed energy by the resources using a two-way strategy. In Fig. 1, the working strategy of this algorithm is depicted. This algorithm considers makespan and energy consumption of a task for designing the scheduling policy. Energy consumption of a task is estimated using Dynamic Voltage Scaling (DVS). This algorithm considers small sets of incoming tasks as chromosomes. These chromosomes are initialized, and fitness functions for them are calculated. The selection operation is then applied on these functions. Fitness calculation and selection processes are carried out on the basis of two strategies namely, Energy consumption Time Unify Genetic Algorithm (ETU-GA) and Energy consumption Time Double Fitness Genetic Algorithm (ETDF-GA). EUA-GA converts energy parameter to time parameter to calculate fitness function. ETDF-GA considers two fitness functions, one is the reciprocal of makespan, and the other is the reciprocal of total energy consumption. Single-point crossover and mutation operation is then applied on these functions. In this algorithm, only energy parameter and makespan are considered for scheduling a task. However, other QOS parameters like cost,



**Fig. 1** Task and resource mapping with Energy Aware Genetic Algorithm

deadline, waiting time, and arrival time are not considered. Thus, the algorithm may not be effective in a versatile situation.

## 2.2 Energy Saving Task Scheduling Depending on Vacation Queue

The task scheduling scheme in [10] proposes an energy saving algorithm which depends on vacation queue. The resources in virtual machine of cloud system are not always in use. So when there is no task in queue, the resources go to idle state, and finally, it goes off to sleep mode where energy consumption is minimum. Again when tasks arrive, it resumes back to alive state. This span is called vacation. Figure 2 shows the processors steps for this scheduling policy. This scheme considers that the time required for executing a task is proportional to total energy consumed by that task because energy consumption of a task depends on how long that task is operating on CPU.

This task scheduler operates on Dynamic Voltage Frequency Scaling (DVFS) and Dynamic Power Management (DPM) technologies to reduce total power consumption by a task. Every server's power consumption is reduced by Dynamic Power Management voltage, and frequency is controlled by Dynamic Voltage Frequency Scaling. The vacation queuing model is suitable for heterogeneous cloud computing environment. However, this algorithm considers only the energy consumption of a task while taking scheduling decisions other factors such as cost and deadline of the task are not considered. So adaptability of this algorithm from user perspective is not very good.

## 2.3 Toward Energy Efficient Scheduling for Online Tasks in Cloud Data Centers Based on DVFS

In [11], authors have assumed a cloud cluster, which is a collection of multicore homogeneous servers. Each processor has two types of power consumption. One is static power consumption which indicates constant power consumption when processor is turned on. The other is dynamic power consumption which is

**Fig. 2** Processor states in Vacation Queue Strategy

dependent on frequency of the processor in which it is executing. This dynamic power can be controlled by the system.

Power consumption is proportional to frequency of the processor in which it is operating. From this, it can be derived that when a processor is running in high frequency then energy consumption by the processor is high. Here author has used Linux kernel 2.6.0 subsystem CPUFREQ to control CPU frequency by software application. Now the online task scheduling problem is solved by two different DVFS models. In performance model, they have established an equation where they found a relation between CPU utilization and frequency of the processor. CPU utilization is inversely proportional to frequency at that time. The product of CPU utilization and frequency is equated to a variable termed as Conversion Parameter. When this conversion parameter is greater than 1, then lowering the frequency value will save more power by stabilizing CPU utilization in turn reducing the dynamic power consumption of a task scheduler.

However, there are some drawbacks with this algorithm they assumed that in a multicore server system, each core will run in same frequency, which is always not true. It is difficult to segregate the cores of a multicore processor and analyze their frequency. They also assumed that deadline of a task will re-evaluate itself when the frequency of the processor is reduced to save power consumption.

## 2.4 An Innovative Energy-Aware Cloud Task Scheduling Framework

In [12], the authors have proposed an energy aware task scheduling (EATS) skeleton for CSPs' data centers. Here DVFS and VM reuse techniques are combined to control dynamic energy consumption, and a task's deadline will not be affected by their framework. Here authors have collected number of physical elements connected with a cloud data center, and from them, number of VMs can be launched. Again lists of acceptable discrete frequencies are created for a PE in which it can run. Upon these DVFS and VM reuse policies, First Fit Decreasing (FFD) or Weighted Round Robin (WRR) is applied. VM reuse mainly focuses on reusing those working VMs rather than launching new ones which will save energy otherwise required for launching a new VM. In VM reuse step, when new tasks arrive, then scheduler checks if that new task can be executed with an older working virtual machine or there is need of launching new VM. To use older VM, scheduler checks if that particular VM has the capacity to execute that task request. If all the ongoing VMs are checked and no reusable VM is figured out, then a new VM is launched for that task. After a VM is allocated, DVFS of the scheduler is calculated to adjust the CPU frequency to reduce energy consumption for a processor. On deciding the CPU frequency using the EATS-FFD technique, more demanding requests are mapped with first fit PE. In EATS-WRR as it is known that RR works in time sharing domain, so every request is mapped to PE in time sharing manner.

This scheduler framework is very easily adoptable and innovative one. Energy consumption of data centers can be minimized by this scheduler. However, in this paper the author has conducted the experiment for two basic scheduling algorithms. Here newer more complicated well-performing algorithms need to be executed in place of FFD or WRR to get more satisfying outcomes.

## 2.5 Energy Aware Scheduling of HPC Tasks in Decentralized Cloud Systems

In [13], the authors designed an energy aware global task scheduling algorithm for high-performance computational tasks. Like the above-discussed algorithms, dynamic energy consumption by processor is controlled by DVFS. This algorithm focuses to work in multicloud systems and where the data centers are decentralized. It takes its mapping decision by considering the status of resource. Tasks are arranged in Directed Acyclic Graph (DAG) where a task request may have some successors and predecessors. These successors are dependent on its predecessor. For a task, five parameters like number of VMs, computing power per VM, volume of disk image, output data volume, and deadline are defined with the help of these parameters decentralized multicloud scheduling framework is applied to those tasks. Total execution time period is determined by this algorithm, and task allocation delays are also considered. The static energy and dynamic energy consumed by a processor is the total energy consumption of the processor. Dynamic energy
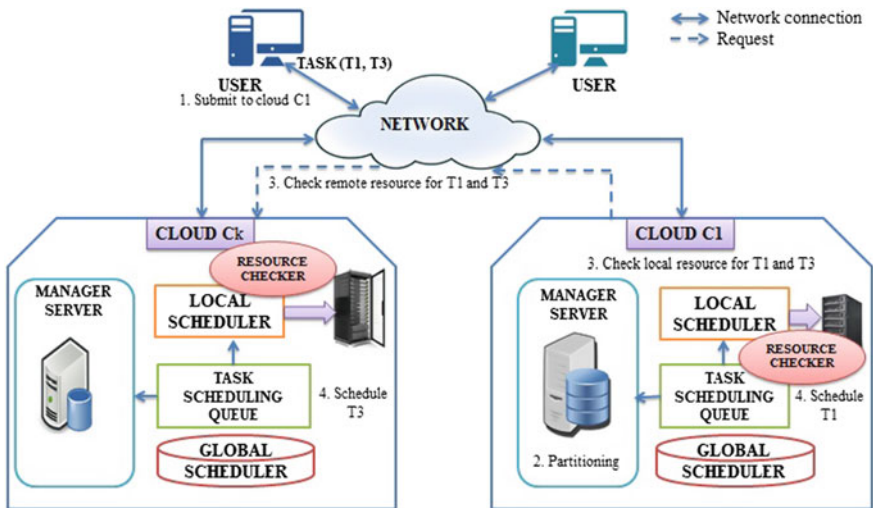


**Fig. 3** Multicloud system structure

consumption is dependent on frequency on which the scheduler is running. In multicore system for all the cores, energy consumption is calculated.

Then, tasks are mapped with VM by checking if their deadline can be met. Otherwise, that task gets rejected. In the above Fig. 3, multicloud system architecture is shown. This process minimizes energy consumption with high user satisfaction, and dependencies of various tasks are considered. However, though deadline is considered for scheduling it only checks if a task will met its deadline during task allocation or not. It does not give higher precedence for tasks with lesser deadline value.

## 3 Comparative Analysis

Table 1 shows a comparative analysis of these existing algorithms. From the above comparative study, some future scope on those existing algorithms is derived. In the paper of Changtian et al. [9], a scheduler is designed that calculates DVFS value, and then, GA is applied on the acquired DVFS value. However, the disadvantage of GA is it does not guarantee to find the global maxima of the problem. To overcome this problem and increase the efficiency of scheduler, simulated annealing (SA) can be applied to the scheduling policy. Cost of executing a task and deadline of a task, which decide the lifetime for every task, are two important parameters, which determine the Quality of Service and Quality of Experience of every user while executing task in a schedulers. Changtian et al. [9], Cheng et al. [10] does not consider cost and deadline while deciding every task's priority, and Alsughayyir et al. [13] does not consider cost while Huai et al. [11], and Alahmadi et al. [12] does not considered deadline parameter while taking scheduling decision for every task. Makespan is a time period that indicates in how much time all allocated tasks will be executed. In schedulers of Changtian et al. [9], Huai et al. [11], Alahmadi et al. [12], and Alsughayyir et al. [13], makespan is not considered for task scheduling. All the above-mentioned schedulers except scheduler of Alahmadi et al. [12] aim to reduce dynamic energy consumption by applying DVFS level. In [12], authors aim to reduce energy by VM reuse policy applying this technique to other scheduling policy will further result in lowering the energy consumption. All above-mentioned schedulers except scheduler of Alsughayyir et al. [13] are designed for single cloud system. In real-life scenario, cloud systems are multicloud structure working in collaboration. So these schedulers need to be modified so that they can accommodate with multicloud system and implemented in real-life scenario.

A well-designed task scheduling algorithm always sticks its focus to both user and service provider satisfaction. From the above study, this can be stated that most algorithms have certain drawbacks. These drawbacks cannot be fully overcome; however, they can be minimized to deliver a better task scheduling policy with high QoS for end users.

**Table 1** Comparison between above-discussed algorithms

| Scheduling name | Comparison parameter | Advantage | Disadvantage |
|---|---|---|---|
| Energy aware GA [9] | Energy consumption | 1. Power consumption of task schedulers is reduced | 1. Other factors like deadline, cost makespan are not considered for better QoS to its users<br>2. GA's solution is bounded to local maxima |
| Energy saving task scheduling on vacation queuing [10] | Energy consumption | 1. Energy consumption is reduced as processors goes to sleep mode when no task is available<br>2. Makespan is reduced in this scheduling algorithm thus giving better efficiency for the system | 1. Task performance is not taken care of<br>2. Adaptability is low<br>3. Cost and implementation policies need to be developed for better QoE value |
| Toward energy efficient scheduling for online tasks in cloud data centers based on DVFS [11] | Frequency, CPU utilization | 1. Energy consumption is low | 1. In real applications, it is difficult to get homogenous servers<br>2. Deadline of a task is assumed to be adjusted, which is not possible in real-time applications |
| Innovative Energy-Aware Cloud Task Scheduling Framework [12] | Frequency | 1. An easy to implement framework<br>2. Cost for task execution is minimized | 1. The framework is not scalable to various scheduling algorithm<br>2. Deadline parameter is not considered; thus, task reaching deadline cannot be handled efficiently |
| Energy aware scheduling of HPC tasks in decentralized cloud systems [13] | Frequency, deadline | 1. Energy consumption is reduced<br>2. Tasks are mapped with resources by checking deadline constraint, so deadline miss is partly decreased | 1. Does not follow prioritized scheduling |

## 4 Conclusion

In cloud computing, resources work collectively or singly to create a VM. These VMs are mapped with tasks virtually based on the need of users. In order to manage these VMs, good scheduling algorithms are required. Energy consumption by different computing resources is huge, and a lion's share of this consumed energy is used for tasks execution. Tasks schedulers take the responsibility to map tasks to its corresponding resource. This mapping is desired to be carried out in an energy

efficient manner. So, that the total energy consumption by cloud systems reduces considerably. Throughout this paper, a brief study is done on different existing energy aware task scheduling algorithms. Most of the existing energy aware algorithms do not consider deadline of a task and cost of executing a task while taking the scheduling decisions. So these algorithms could be further developed to get better result considering important parameters like reducing the number of deadline miss for a task, the cost to complete a task. Comparative study included in this paper will help the CSP to choose an algorithm according to the kind of services they are going to provide.

# References

1. Mathew, T., Sekaran, K.C., Jose, J.: Study and analysis of various task scheduling algorithms in the cloud computing environment. In: International Conference on Advances in Computing, Communications and Informatics (ICACCI), pp. 658–664 (2014)
2. Salot, P.: A survey of various scheduling algorithm in cloud computing environment. Int. J. Res. Eng. Technol. **2**(2), 131–135 (2013)
3. Arya, L.K., Verma, A.: Workflow scheduling algorithms in cloud environment—a survey. In: Recent Advances in Engineering and Computational Sciences, pp. 1–4 (2014)
4. Dave, Y.P., Shelat, A.S., Patel, D.S., Jhaveri, R.H.: Various job scheduling algorithms in cloud computing: a survey. In: International Conference in Information Communication and Embedded Systems, pp. 1–5 (2014)
5. Fakhfakh, F., Kacem, H.H., Kacem, A.H.: Workflow scheduling in cloud computing: a survey. In: 18th IEEE International Enterprise on Distributed Object Computing Conference Workshops and Demonstrations, pp. 372–378 (2014)
6. Patil, S., Kulkarni, R.A., Patil, S.H., Balaji, N.: Performance improvement in cloud computing through dynamic task scheduling algorithm. In: 1st International Conference on Next Generation Computing Technologies, pp. 96–100 (2015)
7. Nagadevi, S., Satyapriya, K., Malathy, D.: A survey on economic cloud schedulers for optimized task scheduling. Int. J. Adv. Eng. Technol. **4(1)**, 58–62 (2013)
8. Awada, U., Li, K., Shen, Y.: Energy consumption in cloud computing data centers. Int. J. Cloud Comput. Serv. Sci. **3**(3), 145 (2014)
9. Changtian, Y., Jiong, Y.: Energy-aware genetic algorithms for task scheduling in cloud computing. In: Seventh China Grid Annual Conference, pp. 43–48 (2012)
10. Cheng, C., Li, J., Wang, Y.: An energy-saving task scheduling strategy based on vacation queuing theory in cloud computing. Tsinghua Sci. Technol. **20**(1), 28–39 (2015)
11. Huai, W., Huang, W., Jin, S., Qian, Z.: Towards energy efficient scheduling for online tasks in cloud data centers based on DVFS. In: 9th International Conference on Innovative Mobile and Internet Services in Ubiquitous Computing (IMIS), pp. 225–232 (2015)
12. Alahmadi, A., Che, D., Khaleel, M., Zhu, M.M., Ghodous, P.: An innovative energy-aware cloud task scheduling framework. In: 8th IEEE International Conference on Cloud Computing (ICCC), pp. 493–500 (2015)
13. Alsughayyir, A., Erlebach, T.: Energy aware scheduling of HPC tasks in decentralized cloud systems. In: 24th Euromicro International Conference on Parallel, Distributed, and Network-Based Processing (PDP), pp. 617–621 (2016)

# Design of Jitter Spectral Shaping as Robust with Various Oversampling Techniques in OFDM

**Syed Umar, N. Priya, P. Gayathri, T. Subba Reddy**
**and Azeem Mohammed Abdul**

**Abstract** Deformation caused by the jitter acts as a limiting factor for the performance of the OFDM system in high data rate. The letter says that oversampling is used to reduce jitter noise. Two types of techniques fractional oversampling and integral oversampling are considered. The simulation results are compared with the theoretical results of phase jitter analysis showing very precise obligations. Oversampling results in 3 dB reduction of jitter noise for every doubling of the sampling frequency.

**Keywords** OFDM · Jitter · Fractional oversampling

## 1 Introduction

Orthogonal frequency division multiplexing transmitters (OFDM) are used in many radio communications because it is simple and scalable solution for inter-symbol interference caused by multichannel. Increasing interest of the recently optical OFDM tight system (see [1] and references). In the fiber optic system, data transfer speed is much higher compared to RF wireless systems in general. At these high

S. Umar (✉) · N. Priya
Department of Computer Science Engineering, MLR Institute of Technology, Hyderabad, India
e-mail: umar332@gmail.com

P. Gayathri
Department of Information Technology, Gokaraju Rangaraju Institute of Engineering and Technology, Hyderabad, India

T. Subba Reddy
Department of Computer Science Engineering, Narsaraopet Engineering College, Guntur, India

A. M. Abdul
Department of Electronics and Communication Engineering, K L University, Vaddeswaram, India

speeds, timing jitter is a serious obstacle to the implementation of the OFDM system. The main source of sampling jitter of the fast analog to digital clock (ADC) is required in the system. Timing jitter is becoming a problem in high sampling radio frequency band OFDM [2]. The analysis of the timing jitter effects [3, 4]. The documents focus on specific jitter fault, typical colors with phase ring in Phase Locked Loop systems (PLL). They believe that the core samples. OFDM, fractional resampling is used to make such a connection sub-carriers. In the letter, both the expansion oversampling. We fractional matrix and supportive jitter proposed [5] the analysis of the data between the carrier (ICI), a scalable system. The high-speed ADC typically use the parallel architecture of pipelines PLL [6] and nervous white, is the topic of this article, the model is more appropriate.

## 2   Determining Jitter Matrix and Design of System Model

Consider the system is given in Fig. 1. The data in the cyclic prefix OFDM, the transmitter T at any time this mark values complex values $N$ refers to the time required to reduce jitter $N$ placed in different places variable OFDM the opinion writing process, but the jitter of the ADC sample.

Figure 2 shows that the timing jitter. Ideally, the OFDM periodically sampled $T/N$ receives the dotted line in Fig. 2A the same sampling intervalsub real number. The results of the phase jitter caused by other branches $n$ $\tau$ and examples given time interval. Figure 2b $\tau$ $n$ standard dithering soft form. OFDM period jitter degrades system performance, and changes the color of the equalizer numbers called "ideal."



**Fig. 1**  Block diagram of OFDM

[5] It found that the jitter performance no time to express themselves in a matrix jitter variable matrix OFDM time jitter system

$$\mathbf{Y} = \mathbf{WHX^T} + \mathbf{N} \tag{1}$$

X, Y, and n are sent, received, and an additive Gaussian noise (AWGN) vector White, H is the channel response matrix

$$
\begin{aligned}
\mathbf{Y} &= \begin{bmatrix} Y_{-N/2+1} & \cdots & Y_0 & \cdots & Y_{N/2} \end{bmatrix}^{\mathrm{T}} \\
\mathbf{H} &= \mathrm{diag}\begin{pmatrix} H_{-N/2+1} & \cdots & H_0 & \cdots & H_{N/2} \end{pmatrix} \\
\mathbf{X^T} &= \begin{bmatrix} X_{-N/2+1} & \cdots & X_0 & \cdots & X_{N/2} \end{bmatrix}^{\mathrm{T}}
\end{aligned}
$$

$$
W = \begin{bmatrix}
w_{-\frac{N}{2}+1,-\frac{N}{2}+1} & \cdots & w_{-\frac{N}{2}+1,0} & \cdots & w_{-\frac{N}{2}+1,\frac{N}{2}} \\
\vdots & \ddots & \vdots & \iddots & \vdots \\
w_{0,-\frac{N}{2}+1} & \cdots & w_{0,0} & \cdots & w_{0,N/2} \\
\vdots & \iddots & \vdots & \ddots & \vdots \\
w_{\frac{N}{2},-\frac{N}{2}+1} & \cdots & w_{\frac{N}{2},0} & \cdots & w_{\frac{N}{2},\frac{N}{2}}
\end{bmatrix} \tag{2}
$$

$$Y = HX^T + (W - I)HX^T + N \tag{3}$$

$n \times n$ matrix, where the first (3) of the sector in the second period of jitter and noise [5]. Reported that the purpose of the matrix jitter period W

$$w_{l,k} = \frac{1}{N} \sum_{n=-N/2+1}^{N/2} e^{j2\pi k \frac{\tau_n}{T}} e^{j\frac{2\pi}{N}(k-l)n} \tag{4}$$



**Fig. 2**  Jitter working principle

# 3   Theoretical Calculation of Jitter

Now, two fractional studies and the figures show the number of OFDM can be used one or two times higher decay of reduce jitter. As oversampling component is the sampling frequency $MN/T$, where $m$ is an integer small fractional black tape taste used for signal transmission? When the $N$-modulated, the bandwidth of the in-band OFDM signal $N/2T$, for example, $T/M$, as shown. 2. Nyquist sampling frequency. Otherwise, only the frame index $N$ and $L + n$, band the signal $(N_L + N_U)/2T$. In this case, the sampling interval by $T/N$, the Nyquist frequency. Why $/N$. Usually, if the level of oversampling $(N_L+N_U)$ candidate fractional resampling function after the ADC buy test.

$$
\begin{aligned}
y_{n_M} &= y\left(\frac{n_M T}{NM}\right) \\
&= \frac{1}{\sqrt{N}} \sum_{k=-N_L}^{N_U} H_k X_k e^{\left(\frac{j2\pi k}{T} \times \frac{n_M T}{NM}\right)} + \eta\left(\frac{n_M T}{NM}\right),
\end{aligned}
\tag{5}
$$

where $n_m$ is the index of $\eta\, M$ oversampled attention AWGN. Built oversampling instead of FFT $N$-points Receiver "Big" $M$ of $N$ FFT points. FFT output vector of length MN.

$$
y_{l_M} = \frac{1}{\sqrt{M}} \frac{1}{NM} \sum_{n_M=-NM/2+1}^{NM/2} y_{n_M} e\left(\frac{-j2\pi n_M l_M}{NM}\right).
\tag{6}
$$

If the index signals $L\,M\,N$ FFT and $M$. coupling (4) (5) and (6), we can change the weighting for the oversampling

$$
w_{l_M,k} = \frac{1}{NM} \sum_{n_M=-NM/2+1}^{NM/2} e^{j2\pi k \frac{n_M}{T}} e^{j\frac{2\pi}{NM}(k-l_M)n_M}.
\tag{7}
$$

With approximation of $e^j = 1+j$ for small , Eqs. (5) and (7) become

$$
w_{l_M,k} \approx \frac{1}{NM} \sum_{n_M=-NM/2+1}^{NM/2} \left(1 + \frac{j2\pi k \tau_{n_M}}{T}\right) e^{j\frac{2\pi}{NM}(k-l_M)n_M}
\tag{8}
$$

$$
w_{l_M,k} \approx
\begin{cases}
\dfrac{1}{NM} \displaystyle\sum_{n_M=-NM/2+1}^{NM/2} \dfrac{j2\pi\tau_{n_M}}{T} e^{j\frac{2\pi}{NM}(k-l_M)n_M} & k \neq l_M \\[4mm]
1 + \dfrac{1}{NM} \displaystyle\sum_{n_M=-NM/2+1}^{NM/2} \dfrac{j2\pi\tau_{n_M}}{T} & k = l_M
\end{cases}
\tag{9}
$$

$$E\left\{\left|w_{l_M,l_M}\right|^2\right\} \approx 1 + \left(\frac{1}{NM}\right)\left(\frac{2\pi k}{T}\right)^2 E\left\{\tau_{n_M}^2\right\} k = l_M \tag{10}$$

$$y_{l_M} = H_{l_M} X_{l_M} + \sum_{k=-N_L}^{N_U} \left(w_{l_M,k} - I_{l_M,k}\right) H_k X_k + \mathrm{N}(l) \tag{11}$$

$$\frac{P_j(l)}{\sigma_s^2} = \frac{E\left\{\left|\sum_{k=-N_l}^{N_u} \left(w_{l,k} - I_{l,k}\right) X_k\right|^2\right\}}{\sigma_s^2} \tag{12}$$

$$= \sum_{k=-N_l}^{N_U} E\left\{\left|w_{l,k} - I_{l,k}\right|^2\right\}.$$

This is the second noise dither phase. Below we consider the channel flat $k$ $H = 1$, and if the signal strength for each subcarrier fair hand with the help of comparison of average performance {}, so a little 'noise, PJ $(L)$ receives the sub-carrier signal difference at the page level

$$\frac{P_j(l)}{\sigma_s^2} = \frac{\pi^2}{3M}\left(\frac{N_v N}{T_N^2}\right) E\left\{\tau_n^2\right\} \tag{13}$$

when $M=1$ and $Nv = N$ the Eq. (3) becomes as follows

$$\frac{P_j(l)}{\sigma_s^2} = \frac{\pi^2}{3}\left(\frac{N^2}{T_N^2}\right) E\left\{\tau_n^2\right\} \tag{14}$$

When we compare Eqs. (13) and (14), it can be expressed that the combination of these integral and fractional oversampling reduces the jitter noise power by the factor of $Nv/NM$.

## 4   Simulation

It should be noted that various changes jitter, if used, oversampling is that most of the examples of the period jitter of the sample. From Fig. 3 shows the theoretical results and simulation on the average jitter noise on the basis of sampling. The theoretical possibility to be close to agreement. Elevation model 10log10 reduction factor of the jitter noise power, so that every doubling of the sampling frequency, 3 dB reduction in jitter noise. Increase of the upper part of the sub-band with

**Fig. 3** Sampling factor versus avg jitter noise

organic oversampling and run through the extraction tube. This time jitter is dependent on event process leads to a reduction of noise jitter linear current sample rate. 3 offers oversampling reduction of noise performance jitter dB for each doubling of the sampling frequency. It also shows that the current jitter, high frequency, that ICI multilayer, but equally the future.

# 5   Conclusion

To reduce jitter on the sampling theory and simulation of OFDM systems also decrease. Two methods were used Rift oversampling over sampling is achieved by letting the parties Increase of the upper part of the sub-band with organic over-sampling and run through the extraction tube. This time jitter is dependent on event process leads to a reduction of noise jitter linear current sample rate. 3 offers oversampling reduction of noise performance jitter dB for each doubling of the sampling frequency. It also shows that about more than low frequencies, but the presence of time jitter, high-frequency side ICI equally here on all.

# References

1. Armstrong, J.: OFDM for optical communications. J. Lightwave Technol. **27**(1), 189–204 (2011)
2. Syrjala, V., Valkama, M.: Jitter mitigation in high-frequency bandpass- sampling OFDM radios. In: Proceedings of WCNC, pp. 1–6 (2012)
3. Manoj, K.N., Thiagarajan, G.: The effect of sampling jitter in OFDM systems. Proc. IEEE Int. Conf. Commun. **3**, 2061–2065 (2003)

4. Onunkwo, U., Li, Y., Swami, A.: Effect of timing jitter on OFDM based UWB systems. IEEE J. Sel. Areas Commun. **24**, 787–793 (2006)
5. Yang, L., Fitzpatrick, P., Armstrong, J.: The effect of timing jitter on high-speed OFDM systems. In: Proceedings of AusCTW, pp. 12–16 (2009)
6. Sumanen, L., Waltari, M., Halonen, K.A.I.: A 10-bit 200-MS/s CMOS parallel pipeline A/D converter. IEEE J. Solid-State Circuits **36**, 1048–1055 (2010)

# Architecture for the Strategy-Planning Techniques Using Big Data Analytics

**Puja Shrivastava, Laxman Sahoo and Manjusha Pandey**

**Abstract** The rapid growth in technology and market has posed a throat-cut competition among the service providers. Retaining of existing customers in place of catching new one is 90% cheaper, but it needs to know the customer very well, which is possible by analyzing the customer data. To analyze customer data and provide customer-oriented services, this paper recommends an architecture for the development of techniques which would further be able to design strategies, such as tariff plan, on the basis of available information from the customer relationship management system of any service-based company, which is finally known as customer-oriented data product. This architecture has five phases with data source, data collection, data refinement, analysis of collected data, and generation of data product. This paper summarizes the requirement of data sets and systems to develop strategy-planning techniques with the study of available architectures in the big data and CRM environment.

**Keywords** BSS · OSS tariff plan · CRM · Big data · Data mining
Clustering · Frequent-pattern

## 1 Introduction

The current trend of market has posed very tough competition among the service providers, and to win this war, companies are designing and adapting new business strategies which are beneficial to both customer plus company with the help of big

P. Shrivastava (✉) · L. Sahoo · M. Pandey
School of Computer Engineering, Kalinga Institute of Industrial Technology University,
Bhubaneswar, Odisha, India
e-mail: pujashri@gmail.com

L. Sahoo
e-mail: laxmansahoo@yahoo.com

M. Pandey
e-mail: manjushapandey82@gmail.com

data analytics. In present scenario, big data analytics seems most promising equipment, due to the amount of data with variety and velocity—real-time data. Data of customer relationship management (CRM) system of any service-based company are a pile of worth, if analyzing techniques are applied on it, to find some patterns which reveal the behavior of customer that motivates, to management and strategy planners, for designing of customer-oriented schemes. Business models have been moved from product-oriented concepts to customer-oriented products, and due to the technology, it is possible to know customer in a fast and better manner. The traditional customer relationship management systems were mainly concerned with the customer details, billing operation, service usage, etc., and the system architecture was the conventional one. Big data technology demands upgradation in the architecture, software, and analytic techniques, since use of big data technology in customer relationship management systems provides a complete picture of the customer. This picture includes details of customer and behavior of customer. An architecture proposed in [1] consists of three modules: first is concurrent analysis, second is storage apparatus, and third is dealing out tools in the big data ecosystem including collaborative technologies, analytical technologies, and operational technologies. Hadoop, Hbase, and MongoDB are available tools to be used in the new CRM model to provide solutions for processing and storage of data. Organization of this paper discusses about telecom CRM in second part since as a case of CRM is taken in this research paper, third section discusses the state of the art, fourth part describes the recommended architecture, and fifth segments include brief review of big data mining techniques applied to be on the customer relationship data, and last section is the conclusion and future scope.

## 2 Customer Relationship Management System of Telecom Industry

Telecom market is one of the highest varying and rising markets, and further for being more productive and get better customer retention, operators are focusing on the creation of multiway customer understanding with reliable business processes. To make honest relationships with customers and craft a good management environment, companies lean to implement the CRM solutions. CRM in telecommunication ties together a broad variety of tools—from usual sales promotions to generate feedback on customers. The goals of telecom CRM can be named as (1) customer data management, (2) sales opportunities improvement, (3) sales increase, (4) cross-sales augmentation, (5) network profitability enhancement, (6) delivery of true customer service, (7) building of customer loyalty, (8) manage documents, (9) boost productivity, and (10) elevate contact center efficiency [2]. Customer experience is one of the few areas where companies can differentiate themselves, but the backend system frustrates both customer and customer-facing staff. The objective of telecom CRM is to integrate billing support system (BSS)/

operation support system (OSS) to attain a holistic view and understanding of the customer base [3]. CRM forms core of any telecommunications service provider's enterprise architecture due to their focus on current market scenario, customer experience, acceleration of customer acquisition, increasing customer loyalty, and improving partner management. It is a transformation from a product-centric to a customer-centric organization [4, 5]. Customer segmentation is one of the really significant issues in the customer relationship management which can be obtained by different classification and clustering techniques of data mining, neural network, and machine learning to improve economic efficiency of enterprises [6]. Benefits of CRM are as follows [7]:

(1) Identification: clean customer data. Single customer view helps sales force and cross-selling.
(2) Differentiation: recognize customer, cost effectual promotion campaign, and reduction of direct mailing cost.
(3) Interaction: customer satisfaction and loyalty and cost-effective customer service.
(4) Customization: lower cost of acquisition and retention of customer and maximize share of wallet.

It can be concluded that CRM is playing very important role in telecommunication companies by providing insights into customer data through different segmentations and mining techniques. Next section provides the state of the art.

## 3 State of the art

Very few academic works are available in the context of architectures for the big data environment and attached systems. Evolution of big data technology has changed the perception of computer scientists with the shift of technology toward the data-centric architecture and operational models. An architecture called as big data architecture framework (BDFA) is proposed in [8] to deal with all portions of big data environment such as infrastructure, analytics, data structures, models, life cycle of data, and data security. Real-time big data is the current and upcoming challenge, and an architecture is proposed in [9] for the real-time big data generated from the different sensors. The three main elements of this architecture include remote sensing big data acquisition unit (RSDU), data processing unit (DPU), and data analysis decision unit (DADU) where the first unit obtains data from the satellite and sends it to the base station, second unit compiles the data, and third unit provides decisions based on the results received from the DPU. Architectures for customer relationship management system in the context of big data technology is studied in [10] proposed by different service-based companies. It concludes in the three sections: collaborative technologies, operational technologies, and analytical technologies in the ecosystem of big data. Design of a new architecture for the

customer relationship management system for the tourism industry in the context of big data environment is proposed in the [11]. On hand architectures incorporate data source, data compilation, and conclusion making as major sections, and internal operation of these parts get modified according to the data set.

## 4 Projected Architecture

Projected architecture for the development of customer-oriented data product, tariff plan, is given in Fig. 1. This architecture consists of five phases. Feedback generation is considered as first phase where the output of business support systems (BSS) and operation support systems (OSS) provides the input to the second phase which is the data of customer and its activities, system constraints, and the feasibility of the system. Third phase is a detailed study of all three sections of second phase. The fourth phase is analytical engine with accumulated big data analytics techniques to be applied on output of the third phase. Finally, fifth phase generates the desired result as data product which is tariff plan.

**Phase 1**: Feedback generation phase is the first phase of the proposed architecture; it provides the data of the existing telecom CRM system which can be further classified into customer behavior, system constraints, and feasibility study. BSS and OSS are two important components of telecommunication systems for the maintenance and support of customers, abbreviated as BSS/OSS, where business support systems (BSS) contracts with consumers and upholds processes like order taking, invoice processing, fee collection, sales and advertising, sustaining customer care agents in answer to service needs, problem reporting and billing investigation, etc. It is an S/W appliance. Operation support system (OSS) is a network system deals through communications network and processes like preserving network catalog, provisioning services, configuring network apparatus, and supervision of
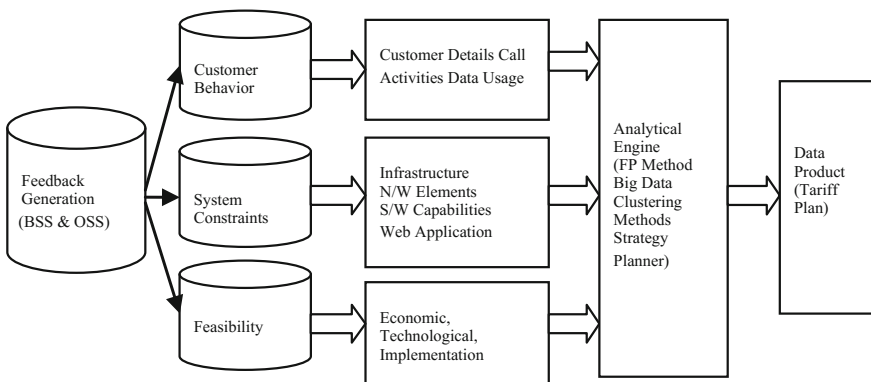


**Fig. 1** Proposed architecture for customer-oriented data product

faults. OSS is combination of hardware and software to hold back-office activities of telecommunication network functions, and maintains customer services. OSS is conventionally used by network planners, service designers, architects, and engineering teams of service provider company [12].

**Phase 2**: The second phase consists of three sections: customer behavior, system constraints, and feasibility study.

**Customer Behavior** includes both the customer details and the service usage. Customer details mean the name, age, gender, address, id proof, profession, etc. The service usage is one of the most important attributes of it because the type of service user availing, call activities, and data usage are big data that will help in discovering the habit of user, so the new strategy can be designed and offered.

**System Constraints** are the limitations of physical infrastructure, software applications, monetary, strategic planning, skill development, and vision.

**Feasibility** discusses the practicability of proposal in light of economics, technology, and implementation.

**Phase 3**: This phase of the proposed architecture goes in detail of each of its three sections. The first section is customer behavior: With the details of customer identification, it keeps records of customer activities such as what kind of service is being used like mobile phone or landline phone, net connection which further can be classified in cable, Wi-fi, or data card plus the data pack recharge in mobile phone. Frequency of calls, duration of calls, type of call like local, national, or international, time of call such as week days plus working hours and nonworking hours plus weekend, and all such call activities are the call activities record of customer. Details of local, national, and international SMS and MMS, their frequencies, daily routine, or festive time are also kept as record. The data usage like user is using apps in place of call or not using apps due to lack of knowledge or expertness; net activities such as social networking or e-commerce; and type of tariff plan and data pack, prepaid, and postpaid; all such records come under the customer behavior. Further, these details are used to classify customers in no. of groups and give some insight into user behavior so the customer-oriented strategies know how to be designed and marketed to retain and magnetize customers.

The second section is system constraints that show the limitation of current system like the lack of physical infrastructure such as the fewer networks of cables, towers, transmitters, operating stations, no adaptation of new technologies, and so on. Deficiency of expertise work force in expertise, management intention, and immature service planning; employee–employer relationship also plays an important role in the growth of any company; bad working environment and unsatisfied employee plus operators can also be considered in system constraints. The monetary gain and expenditure also play important roles.

The third section is the study of economic feasibility, technological feasibility, and implementation feasibility. The economic feasibility is the study of financial condition of both company and customer. For example, company is launching a very good scheme, but the customer is not able to avail that facility, so the initiative will be waste; in the same way, if market is demanding a new launch according to the new technology advancement, but the financially the technology is very costly,

it will not be feasible to commence it. The economic feasibility studies the monetary possibility of the system in the context of both market and customer. The technological feasibility checks the availability of technology; such new device or plan is developed, but the compatible technology is not present in the market; for example, new mobile apps are developed, but the existing range of mobiles are not able to install that apps; another example is of tariff plan with no roaming charges, but the company's network is not all over; so, either there will not be network in company's SIM, or company will have to be overburdened financially by paying other company, so neither technically nor economically it will be possible. The implementation feasibility looks for the execution scenario of the system. Might be there is no shortage of funds and technology, but there is no market, so it will be a big waste. So, the third section of the third phase searches about the feasibility of the system.

**Phase 4**: The fourth phase is the most important phase since it is the analytical engine of the system. It receives input in the form of details of customer, call activities, and log records from the third phase and applies the techniques of big data analytics on the received input which may be either structured, semi-structured, or unstructured data or combination of all three types of data. Here, the proposal is to first refine data by fixing the no. of attributes to be used for analysis, then apply some pattern identification big data algorithms on data to identify patterns, and then use big data clustering techniques on the obtained patterns, so the clusters of customers with similar choice of habits are recognized. Analysis of the shaped clustered provide insight into data, so the customer-oriented data products, tariff plan, be devised to further help management and market people. Figure 2 shows the no. of steps inside the analytical engine for the generation of tariff plan, the subscriber information, call detail records (CDRs), and BSS/OSS data are supplied as input to pattern recognition model; the obtained patterns are supplied as input to clustering methods, and the acquired clusters are input for the decision-making model which is the **Phase** 5 to finally provide data product—tariff plan as output. Memory scalability, work partitioning, and load balancing are three challenges of conventional data mining frequent pattern growth algorithm to be applied in big data [13].

Analytical engine is completely made up of the data analytic techniques. First is frequent pattern growth algorithm based on MapReduce to obtain the behavior of customer as patterns, such as a customer makes daily five international calls of 30 min each in working hours, or customer daily sends more than 100 SMS nations wide for business purposes. Once the patterns are identified, MapReduce based K-Means clustering algorithm is applied, and clusters of customers with similar behavior are obtained. From the density of clusters, revenue is determined, and according to that, tariff plans are generated.

Fig. 2 Proposed no. of steps
inside the analytical engine



# 5 Data Mining Techniques

Frequent pattern mining is a significant job in the discovery of knowledge hidden in the form of repeated blueprints. In the circumstances of big data technology, frequent pattern mining algorithms suffer from the problems of memory scalability, work partitioning, and load balancing. These problems of conventional pattern mining algorithms can be solved with the MapReduce distributed systems where the put in data ought to fit in the comprehensive system memory and available to every process. The HDFS (Hadoop distributed file system) ensures ample space for large input data set. Big data clustering techniques include K-Means algorithm in the MapReduce program structure. Available big data clustering techniques are represented in Fig. 3 with the single-machine clustering techniques and multi-machine clustering techniques [14].

A parallelize frequent pattern growth algorithm (PFP) is to be applied on distributed machines where the large volume of data is stored in distributed file system. Implementation of PFP-growth algorithm in MapReduce framework first converts the database of customers into database dependent on call activities and customer characteristics, and then, construction of FP-tree and FP-growth are done recursively in reduce phase [15]. MapReduce-based K-Means and Fuzzy C-Means algorithm to obtain the clusters of similar patterns (behavior) [16] will be more suitable [17].

**Fig. 3** Big data clustering techniques

## 6   Conclusion and Future Scope

Current era is of customers. Service providers are fighting hard to retain their customers by offering customer-oriented services, since marketing and attracting new customer is 90% costlier than retaining the existing customers. This paper discusses the generation of, analysis based, tariff plans for telecom company by introducing a model which takes customer information, call detail records plus the data of BSS/OSS, and finds out the patterns which reveal the behavior of customers; after that, it discovers the clusters of customers who have similar behavior in their call activities and data usage. Once the clusters are obtained, tariff plans are designed for the customers of similar manners such as less charges on international calls for customer who is making more than 5 international calls everyday. State of the art is also considered. Discussed model is applicable in most of the service-oriented businesses, in understanding their customers and offering services according to their requirements, with little modifications in the no. of attributes and the algorithms used, since according to the data set, analysis algorithms get changed. Still, the requirement of standard generalized model is alive.

# References

1. Daif, A., et al.: Review current CRM architectures and introducing new adapted architecture to big data. In: IEEE 2015 (2015)
2. https://www.bpmonline.com/telecom
3. https://www.ericsson.com/ourportfolio/products/telecom-crm
4. http://www.tcs.com/enterprise-solutions/Pages/Telecom-CRM.aspx
5. http://www.crmnext.com/industries/telecom/
6. Zhang, T., et al.: Case study on cluster analysis of telecom customers based on consumers' behavior. In: IEEE 2011, pp. 1358–1362 (2011)
7. Haridasan, V., Venkatesh, S.: CRM implementation in Indian telecom industry-evaluating the effectiveness of mobile service providers using data envelopment analysis. Int. J. Bus. Res. Manag. (IJBRM) **2**(3), 110–127 (2011)
8. Demchenko, Y., et al.: Defining architecture components of the big data ecosystem. In: IEEE 2014, pp. 104–112 (2014)
9. Rathore, M.M.U., et al.: Real-time big data analytical architecture for remote sensing applications. IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens. **8**(10), 4610–4621 (2015)
10. Dair, A., et al.: Review current CRM architectures and introducing new adapted architecture to big data. In: IEEE 2015 (2015)
11. Fardoie, S.R., et al.: A new design architecture for e-CRM systems (case study: tour package choice in tourism industry). In: ICMIT, IEEE 2008, pp. 463–468 (2008)
12. http://www.ossline.com/2010/12/definition-oss-bss.html
13. David, C., et al.: Big data frequent pattern mining. glaros.dtc.umn.edu/gkhome/fetch/papers/pfpm14TR.pdf
14. Zerhari, B., et al.: Big data clustering: algorithms and challenges. https://www.researchgate.net/publication/276934256
15. Zhou, L., et al.: Balanced parallel FP-growth with MapReduce. In: IEEE, 2010, pp. 243–246 (2010)
16. Eren, B., et al.: A K-means algorithm application on big data. In: WCECS 2015 (2015)
17. Fahad, A., et al.: A survey of clustering for big data: taxonomy and empirical analysis. IEEE Trans. Emerg. Top. Comput. **2**, 267–279 (2014)

# Indian Classical Dance Mudra Classification Using HOG Features and SVM Classifier

**K. V. V. Kumar and P. V. V. Kishore**

**Abstract** Digital understanding of Indian classical dance is least studied work, though it has been a part of Indian Culture from around 200 BC. This work explores the possibilities of recognizing classical dance mudras in various dance forms in India. The images of hand mudras of various classical dances are collected from the internet, and a database is created for this job. Histogram of oriented (HOG) features of hand mudras input the classifier. Support vector machine (SVM) classifies the HOG features into mudras as text messages. The mudra recognition frequency (MRF) is calculated for each mudra using graphical user interface (GUI) developed from the model. Popular feature vectors such as SIFT, SURF, LBP, and HAAR are tested against HOG for precision and swiftness. This work helps new learners and dance enthusiastic people to learn and understand dance forms and related information on their mobile devices.

**Keywords** Indian classical dance mudras · HOG features · SVM classifier
Mudra recognition frequency · Scale invariant feature transform

## 1 Introduction

Indian dance forms are a mirror to rich cultural heritage that existed from the past 5000 years. The name for these classical dance forms is called 'Natya Rasa' as portrayed in the bible of Indian dance 'Natya Shastra.' According to Natya Shastra, there are 108 karanas [1] meaning action of hands, feet, and body. These poses symbolize various physical meaning related to nature, god, and actions. A few hasta mudras are in reference [2].

K. V. V. Kumar (✉) · P. V. V. Kishore
Department of Electronics and Communication Engineering, K L University, Green Fields, Vaddeswaram, Guntur District, Andhra Pradesh, India
e-mail: kumarece405@gmail.com

P. V. V. Kishore
e-mail: pvvkishore@kluniversity.in

In this work, we seek to classify hasta mudras used in various Indian classical dance forms. Creation of dataset for this task is a complex process. For the initial phase of testing, we work on images available on internet dance sites. A careful section of precisely captured images is grouped together to form our training and testing datasets for the classifier. Our dataset is having five different sets of 24 hasta mudras from various dance forms.

The regular image segmentation algorithms such as thresholding [3], edge [4], and color clustering [5] fail to extract full features of the mudras. This is due to occlusions of fingers during capture and coloring used for fingers during performances. General feature extraction models in literature are used to represent the dance mudras. These features are histogram of oriented gradients (HOG), speed up robust features (SURF), scale invariant feature transform (SIFT), local binary patterns (LBP), and haar wavelet features (HAAR).

A comparison of these features on images of dance mudras from web and self-captured data is performed. Support vector machines (SVM) are used as classifier of these features to identify a mudra class in kuchipudi dance mudras. We found that these 24 mudras are the basis for all eight classical dance forms in India. Hence classification in kuchipudi can be extended to other dance forms as well.

## 2 Literature

The idea is to represent Indian classical dance on a digital platform. Indian cultural dance forms are most complex human gestures to be represented in digital format. Feature extraction is most complicated task as the images are full of color, occlusions, and finger closeness. This can be observed in images in Fig. 1. The regular image processing segmentation models such as thresholding and edge detection fail to represent the correct shapes as found in the original images.

The most vibrant and highly used models of segmentation in recent times are active contours. Many models have been proposed in literature [6, 7]. But the basic model sufferers from many drawbacks such as illumination, position of the mask, and number of iterations. We believe focused active contour models with more spatial information using color, texture, and shape have profound effect on extracting the correct segments [8].

A number of methods in literature help in determine shape features. This paper tests five such features and provides an indicator telling the dance mudra classification algorithm for best matching score. A two-decade long challenge for producing an imaging feature that is immune to illumination, noise, scale, orientation, partial occlusions giving accuracy and speed is coming good.

Liang and Juang [9] is proposed moving object classification likes: cars, motorcycles, pedestrians, and bicycle by using local shape and wavelet transform HOG features with hierarchical SVM classification. The proposed method used to test in six video sequences for classification. The computer processing times of the

**Fig. 1** Setup used for capturing dance images of kuchipudi mudras

object segmentation in 79 ms, object tracking in 211 ms, feature extraction and classification in 0.01 ms, respectively.

In Recent years, support vector machine (SVM) classifier with Histogram of oriented gradients (HOG) features are the most popular techniques for vehicle detection [10]. In real time implementation which is important for advanced driver assistance systems applications. To reduce the complexity of the SVM is to reduce the dimensions of HOG features. The proposed method SVM classification for vehicle detection is three times speed up in other detection performance.

The rest of the paper is organized as: Sect. 3 describes the followed methodology for mudra classification. Results and discussion are presented in Sect. 4 with conclusions in Sect. 5.

# 3 Methodology

The experiment involves only dance mudras from kuchipudi dance form as they are the basic structures for formation of any dance. Methodology involves two phases: training phase and testing phase. During training phase 24 dance mudras are used to the train the SVM classifier. The capabilities of SVM classifier are mapped to multiple classes to form a Multi-Class SVM.

## 3.1 Dataset for Kuchipudi Dance Form

The dataset is made from a combination of lab dance mudras and dance mudras images on the websites of Indian art and culture. For each mudra, we have made a set of five images from five different artists. Figure 1 shows the set used for capturing the dataset at university multimedia center.

A mixture of $5 \times 24$ images is used for training and testing. We have two sets of images from our dancers, two sets from dance websites and one set from YouTube video frames.

## 3.2 Banalization of Images and Feature Extraction

Processing easiness for feature extraction calls for this step. The dimensionality is reduced to red plane and local maxima are computed. The local maxima in a $16 \times 16$ block are used as a threshold for that particular block making the process invariant to brightness and contrast. A set of binary sign images is coupled in Fig. 1 in the training side. Shape features are modeled from these binary images. Five feature vectors and their combination are used to extract features from the mudras.

## 3.3 Support Vector Machines

SVM's analyze data and produces binary responses for classification come under a class of supervised learning. The basic SVM classifies only two class problems by projecting a hyper plane between data during training phase. The hyper plane is characterized by a subset of data points acting as support vectors. During training the SVM is presented with example vectors $x_i \in \Re^n, i = 1\ldots l$; $l$ training sample, to label each data sample as either as class label +1 or −1 which forms the indicator vector $y_i \in \{+1, -1\}$. SVM formulates the optimization problem as a decision boundary $D(x)$ such that

$$D(x) = \min_{w,b,\lambda} \left( \frac{1}{2} w^T w + C \sum_{i=1}^{l} \lambda_i \right)$$

$$\text{Subject to } y_i \{ w^T \phi(x_i) + b \} \geq 1 - \lambda_i \text{ with } \lambda_i \geq 0, i = 1, 2, \ldots, l; \tag{1}$$

Where $C$ is a positive constant defining regularization. The terms $w$ and $b$ are weight and bias. $\lambda$ is the misclassification handler. The function $m(x) : x \to \phi(x)$ maps feature vector $x$ to a higher dimensional space. The mapping function $m$ $(x)$ maps $x$ into a dot product of feature space that satisfies $m(x_{i-1}, x_i) = \phi^T(x_{i-1})\phi(x_i)$.

## 3.4   Multi-Class SVM

The most widely used multi-class SVM models are one versus all (OVA), one versus one (OVO) [11], directed acyclic graph (DAG) [12], and error correcting output codes (ECOC) [13]. OVA create $N$ binary SVM's for all categories where $N$ is class number. For a $n$th SVM, only examples in that class is positive and reaming is negative. The computation time is less but at a compromised efficiency. OVO creates a pairwise $0.5N(N-1)$ SVM's and pairwise voting to accommodate new samples for solving multi-class problems. DAG training is from OVO model and testing is from binary acyclic graph model. ECOC disambiguates output binary codes to construct a code word matrix which is compared with generated bit vectors by selecting row as a class having minimum hamming distance. This method gives good classification rates compared to other four at the cost of speed of execution. The slower speed is due to the increased length of code words to disambiguate N classes. The minimum code words in ECOC are $\log_2 N$ to a maximum of $2^{N-1} - 1$ bits. Comparing the multi-class SVM methods from MATLAB implementation, we found ECOC is performs better at optimum speeds. The similarity measure for 24 different kuchipudi dance mudras using computer vision model and machine learning algorithm is executed.

## 4   Results and Discussion

The dataset consists of $24 \times 5$ images of Indian dance form kuchipudi are collected from various sources. These 120 images are contrast enhanced by 20% to smooth pixel values. Feature extraction module is initiated to extract features for the 120 mudras. Each mudra feature is labeled to identify them with a particular class. The class labels are the names of the mudras in kuchipudi dance form. We came to understand that these basic mudras are common to all Indian dance forms listed in [1].

Figure 2 shows the five features used on some mudras. Visual observations of the Fig. 5 provide a better performing feature vector in the set {HOG, SURF, SIFT, LBP, and HAAR}. This particular application of computer vision on human hands indicates HOG as the better performer compared to other four feature extraction models.

The other problem encountered is in the number of feature vectors produced per image mudra. Depending on the hand density in the image frame, the number of features in each feature vector is of different size. To tackle this problem, we modified the feature vectors of all images to a normalized feature size. The normalization process involved selection of important features based on magnitude of feature vectors.

The max pooling algorithm is used to extract useful features from a set of features. The support vector machine is supplied with the normalized feature

**Fig. 2** Feature vector generations from five feature extraction models

vectors from a set of mudras. Training on these set of mudras is given to the SVM. A set of 24 multiple classes of kuchipudi dance mudras is the target vector. Training vector consists of only one mudra set that was perfectly captured during lab trials.

Testing SVM by using the same set of mudras resulted in a 100% match to the class. The class labels are [Pataka, Tripataka, Ardhapataka, Kartarimukha, Mayura, Ardhachandra, Arala, Shukatunda, Mushthi, Shikhara, Chandrakala, Sarpashirsha, Simhamukha, Trishula, Swastikam, Matsya, Kurma, Varaha, Garuda, Shivalingam, Pushpaputam, Karkatam, Kapotam, and Bherunda]. For HOG features, the confusion matrix is shown in Fig. 3.

| | Pat | Trip | Ardp | Kar | May | Ardc | Ara | Shu | Mus | Shi | Cha | Sar | Sim | Tris | Swa | Mat | Kur | Var | Gar | Shi | Push | Kark | Kap | Bher |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Pataka | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Tripataka | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Ardhapataka | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Kartarimukha | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Mayura | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Ardhachandra | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Arala | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Shukatunda | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Mushthi | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Shikhara | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Chandrakala | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Sarpashirsha | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Simhamukha | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Trishula | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Swastikam | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Matsya | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Kurma | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| Varaha | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| Garuda | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| Shivalingam | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| Pushpaputam | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| Karkatam | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| Kapotam | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Bherunda | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |

**Fig. 3** Confusion matrix between lab dataset and YouTube frame dataset for HOG features

The confusion matrix gives the matching metrics between the trained samples and the testing samples. The training samples for SVM and the testing samples were different in Fig. 3. Training was given only once using the lab captured dance images. Testing is done with lab images, web images, and YouTube images. The confusion matrix indicates a 75% match between the mudras of unseen by the SVM.

The mudra recognition frequency (MRF) is defined as number of correctly recognized mudras to total number of mudras used for classification. MRF is improved by increasing the training vector from one mudra set to two mudra sets. For a two mudra set training involving both labs captured and web images the classification of unseen mudras in improved by 14%. The MRF in this case is around 89%.

The reason for misclassification between similar mudras can be attributed to the hand shapes, hand colors, hand orientations, and hand textures. The next best feature vector that has shown maximum MRF is SIFT. For same training and testing vectors, SIFT has a MRF of 96%. However, the MRF is down to 71% for different training and testing samples. In case of multiple training vectors, SIFT showed an increase in MRF by 9%.

The overall MRF's for all five feature vectors used is plotted in Fig. 4. Plot shows MRF values for single same set training and testing (1STTV), single different set training and testing (1DTTV), two same set training and testing (2STTV), and two different set training and testing (2DTTV). Along with the proposed feature extraction models, we also tested SVM with multi-feature extraction models with HOG combination. The plots show an increase in the MRF value for HOG combined features. Other combinations of feature vectors are also tested but the results were not encouraging and were discarded from reporting in the plot of Fig. 4
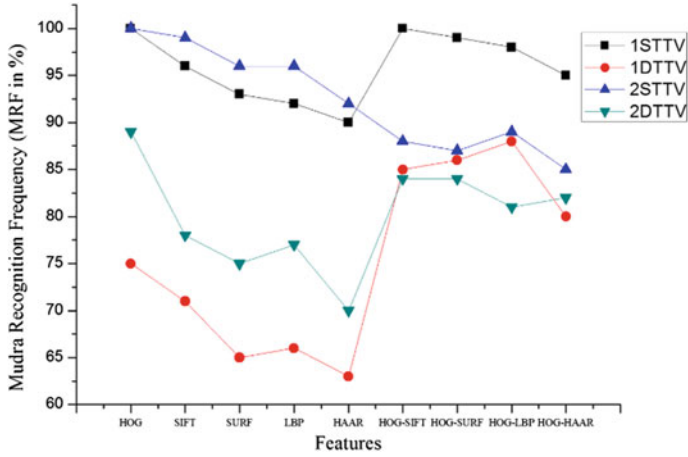
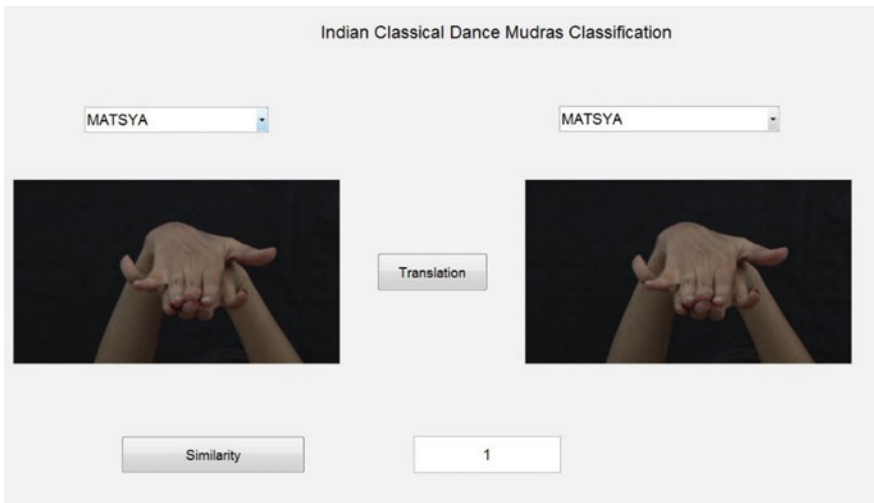**Fig. 4** Plot comparing various feature vectors with respect to MRF



**Fig. 5** GUI showing SSIM for same mudra image of a class for training and testing

The results are validated by measuring a parameter called structural similarity index (SSIM) between the mudras using a graphical user interface (GUI). The GUI has on one side a query mudra and the other side in a mudra of different image along with their names. Figures 5 and 6 shows are the operation of GUI for same mudra images and different mudra images for the same class.
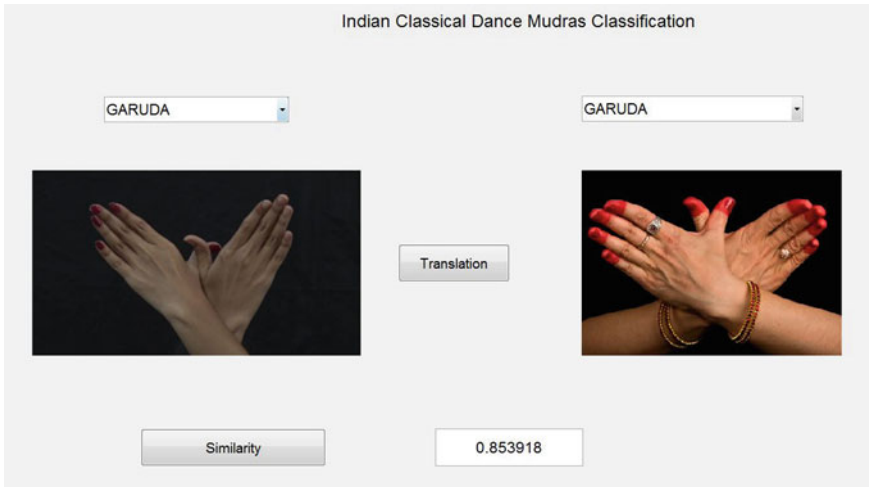
**Fig. 6** GUI showing SSIM for different mudra image of a class for training and testing

## 5 Conclusion

An attempt is made to find similarity between dance mudras of Indian classical dance form kuchipudi based on image processing models and pattern classifiers. Five feature extraction techniques are compared for this work. Multi-class support vector machine classified these features and the performance of the classifier with respect to a particular feature is measured. Visual verification and structural verification using SSIM are preformed to check the classifiers performance. The SVM classifier registered an average MRF of 90% with HOG feature vector and the remaining feature vectors produced less than 80% matching. A GUI is built to validate the results produced by feature vectors using SSIM indicator. The SSIM indicator has closely related the feature models of HOG and SIFT for same and different data sets. Most of the mudras with two hands produced occlusions that induced bottlenecks during feature extraction stage. This model of mudra classification will help enhance the learning capacity of a first time learner.

## References

1. Naidu, N., Naidu, B.V., Pantulu, P.R.: Tandava Lakshanam: the fundamentals of ancient hindu dancing. A Karana in dance is defined as "the coordination of the movements of the hands and feet". New Delhi, Munshiram Manoharlal, 1971, p. 19 (1980)
2. http://natyanjali.blogspot.in/
3. Chatterjee, S.: Matrix estimation by universal singular value thresholding. Ann. Stat. **43**(1), 177–214 (2015)

4. Davis, L.S.: A survey of edge detection techniques. Comput. Graph. Image Process **4**(3), 248–270 (1975)
5. Kishore, P.V.V., Anil Kumar, D., Goutham, E.N.D., Manikanta, M.: Continuous sign language recognition from tracking and shape features using fuzzy inference engine. In: International Conference on Wireless Communications, Signal Processing and Networking (WiSPNET), pp. 2165–2170. IEEE (2016)
6. Kishore, P.V.V., Kishore, S.R.C., Prasad, M.V.D.: Conglomeration of hand shapes and texture information for recognizing gestures of Indian sign language using feed forward neural networks. Int. J. Eng. Technol. (IJET), ISSN: 0975-4024 (2013)
7. Anandh, A., Mala, K., Suganya, S.: Content based image retrieval system based on semantic information using color, texture and shape features. In: International Conference on Computing Technologies and Intelligent Data Engineering (ICCTIDE), pp. 1–8. IEEE (2016)
8. Lowe, D.G.: Distinctive image features from scale-invariant keypoints. Int. J. Comput. Vis. **60**(2), 91–110 (2004)
9. Liang, C.-W., Juang, C.-F.: Moving object classification using local shape and HOG features in wavelet-transformed space with hierarchical SVM classifiers. Appl. Soft Comput. **28**, 483–497 (2015)
10. Lee, S.-H., et al.: An efficient selection of HOG feature for SVM classification of vehicle. In: 2015 International Symposium on Consumer Electronics (ISCE). IEEE (2015)
11. Galar, M., et al.: An overview of ensemble methods for binary classifiers in multi-class problems: experimental study on one-vs-one and one-vs-all schemes. Pattern Recogn. **44**(8), 1761–1776 (2011)
12. Zhang, X., Ding, S., Sun, T.: Multi-class LSTMSVM based on optimal directed acyclic graph and shuffled frog leaping algorithm. Int. J. Mach. Learn. Cybern. **7**(2), 241–251 (2016)
13. Bai, X., et al.: Learning ECOC code matrix for multiclass classification with application to glaucoma diagnosis. J. Med. Syst. **40**(4), 1–10 (2016)

# Feature Extraction Model for Social Images

Seema Wazarkar and Bettahally N. Keshavamurthy

**Abstract** Extraction of appropriate features is a difficult task because it mainly depends on a specific application domain. In this paper, we presented a 5-layered feature extraction model for social images. This model extracts color, texture, geometric, and regional features from given image and also checks presence or absence of people in an image by face detection. Then, normalization of the feature vector is done with the help of priority element. Proposed model is able to deal with the heterogeneous nature of social images. It is useful to get good results in the field of social data analytics.

**Keywords** Feature extraction · Social images · Social data analytics

## 1 Introduction

Rise in number of users of social networks like Facebook, Twitter, etc. makes huge amount of data sprinkle through the social network. This social data need to be analyzed for further use. It is very useful in various sectors such as marketing, security, public health for decision-making, suspicious activity detection and crisis management respectively and many more. It is challenging to analyze social data because it is large in volume, unstructured and having heterogeneous nature. Basically it contains two types of data, content and linkage data. Content data is present in forms like numeric (number of likes, tags), text (comments), images (profile picture, posts), audio, video, etc. Linkage data is about a relation between different users [1].

As image data is one of the most expressive and interesting data type, it plays an important role in content data analysis. Hence, it is important to analyze

S. Wazarkar (✉) · B. N. Keshavamurthy
National Institute of Technology Goa, Ponda, India
e-mail: wazarkarseema@nitgoa.ac.in

B. N. Keshavamurthy
e-mail: bnkeshav.fcse@nitgoa.ac.in

the image data. Information present in image data can reveal many important things to accomplish a given task. For example, if any user is sharing images with negative thoughts on specific topic consistently, then he can be identified as a depressed person. By verifying and taking necessary steps like counseling, one can avoid a severe situation in future.

In image data analytics, image feature extraction is a basic and most important step. Generally, features are represented in the form of vectors. Feature vector is generated by extracting only important representative information from given image. It can be low-level features such as color, texture, edge detection, etc. or high-level features like semantic features. Low-level features are extracted directly from the image, but high-level features are generated with the help of low-level features. These feature vectors are reduced version of the image which is further used as a representative of image [2]. Quality of final results of the data analytics is mainly depends on the features extracted as an initial step. Therefore, feature extraction needs to be done very carefully. Features extracted for one application may not be useful or unable to provide good results for another application. For example, suppose we want to apply clustering on two data sets having (1) scenery images and (2) person images as shown in Fig. 1. If we select color as type of feature for this task, then it will be able to provide good results for scenery image data set but not for the people image data set. We can perform better with the help of semantic features for people image dataset. Computations in feature extraction are not a deterministic, it depends on specific application. Hence, in this paper feature extraction model for operations on the social images is provided. It will be useful for various applications which need to deal with the social images such as sentimental analysis, activity detection, trend analysis, recommendation and so on.

The rest of this paper is organized as follows: Related work is discussed in Sect. 2. Proposed work is presented in Sect. 3. In Sect. 4, experimental results are provided. Paper is concluded, and future directions are given in Sect. 5.



**Fig. 1** Sample of (**a**) scenery image [3] (**b**) person image

## 2 Related Work

Color is an elementary feature but widely used because it is stable and insensitive to the rotation and zooming of an image. Histogram intersection, dominant color, and color correlogram are popular color-based image descriptors [4]. Mabrouk et al. [5] proposed to speed up robust feature (SURF) detector based color extraction method which is further used for the flower image classification by using support vector machine. Experimental work is carried out on dataset provided by the University of Oxford where better results are obtained in terms of execution time and classification rate as compared to method proposed by Nilsback and Zisserman in 2008 [6].

Texture feature provides spatial information of image. Texture feature extraction methods are classified as statistical, structural, and spectral approaches [7]. Stochastic properties of the spatial distribution of gray levels in the image are characterized by spatial method. Gray-level co-occurrence matrix [8], color co-occurrence matrix, and Gabor filter [9] are examples of spatial approach. Structural approach works on the basis of set of primitive texture elements i.e., micro-texture having relationship among them. It is useful for analysis of artificial textures. Spatial domain of image is converted into the frequency domain and vice versa in spectral method. Fourier spectral transform [10] and wavelet transform [11] are the examples of spectral approach.

Shape of object present in given image also provides important information. Hence, it is used to describe the image contents. In 2008, Liu et al. [12] extracted shapes from image on the basis of edge detection and Zhang et al. [13] extracts based on Fourier descriptors with brightness. Pedrosa et al. [14] proposed salient/corner point based shape feature extraction method in 2013. Then, in 2015, Hong et al. [15] presented integral kernel based shape extractor.

Other high-level feature extraction methods such as bag of words [16], deep learning [17, 18], etc. are popular. But, these methods are working on the basis of low-level feature extraction methods. Individual color, texture or shape features are not capable to get good results in many cases. Therefore, layered feature extraction model is proposed in this paper.

## 3 Proposed Work

In this section, proposed 5-layered feature extraction model for social image is presented as shown in Fig. 2 and each layer of the model is discussed. Initially, pre-processing is carried. Then, color, texture, regional, and geometric features are extracted and face detection is done by using five layers. Feature vector generated by earlier step is normalized and updated feature vector is provided as an output of the proposed model.
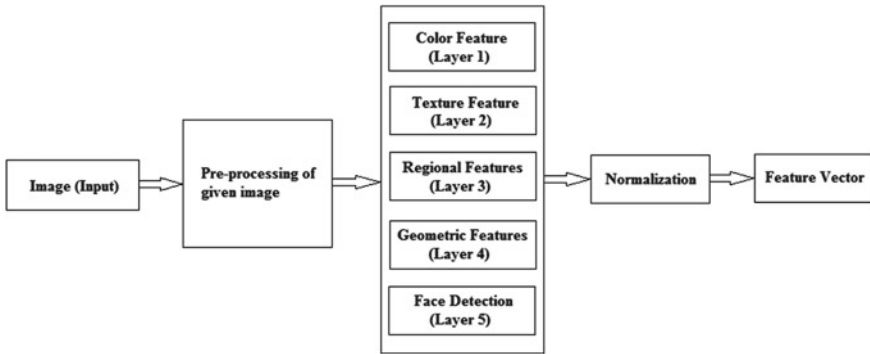
**Fig. 2** Proposed feature extraction model

## 3.1 Pre-processing

Prior to extract features, it is very important to pre-process given data in order to get the good results. In proposed model before providing image to the layer 3 and 4, it is checked that whether given image is RGB image or not. If given image is RGB image, convert it into binary image. Further skeletonization is performed on a binary image and obtained image is zoned into 9 zones as shown in Fig. 3.

## 3.2 Feature Extraction

In this feature extraction model, color, texture, regional, geometric features and face detection are taken into consideration in order to deal with heterogeneity in social images. Details about each layer from the proposed model are discussed below:

**Layer 1**: In layer 1, color features are extracted directly from the original image with the help of histogram. Three components of RGB image i.e., red, green, and blue are separated. In this process, three-dimensional RGB images get converted into three-two-dimensional components. Histogram of those components is illustrated in Fig. 4.

**Layer 2**: Texture features are extracted by using gray-level co-occurrence matrix. It is a statistical approach which represents the spatial relationships among pixels. Statistics like contrast, correlation, energy, and homogeneity are computed by using obtained gray-level co-occurrence matrix. Local variations in that matrix are measured by contrast. Correlation provides joint probability of the pairs of specified pixel. Energy is a sum of squared elements from gray-level co-occurrence matrix as given in Eq. (1). Homogeneity represents the closeness of distribution of elements from gray-level co-occurrence matrix to the diagonal elements of gray-level co-occurrence matrix [19]
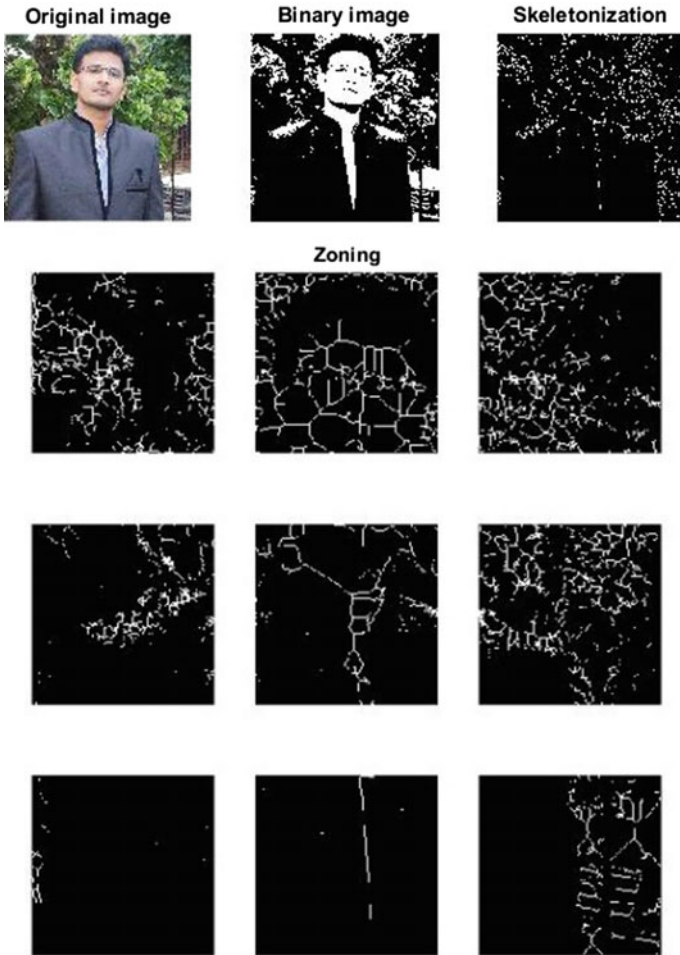
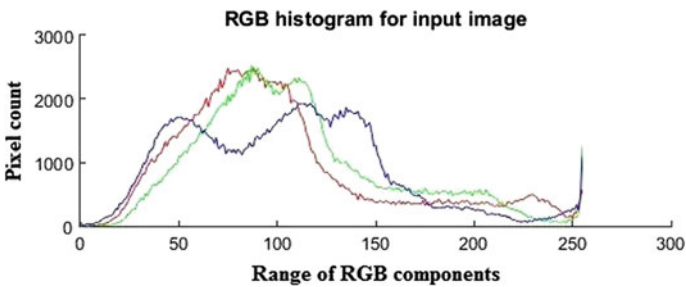**Fig. 3** Pre-processing results of original input image



**Fig. 4** Histogram of original input image

$$E = \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} p_{ij}^2 \tag{1}$$

Where,

$E$  is energy,
$N$  is gray tone,
$P$  is relative frequency of pixels.

**Layer 3**: Regional Features like eccentricity, extent, and orientation are computed in this layer. Eccentricity is the distance between the foci of the ellipse per unit lengthy of major axis. Its value ranges between 0 and 1. Extent is division of number of pixels in the region and number of pixels in the total bounding box. Orientation is the angle suspended by major axis of ellipse on the $x$ axis.

**Layer 4**: In layer 4, to get geometric features, Euler number is computed. Euler number is difference between number of objects and number of holes present in the given image.

$$En = O - H \tag{2}$$

Where,

En  is Euler number,
$O$  is number of objects,
$H$  is number of holes.

**Layer 5**: Face detection is carried out to check the presence of a person in image. Face is detected with the help of extracted facial features. Result of this layer is represented in Fig. 5.



**Fig. 5** Face detection of original input image

### 3.3 Normalization

Normalization of obtained feature vector is done to equalize the effect of all the features. Generally, one of the retrieved features should not dominate the other features. But, here we used priority element to improve the effect of important feature in order to get the good results.

## 4 Experimental Results

Proposed feature extraction model is experimented on the social images taken from three datasets—synthetic, standard, and real-world dataset. Standard dataset titled as MIRFLICKR-25000 [20] is used. For this dataset, images are taken from Flickr. Real-world dataset is generated by collecting images from Facebook [21]. Results obtained at each stage of the model are given below. As a final output of proposed model, feature vector for given image is obtained.

Results of pre-processing at each step such as RGB image to binary image transformation, skeletonization and zoning are shown in Fig. 3. Size of original image is $476 \times 476$. Color features are directly extracted from the original image by using histogram. It contributes in order to bring similar colored objects together. In Fig. 4, histogram for red, green, and blue is represented together. Properties of image like contrast, correlation, energy, and homogeneity are computed by using texture features as discussed in Subsect. 3.2. Obtained values for original image are given in Table 1. Regional features i.e., eccentricity, extent, and orientation of original image are given in Table 2. Face detection in given image is represented in Fig. 5. Priority element with value 2 is used for normalization. Here, priority element is used to improve the importance of a particular feature as per the requirement. Obtained feature vector is 916-dimensional.

As an output of geometric feature, Euler number is computed for original image. Euler number = 273.

To evaluate the proposed model, obtained feature vector is used for the clustering using $k$-means and mean shift clustering algorithm. Accuracy is computed as given in Eq. (3) across all three datasets and shown in Fig. 6.

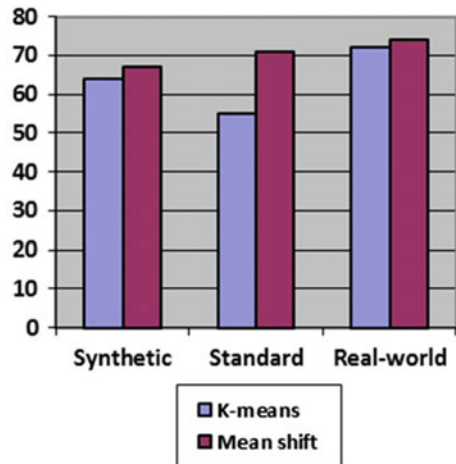$$\text{Accuracy} = \frac{\text{Correctly clustered objects}}{\text{Total objects}} \tag{3}$$

**Table 1** Texture features

| Property | Obtained values |
|---|---|
| Contrast | 0.6394 |
| Correlation | 0.8609 |
| Energy | 0.1113 |
| Homogeneity | 0.8247 |

**Table 2** Regional features

| Property | Obtained values |
|---|---|
| Eccentricity | 0.9785 |
| Extent | 0.1961 |
| Orientation | 80.8737 |

**Fig. 6** Illustration of obtained accuracy (in %) for clustering results



## 5   Conclusion

5-layered feature extraction model is provided for the social images. Multiple features are taken into account by this model which makes it capable to deal with the heterogeneous nature of social images. To evaluate proposed feature extraction model, accuracy of the clustering task is computed and analyzed. K-mean and mean shift algorithms are used for clustering where better accuracy is obtained for the social images. In future, this model can be applied to optimize the results of social data analytics.

**Declaration**   Personal images used in this paper are taken with due permission from the concerned person/authority.

## References

1. Aggarwal, C.C.: An introduction to social network data analytics. In: Social Network Data Analytics. Springer, Berlin (2011)

2. Theodoridis, S., Koutroumbas, K.: Pattern Recognition, 4th edn. Academic Press, Elsevier, London (2008)
3. You, J.-Y., Chien, S.-I.: Saturation enhancement of blue sky for increasing preference of scenery images. IEEE Trans. Consum. Electron. **54**(2), 762–768 (2008)
4. Fierro-Radilla, A.N., Nakano-Miyatake, M., Pérez-Meana, H., Cedillo-Hernandez, M., Garcia-Ugalde, F.: An efficient color descriptor based on global and local color features for image retrieval. In: Proceedings of the 10th IEEE International Conference on Electrical Engineering, Computing Science and Automatic Control, pp. 233–238 (2013)
5. Mabrouk, A.B., Najjar, A., Zagrouba, E.: Image flower recognition based on a new method for color feature extraction. In: Proceedings of the IEEE International Conference on Computer Vision Theory and Applications, vol. 2, pp. 201–206 (2014)
6. Nilsback, M.-E., Zisserman, A.: Automated flower classification over a large number of classes. In: Indian Conference on Computer Vision, Graphics and Image Processing, pp. 722–729 (2008)
7. Vel Murugan, M., Jeyanthi, P.: Content based image retrieval using color and texture feature extraction in Android. In: Proceedings of the IEEE International Conference on Information Communication and Embedded Systems, pp. 1–7 (2014)
8. Kong, F.-H.: Image retrieval using both color and texture features. In: Proceedings of the IEEE International Conference on Machine Learning and Cybernetics, vol. 4. pp. 2228–2232 (2009)
9. Fogel, I., Sagi, D.: Gabor filters as texture discriminator. In: van Hemmen, J.L. (ed.) Biological Cybernetics, vol. 61(2), pp. 103–113. Springer, Berlin (1989)
10. Bama, S.B., Raju, S.: Fourier based rotation invariant texture features for content based image retrieval. In: Proceedings of the IEEE National Conference on Communications, pp. 1–5 (2010)
11. Chang, T., Jay Kuo, C.C.: Texture analysis and classification with tree-structured wavelet transform. IEEE Trans. Image Process. **2**(4), 429–441 (1993)
12. Liu, H., He, G.: Shape feature extraction of high resolution remote sensing image based on susan and moment invariant. Congr. Image Signal Process. **2**, 801–807 (2008)
13. Zhang, G., Ma, Z.M., Tong, Q., He, Y., Zhao, T.: Shape feature extraction using Fourier descriptors with brightness in content-based medical image retrieval. In: Proceedings of the IEEE International Conference on Intelligent Information Hiding and Multimedia Signal Processing, pp. 71–74 (2008)
14. Pedrosa, G.V., Batista, M.A., Barcelos, C.A.Z.: Image feature descriptor based on shape salience points. Neurocomputing **120**, 156–163 (2013)
15. Hong, B.-W., Soatto, S.: Shape matching using multiscale integral invariants. IEEE Trans. Pattern Anal. Mach. Intell. **37**(1), 151–160 (2015)
16. Li, T., Mei, T., Kweon, I.S., Hua, X.S.: Contextual bag-of-words for visual categorization. IEEE Trans. Circuits Syst. Video Technol. **21**(4), 381–392 (2011)
17. Deng, L., Yu, D.: Deep learning: methods and applications. J. Found. Trends Signal Process. **7**(4), 197–387 (2014)
18. Zhu, Z., Wang, X., Bai, S., Yao, C., Bai, X.: Deep learning representation using autoencoder for 3D shape retrieval. Neurocomputing **204**, 41–50 (2016)
19. Mohanaiah, P., Sathyanarayana, P., GuruKumar, L.: Image texture feature extraction using GLCM approach. Int. J. Sci. Res. Publ. **3**(5), 1–5 (2013)
20. Huiskes, M.J., Lew, M.S.: The MIR Flickr retrieval evaluation. In: Proceedings of the ACM International Conference on Multimedia Information Retrieval, Vancouver, Canada (2008)
21. Facebook Profiles. https://www.facebook.com/

# Study of Dimensionality Reduction Techniques for Effective Investment Portfolio Data Management

**Swapnaja Gadre-Patwardhan, Vivek Katdare and Manish Joshi**

**Abstract** The aim of dimensionality reduction is to depict meaningful low-dimensional data of high-dimensional data set. Several new nonlinear methods have been proposed for last many years. But the question of their assessment is still open for the study. Dimensionality reduction is the vital problem in supervised and unsupervised learning. For high-dimensional data, computation becomes heavy if no pre-processing is done before supplying it to any of the classifiers. Because of the constraints like memory and speed, it is not suitable for certain practical applications. As per the method of attribute selection process, attribute sets are provided as an input to the classifier. The attributes that incorrectly classified are supposed to be irrelevant and are removed by obtaining the subset of selected attributes. Thus, accuracy of the classifier is improved, and time is also reduced. Attribute evaluators such as cfsSubset evaluator, information gain ranking filter, chi-squared ranking filter and gain ration feature evaluator are used for the classifiers viz. decision table, decision stump, J48, random forest. Individual investor's investment portfolio data is used for the present study. Twenty-six attributes are obtained from the questionnaire. By applying dimensionality reduction techniques, five major attributes are obtained using information gain ranking filter, chi-squared ranking filter, gain ratio feature evaluation and seven attributes using cfsSubset evaluator. Around 70.7692% accuracy is obtained using three attribute evaluators for all five classification algorithms, whereas cfsSubset evaluator along with random forest classifier gives 81.5385% accuracy. It has been observed that cfsSubset

S. Gadre-Patwardhan (✉)
MES's Institute of Management and Career Courses, Pune, India
e-mail: swapnaja.gadre24@gmail.com

V. Katdare
KCESs Institute of Management and Research, North Maharashtra University,
Jalgaon, India
e-mail: vvkatdare@rediffmail.com

M. Joshi
School of Computer Sciences, North Maharashtra University, Jalgaon, India
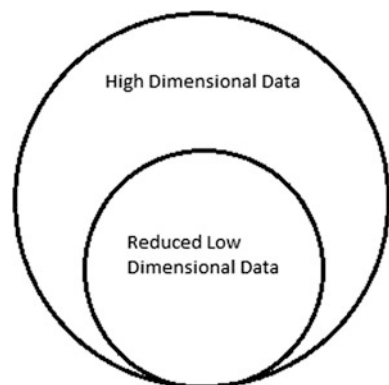e-mail: joshmanish@gmail.com

evaluator with partition membership as a pre-processing technique and random forest as classification algorithm performs reasonably better in terms of accuracy.

**Keywords** Dimensionality reduction · Classification · cfsSubset evaluator Random forest · Chi-squared ranking filter · Information gain ranking filter

# 1 Introduction

With the technological advancement, information technology is used in almost all aspects of daily lives. A huge amount of data is generated such as financial and commercial data, medical data, climate patterns, agricultural and much more. The data is accumulated in several databases and data warehouses. Most of these data have various attributes which are spread over the high-dimensional space. People working with these regularly face the problem of dimensionality reduction which is the process of random variable reduction under consideration for obtaining principal variables [3]. Study of dimensionality reduction depicts meaningful low-dimensional data of high-dimensional data set as shown in Fig. 1. There are two aspects of dimensionality reduction: (1) feature selection and (2) feature extraction. Feature selection is the active and promising research areas in the fields like statistics, data mining, machine learning and pattern recognition. The objective of the feature selection is to determine a subset of input data items by eliminating irrelevant data elements. Feature extraction transforms the data of high-dimensional space to low-dimensional space, thereby reducing the number of attributes [6]. For high-dimensional data, computation becomes heavy if no pre-processing is done before supplying it to any of the classifier. Because of the constraints like memory and speed, it is not suitable for certain practical applications. As per the results of attribute selection methods, attribute sets are provided as an input to the classifier. The attributes that incorrectly classified are supposed to be irrelevant and are removed by obtaining the subset of selected attributes. Attribute evaluators such as

**Fig. 1** Dimensionality reduction

cfsSubset evaluator, information gain ranking filter, chi-squared ranking filter and gain ratio feature evaluator are used for reduction of dimension.

Individual investor's portfolio data is used for the present study. Decision of the investment is dependent on some demographic characteristics. The total risk associated with the investment portfolio can be calculated or predicted on the basis of the investment amount, the purpose of the investment returns expected from the investment and the period of the investment [8, 9].

The paper is organized as—Sect. 2 deals with related work. Literature related to dimensionality reduction is studied under this section. Section 3 contains the information about the data set generated and used. Section 4 is based on experimental work. Comparative analysis of attribute selection methods is performed in this. In Sect. 5, results and observations are discussed, and Sect. 6 talks about the conclusion of the present dimensionality reduction study.

## 2 Related Work

Dimensionality reduction techniques are classified as supervised and unsupervised techniques on the basis of learning process. Supervised algorithms require a training data set with the class label data to learn the lower dimensional representation of data based on some criteria and then predict the class labels for unknown data. Unsupervised approaches project the original data set to a new lower dimensional space without using the label information.

Dimensionality reduction techniques operate either by selecting the subset of the existing data or by transforming the existing features to a reduced data set. Several new nonlinear methods have been proposed for last many years. But the question of their assessment is still open for the study. Since from a long time, researchers are exploring the area of dimensionality reduction in different domains for many years. Stating below the study of some researchers.

Fu et al. [2] have applied dimensionality reduction for simplifying radial basis function (RBF) network structure to improve the classification performance. They have also proposed a novel method, i.e. separability correlation measure (SCM) for ranking the importance of the attributes. In comparison of existing methods such as SUD and Relief-F methods, SCM points to smaller attribute subsets and higher classification accuracies in simulation. They also proposed a better modification for RBF network construction and training by allowing cluster overlap of the same class.

Wang et al. [13] used feature extraction and dimensionality reduction techniques for vowel recognition. They investigated minimum classification error (MCE) for feature extraction. A generalized MCE (GMCE) has been proposed to overcome the shortcoming of MCE algorithm. Linear discriminant analysis (LDA), principal component analysis (PCA), MCE and GMCE are applied for feature extraction. They also investigated support vector machine (SVM) and compared it with linear feature extraction method. They observed that linearly extracted models are fit for training data, whereas SVM has improved generalized properties.

Esser et al. [1] have demonstrated a convex model for non-negative matrix factorization and dimensionality reduction on physical space. They have successfully applied it to hyper-spectral end-member detection and separation of blind source in NMR. A model that can handle outliers is also proposed by them.

Joseph et al. [7] have performed dimensionality reduction and classification on hyper-spectral image. The objective was to find and classify the constituent materials for every pixel in a scene and reduce the dimensionality or volume without losing any critical information. Joseph et al. have described a technique that reduces data dimensions simultaneously and suppresses unwanted signature of interest. The main concept was to project each pixel vector on an orthogonal subspace to the undesired signatures. This is an optimal interference suppression process. Using this process signal-to-noise ratio was maximized and resulted in single component image. Orthogonal subspace projection (OSP) operator can handle k signatures, thereby simultaneously classifying and reducing dimensionality.

Kumar [11] studied unsupervised dimensionality reduction techniques for text data retrieval. The analysis was done on the basis of retrieval quality, complexity and approximation error. He concluded that semantic space obtained by singular value decomposition (SVD) and fuzzy k-means (FKM) produced better results as compared to other dimensionality techniques.

Kaushik et al. [10] have performed dimensionality reduction for indexing large time series database. They studied various dimensionality reduction techniques including singular value decomposition (SVD), discrete wavelet transform (DWT) and discrete Fourier transform (DFT). A new technique viz. adaptive piecewise constant approximation (APCA) is proposed by them. They have used two distant measures for fast searching: (1) a lower bounding Euclidean distance approximation and (2) a non-lower bounding and found that APCA is superior to other techniques.

John et al. [12] examined the quality of dimensionality reduction using rank-based criteria. They reviewed some quality measures which were based on k-ray neighbourhood and distance ranking. According to their observation co-ranking matrix can be used for rank comparison for the initial data set and also for some low-dimensional embedding. After that, rank errors and concepts can be associated with various blocks of the co-ranking matrix. They used real and synthetic data for experiments. They concluded that co-ranking matrix is better tool for simple and discriminatory quality criteria with the help of straightforward interpretation.

Amir et al. [4] introduced an information theoretic nonlinear method for informative dimensionality reduction. They selected continuous feature functions from the co-occurrence matrix. The algorithm proposed by them is analogous to the association analysis used in statistics with the help of feature extraction. They used synthetic co-occurrence data for experimental analysis and observed performance improvement in original feature set.

Zhang et al. [15] applied dimensionality reduction technique on hyper-spectral imagery based on clonal selection. The clonal selection was used to describe the basic features. They developed feature weighting algorithm, clonal selection feature selection (CSFS) algorithm, a feature subset search algorithm and clonal selection

feature weighting (CSFW) algorithm. Experimental analysis demonstrated that CSFS and CSFW are much better than other algorithms, and hence found effective for dimensionality reduction of hyper-spectral remote sensing imagery.

Kilian et al. [14] introduced nonlinear dimensionality reduction using maximum variance unfolding. They reviewed maximum variance unfolding algorithm for low-dimensional representation of high-dimensional data. They concluded that algorithm is based on modern tools and found applicable in the field of machine learning.

## 3 Data Set

The present study is an attempt to understand the behaviour of individual investor in various investment avenues. The investment decision-making depends upon many factors. The objective of the present study is to identify the factors that are affecting on individual investors decision-making process.

### 3.1 Study Sample

The population for the present study comprises of the investors from Pune (Maharashtra, India) region. Investors from all age group, education level, occupation and income group are considered.

### 3.2 Data Collection Method

Questionnaire is designed to collect the data from the investors. Overall, there are 26 questions in the questionnaire that are divided in two sections. It consists of two sections as—Section I personal information and section II includes the questions about the portfolio composition and risk-bearing capacity of an individual investor. Section II is divided into six axes viz. first axis: investment objectives, second axis: time horizon, third axis: level of satisfaction, fourth axis: factors influencing investments, fifth axis: knowledge and sixth axis: risk tolerance.

## 4 Experimental Work

Due to a large amount of data elements in data sets, we applied some filtering approach for feature selection. Attribute selection methods allow us to extract the subset of data elements from the high-dimensional data set, thereby improving the

computation accuracy and efficiency. Principal component analysis, factor analysis, backward feature elimination are some of the popular techniques of dimensionality reduction. For our study, we have used cfsSubset evaluator, information gain ranking filter, chi-squared ranking filter and gain ratio for dimensionality reduction.

Information gain ratio is the ratio of information gain to the intrinsic information. When an attribute $X$ splits the set $S$ into subsets $S_i$, average entropy is calculated, and the sum of the entropy is compared with the original set $S$. The attributes that maximize the difference are selected.

Chi-squared ranking filter is the statistical method that measures the closeness of actual and expected result. In this method, it is assumed that variables are random and drawn from sufficient sample of independent variable. The result indicates the difference between the actual and expected outcomes. Each feature is independently evaluated with reference to the class label. The larger the value of chi-square, features are more relevant to the class model.

Gain ratio is the modified version of information gain. Information gain selects the features with large number of values, whereas gain ratio maximizes the features information gain by minimizing its value.

Correlation-based feature selection (CFS) subset evaluator uses the search algorithm together with a function to evaluate the merits of the feature subset. The method by which CFS measures the usefulness of feature subset considers the importance of the individual feature for class label prediction along with the level of inter-correlation between them [5].

Diagrammatic representation of the dimensionality reduction process used in the study is shown in Fig. 2.

For dimensionality reduction process, WEKA attribute selection techniques such as information gain ranking filter, chi-squared ranking filter, gain ratio feature evaluation and cfsSubset evaluator are demonstrated. After applying dimensionality reduction techniques, we obtained five attributes in information gain ranking filter, chi-squared ranking filter, gain ratio feature evaluation and seven attributes in cfsSubset evaluator. These techniques are applied on the classifiers like decision tree, decision stump, J48, random forest and logistic model tree (LMT). Comparative study is done for all these classifiers. The accuracy of the classifier can be calculated as follows:

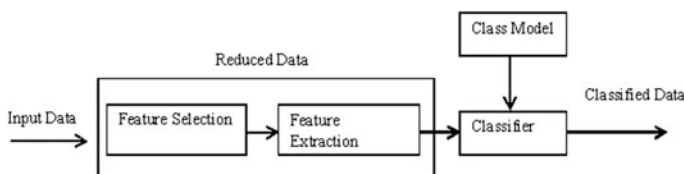$$\text{accuracy} = (\text{tp} + \text{tn})/(\text{tp} + \text{tn} + \text{fp} + \text{fn}) \qquad (1)$$



**Fig. 2** Dimensionality reduction process

and precision by following formula

$$precision = tp/(tp + fp) \qquad (2)$$

## 5  Result and Observations

On the selected five attributes using information gain ranking filter, chi-squared ranking filter, gain ratio feature evaluation techniques and seven attributes using cfsSubset evaluator some pre-processing is done. Pre-processing techniques as follows:

**Class-order** alters the order of the classes. The values are arranged in the order specified by the user, i.e. either ascending or descending or in random order.

**Discretise** discretises the numeric attribute data set to nominal attributes.

**Partition Generator** partitions the membership value.

**Class condition** converts the values of numeric and/or nominal attributes to class conditional probabilities.

Analysis of all these dimensionality reduction methods along with the pre-processing techniques is as below.

Figure 3 represents the accuracy obtained by using decision tree, decision stump, LMT, J48 and random forest classifiers. By applying these classifiers on selected five attributes, we can observe that LMT gives 69.2308% accuracy, i.e. 45 instances are correctly classified.

Figure 4 represents the graphical view of correctly classified and incorrectly classified instances using decision table, decision stump, LMT, J48 and random forest classifiers with the help of pre-processing technique as class conditional probability. From the graphical representation, it can be seen that LMT gives the 69.2308% accuracy. Thus, we can say that with pre-processing technique as class conditional probabilities min values, maximum obtained accuracy remains the same.



**Fig. 3** Five attributes using class conditional probabilities min values
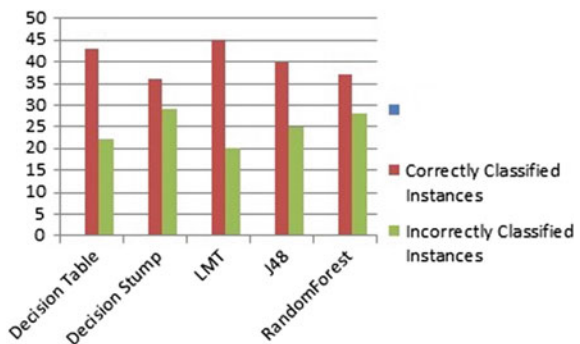
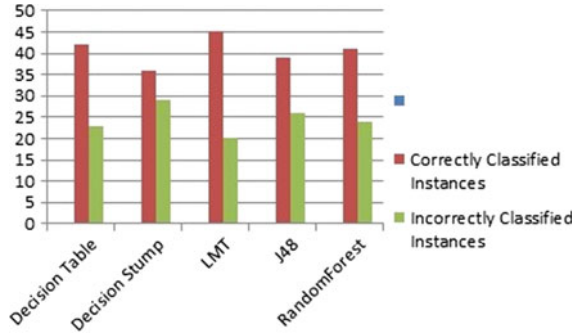**Fig. 4** Seven attributes using class conditional probabilities min values



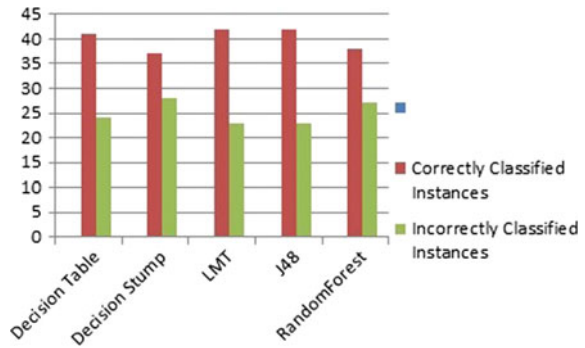**Fig. 5** Five attributes using class-order



**Fig. 6** Seven attributes using class-order

With the help of class-order pre-processing technique along with LMT classifier, 64.6154% accuracy is obtained, i.e. 42 instances are correctly classified as shown in Fig. 5 and whereas from Fig. 6 we can see that Random Forest classifier gives 66.1538% accuracy which means 43 instances are correctly classified.

By applying discretise pre-processing technique, LMT classifier performs better. Using all the techniques, we obtained 67.6923% accuracy which maximum of all. The results are illustrated in Tables 1 and 2.

**Table 1** Five attributes using discretise

| Classifier used | Decision table | Decision stump | LMT | J48 | Random forest |
|---|---|---|---|---|---|
| Correctly classified instances | 66.1538% | 60% | 67.6923% | 66.1538% | 66.1538% |
| Incorrectly classified instances | 33.8462% | 40% | 32.3077% | 33.8462% | 33.8462% |
| Kappa statistics | 0.4758 | 0.3749 | 0.5051 | 0.4758 | 0.4838 |
| MAE | 0.1904 | 0.1813 | 0.1682 | 0.1631 | 0.1632 |
| RMSE | 0.2979 | 0.307 | 0.2948 | 0.2999 | 0.3003 |
| RAE | 83.0547 | 79.0818 | 73.3937 | 71.1426 | 71.1870 |
| RESE | 88.7312 | 91.4458 | 87.7988 | 89.3214 | 89.4303 |
| Cov of clas | 93.8462 | 90.7692 | 92.3077 | 87.6923 | 87.6923 |
| MRRS | 85.1282 | 60.5128 | 65.6410 | 45.8974 | 46.4103 |
| Total inst. | 65 | 65 | 65 | 65 | 65 |

**Table 2** Seven using discretise

| Classifier used | Decision table | Decision stump | LMT | J48 | Random forest |
|---|---|---|---|---|---|
| Correctly classified instances | 66.1538% | 60% | 67.6923% | 66.1538% | 66.1538% |
| Incorrectly classified instances | 33.8462% | 40% | 32.3077% | 33.8462% | 33.8462% |
| Kappa statistics | 0.4758 | 0.3749 | 0.5051 | 0.4758 | 0.4838 |
| MAE | 0.1904 | 0.1813 | 0.1682 | 0.1631 | 0.1631 |
| RMSE | 0.2979 | 0.307 | 0.2948 | 0.2999 | 0.3005 |
| RAE | 83.0547 | 79.0818 | 73.3937 | 71.1426 | 71.1688 |
| RRSE | 88.7312 | 91.4458 | 87.7988 | 89.3214 | 89.4917 |
| Cov of clas | 93.8462 | 90.7692 | 92.3077 | 87.6923 | 87.6923 |
| MRRS | 85.1282 | 60.5128 | 65.6410 | 45.8974 | 46.1538 |
| Total ins | 65 | 65 | 65 | 65 | 65 |

With the help of partition membership pre-processing technique and by applying information gain ranking filter, chi-squared ranking filter and gain ratio feature evaluation techniques, maximum 70.7692% accuracy is obtained as shown in Table 3.

On the other hand with the help of partition membership pre-processing technique along with cfsSubset evaluator as a dimensionality reduction technique, 81.5385% accuracy is obtained which is the maximum and is shown in Table 4.

**Table 3** Five using partition membership

| Classifier used | Decision table | Decision stump | LMT | J48 | Random forest |
|---|---|---|---|---|---|
| Correctly classified instances | 70.7692% | 56.9231% | 70.7692% | 70.7692% | 70.7692% |
| Incorrectly classified instances | 29.2308% | 43.0769% | 29.2308% | 29.2308% | 29.2308% |
| Kappa statistics | 0.5489 | 0.3211 | 0.5489 | 0.5489 | 0.5489 |
| MAE | 0.1871 | 0.187 | 0.162 | 0.1604 | 0.1611 |
| RMSE | 0.2936 | 0.317 | 0.2862 | 0.2924 | 0.2925 |
| RAE | 81.6124 | 81.5680 | 70.6509 | 69.9638 | 70.2857 |
| RESE | 87.4384 | 94.4045 | 85.2572 | 87.0882 | 87.1167 |
| Cov of clas | 95.3846 | 90.7692 | 96.9231 | 89.2308 | 89.2308 |
| MRRS | 89.4872 | 59.4872 | 83.3333 | 61.2821 | 61.2821 |
| Total inst | 65 | 65 | 65 | 65 | 65 |

**Table 4** Seven attributes using partition membership

| Classifier used | Decision table | Decision stump | LMT | J48 | Random forest |
|---|---|---|---|---|---|
| Correctly classified instances | 63.0769% | 56.9231% | 61.5385% | 75.3846% | 81.5385% |
| Incorrectly classified instances | 36.9231% | 43.0769% | 38.4615% | 24.6154% | 18.4615% |
| Kappa statistics | 0.4413 | 0.3211 | 0.4153 | 0.6077 | 0.7143 |
| MAE | 0.1846 | 0.187 | 0.1662 | 0.1074 | 0.096 |
| RMSE | 0.2933 | 0.317 | 0.2928 | 0.2618 | 0.236 |
| RAE | 80.5680 | 81.5680 | 72.5043 | 46.8648 | 41.8794 |
| RESE | 87.3479 | 94.4045 | 87.2007 | 77.9908 | 70.2879 |
| Cov of clas | 100 | 90.7692 | 96.9231 | 93.8462 | 90.7692 |
| MRRS | 85.6410 | 59.4872 | 51.7949 | 38.2051 | 31.7949 |
| Total inst | 65 | 65 | 65 | 65 | 65 |

# 6   Conclusion

Dimensionality reduction is a vital problem in supervised and unsupervised learning. In this paper, we studied feature selection, feature extraction and classification algorithms for dimensionality reduction. Five classification algorithms are taken into consideration. Individual investor's investment portfolio data is used for the present study. Data is collected through questionnaire, and it is observed that some demographic characteristics affect upon the individual investors investment decision-making process. In all, we have got 26 attributes through the questionnaire. By applying dimensionality reduction techniques, we obtained five major attributes using information gain ranking filter, chi-squared ranking filter, gain ratio

feature evaluation and seven attributes using cfsSubset evaluator. Experimental results show that there is no significant difference in the results obtained by information gain ranking filter, chi-squared ranking filter, gain ratio feature evaluation attribute evaluators. Around 70.7692% accuracy is obtained using three attribute evaluators for all five classification algorithms, whereas cfsSubset evaluator along with random forest classifier gives 81.5385% accuracy. It has been observed that cfsSubset evaluator with partition membership as a pre-processing technique and random forest as classification algorithm performs reasonably better in terms of accuracy.

**Declaration** Authors have obtained permission to use the data from the investors. Authors take full responsibilities to bear any consequences if any issues arise due to this. Publisher or Editors are not responsible for this.

# References

1. Esser, E., Moller, M., Osher, S., Sapiro, G., Xin, J.: A convex model for nonnegative matrix factorization and dimensionality reduction on physical space. IEEE Trans. Image Process. **21** (7), 3239–3252 (2012)
2. Fu, X., Wang, L.: Data dimensionality reduction with application to simplifying RBF network structure and improving classification performance. IEEE Trans. Syst. Man Cybern. B (Cybern.) **33**(3), 399–409 (2003)
3. Geng, X., Zhan, D.-C., Zhou, Z.-H.: Supervised nonlinear dimensionality reduction for visualization and classification. IEEE Trans. Syst. Man Cybern. B (Cybern.) **35**(6), 1098–1107 (2005)
4. Globerson, A., Tishby, N.: Sufficient dimensionality reduction. J. Mach. Learn. Res. **3**, 1307–1331 (2003)
5. Hall, M.A., Smith, L.A.: Feature subset selection: a correlation based filter approach (1997)
6. Han, J., Pei, J., Kamber, M.: Data Mining: Concepts and Techniques. Elsevier, Amsterdam (2011)
7. Harsanyi, J.C., Chang, C.-I.: Hyperspectral image classification and dimensionality reduction: an orthogonal subspace projection approach. IEEE Trans. Geosci. Remote Sens. **32**(4), 779–785 (1994)
8. Hoffmann, A.O.I., Post, T., Pennings, J.M.E.: Individual investor perceptions and behavior during the financial crisis. J. Bank. Finance **37**(1), 60–74 (2013)
9. Kaur, M., Vohra, T.: Understanding individual investor's behavior: a review of empirical evidences. Pac. Bus. Int. **5**(6), 10 (2012)
10. Keogh, E., Chakrabarti, K., Pazzani, M., Mehrotra, S.: Locally adaptive dimensionality reduction for indexing large time series databases. ACM SIGMOD Rec. **30**(2), 151–162 (2001)
11. Kumar, A.C.: Analysis of unsupervised dimensionality reduction techniques. Comput. Sci. Inf. Syst. **6**(2), 217–227 (2009)
12. Lee, J.A., Verleysen, M.: Quality assessment of dimensionality reduction: rank-based criteria. Neurocomputing **72**(7), 1431–1443 (2009)
13. Wang, X., Paliwal, K.K.: Feature extraction and dimensionality reduction algorithms and their applications in vowel recognition. Pattern Recogn. **36**(10), 2429–2439 (2003)
14. Weinberger, K.Q., Saul, L.K.: An introduction to nonlinear dimensionality reduction by maximum variance unfolding. In: AAAI, vol. 6, pp. 1683–1686 (2006)
15. Zhang, L., Zhong, Y., Huang, B., Gong, J., Li, P.: Dimensionality reduction based on clonal selection for hyperspectral imagery. IEEE Trans. Geosci. Remote Sens. **45**(12), 4172–4186 (2007)

# Automatic Text Recognition Using Difference Ratio

Shamama Anwar

**Abstract** With the rapid advancement in technology, digitization of documents is gaining popularity. For digitization, the printed or handwritten text needs to be converted to a computer-readable form. For this, the document has to go through line detection, character extraction, recognition and finally conversion to a computer-readable form. A variety of methods have been proposed for the same. The paper proposes a method for text extraction and recognition which is based on a data set called as a learn file which is a vector representation of the images in the data set. Recognition is achieved by using the difference ratio between the input image and the learn file. The paper also presents two applications of the proposed method: text extraction from printed document and automatic number plate recognition. After recognition, the identified characters are written on to a text file.

**Keywords** Text recognition · Text extraction · Absolute difference
Difference ratio · Automatic number plate recognition

## 1 Introduction

Text extraction is a process by which a handwritten page/printed document/scanned page or an image in which text is available is converted to ASCII characters that a computer can recognize. With advancement in technology, digitization of documents is encouraged for a paperless office environment. Lots of resources are now easily available in electronic medium. A challenge for this document digitization is converting the text into digital form.

 A system is thus needed that involves the automatic conversion of text from paper documents to a computer-readable form. These documents can be scanned and treated as image files before text extraction. Therefore, the aim of such systems

S. Anwar (✉)
Department of Computer Science and Engineering, Birla Institute of Technology, Mesra, Ranchi, India
e-mail: shamama@bitmesra.ac.in

developed is to locate, recognize and convert text present in an image to a form that is usable by a computer and text-processing applications. The aim of this paper is to detect, extract and recognize text from images and write the same to an output text file for further use. The paper also presents two real application of the proposed system. In the first application a printed or handwritten text is recognized and converted to a text file. The second application is the automatic identification of vehicle number plate and its conversion to a text file. This can make number plate matching a very easy job. Since matching number plates as images from an image database is a tedious task as compared to matching text.

## 2  Related Work

Some of the earlier work in text identification used clustering algorithms. One such method combined the top-down approach for image splitting and the bottom-up approach for region growing which gave good results for text identification on books and journal covers [1]. To reduce the effect of colour variations, the method used a colour clustering algorithm in the pre-processing stage. An approach using adaptive local connectivity map was proposed to locate lines in a text document, and the method was tested on a collection of 30 historical letters [2]. A generalized method for text segmentation by [3] included four steps. Firstly, the input image was enhanced to make text block extraction easier. Then the image was split into connected components, and component filtering was applied to discard the useless components. Next, text layers were formed by merging the remaining components, and finally the text layer was used to generate a binary output. Another hybrid approach was proposed by [4]. This method used Canny edge detector to detect all possible text edge pixels and then used connected component analysis to identify candidate text regions. An extraction method using wavelet transform and a dynamic threshold has been proposed by [5]. The accuracy reported for the method was around 91.2%. Other significant approaches for text extraction and recognition have been done by many researchers [9, 10]. A survey of the different approached for text extraction can be found in [6].

An application of text extraction method discussed in the paper includes automatic vehicle number plate recognition. A technique for recognition using morphological operations, histogram manipulation and edge detection was given by [7]. They also used artificial neural network for character classification and recognition [8] proposed a feature-based number localization method which comprises of edge finding and window filtering method. Some other significant work proposed in this field is given in [11, 12]. A more detailed study of the methods can be found in [13].

The method proposed in this paper does not need any classifier and nor is any training required. Hence, making it easier to implement. The implementation has been done in MATLAB.

# 3 Proposed Methodology

The proposed text recognition and conversion system relies on a database created from a certain set of characters. It compares the characters in the scanned image file to the characters in the learned set. Generating the learned set is quite simple. It requires that an image file with the desired characters in the desired font be created.

Once the learned set has been read in from the image file, it is converted to a vector and stored in a learn file. Vectorization is done by reshaping the image matrix. The advantage of the learn file is that it eliminates the need for retaining the images of the learned characters and can be read in very quickly. It also makes the conversion system easily adaptable to different types of character sets and makes it adaptable to a variety of languages. The paper presents two applications for the system which were implemented in MATLAB. For printed document recognition, the data set consists of 91 characters, 26 capital letters, 26 small letters, 26 small cursive letters, 10 numbers and 3 characters to represent comma, full stop and blank space. Since it is not possible to include the entire data set considering the length of the paper, a sample of the data set can be seen in Fig. 1. These images are converted to a vector form and stored in a learn file. Once the learn file is created, the images are discarded. The data set can be easily expanded to include more characters per the application area. Likewise, for automatic number plate recognition, the data set would simply comprise of 26 capital letters, 10 numbers and 1 character to represent a space.

After the database creation, the extraction of characters from the given input image is done followed by recognition of extracted characters and then finally conversion to text document. The steps are represented in Fig. 2.

## 3.1 Line Detection

For recognition, it is vital to extract characters from the input image. Given an entire page of the document, the first aim is to detect lines, i.e. to break the page into lines. To detect line breaks, blank spaces are detected along vertical spacing of characters. The page is scanned pixel by pixel vertically (i.e. from top to bottom), and blank spaces are marked between characters. The process is repeated through the entire image thereby marking all blank spaces. Next the marked points are now joined so
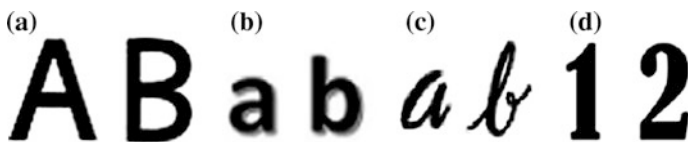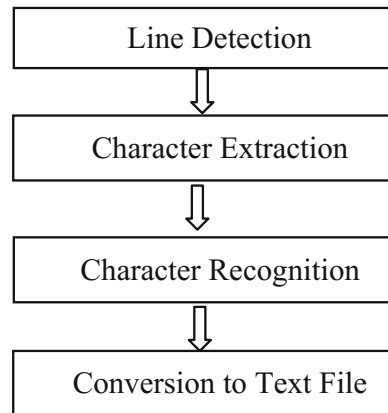


**Fig. 1** Sample data set for English alphabets; **a** Capital letters; **b** Small letters; **c** Small cursive letters; **d** Numbers

**Fig. 2** Automatic text
recognition flow chart

```
┌─────────────────────────────────────┐
│          Line Detection              │
└─────────────────────────────────────┘
                  ⇓
┌─────────────────────────────────────┐
│         Character Extraction         │
└─────────────────────────────────────┘
                  ⇓
┌─────────────────────────────────────┐
│        Character Recognition         │
└─────────────────────────────────────┘
                  ⇓
┌─────────────────────────────────────┐
│        Conversion to Text File       │
└─────────────────────────────────────┘
```

that it makes a line to represent the line break. Any marked points that do not fall in the joint lines are discarded. This method of line detection is simple to implement and also works if the handwritten text is written slanting instead of a straight line.

## 3.2 Character Extraction

Next comes the process of locating regions of printed text. The process of extraction isolates characters. The most common problem in extraction is: it causes confusion between text and graphics in case of joined and split characters. For character extraction, the image is scanned from left to right to find a dark pixel. If a dark pixel is found all connected dark pixels are grouped as a block. When successive light pixels occur, then the group or blocks end. When another dark pixel is found, then it forms another group or block. A box is marked along the groups detected by finding the distance between the extreme points in the block. This process is repeated along the lines detected from left to right. Simultaneous pixel tracking can be done across a line to make it faster. The output of the extraction process marks block of characters across the page.

## 3.3 Character Recognition

After extracting the individual characters in a document, it has to be compared to the data set. The comparison algorithm is simple: it uses the absolute difference between the characters and the data set, thereby calculating the difference ratio. The image is first converted to a vector. Then the absolute difference is calculated for the input image and different images in the data set. Since the data set is stored as a vector, calculating the absolute difference is easy. The absolute difference for

different images is nonzero, whereas similar images have more number of zeros in their difference. The absolute difference is calculated as [14]:

$$AD(i) = A(:) - B_i(:)$$ (1)

where $A(:)$ is the vector representation of the input image, $B_i(:)$ is the vector representation of characters in the data set and $1 \leq i \leq n$, where $n$ is the total number of elements in the data set.

Next a ratio between the number of uncommon pixels and the total number of pixels is then found which is termed as the difference ratio as in Eq. (2). The lesser the difference ratio the similar the images are.

$$DR = \frac{\sum_{i=1}^{n} AD(i)}{n}$$ (2)

The difference ratio of the input character and the data set is stored as a vector. The minimum value from the vector is searched. The position of the minimum value in the vector helps in recognition as each position represents a character. For example, for printed text recognition, the vector would contain 91 entries, the first entry would correspond to the letter A, next would correspond to letter B and so on. The recognized characters are then sent down to be written down in a text file.

The algorithm for the process is given in Fig. 3.

## 4    Results

The proposed method is implemented for printed text recognition and for automatic number plate detection.

---

Input: Image file
Output: Recognized character in text file
Steps:    1. (a) Scan the image vertically.
              (b) Detect and mark white spaces.
              (c) Divide the image into lines based on the marked white lines.
          2. (a) Scan resulting image from left to right.
              (b) Find and mark connected dark pixels as one block.
          3. Convert each block into vector representation.
          4. Find the Difference Ratio (DR) between input character and the dataset characters and store it in a vector.
          5. The recognized character is the one having the least DR value.

---

**Fig. 3**  Algorithm for character recognition

## 4.1 Printed Text Recognition and Conversion

The data set for a printed text recognition system has been discussed earlier, and a sample has been included in Fig. 1. The input image is then taken, and it has to go through the first step of line detection followed by character extraction as depicted in Fig. 4. The extracted characters are now individually represented in a vector form. The difference ratio is found between the characters and the data set. A portion of the list of the difference ratio is included in Table 1 for the first extracted character, i.e. 'B'. It is not possible to include the entire 91 entries in the table given the consideration of the paper length.

Table 1 clearly shows that the lowest difference ratio is for the character 'B', hence the recognition is done correctly. The DR value for the entire data set is stored as a 91 element vector. The lowest value is picked from the vector. The position of the vector represents the character 'B', and hence, this character is written on to a text file. Similarly, the entire text document is converted and represented in a text file as in Fig. 5.

## 4.2 Automatic Number Plate Recognition

A similar application is for automatic number plate recognition for vehicles. The data set just comprises of 37 images, and hence, the learn file consists of 37 vectors. The input image is first converted to grey or binary image. As this application does not need line detection, this step can be missed. The process begins by extracting characters by identifying connected dark pixels and forming bounding boxes around them to isolate individual characters. Next, the difference ratio between the input character and the data set is calculated as explained above. This yields a vector of length 37. The minimum value if found in the vector and the letter to which it corresponds is the recognized letter. A example is shown in Fig. 6.
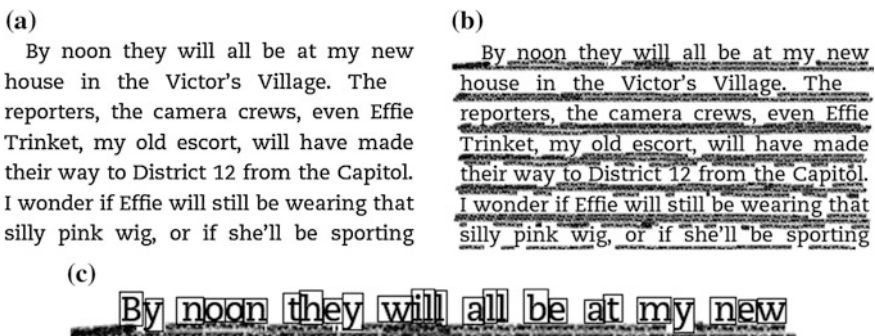


**Fig. 4** **a** Sample input text document; **b** Text document after line detection; **c** Output of character extraction

**Table 1** Difference ratio

| DR | Data set | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | A | B | a | b | *a* | *b* | 1 | 2 |
| | 0.045 | 0.026 | 0.057 | 0.031 | 0.052 | 0.035 | 0.041 | 0.056 |



**Fig. 5** Text recognized and written to a text file



**Fig. 6  a** Sample input image; **b** Character extracted as blocks; **c** Output text file

## 5   Conclusion

In the present research, text extraction and recognition has been done. The advantage of the method is that the size of the data set has been considerably reduced by storing the images in the data set as a vector for recognition purpose. Finding the absolute difference and hence the difference ratio aids the recognition process and reduces the computational time considerably. The application of the proposed system has also been discussed. The algorithm proposed has been tested for automatic text recognition on a set of 40 input images from different sources, and the accuracy reported is 95%. For automatic number plate recognition, a set of 25 images was considered, and the accuracy reported is about 94%. The accuracy is calculated as the ratio between the number of correct recognition instances and the total number of input images used. Further, the method can also be tested for data set comprising of other regional languages.

## References

1. Sobottka, K., Bunke, H., Kronenberg, H.: Identification of text on colored book and journal covers. In: International Conference on Document Analysis and Recognition, p. 57 (1999)
2. Shi, Z., Setlur, S., Govindaraju, V.: Text extraction from gray scale historical document images using adaptive local connectivity map. In: 8th International Conference on Document Analysis and Recognition, pp. 794–798 (2005)
3. Zhan, Y., Wang, W., Gao, W.: A robust split- and -merge text segmentation approach for images. In: 18th International Conference on Pattern Recognition, vol. 2, pp. 1002–1005 (2006)
4. Nagabhushan, P., Nirmala, S.: Text extraction in complex color document images for enhanced readability. Intell. Inf. Manag. **2**(2), 120–133 (2010)
5. Zaravi, D., Rostami, H., Malahzaheh, A., Mortazavi, S.S.: Journals subheadlines text extraction using wavelet thresholding and new projection profile. Int. J.Comput. Electr. Autom. Control Inf. Eng. **5**(1), 33–36 (2011)
6. Sumathi, C.P., Santhanam, T., Devi, G.G.: A survey on various approached of text extraction in images. Int. J. Comput. Sci. Eng. Surv. **3**(4), 27–42 (2012)
7. Badr, A., Abdelwahab, M.M., Thabet, A.M., Abdelsadek, A.M.: Automatic number plate recognition system. Math. Comput. Sci. Ser. **38**(1), 62–71 (2011)
8. Kranthi, S., Pranathi, K., Srisaila, A.: Automatic number plate recognition. Int. J. Adv. Technol. **2**(3), 408–422 (2011)
9. Hoang, T.V., Tabbone, S.: Text extraction from graphical document images using sparse representation. In: 9th International Workshop on Document Analysis System, pp. 143–150 (2010)
10. Grover, S., Arora, K., Mitra, S.K.: Text extraction from documnet images using edge information. In: Annual IEEE India Conference, pp. 1–4 (2009)
11. Kocer, H.E., Cevik, K.K.: Artificial neural networks based vehicle license plate recognition. Procedia Comput. Sci. **3**, 1033–1037 (2011)

12. Roy, A., Ghoshal, D.P.: Number plate recognition for use in different countries using an improved segmentation. In: 2nd National Conference on Emerging Trends and Applications in Computer Science, pp. 1–5 (2011)
13. Patel, C., Shah, D., Patel, A.: Automatic number plate recognition system (ANPR): a survey. Int. J. Comput. Appl. **69**(9), 21–33 (2013)
14. Singla, N.: Motion detection based on frame difference method. Int. J. Inf. Comput. Technol. **4**(15), 1559–1565 (2014)

# Sign Language Conversion Tool (SLCTooL) Between 30 World Sign Languages

**A. S. C. S. Sastry, P. V. V. Kishore, D. Anil Kumar and E. Kiran Kumar**

**Abstract** This paper proposes to find similarity between sign language finger spellings of alphabets from 30 countries with computer vision and support vector machine classifier. A database of 30 countries sign language alphabets is created in laboratory conditions with nine test subjects per country. Binarization of sign images and subsequent feature extraction with histogram of oriented gradients gives a feature vector. Classification with support vector machine provides insight into the similarity between world sign languages. The results show a similarity of 61% between Indian sign language and Bangladesh sign language belonging to the same continent, whereas the similarity is 11 and 7% with American and French sign languages in different continents. The overall classification rate of multiclass support vector machine is 95% with histogram of oriented gradient features when compared to other feature types. Cross-validation of the classifier is performed by finding an image structural similarity measure with Structural Similarity Index Measure.

**Keywords** Sign language recognition · World sign languages comparison
Feature extraction · Support vector machines · Sign−to−sign translator

A. S. C. S. Sastry (✉) · P. V. V. Kishore · D. Anil Kumar · E. Kiran Kumar
Department of Electronics and Communication Engineering, K L University, Green Fields, Vaddeswaram, Guntur District, Andhra Pradesh, India
e-mail: ascssastry@kluniversity.in

P. V. V. Kishore
e-mail: pvvkishore@kluniversity.in

D. Anil Kumar
e-mail: danilmurali@kluniversity.in

E. Kiran Kumar
e-mail: kiraneepuri@kluniversity.in

# 1   Introduction

Language translator from Google [1] is helping 200 million people to communicate from all over the word. Although there are many such language translators [2], the primary goal is translation of words and sentences from one language to another language. The program compares language structures instead of word or sentence features in both languages. The language is modelled through vector spaces, and the transformations happen by vector space mapping between different languages. The rate of accuracy for a 5-word conversion is around 90%. There are many such models for language converters in speech and text [3], but this paper articulates a sign language translator between multiple countries.

Vocal languages are produced by voice, and basic structure is decided by the alphabets. Every language around the world is represented by a set of alphabets, and their infinite combination produces words that convey information. But for hearing-impaired people, this is of no use. Their alternative is sign language. Sign languages are produced by finger shapes, hands location with respect to head, face and body along with facial expressions. The alphabets in sign languages are finger mapped. Each English alphabet is mapped into either five fingers (single hand) or ten fingers (double hand). The structural representation of fingers forms alphabets for sign languages.

The Ethnologue—language encyclopaedia of the world lists 6909 living languages from which only 130 are deaf sign languages. Before exploring the possibility of a Sign−to−sign translator that transforms one country sign language into another, this work focuses on identifying a similarity between these visual languages. We have carefully chosen 30 countries whose sign languages are popular, and extensive research is going on in developing machine translation of these sign languages with non-visual (glove based) and visual (video camera based) techniques. The countries are America, Mexico, India, Bangladesh, Pakistan, Sri Lanka, China, Philippines, Indonesia, Britain, France, Ireland, Spain, Czech, Estonia, Finland, Germany, Hungary, Netherland, Norway, Poland, Chile, Australia, New Zealand, Iceland, Brazil, Kenya, South Africa, Uganda and Zambia.

Visually the structural similarity between the letters can be decoded by the human brain with some efforts, but it is quite a challenge for the computer. In an experiment at our laboratory, even the humans who learned one sign language found it difficult to follow signs from another sign language. Their failure rate was 60% for other sign languages, but again this is a subjective evaluation. This visual decoding and mapping of signs to text or speech is challenging researchers for around two and half decades. For an efficient Sign−to−sign translation between countries, the following are important factors for evaluation.

1. The first part is to find a similarity between 30 world sign languages using histogram of oriented gradients (HOG) features and support vector machine (SVM).

2. To draw a confusion matrix for these 30 countries and to evaluate the performance of the classifier.
3. The third part we used various feature extractors to test the robustness of the HOG as it maps nine bin gradient orientations into histograms making it rotation and scale invariant for small variations.
4. Lastly, we plot the conversion efficiency of one sign language into another and also measure the relativity between sign languages geographically.

Liang [4] proposed moving object classification like cars, motorcycles, pedestrians and bicycle by using local shape from wavelet transform and HOG features with hierarchical SVM classification. The proposed method is tested on six video sequences for classification. The average computer processing times of the object segmentation is 79 ms, object tracking is 211 ms, and classification is 0.01 ms, respectively.

In recent years, SVM classifier with histogram of oriented gradients (HOG) features is the most popular technique for vehicle detection [5]. In real-time implementation, this is important for advanced driver assistance system applications. To reduce the complexity of the SVM, the dimensions of HOG features are to be reduced. The proposed method in [5] using SVM classifies for vehicle detection is three times faster than other algorithm in the area.

The rest of the paper is organized as follows: Sect. 2 describes the followed methodology in determining the sign similarity. Results and discussion is presented in Sect. 3 with conclusion in Sect. 4.

## 2 Methodology: Inter-Country Sign Language Classification

Figure 1 shows the procedure followed in this paper to investigate the similarity between basic structures of world sign languages. The experiment involves only alphabets as they are the basic structures for formation of any language. Methodology involves two phases: training phase and testing phase.

### 2.1 Support Vector Machines

SVMs analyse data and produce binary responses for classification problem, which come under a class of supervised learning classifier models. The basic SVM classifies a two-class problem by projecting a hyperplane between data during training phase. The hyperplane is characterized by a subset of data points acting as support vectors. During training, the SVM is presented with example vectors $x_i \in \Re^n, i = 1 \ldots l; l$ training samples, to label each data sample as either +1 or -1 class
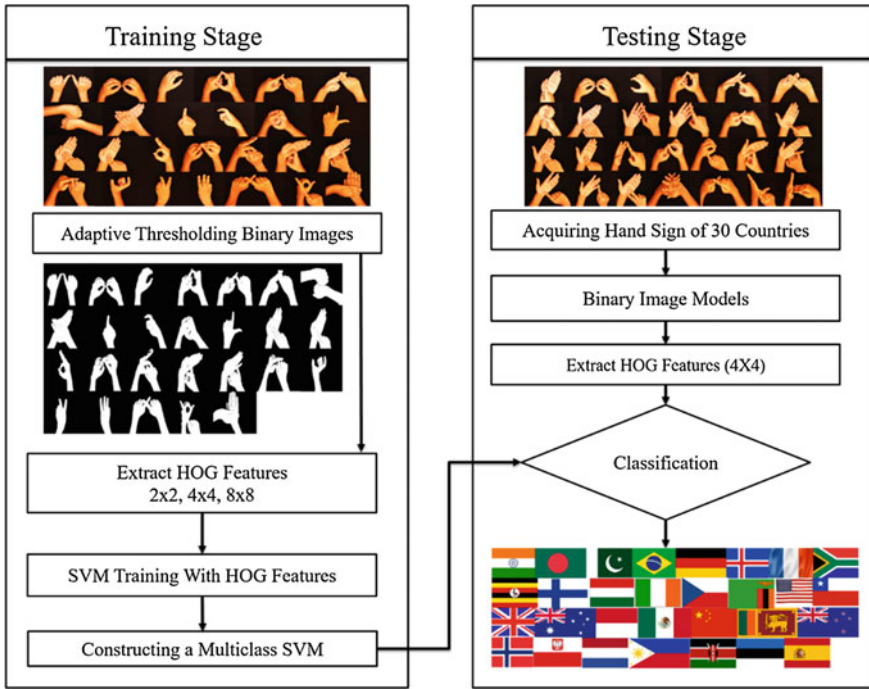
**Fig. 1** Algorithm for decoding relativity among world sign languages

label which forms the indicator vector $y_i \in \{+1, -1\}$. SVM formulates the optimization problem as a decision boundary $D(x)$ such that

$$D(x) = \min_{w,b,\lambda} \left( \frac{1}{2} w^{\mathrm{T}} w + C \sum_{i=1}^{l} \lambda_i \right) \tag{1}$$

$$\text{Subjected to } y_i \{ w^{\mathrm{T}} \phi(x_i) + b \} \geq 1 - \lambda_i \text{ with } \lambda_i \geq 0, i = 1, 2, \ldots, l;$$

where $C$ is a positive constant defining regularization. The terms $w$ and $b$ are weight and bias. $\lambda$ is the misclassification handler. The function $m(x) : x \rightarrow \phi(x)$ maps feature vector $x$ to a higher dimensional space. The mapping function $m(x)$ maps $x$ into a dot product of feature space that satisfies $m(x_{i-1}, x_i) = \phi^T(x_{i-1})\phi(x_i)$.

## 2.2 Multiclass SVM

The most widely used multiclass SVM models are one vs. all (OVA), one vs. one (OVO) [6], directed acyclic graph (DAG) [7] and error correcting output codes (ECOC) [8]. OVA creates $N$ binary SVMs for all categories where $N$ is class

number. For a $n$th SVM, only examples in that class are positive and remaining are negative. The computation time is less but at a compromised efficiency. OVO creates a pairwise $0.5N(N - 1)$ SVMs and pairwise voting to accommodate new samples for solving multiclass problems. DAG training is from OVO model, and testing is from binary acyclic graph model. ECOC disambiguates output binary codes to construct a code word matrix which is compared with generated bit vectors by selecting a row as a class having minimum hamming distance. This method gives good classification rates compared to other four at the cost of execution speed. The slower speed is due to the increased length of code words to disambiguate N classes. The minimum code word in ECOC is $\log_2 N$ to a maximum of $2^{N-1} - 1$ bits. Comparing the multiclass SVM methods from MATLAB implementation, we found ECOC performs better at optimum speeds.

The similarity measure for 30 different world sign language alphabets using computer vision model and machine learning algorithms is proposed. Experimental results show the sign language relativity between countries and continents. Validation is through human expert identification and Structural Similarity Index Measure (SSIM).

## 3 Results and Discussion

Experimentation with the proposed methodology aims to answer the following questions.

1. How much similarity is observed between sign languages of the 30 countries?
2. Does countries of the same continent exhibit more similarity than others?
3. What is the overall similarity in sign language between continents of the world?
4. Can a Sign−to−sign converter is possible at the image level between different sign languages of the world?

The captured sign images are large, and cubic interpolations trimmed their size to $64 \times 64$. The RGB colour images have large R (red) content and hence R plane is extracted for processing. Block thresholding within a 16-pixel block separates foreground hand regions from background. Ten features are extracted from these binary images. For each country, a feature matrix is build. The size of each feature matrix is $\boldsymbol{m}^f \times \boldsymbol{n}^f$, where $m = 26$, i.e. the number of alphabets and $\boldsymbol{n}$ is variable column vector that captures feature values. $\boldsymbol{f}$—consists of country and test subject indicator. The first problem encountered during feature matrix creation is the inability of our algorithm to control the length of $\boldsymbol{n}$, where n is initial length of the feature vector. For each image, the length of the feature vector changes due to number of feature points detected during the feature extraction phage. For 26 different images, we have 26 different feature lengths. Feature length normalization has been challenging, as it is difficult to decide on the number of features required to produce good classification rate.

Figure 2 shows variational feature plots of each alphabet in Indian sign language. The plots also show that the feature variations are almost constant cross-features even though the number of features per sign per country changed marginally. Normalization of $n$ through maximum feature size is done to preserve the actuals, and the remaining features are zero padded to design a constant size feature matrix. This procedure gives a fixed feature size matrix of size $m \times \max(n)$.

The first part is to find the similarity between sign languages from 30 different countries. For this, the feature matrices of all countries from all feature vector models are prepared. A multiclass SVM with ECOC model is trained with one country and tested with all other countries for each feature type.

Testing results in a classification matrix or a confusion matrix between two countries. All countries' sign languages are tested against one trained country, and cross-verification is done by testing the multiclass SVM for all other countries. The SVM is trained with single sample and tested with a different sample from our database. Multiple testing of this kind produced more or less similar results with a deviation of $\pm 3\%$.

Misclassifications between the Indian signs (ISL) and Bangladesh signs (BanSL) are projected from the confusion matrix in Fig. 3. The green is Bangladesh and saffron is India. From the confusion matrix, the Bangladesh 'E' is classified as Indian 'D'. A total of ten signs are misclassified using our proposed method of classification. Total 16 signs match between the two countries.

From the following observations, the similarity of world sign languages is formulated as

1. Spain and German sign languages are 96% similar with 25 signs being matched in two-way training and testing.
2. Mexican–Spain, Mexican–German and Kenya–South Africa are next with 24 sign matches having 92.3% similarity.
3. The lowest similarity set countries are (Australia, American SL), (American, Indian SL), (Netherlands, Australia SL), (Sri Lanka, French SL), (Estonian, French SL), (Netherlands, New Zealand SL) and (Polish, Sri Lanka SL) where the matching signs in both directions range between 0 and 1. Visual verification can be made using Fig. 4 for a set of two sign alphabets 'C' and 'N'.
4. The reason interpreted by us for lowest and highest similarity match among sign languages of different countries depends on the geographical regions in which the country is located.
5. The continentwise similarity measure is checked, and the results for one continent, i.e. Asia, is projected in the plot in Fig. 5.

Figure 5 has seven Asian countries, namely India (IN), Bangladesh (BA), Pakistan (PA), Sri Lanka (SR), China (CN), Indonesia (IA) and Philippines (PH). The plots show histogram of matching signs with ten different types of features. Each feature representing a particular colour; red-HOG, green-SIFT, blue-SURF, cyan-MESR, magenta-BRISK, yellow-LBP, dark yellow-LSS, navy-HAARS, purple-HCORNERS, wine-FAST.
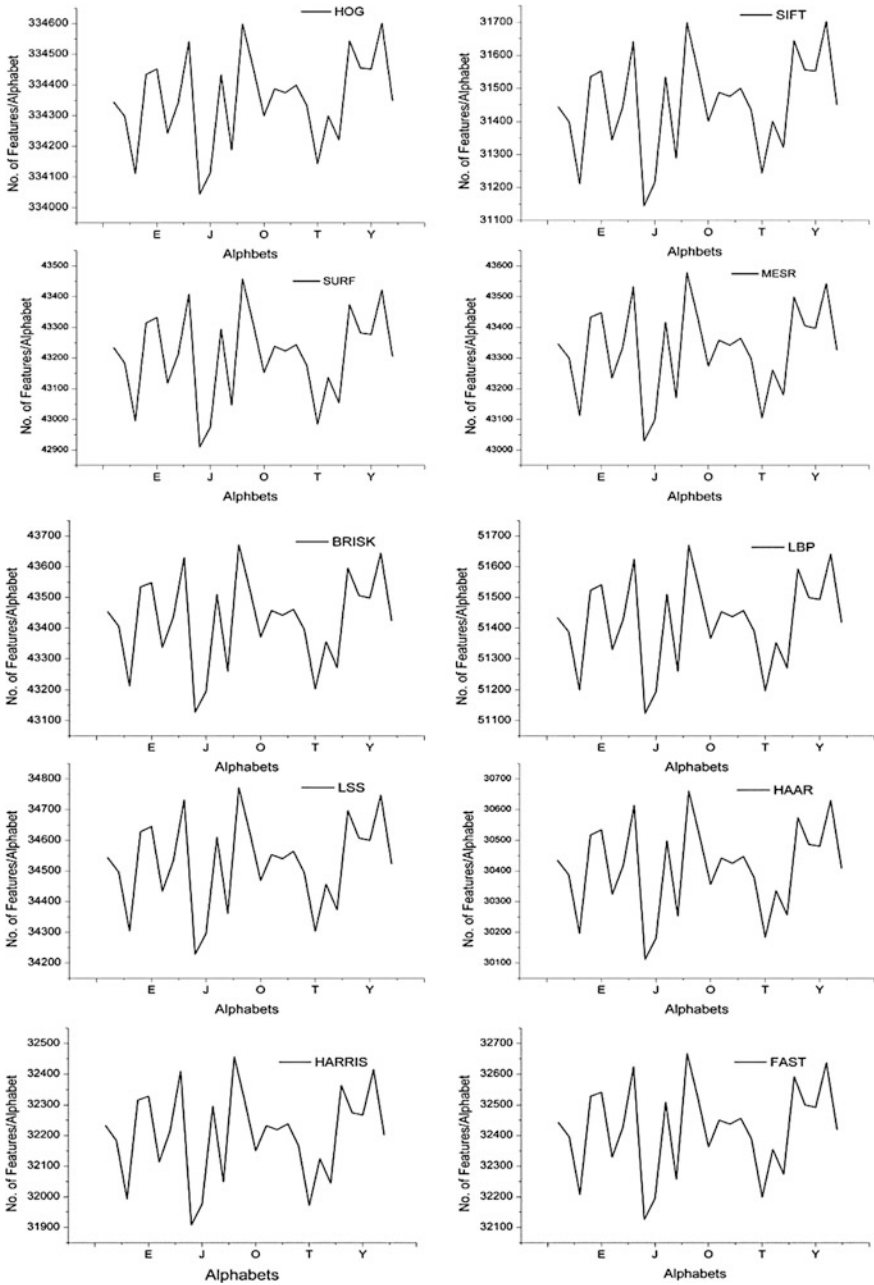
Fig. 2 Feature number variations of alphabets from ISL

|   | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U | V | W | X | Y | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| B | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| C | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| D | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| E | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| F | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| G | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| H | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| I | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| J | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| K | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| L | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| M | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| N | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| O | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| P | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Q | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| R | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| S | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| U | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| V | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| W | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| X | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| Y | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| Z | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

**Fig. 3** Confusion matrix between Indian sign alphabets and Bangla sign alphabets with SVM classifier



**Fig. 4** For visual verification between sign languages of five different countries

Except China and Philippines, all other countries' sign languages show a high range of similarity of around 50–60%. China and Philippines have a high range of similarity due to their cultural influences on each other. HOG features give a high range to classifier performance compared to other features in the list during multiple instances of testing as shown in Fig. 5.

6. There is high similarity between countries from same continent compared to that of countries from different continents as can be analysed.

We also explored the idea of Sign−to−sign translation as in case of spoken language translators [1]. HOG features and SVM are used for training and testing. But cross-verification of the feature vector is checked using a known image
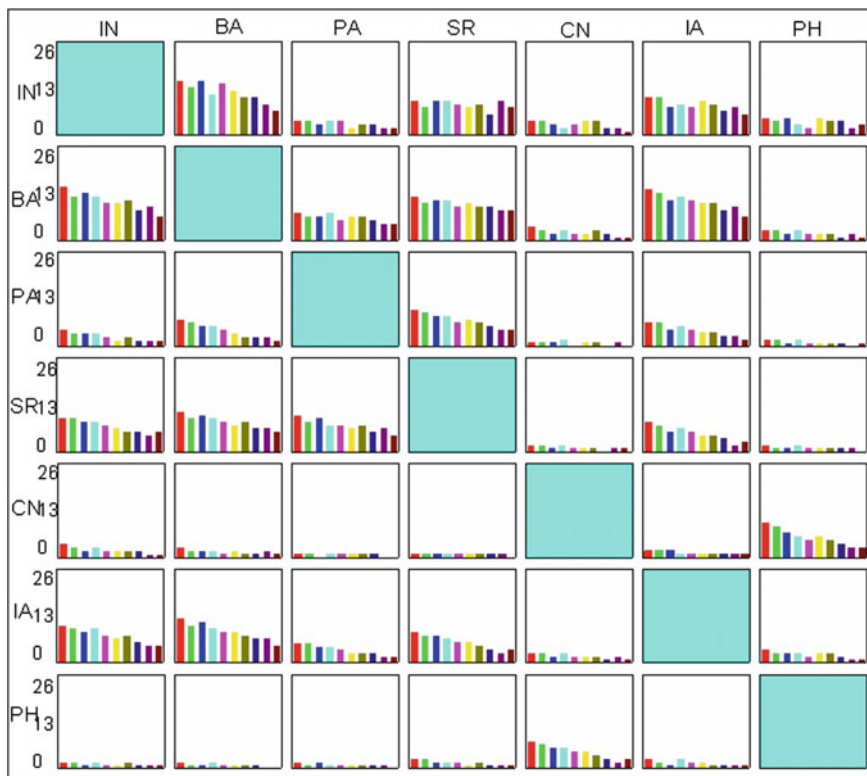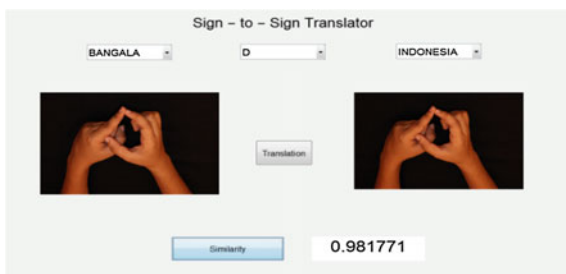
**Fig. 5** Sign language similarity measure for Asian countries

**Fig. 6** Sign−to−sign translator between Bangladesh and Indonesia for sign D



structure measurement parameter called Structural Similarity Index Measure (SSIM). A graphical user interface (GUI) is built in MATLAB to do the job. The user of the GUI can translate sign language alphabets between countries and check the similarity index (SSIM) value. The translator uses HOG features and SVM classifier for the recalling the corresponding signs. Snapshots of GUI testing are in Figs. 6 and 7.

**Fig. 7** Sign−to−sign translator between Sri Lanka and Ireland for sign B



**Fig. 8** Performance of SVM with features used and cross-verification with SSIM

Matching the performance of HOG+SVM with SSIM has a deviation of ±3%. The performance of the best feature for a Sign−to−sign translator with respect to structural similarity of signs is computed rigorously with nine different sets of data from 30 different sign languages for six continents around the world is shown in Fig. 8.

# 4 Conclusion

An attempt is made to find similarity between sign languages from 30 different countries based on image processing models and pattern classifiers. Ten feature extraction techniques are compared for this work. Multiclass support vector machine classified these features, and the performance of the classifier with respect to each feature is measured. Visual verification and structural verification using SSIM are preformed to validate the classifiers performance. Overall the SVM classifier registered a 95% matching with HOG feature vector and the remaining feature vectors produced less than 90% matching. A high similarity in sign languages is found in countries of same continent which are geographically close to each other. Cultural variation is also a cause for large variations in neighbouring countries having different sign languages, e.g. India and China. A Sign−to−sign translator between alphabets of 30 countries with their similarity is created and tested. This translator can be made dynamic to accept signs from various countries online and use the translator to communicate effectively by sign language users of different countries without learning other countries' sign languages.

# References

1. Leite, F.O., et al.: Using Google Translate© in the hospital: a case report. Technology and Health Care (Preprint), pp. 1–4 (2016)
2. Cheriton, D.R.: Interpreter-based program language translator using embedded interpreter types and variables. Google Patents (2016)
3. Huang, Y.-M., Shadiev, R., Hwang, W.-Y.: Investigating the effectiveness of speech-to-text recognition applications on learning performance and cognitive load. Comput. Educ. **101**, 15–28 (2016)
4. Liang, C.-W., Juang, C.-F.: Moving object classification using local shape and HOG features in wavelet-transformed space with hierarchical SVM classifiers. Appl. Soft Comput. **28**, 483–497 (2015)
5. Lee, S.-H., et al.: An efficient selection of HOG feature for SVM classification of vehicle. In: 2015 International Symposium on Consumer Electronics (ISCE). IEEE (2015)
6. Galar, M., et al.: An overview of ensemble methods for binary classifiers in multi-class problems: experimental study on one-vs-one and one-vs-all schemes. Pattern Recognit. **44**(8), 1761–1776 (2011)
7. Zhang, X., Ding, S., Sun, T.: Multi-class LSTMSVM based on optimal directed acyclic graph and shuffled frog leaping algorithm. Int. J. Mach. Learn. Cybern. **7**(2), 241–251 (2016)
8. Bai, X., et al.: Learning ECOC code matrix for multiclass classification with application to glaucoma diagnosis. J. Med. Syst. **40**(4), 1–10 (2016)

# Hybrid ITLBO-DE Optimized Fuzzy PI Controller for Multi-area Automatic Generation Control with Generation Rate Constraint

**Aurobindo Behera, Tapas Ku Panigrahi, Arun Ku Sahoo and Prakash Ku Ray**

**Abstract** The paper projects the gains of a fuzzy controller with its parameter being tuned by the hybrid improved teaching learning based optimization and differential evolution (hITLBO-DE). The foremost apprehension with the operation of AGC is satisfying equivalence of generation and gross demand with reference to a system. The frequency and the interline exchange have to be maintained for a stable and reliable operation of the system. The prime motive addressed in this chapter is to scheme a profligate and accurate controller with ability to sustain the frequency for the power system within nominal operating limits. A two-area reheat thermal system with generation rate constraint is considered, and a fuzzy logic with proportional integral controller is included for the enhanced operation in control of the governor and system response. The comparison of the obtained response for the hITLBO-DE to particle swarm optimization (PSO), pattern search (PS) and recently published results with hPSO-PS technique gives a clear view of the improvement in the system response.

A. Behera (✉) · T. K. Panigrahi · A. K. Sahoo · P. K. Ray
Department of Electrical and Electronics Engineering, International Institute of Information Technology Bhubaneswar (IIIT BBSR), Bhubaneswar, India
e-mail: C115002@iiit-bh.ac.in

T. K. Panigrahi
e-mail: tapas@iiit-bh.ac.in

A. K. Sahoo
e-mail: C116003@iiit-bh.ac.in

P. K. Ray
e-mail: prakash@iiit-bh.ac.in

# 1   Introduction

In current scenario power system automatic generation control (AGC) introduced a vital role to sustain the steadiness between generation and demand of load with reducing the variation of frequency [1, 2]. Tie-line power exchange is used here to regulate and facilitate contracts between covering all control areas for confirming reliable and quality process of the interconnected transmission system. AGC observes and controls frequency of the system and tie-line power flows, try to match demand variation and required generation considering average time of the area control error (ACE) with short value. Various control and optimization processes have been suggested for AGC of power system in the last decade such as genetic algorithm (GA) and nature-inspired techniques such as particle swarm optimization (PSO), pattern search (PS), bacteria foraging optimization algorithm (BFOA), teaching learning based optimization (TLBO), improved teaching learning based optimization (ITLBO), differential evolution (DE) and fuzzy logic controller (FLC). Fuzzy logic controller (FLC) is used to increase the ability of PI controller which is able to operate even with the alterations in functional point. It applies a dynamic modification of the parameters involved with the functioning of the controller. Fuzzy logic centred PI controller was used effectively for nonlinear system using precise mathematical formulation to choose suitable fuzzy parameters (i.e. rule base, input, membership functions, scaling factors, etc.). Certain pragmatic rules are applied for selection of fuzzy parameters. Proper choice of input and output scaling factors has been taken so as not to affect FLC performance.

The performance of AGC depends on the heuristic techniques as well as the controller structure. The adequate regulation by the applied technique for exact factors can escalate computational work or produce results with optimum values. Any algorithm-specific parameter is not required by TLBO but only involves mutual monitoring parameters (i.e. size of the population and generators number), which are common in running any population-based optimization algorithms. Certain developments of elementary TLBO technique are presented, which improve its search and exploitation capacities.

# 2   System Specification and Designing

A system rating of 2000 MW and nominal load of 1000 MW in each area are considered for the analysis. A load demand variation $\Delta P_{D1}$, $\Delta P_{D2}$ is applied to a system represented by system gain $K_{P1}$, $K_{P2}$ and time constant $T_{P1}$, $T_{P2}$. The governor droop regulator $R_1$, $R_2$ along with governor time constant $T_{G1}$, $T_{G2}$ in seconds controls the input to the turbine (i.e. $\Delta P_{G1}$, $\Delta P_{G2}$). The turbine characteristics are determined by time constant $T_{T1}$, $T_{T2}$ in seconds. The frequency deviation is given by $\Delta F_1$ and $\Delta F_2$ and the frequency bias parameter $B_1$ and $B_2$,

whereas the incremental tie-line power is given by $\Delta P_{\text{Tie}}$ and the synchronization coefficient by $T_{12}$. The significant values of the parameters are taken from [3].

The design of the system uses nonlinearities such as the generation rate constraint (GRC) and the reheat thermal system. The thermal unit has a GRC of 3%, and the reheat system is designed by $T_{\text{Re1}}$ and $T_{\text{Re2}}$. Figure 1 describes the power system models with the above parameters.

## 3  Structural Design of Controller

The proportional integral (PI) controller is used extensively due to its ability of improving steady state stability. They are operationally simple yet exhibit robust response for a widespread variety of working circumstances. The basic mathematical description in time domain for a PI controller is:

$$u(t) = \frac{c(t)}{r(t)} = K_{\text{P}} e(t) + K_i \int_0^t e(t).\,\mathrm{d}t \tag{1}$$

In Laplace domain:

$$U(s) = \frac{C(s)}{R(s)} = K_{\text{P}} + \frac{K_i}{s} \tag{2}$$

The above mathematical equation has been modified for the modelling of the fuzzy PI controller, described in Fig. 2, where $c(t)$ and $r(t)$ are the response signal



**Fig. 1** Simulink diagram for 2 area system studied with Fuzzy PI controller

and input to the PI controller and $K_p$, $K_i$ are controller control parameters. The Eqs. (1) and (2) are used for a system design with simple PI controller application.

# 4 Fuzzy Logic Controller (FLC) Modelling

The FLC provides a crisp output for a range of input depending on the values of the rule base and the membership function provided. The model is given in Fig. 2, and the area control error (ACE) is the input to the system as described by the Eqs. (3) and (4).

$$ACE_1 = B_1 \Delta f_1 + \Delta P_{tie} \tag{3}$$

$$ACE_2 = B_2 \Delta f_2 - \Delta P_{tie} \tag{4}$$

The "error (e) and derivative of error (e′)", [4] is taken as control signal to fuzzy logic controller, and signal at the FLC terminal is the control or reference signal for the two-area power system. The input tuning factors $K_1$, $K_2$, $K_p$ and $K_i$ and the parameters of PI controller are enhanced.

The FLC has five membership functions (i.e. "NB (negative big), NS (negative small), Z (zero), PS (positive small) and PB (positive big)") for each of the inputs and output [4, 5]. The computational proficiency and simple procedure of application of the triangular membership function is the reason of such extensive application. Mamdani selected as the fuzzy interface system and for defuzzification centre of gravity method is implemented.

# 5 Objective Function Formulation

The problematic consideration of the paper is making the system quicker and lesser liable to uncertainty through perturbation. To boost the system response, parameters such as rise time ($T_r$), peak time ($T_p$), overshoot/undershoot (OS/US), settling time ($T_s$) and steady state error ($e_{ss}$) are considered for the formulation of the objective



**Fig. 2** Simulink diagram for the fuzzy PI controller

function. Diverse objective function had been realised by researchers over the decade for optimization of the problem. The recurrently used functions are "integral of time and squared error multiplied (ITSE), integral of squared error (ISE), integral of time and absolute value of error multiplied (ITAE) and integral of absolute value of error (IAE)" as revealed by research work. The mathematical representation for the objective function used in the analysis is:

$$ITAE = \int_0^t t.(\Delta f_1 + \Delta f_2 + \Delta P_{tie})dt \tag{5}$$

## 6 Hybrid Improved Teaching Learning Based Optimization-Differential Evolution (HITLBO-DE)

The fundamental asset of differential evolution and improved teaching learning based optimization has been pooled so as to offer a methodology with improved optimization fitness [6, 7]. Under the status of learner not attaining enhanced value, a location different to the prior and lying inside a narrow domain is considered. Here, the differential evolution technique [8] has been inserted for realizing a new search area where more credible solution is obtainable. Figure 3 shows the flow chart for the hITLBO-DE technique. The step considered from individual algorithm is discussed herein; the productiveness of the approach is verified in two different models of a two equal area system and two unequal area systems as discussed in the result analysis section.

So the following equations are applied for the modification and faster learning of the students.

$$X_{mod}^i = X^i + \text{rand}(0, 1)[X_{teacher} - (T_f^i.X_{mean})] \tag{6}$$

$$X_{teacher} = X^i + M(X^m - X^n) \tag{7}$$

So three students are selected such as $X^i \neq X^k \neq X^j \neq X^l$.

$$X_{Gr} = X^j + F(X^k - X^l) \tag{8}$$

$$X_{mod}^i = X^i + \text{rand}(0, 1)[X^i - X_{Gr}] \tag{9}$$

The crossover operation in DE provides it the advantage of not being trapped by local minima and thereby converging to a global minimum point [9].

**Fig. 3** Flow chart for the hITLBO-DE technique

# 7 Result and Analysis of the System with Fuzzy PI Controller

The application of fuzzy logic to the tuning process of PI controller involves the hITLBO-DE optimization technique. The controller parameters $K_1$, $K_2$, $K_P$ and $K_I$ are tuned by the optimization technique based on the fuzzy logic interface. The tuned values for the controller with various techniques are presented in Table 1. The response is shown in Fig. 4, and the mathematical analysis is presented in Table 2 for case-1 where a 10% load disturbance is applied for both the areas. For case-2, the load disturbance is assumed to be 10% for area-1 and 20% for area-2 which is presented in Fig. 5 with mathematical analysis in Table 3.

**Fig. 4** **a** Deviation in frequency for area-1, **b** for area-2, **c** tie-line power with 10% step load change in both areas

**Fig. 5 a** Deviation in frequency for area-1, **b** frequency for area-2, **c** inter area exchange power with 10% step load change in area-1 and 20% step load change in area-2

**Table 1** The fuzzy PI controller parameter tuned by hITLBO-DE

| Techniques | Optimum controller parameter | | | |
|---|---|---|---|---|
| | $K_p$ | $K_i$ | $K_1$ | $K_2$ |
| PSO | 0.8176 | 0.7948 | 0.5085 | 0.5108 |
| PS | 0.4509 | 0.5470 | 0.7317 | 0.6477 |
| hPSO-PS [3] | 0.9336 | 0.7203 | 0.9852 | 0.5595 |
| hITLBO-DE | 1.9001 | 2.0709 | 1.7593 | 1.1398 |

**Table 2** System response parameters for 10% load variation in both area-1 and area-2

| Techniques | ITAE | Settling time (2%) $T_s$ (s) | | | Overshoot (OS) ($\times 10^{-3}$) | | |
|---|---|---|---|---|---|---|---|
| | | $\Delta f_1$ | $\Delta f_2$ | $\Delta P_{tie}$ | $\Delta f_1$ | $\Delta f_2$ | $\Delta P_{tie}$ |
| PSO | 0.4470 | 3.95 | 1.19 | 8.22 | −48.61 | −35.04 | −47.40 |
| PS | 0.6334 | 3.11 | 3.11 | 6.54 | −65.29 | −44.24 | −42.61 |
| hPSO-PS [3] | 0.1438 | 2.62 | 2.54 | 4.29 | −43.66 | −30.53 | −58.51 |
| hITLBO-DE | 0.0874 | 1.47 | 2.48 | 3.18 | −22.59 | −16.51 | −15.12 |

**Table 3** System response parameters for 10% load variation in area-1 & 20% in area-2

| Techniques | (2%) Settling time ($T_s$ (s)) | | | Overshoot (OS) ($\times 10^{-3}$) | | |
|---|---|---|---|---|---|---|
| | $\Delta f_1$ | $\Delta f_2$ | $\Delta P_{tie}$ | $\Delta f_1$ | $\Delta f_2$ | $\Delta P_{tie}$ |
| PSO | 3.81 | 3.59 | 7.07 | −151.72 | −193.0 | 125.5 |
| PS | 3.95 | 3.18 | 5.22 | −190.55 | −156.0 | 122.3 |
| hPSO-PS [3] | 2.78 | 3.81 | 3.83 | −128.59 | −135.2 | 87.87 |
| hITLBO-DE | 1.25 | 3.13 | 1.46 | −62.78 | −136.2 | 58.21 |

## 8 Conclusion

This chapter proposes a hybrid improved teaching learning based optimization with that of differential evolution (hITLBO-DE) for tuning of a fuzzy PI controller applied to a two-area reheat thermal system with generation rate constraint. The improved result is compared to recently published paper with hPSO-PS; thus, the superiority of the hITLBO-DE can be observed from the above graphs and tables in two different cases. In case-1, a 10% load disturbance is considered to be occurring simultaneously in both the areas, whereas in case-2 a 10% load disturbance is considered to be occurring in area-1 and a 20% load disturbance occurring in area-2 simultaneously. The complexity of the system is handled well by the proposed hITLBO-DE, which has the advantages of both ITLBO and DE algorithm and can be instrumental in providing accurate result in the field of automatic generation control.

# References

1. Bevrani, H.: Robust Power System Frequency Control. Springer, Berlin (2009)
2. Elgerd, O.I.: Electric Energy Systems Theory—An Introduction, 2nd edn. Tata McGraw Hill, New York (2000)
3. Sahu, R.K., Panda, S., Chandra Sekhar, G.T.: A novel hybrid PSO-PS optimized fuzzy PI controller for AGC in multi area interconnected power systems. Electr. Power Energy Syst. **64**, 880–893 (2015)
4. Chandrakala, K.R.M.V., Balamurugan, S., Sankaranarayanan, K.: Variable structure fuzzy gain scheduling based load frequency controller for multi-source multi area hydro thermal system. Int. J. Electr. Power Energy Syst. **53**, 375–381 (2013)
5. Sahu, B.K., Pati, S., Mohanty, P.K., Panda, S.: Teaching–learning based optimization algorithm based fuzzy-PID controller for automatic generation control of multi-area power system. Appl. Soft Comput. **27**, 240–249 (2015)
6. Chen, D., Zou, F., Li, Z., Wang, J., Li, S.: An improved teaching–learning-based optimization algorithm for solving global optimization problem. Inf. Sci. **297**, 171–190 (2015)
7. Venkata Rao, R., Patel, V.: An improved teaching-learning-based optimization algorithm for solving unconstrained optimization problems. Sc. Iran. D **20**(3), 710–720 (2013)
8. Mohanty, B., Panda, S., Hota, P.K.: Controller parameters tuning of differential evolution algorithm and its application to load frequency control of multi-source power system. Int. J. Electr. Power Energy Syst. **54**, 77–85 (2014)
9. Rout, U.K., Sahu, R.K., Panda, S.: Design and analysis of differential evolution algorithm based automatic generation control for interconnected power system. Ain Shams Eng. J. **4**(3), 409–421 (2013)

# Data Analysis of Weather Data Using Hadoop Technology

**Diallo Thierno Mamadou Oury and Archana Singh**

**Abstract** In the present age, all the real business forms comprise of colossal or huge information. As per the reports in the most recent 3 years, 90% of information was made from the utilitarian modules. The information is advancing exponentially, and the IT databases are moving to the capacity in terabytes along with these lines venturing into the period of huge information. The paper, investigates the progressions in different parameters of climate conditions like temperature, precipitation, snowfall, and so forth in the Pannonian Basin (focal Europe) and other comparative districts, from 1900 to 2014 with a base dataset containing day by day estimations by the climate stations arranged close to our purpose of examination. We have utilized the Hadoop innovation to actualize the weather data. For map reduction, we have utilized Apache PIG, and information is pictured by utilizing Python technology. The data results show the visualization Flot of weather data.

**Keywords** Big Data · Hadoop · Weather data · Map reduction

## 1 Introduction

The era of data and information at monstrous rate in each business issue needs support and learning extraction out of it. The conversion of this information into learning by gathering valuable inputs and significant yields in associations and connections existing in the data is what is really required. The accessible option strategy for prescient investigation of unstructured or semi-organized information is utilizing conventional strategies like RDBMS. Be that as it may, we have utilized Hadoop system for investigation, which is still creating an innovation. The favorable position is that Hadoop has substantial scale handling force and capacity with ease when contrasted with RDBMS. Aside from this Hadoop, likewise offers huge parallel handling capacities where one can run occupations in parallel to process

D. T. M. Oury · A. Singh (✉)
Amity University, Noida, Uttar Pradesh, India
e-mail: archana.elina@gmail.com

vast volumes of information. For investigation of organized or semi-organized information, we have been long utilizing the conventional database approaches like Relational Database Management Systems. Be that as it may, information these days is accessible in broadly substantial volumes (terabytes and petabytes) and in assortment of arrangements (organized, unstructured, and semi-organized). The current arrangement of conventional RDBMS is obviously not able to take into account this continually expanding information blast administration request. Therefore, an innovation like Hadoop comes into the picture and gives practical answers for predictive data analysis.

Hadoop answers a significant number of the enormous information challenges that we see today. The way to store terabytes of information and recover the data rapidly, how would we work with information that is in an assortment of various organizations, this question still remains? The diverse dataset arrangements being organized record frameworks like tables, semi-organized document frameworks like .xml records or any tab or comma isolated documents or unstructured record frameworks like messages, pdfs, and so forth. Hadoop can give solutions in an adaptable, modifiable, and flaw safe way. A typical example of Big Data is the climate information data looked over a great many years and managed by a substantial number of climate stations. We generally ponder in the matter of what the enormous organizations like Google, Facebook, Twitter, and so on do with this huge information. Google, back in the days when they were beginning, envisioned with respect to how they can make the web searchable? They expected to file a billion pages for that reason. So they assembled the innovation called map reduce alongside google file system (GFS) and that is truly what Hadoop depends on. Google recorded those billion pages taking into account this very innovation. Presently today, 60 billion pages are the thing that Google records to make web searchable. Another case in huge information business sector is that of Facebook which has the biggest Hadoop cluster on the planet.

They have a 100 PB of information. Also, on top of that, they create a large portion of a petabyte of information each and every day. All that anyone does on Facebook, for example, from signing in, to clicking, to preferring or remarking on something, is followed on Facebook, and this is the way they do communitarian separating procedures of advertisement focusing on as well. Twitter produces 400 million tweets a day, or we can say that it is roughly 100,000 tweets a moment. This is equivalent to around 15 terabytes of information every day. We ponder what they do with that information of theirs. Indeed, they have a considerable measure of Hadoop clusters where they run a huge number of MapReduce occupations toward the end of every day, breaking down that information in an assortment of approaches to find patterns. They know all the most recent and most prominent patterns thusly, and they likely offer that data to individuals who make items so they can focus on those patterns. Taking a shot at this information as base information for our venture, we first gather it and afterward import it into the Hadoop environment. We then examinations this information by utilizing PIG scripting inquiries and in the long run total the outcomes for information perception utilizing plot

graphs to give a prescient investigation for an unnatural weather change patterns over the years.

The paper investigated the utilization of Hadoop improvement environment and mines and examinations information by running questions utilizing PIG alongside a fundamental MapReduce operation. We significantly underlined on intuitive information perception by composing Python script for sorting, as it is of most extreme significance to make the appreciation of data simpler. Thirdly, a near examination between performing ETL utilizing conventional RDBMS and Hadoop is likewise done.

The paper is organized as follows in Sect. 2 related work done is explained, Sect. 3 explored the research methodology used. In Sect. 4, data analysis and results are explained. The Sect. 5 closes with discussion of results. Finally, the paper winds up with conclusion.

## 2 Related Work Done

The application of Big Data technology is applied in infrastructure data mining and applied analytics on Twitter [1]. The paper [2] explores the technology in heterogeneous information networks by using the structural analysis approach. This approach used relational database and semi-structured data to develop a model with nodes and links. The Big graphs and mining are demonstrated in the paper [3–5]. It explored the big web graphs and Twitter social data. Various websites mentioned the application of Big Data analytics in the weather data [6]. The data visualization technique and scripts were referred from the websites [5, 7, 8].

## 3 Research Methodology

In Fig. 1, the research model shows the flow of information throughout the work.

Step 1: The climate data and information which requires picking up a considerable measure of foundational data about the historical backdrop of and in addition to the present climate patterns and changes.

Step 2: This stage incorporates check and approval of information sources and settling on an appropriate toolset supporting the undertaking.

Step 3: Historical climate information is accessible and effectively downloadable from the National Climatic Data Center (NCDC).

Step 4: Construction of Metadata for the base information and focussing on the parameter to be worked upon like temperature or precipitation or snowfall and so forth.

Fig. 1 Research Model

Step 5: After utilizing the Hadoop structure introduced on the virtual machine and utilizing the abnormal state inquiry dialect, that is, PIG, for performing ETL operations.

Step 6: After conglomerating the outcomes from PIG ETL script, the data was assessed, made good with the organization reasonable for representation (utilizing SED summons) lastly approved by contrasting the outcomes and another comparable toolset and methodology like the customary RDBMS.

Step 7: It includes the visualization or interpretation of data, the step is fundamental to improve the presentation of information and its correspondence to the client's effectively justifiable utilizing Flot javascript plotting library.

## 4 Data Analysis and Results

The different modules of the framework the accompanying programming apparatuses and products were utilized. The Cloudera VMware workstation empowers the arrangement and administration of undertakings and toolsets like Apache Hadoop and other comparable activities. It likewise gives a coordinated and adaptable stage that encourages effortlessly reasonable and alarmingly expanding volumes and assortments of information. It additionally controls and breaks down information, and keeps it secure and ensured also. Hadoop is the structure backings the preparing of substantial information sets in a circulated processing environment. Utilizing Hadoop, we can run applications on the frameworks having a great many hubs including information of a large number of terabytes. HDFS is an appropriated and

elite document framework utilized for capacity and the MapReduce operation framework for parallel information handling.

## 4.1 Scripting Dialects Utilized

PIG script: This dialect is used by experts in SQL. The PIG script is used to perform the ETL operation for the data extraction.

Python script: This Python script language is used to access the library and operations like parsers. It is used here to program. The appropriate fitting in the datasets.

jQuery Flot script: jQuery Flot is a javascript plotting library utilized for information representation. It is used for data visualization using Flot.

## 4.2 Steps Used in Data Analysis

- Writing the PIG script evaluation and ETL operations.
- Linux SED charges to control yield information.
- Visualization of information with intelligent jQuery Flot Library.
- *Phase 1: The PIG script is used to perform the ETL (Extraction, Transform, and Loading) operation.*

Phase 1: Initially, the script loads the.CSV (comma separated file) as variable A, B variable is used to store the schema. In the variable C1, it stores the filtered or extracted data containing NOT NULL values. The parameter D1 is used to group the station data by station code, year, latitude, and longitude for weather data analysis prediction. By using mathematical function data is transformed. Finally, a separate variable is used to store the data results.

*Phase 2: Linux SED commands to manipulate output data.*

The SED (Stream Editor) commands are used to alter the output dataset from PIG and to make it compatible with the JSON format used for data visualization. The following three SED commands were used for manipulation of resultant dataset:

SED command1:

```
sed 's/(\([A-Z ]*\),\([0-9\.]*\),\([0-9.]*\))/[{"Amity_Weather": "\1", '\1: [\3, \2], /g'
rain_orig.csv > SED1.csv && cat SED1.csv
```

This SED command simply checks for the occurrence of any capital characters in the dataset followed by a comma or any number from zero to nine in blocks of two, separated by a comma, and starts changing it to a format suitable for visualization.

SED command 2:

```
sed 's/{((\([0-9]*\),/"\1": /g' SED1.csv > SED2.csv && cat SED2.csv
```

This SED commands parses for any no. in the dataset and encloses the same to double quotes, just another compatibility criteria in accordance with data visualization.

SED command 3:

```
sed 's/)),((\([0-9]*\),/, "\1":/g' SED2.csv > SED3.csv && cat SED3.csv
```

*Phase 3: Visualization of data with interactive jQuery Flot Library.*

## 5  Discussion

The weather data prediction was analyzed for the basin of central Europe and visually represents the trends rain in those regions. The resultant datasets for the stages of measuring and analysis phase are as follows:

In Fig. 2, using extract, transform and load (ETL) operations, all missing and incomplete data were detected and removed.

*Phase 2: Linux SED commands to manipulate output data.*



**Fig. 2** PIG script for analysis and ETL operations

**Fig. 3** Showing the output of SED commands



**Fig. 4** Visualization of data with interactive jQuery Flot Library

After the three SED commands were run, the resultant manipulated data, compatible with the JSON format, are as shown in Fig. 3.

In Fig. 4, the weather data visualization Flot is demonstrated.

## 6   Conclusion and Future Scope

In this paper, the climatic weather data was analyzed to predict rain using Hadoop technology. It used a virtual machine (here, Cloudera VMware). The point by point investigation of enormous information analysis and utilized its structure to break down the climate information. The paper utilized the scripting dialect and to comprehend the information representation component like Flot and Map information is pictured. This is the place, where huge information examination

advancements like Hadoop act the hero and have colossal future extension. Here, weather data was collected and analyzed using Hadoop and Map reduction techniques.

The future scope of our task is that climate data at different areas can be gathered on an every day, hourly or even minutely premise. Along these lines with more information being produced all the time, the undertaking configuration can be utilized for predictive analysis of changing climate drifts and to consider and make cataclysmic danger models for catastrophe inclined areas.

# References

1. Lin, J.: Mapreduce is good enough? if all you have is a hammer, throw away everything that's not a nail! Big Data **1**(1), 28–37 (2012). CoRR, abs/1209.2191
2. Miller, T.W.: Modeling Techniques in Predictive Analytics: Business Problems and Solutions with R, Revised and Expanded 1st edn. Pearson FT Press (2013)
3. Kang, U., Chau, D.H., Faloutsos, C.: Pegasus: mining billion-scale graphs in the cloud. In: ICASSP, pp. 5341–5344 (2012)
4. Fayyad, U.: Big data analytics: applications and opportunities in on-line predictive modeling (2012). http://big-data-mining.org/keynotes/#fayyad
5. Mayer-Schonberger, V.: Big data: a revolution that will transform how we live, work and think paperback
6. jQuery Flot tutorial. http://www.jqueryflottutorial.com/what-is-jquery-flot.html
7. Weather data analysis and visualization—big data tutorial part 1/9—fundamentals. file:///C:/Open-BigData/Upload/Weather-analysis.htm
8. Python Software Foundation (US). http://www.jqueryflottutorial.com/what-is-jquery-flot.htmlSub
9. Hadoop tutorial. http://hortonworks.com/hadoop-tutorial/hello-world-an-introduction-to-hadoop-hcatalog-hive-and-pig/

# Comparative Analysis of Coherent Routing Using Machine Learning Approach in MANET

Ayushree and G. N. Balaji

**Abstract** Ad hoc network is a network which is dynamic in nature where the mobile nodes form a temporary network in the absence of centralized administration. Due to the absence of centralized administrator in network, routing in mobile ad hoc network (MANET) becomes the fundamental issue which minimizes the selection of an optimal path for routing. Certain performance parameters such as latency, overhead, and packet delivery ratio (PDR) are affected adversely for which numerous techniques are advocated that enhances the selection of efficient and stable path. In the present paper, an attempt is made to select the optimal path and compare the results by varying the number of nodes by using knowledge-based learning algorithm. The optimal path will possess the highest average sum of relay nodes and will be considered as the most optimal and reliable path. We also proposed that analysis of throughput and PDR is better as compared to the traditional methods. The simulation is carried out at NS-2 network simulator, which is employed to implement wired and wireless ad hoc simulation.

Ayushree (✉)
Discipline of Electronics and Communication Engineering, K L University, Guntur District, Andhra Pradesh, India
e-mail: ayushrees@gmail.com

G. N. Balaji
Department of Computer Science and Engineering, K L University, Guntur District, Andhra Pradesh, India
e-mail: balaji.gnb@gmail.com

# 1 Introduction

A MANET is a network system where every node works collectively without any hindrance from the centralized authority. An increase in interest is observed in mobile wireless communication due to their pervasive feature. MANET is such a network which provides the flexibility to the network system by maintaining the fixed network and exchange of information without any access point or base station requirement. Multi hop communication enables the network to achieve this and permits node to approach distant nodes by means of relay nodes or intermediate nodes. A fundamental problem observed in MANET network is the maintenance and selection of stable multilevel hop path. Numerous elements are evolved that would cause asymmetry in the network topology; certain factors can be node mobility, excess power consumption, and signal interference intervention, which leads to the loss of path, and as consequence, the information is lost. An area which is exposed to the natural catastrophe prevents any sort of communication with the outside world [1].

In this paper, our major focus is achievement of stable routing path, as link can be lost due to node mobility. This is attained by the usage of machine learning algorithm. A machine learning algorithm enlists previous studied network and concludes on the basis of obtained knowledge. Further, comparison is made between the cases of varying number of nodes based on certain performance factors. AODV protocol is considered for relay of data from origin to end [2]. AODV is one of the types of reactive protocol. AODV makes use of sequence number rather than routing tables for creation of path so that looping complexity can be ignored [3].

Ample study is done on energy exhaustion, and many techniques are proposed for the attainment of the most reliable and optimal path. In [4], a network topology is proposed where it is centrally coordinated by a leading node. In paper [5], the parameter overall stability is calculated which will compute the stability per node in the network system. An inverse relation is seen in nodes mobility and stability factor.

In [6] and [7], comparison is made between on-demand routing protocols by which we can make decision of finest routing algorithm. Papers [7, 1, 2] advocate efficient methods of estimation and conservation of energy. Followed by these papers, various methodologies are proposed in [3, 8–10] so as to restrain latency and overhead by minimizing the path loss.

An attempt of comparison between malicious attack and black hole attack is made in the presented paper. In proposed work, we have introduced knowledge learning algorithm which would enable the network to gather the information about its initial performance, and on the basis of this gathered information, optimum route is established. Henceforth, comparison is made between the varying number of nodes by considering certain performance parameters: packet delivery ratio (PDR), throughput. One is such a case where nodes are imposed to foreign attack and the other is ideal case (no foreign attack).

The remaining paper is organized as follows: Sect. 2 discusses proposed work; Sect. 3 presents research approach; Sect. 4 comprises simulation, result, and analysis; finally, Sect. 5 presents conclusion and future scope.

## 2 Proposed Work

The pattern recognition allows gathering of the knowledge about mobility of node and by means of these patterns relay numbers are allotted accordingly. Whenever the information is transmitted from one to another node, the power required is inversely proportional to the nth power of the distance ($d$) between these nodes by $1/d{\wedge}n$. Here, $n$ ranges between 2 and 4 on the basis of the distance between the observed nodes. For successful routing, signal to noise ratio (SNR) of node two must be in excess with threshold value. If

$n_i$   origin node
$n_j$   end node
$\Psi$   threshold value



**Fig. 1** Flow chart of proposed work

The SNRj must satisfy the following condition: (Fig. 1)

$$\text{SNR}_j = \frac{P_i G_{i,j}}{\sum_{k \neq i} P_k G_{k,j} + \eta_j} \Psi_j(\text{BER}) \tag{1}$$

In Eq. (1)

$P_i$     Transmitted power of $n_i$,
$G_{i,j}$ Path gain between $n_i$ and $n_j$
$\Psi_j$  Threshold value

$$G_{i,j} = \frac{1}{d_{i,j}^n} \tag{2}$$

M      Node mobility
N      Integer no. 1,2,3…
B      Black hole attack node
Mal.   Malicious attack node
IN     Intermediate/relay node.

$$Q = \{o, d, \text{in}, l\} \tag{3}$$

$$P_{o \to \text{in}} = \left\{ (m, n) \big| \text{power}_{o \to \text{in} \to (m,n)} < \text{power}_{o \to (m,n)} \right\} \tag{4}$$

$Q$   Network system,
$o$   Source node,
$d$   End node,
in   Relay node,
$l$   Link failure,
$P$   Path

The Eq. (4) states that if the data packets are routed in the network system from their origin node to any arbitrary point $(m, n)$, then the power requirement for direct transmission of data packets from origin to intermediate is greater than the power requirement for indirect transmission. With the link failure, the above equation can be modified, and henceforth, the data packets are given as follows:

$$P_{o \to \text{in} \to l} = \left\{ (m, n) \big| \text{power}_{o \to \text{in} \to l \to (m,n)} < \text{power}_{o \to (m,n)} \right\} \tag{5}$$

$$\text{DP}_o > \text{DP}_d \tag{6}$$

$\text{DP}_o$   Data packets at origin node,
$\text{DP}_d$   Data packets at destination node

After the hindrance in the transmission of data packets, knowledge learning algorithm is executed and the equations obtained are as follows:

$$Q' = \{o, d, \text{in}_n\} \tag{7}$$

$$o = \{o_{t1}, o_{t2}, \ldots o_{tm}\} \tag{8}$$

$$d = \{d_{t1}, d_{t2}, \ldots d_{tm}\} \tag{9}$$

$$\text{in}_1 = \{\text{in}_{t11}, \text{in}_{t12}, \ldots .\text{in}_{t1n}\} \tag{10}$$

$$\text{in}_N = \{\text{in}_{tN1}, \text{in}_{tN2}, \ldots .\text{in}_{tNn}\} \tag{11}$$

Equations (7), (8), (9), (10), and (11) present mobility samples of every node at different time interval. By means of these mobility samples, assumptions are done on the basis of which relay number is subsequently allotted to nodes of the examining cluster. If 'R' is the relay number, then the range of relay number is between 1 and 10.

$$1 \leq R \leq 10$$

If $R_s$, $R_d$, $R_i$ are the relay number at source node, destination node, and intermediate node, then the resultant relay number for respective intermediate node is given by following equations.

$$R_1 = \sum (R_s + R_d + R_{\text{in}1}) \tag{12}$$

$$R_2 = \sum (R_s + R_d + R_{\text{in}2}) \tag{13}$$

$$R_N = \sum (R_s + R_d + R_{\text{in}N}) \tag{14}$$

The above Eqs. (12), (13), (14) results the average of relay number linked in the path of relay node i1, i2, up to in.

$$\text{If} \quad R1 > R2$$
$$\text{Than} \quad m1 < m2$$

$$Q'' = \{o, d, in_1\} \tag{15}$$

As stated earlier, there is an inverse relation between relay number and mobility. Increase in relay number would lead to the decrease in mobility of node. The resultant path obtained out of $R_1$ and $R_2$ is $R_1$, and the equation is given as Eq. 15 (Fig. 2).

The above considered case is the ideal case in which no foreign attack is imposed over the network system. In the proposed work, two other cases are

**Fig. 2** Knowledge-based learning algorithm

considered which would involve the foreign attack, malicious attack, and black hole attack over the network system.

The mathematical description of attacks on network system is given as follows:

$$Q_B = \{o, d, \text{in}, B\} \tag{16}$$

Here, B represents the Black hole node.
If

$$o = \{o_{t1}, o_{t2}, \ldots o_{tm}\} \tag{17}$$

$$d = \{d_{t1}, d_{t2}, \ldots d_{tm}\} \tag{18}$$

$$\text{in}_1 = \{\text{in}_{t11}, \text{in}_{t12}, \ldots \text{in}_{t1n}\} \tag{19}$$

$$P_{o \to \text{in}} = \left\{(m, n) \middle| \text{power}_{o \to \text{in} \to (m,n)} < \text{power}_{o \to (m,n)}\right\} \tag{20}$$

$$\text{If} \quad \text{in} \in B$$
$$\text{Than} \quad Q_B \notin \{d\}$$

Therefore the final selected path QB is given by Eq. 21:

$$P_{o \to \text{in}} = \left\{(m, n) \middle| \text{power}_{o \to \text{in} \to (m,n)} < \text{power}_{o \to (m,n)}\right\} \tag{21}$$

Whenever a network possesses black hole attack, loss of route orientation occurs by giving a false reply from the affected node. In black hole node, whenever the affected route receipts route request (RREQ) from the adjacent node, then it reverts

route reply (RREP) to the origin node which leads in the tampering of the network. The reply from tampered path will be considered as the most reliable path which would lead in false network maintenance (Fig. 3).

Another case involves malicious attack which would intercept the communication and inhibit transit of data packets. The mathematical formulation of above case is given as follows:

$$Q_M = \{o, d, \text{in}, M\} \tag{22}$$

M   Malicious node

$$o = \{o_{t1}, o_{t2}, \ldots o_{tm}\} \tag{23}$$

$$d = \{d_{t1}, d_{t2}, \ldots d_{tm}\} \tag{24}$$

$$\text{in}_1 = \{\text{in}_{t11}, \text{in}_{t12}, \ldots \text{in}_{t1n}\} \tag{25}$$

$$P_{o \to \text{in}} = \left\{(m, n) | \text{power}_{o \to \text{in} \to (m,n)} < \text{power}_{o \to (m,n)}\right\} \tag{26}$$

If   in $\in M$

Than   $Q_M \notin \{d\}$

Therefore the final selected path QM is given by Eq. 27.

$$Q_M \in \{o, d, \text{in}\} \tag{27}$$



**Fig. 3**  Induced black hole node

# 3   Research Approach

## 3.1   Implementation

NS2.35 (Network Simulator) is the tool used in Linux (Ubuntu 12.04) operating system for the implementation of proposed work. The parameter values are obtained from awk scripts (Fig. 4).

## 3.2   Performance Parameters

Throughput: It is the measure of how rapidly data is transmitted in a network. Theoretically, it appears that bandwidth and throughput are same, but in reality, both possess a different meaning.

Packet delivery ratio: It is the ratio of summation of total received packets to the summation of sent data packets. PDR is given as:

$$PDR = \frac{\sum \text{Total number of received packets}}{\sum \text{Total number of sent packets}} \tag{28}$$



**Fig. 4** Introduction of malicious node

## 4 Simulation, Result, and Analysis

In our proposed work, the simulation is performed by using NS2 tool (network simulator) in Ubuntu 12.04. Simulation parameters and specifications are as follows in Table 1.

The presented graphs provide the performance comparison of considered performance metrics: throughput and PDR. The results of network systems possessing 12 nodes are shown in Fig. 5a, b which leads to the dynamic increment of throughput and PDR. The results are shown in Fig. 6a, b for the same.

From Tables 2 and 3, it can be concluded that increase in number of nodes can cause increase in throughput in all the cases considered.

## 5 Conclusion and Future Scope

The intention is to evaluate the performance of the network based on throughput and PDR. AODV protocol is made in use for routing of packets. The paper focuses on per node analysis rather than per flow analysis. The comparison made between the varying numbers of nodes in network system shows that with increase in number of nodes throughput as well as PDR is decreasing. The considered knowledge learning algorithm helps us to determine the reliable path by means of relay number which is inversely proportional with mobility. Finally, on the basis of our findings, an approach of finding and selecting optimal route is proposed which would account for improving throughput, PDR, and overhead of routing protocols. Further study is still needed to be done where new reactive protocols can be taken in consideration for the implementation of above discussed method. Moreover, comparison can be made between these cases by variation in the number of nodes.

**Table 1** Simulation parameters

| Simulation parameter | Specification (12 nodes) | Specification (24 nodes) |
|---|---|---|
| Simulation duration-ideal | 1 m 10 s | 1 m 34 s |
| Simulation duration-malicious attack | 1 m 24 s | 1 m 30 s |
| Simulation duration-black hole attack | 1 m 5 s | 1 m |
| Channel specification | Wireless channel | Wireless channel |
| Type of antenna | Omni-directional | Omni-directional |
| Propagation model | Two ray ground | Two ray ground |
| Nodes | 12 | 24 |
| Total malicious node | 1 | 1 |
| Number of black hole node | 1 | 1 |
| Data packets | 50 | 50 |
| Traffic type | Constant bit rate | Constant bit rate |
| Routing protocol | AODV | AODV |

**Fig. 5** **a** Throughput (12 nodes), **b** PDR (12 nodes)



**Fig. 6** **a** Throughput (24 nodes), **b** PDR (24 nodes)

**Table 2** Throughput and PDR for 12 nodes network system

| Performance analysis | Ideal | Malicious attack | Black hole attack |
|---|---|---|---|
| Throughput (kbps) | 274.50 | 176.45 | 145.86 |
| Packet delivery ratio | 0.919235 | 0.69567 | 0.27454 |

**Table 3** Throughput and PDR for 24 nodes network system

| Performance analysis | Ideal | Malicious attack | Black hole attack |
|---|---|---|---|
| Throughput (kbps) | 449.32 | 278.02 | 175.3 |
| Packet delivery ratio | 0.954368 | 0.732756 | 0.35691 |

# References

1. Tseng, Y.-C., Li, Y.-F., Chang, Y.-C.: On route lifetime in multihop mobile ad hoc networks. IEEE Trans. Mobile Comput. **2**(4), 366–376 (2003)
2. Basarkod, P.I., Manvi, S.S., Albur, D.S.: Mobility based estimation of node stability in MANETs. In: International Conference on Emerging Trends in Computing, Communication and Nanotechnology (ICE-CCN), International Conference, India, pp. 126–130 (2014)

3. De Couto, D.S.J., Aguayo, D., Bicket, J., Morris, R.: A high throughput path metric for multi-hop wireless routing. In: Proceedings of the 9th Annual International Conference on Mobile Computing and Networking, pp. 134–146 (2003)

4. Aggarwal, S., Ahuja, A., Singh, J.P., Shorey, R.: Route lifetime assessment based routing protocol for mobile adhoc networks. In: IEEE International Conference on Communications, ICC'00, New Orleans, LA, pp. 1697–1701 (2000)

5. Dube, R., Rais, C.D., Wang, K.-Y., Tripathi, S.K.: Signal stability based adaptive routing (SSA) for ad-hoc mobile networks. IEEE Pers. Commun. **4**(1), 36–45 (1997)

6. Su, W., Lee, S.-J., Gerla, M.: Mobility prediction and routing in ad hoc wireless networks. Int. J. Netw. Manag. **11**(1), 3–30 (2001)

7. Jones, C.E., Sivalingam, K.M., Agrawal, P., Chen, J.C.: A survey of energy efficient network protocols for wireless networks. Wirel. Netw. **7**(4), 343–358 (2001)

8. Muthuramalingam, S., Janani, P., Bavya, B., Rajaram, R.: An energy conserving topology maintenance algorithm for MANET. In: 1st International Conference on Network and Communications, pp. 101–106 (2009)

9. Chiasserini, C.-F., Rao, R.R.: Energy efficient battery management. IEEE J. Sel. Areas Commun. **19**(7), 1235–1245 (2001)

10. Rodoplu, V., Meng, T.H.: Minimum energy mobile wireless network. IEEE J. Sel. Areas Commun. **17**(8), 1333–1344 (1999)

# Constrained Level Validation of Serial Peripheral Interface Protocol

**Avinash Yadlapati and Hari Kishore Kakarla**

**Abstract** The motivation behind this paper is to give a full portrayal of a state-of-the-art SPI master/slave usage. Every single related issue, beginning from the elaboration of introductory details, till the last framework confirmation, is thoroughly examined and justified. In similarity with outline reuse approach, the concerned paper imparts high-grade intellectual properties i.e., IP's that concerns in gathering all essential components that help in achieving presently required and modern ASIC or SoC applications using this SPI master and slave protocol. SPI is a standout among the most usually utilized serial interface protocols.

**Keywords** System on chip (SoC) · Serial to peripheral interface (SPI)
Control registers · Verilog · Constraints · ARM

## 1 Introduction

SPI interface i.e., serial to peripheral interface protocol was introduced with the aim to get fast exchange of data between devices by replacing parallel interface mechanism; as a result, we can avoid parallel routing. SPI has become a standardized way of communication due to its advantages like fast communication, less effort in interfacing thus helping in providing effective communication in present scenario.

As a result, SPI received a good market in the embedded system arena and on chip processing systems including the higher end processors i.e., 32-bit and high.

A. Yadlapati
Department of Electronics and Communication Engineering, K L University, Green Fields, Vaddeswaram, Guntur 522502, Andhra Pradesh, India
e-mail: avinash.amd@gmail.com

H. K. Kakarla (✉)
Department of Electronics and Communication Engineering, K L University, Guntur District, Andhra Pradesh, India
e-mail: kakarla.harikishore@kluniversity.in

For example, those utilizing MIC, power PC, or ARM and even those with different microcontrollers like PIC, AVR. In most of the cases, SPI controllers are integrated and are capable to work as master or slave. SPI is being utilized in Field programmable gate array or chip-based design for efficient communication. Hence, nowadays SPI is a typical protocol utilized for correspondence with Processors or Microcontrollers where information should be exchanged rapidly. ETHERNET, USB, SATA, PCI-EXPRESS [1], controller area network (CAN), UART, USART, etc. are utilized for long distance while protocols like SPI and I2C are utilized for shorter distance communication. SPI is one among the serial interface protocols. It has ease of utilization, less number of pins advantage, and high transmission speed. At least, four interfaces are required by standard SPI protocol for functioning. For the most part, the devices which are based on SPI protocol are partitioned into master and slave to [2] transmit required information. The chip select signal and clock signal have to be engendered by the master when the data is processed. It can be précised that SPI is "Short Distance" communication protocol which can be used for onboard communication [5]. With reference to slower interface protocol definition, there are primarily two protocols i.e., I2C i.e., inter-integrated circuit protocol and the SPI. Both these protocols provide the best way of onboard communication for peripherals and integrated circuits. SPI is one among the most typically utilized serial protocol for low/medium rate information stream transfer on both inter-chip and intra-chip. The SPI design is utilized for communication between various peripherals present in a system on chip (SoC) with that of a processor application and, also utilized as part of real system, in order to provide proper communication, which is necessary for current day applications like Avionics and Defense [1].

## 1.1 SPI Protocol

SPI protocol is a fast, full-duplex, synchronous communication bus. SPI occupies less space when it is integrated on the printed circuit boards (PCB) as it is having just four ports. It has one master device and one slave device when it is in full-duplex mode. The four lines of SPI ports are MISO i.e., master in slave out, MOSI i.e., master out slave in, SCK i.e., slave clock, and SS i.e., slave select. The operation of the SPI is as follows: Master pulls down the SS line low to select slave when it needs to send information to a slave.

The data to the MOSI line is given by the master. SPI transmits the information bit by bit which is nothing but the serial communication. The clock pulse is given by SCLK [2] and MOSI [3]. Based upon on this pulse, MISO transmits the information [6]..At the rising or falling edge of the clock, data is output from the master through the master output line and followed by slave's "read." Hence, minimum eight clock cycles are required for 8-bit data transfer (Fig. 1).

**Fig. 1** SPI schematic symbol

## 2 SPI Block Diagram Descriptions

The figure shows a block diagram of SPI circuitry. Master transmits data to the slave via MOSI and receives data from it through MISO. This illustrates full-duplex mode of transmission with both data-in and data-out synchronized (Fig. 2).

DIVIDER: The clock is the input. It divides the clock frequency by 2, 4, 8, 16, 32, 64, 128, 256, and then it is sent into the selector block (mux).

SELECTOR: The clock from the divider is the input to the selector which acts as 8:1 mux with baud rate registers as select lines (SPPR0–SPPR2), (SPR0–SPR1).

SPI BAUD RATE REGISTER: It consists of SPPR0, SPPR1, SPPR2, SPR0, SPR1, and SPR2. By which we can calculate Baud Rate

Baud Rate Divisor = (SPPR + 1) • $2^{(SPR+1)}$
Baud Rate = Bus clock/Baud Rate Divisor

SPI CLOCK LOGIC: It is set from the control register 1 with MASTER, CPHA, and CPOL and sets the master with the mode the baud clock is received from the baud clock generator and the clock from mux is given as output to the SPI control. Master and slave are operated on clock from the clock logic and from clock logic [7].

SPI CONTROL: The SPI control sets the status register pins SPIF, SPTEF, MODF pins so that the master mode is ready.

SPI CONTROL LOGIC: In this block the master mode and slave modes are set from the inputs pins which come from the clock logic and SPI control block and also takes the inputs from the control register 2 and the data transfer is done from the shift register with the SPI data register from the MISO and MOSI pins the data is transferred the SCLK is set from the clock logic and the ss ∼ is set low until the transmission is done otherwise the transmission is stopped.

SHIFT REGISTER: MSB of shift register data goes out through MOSI pin from master and receives the data into LSB of slave shift register. The data in slave from MSB goes through MISO pin to master, and the data is transferred to LSB of master shift register.

**Fig. 2** Block diagram of SPI

## 2.1 SPI Signal Descriptions

*MISO: Master in Slave out*

There is an input to the master and output of a slave. MISO is one among the two lines which transmits data only in one direction serially, firstly the MSB is sent out. Whenever the slave is deselected, the MISO line is at high-impedance state.

*MOSI: Master out Slave in*

This is an input to a slave device as it is generated by the master device as output. It transfers serial information in one direction, by pushing out MSB first.

*SCK: Serial Clock*

It is used to synchronize data transmission which is done through MOSI and MISO lines. The master and slave are skilled enough to interchange a byte of data within eight clock cycles. SCK is an input to slave as it is generated by the master.

*SS: Slave Select*

SS input is used to select a slave device. It is set to "0" for transforming data and should be low till the end.

## 2.2  Signal Description of SPI

See Table 1.

## 3  Functional Descriptions

The transmission of data in the SPI is serial through a channel. This communication between master and slave is provided serially in the SPI system. To activate the system, the control is taken by the control register by enabling the SPIE bit.

There are four pins that are associated with the SPI function. These four pins are namely.

- Master out slave in pin i.e., MOSI
- Master in slave out pin i.e., MISO
- Serial clock pin i.e., SCK
- Slave select pin i.e., SS

The vital components of the SPI are data transmit and data receive registers. When the transmit register of slave's data is transferred to the master's shift register of width 8, the shift register in the slave is connected by MOSI and MISO pins to form bit registers of widths 16, 32, 64, respectively.

The master and slave [8], which are selected one, complete their data transmission only when a task of serial shifting 16-, 32-, 64-bit registers to 8, 16, 32 positions by the master when a serial clock was generated. The transmission ends with the data transfer on the shift registers to the data receive registers [2], and this

**Table 1**  Signal description

| Pin | Type | Function |
| --- | --- | --- |
| SCLK | Output | Serial clock output in master mode |
| SS | Output | Slave chip select output in master mode |
| MOSI | Output | Serial data output in master mode |
| MISO | Input | Serial data input in master mode |
| DATA IN | Input | Data given to master from peripheral |
| DATA OUT | Input | Data received by the peripheral |

data can be read any time before completion of next transfer. There are four clock formats, and the needed one is selected by the clock phase clock polarity. Clock polarity bit is capable of selecting an inverted or non-inverted clock [4].

### 3.1  Master Code

The master mode is selected when MSTR is set. For the transmission of data, data register has to be operated as a shift register. The data-widths vary as multiples of 8 and the bits start transmitting on the MOSI at SCK. The speed of the transmission 10 is monitored by the prescaler register. The SCK is the SPI clock output. During master mode operation [5], data is send through MISO pin. Slave was selected by using slave-select register, while master's SS pin tries to become active. The entire operation will hinge on the control register pins (CPHA and CPOL) which are characterized by the clock phase and clock polarity (Fig. 3).

### 3.2  Slave Code

SPI is a slave when the MSTR is reset. When the SPI is a slave, the working of the pins will differ from the previous case. SCLK which is the serial clock will act as the clock input from the master, the data output pin will be master in slave out that is MISO, the serial data output pin is nothing but the master out slave in (MOSI), and SS (slave-select) pin is the input. In order to start the transmission, SS pin should be pulled down and must stay low till the end of transmission. SPI will attain



**Fig. 3** Master block

IDLE state when SS is high and the slave is deselected, and hence, the shifting operation stops when SCLK input is ignored. The data created by the odd numbered edges on the serial clock is to be latched when the clock phase is reset.

The shifting of the data from data input pin (serial) into the shift register is caused by the even edges on the serial clock. Data on the data input pin (serial) is shifted and latched when CPHA is set. Serial clock latches the data at the input pin for data on the even numbered edges and shifts the latched data from the data input pin (serial) into the shift register on the odd numbered edges of SCLK, it may be LSB/MSB. Similarly, when CPHA = "0" and SS input is low, the data driven out of MISO pin is first bit of the SPI information. When CPHA is set, the first data bit onto the output pin is delivered at its first edge. The transmission complete flag is set in order to indicate the end of transmission [9] (Fig. 4).

## 3.3 Timing Diagrams

Timing analysis when CPHA = 0 (Fig. 5).
    Timing analysis when CPHA = 1 (Fig. 6).



**Fig. 4** Slave block

**Fig. 5** Timing analysis of MOSI and MISO when CPHA = 0



**Fig. 6** Timing analysis of MOSI and MISO when CPHA = 1

# 4   Simulation Results

Simulation is being carried out on Cadence NC simulator, using Verilog as a hardware description language. The write operation of the proposed method SPI protocol was shown in below Fig. 7. This shows the accurate data transmission rate from master to slave via SPI protocol.

Results are validated as per the specifications for applied clock input to the master. MISO transfers the corresponding data, at this instant SPI will be operating in write mode of operation.

The read operation of the proposed method SPI protocol was shown in below figure. This shows the accurate data transmission rate from master to slave via SPI Protocol.

Read mode of SPI is validated for applied clock input; MOSI is enabled and data is transferred through it (Fig. 8).



**Fig. 7**   Master simulation results



**Fig. 8**   Slave simulation results

## 5   Conclusions

Our work concentrates on the execution and examination of SPI alongside the RTL design work. Design part includes around Motorola 2.07 specification, and codes are written in the Verilog. RTL code thought of all blocks work done using NCSim tool which gives simulation and synthesis validation of specification. This paper delivers that SPI is an efficient protocol for secured, proper data core communication. This work concentrates on how SPI transfers the data in all modes of operations. SPI is designed, verified, and observed the simulation summary.

**Future Work**
Present work can be implemented on a FPGA Kit, and the same thing can be applied for controlling microwave and intelligent SPI Controllers.

## References

1. Pachler, W., Pressel, K., Grosinger, J., Beer, G., Bosch, W., Holweg, G., Zilch, C., Meindl, M.: A novel 3D packaging concept for RF powered sensor grains. In: 2014 IEEE 64th Electronic Components and Technology Conference (ECTC), pp. 1183–1188 (2014)
2. Das, R., Singh, G.K., Mehra, R.M.: Two-phase clocking scheme for low-power and high-speed VLSI. Int. J. Adv. Eng. Sci. Technol. **2**(2), 225–230 (2013)
3. Bais, A., Singh, G.K., Mehra, R.M.: Design of 6-T SRAM cell for enhanced read/write margin. Int. J. Adv. Electr. Electron. Eng. **2**(2), 317–325 (2013)
4. Aditya, K., Sivakumar, M., Noorbasha, F., Praveen Blessington, T.: Design and functional verification of a SPI master slave core using system verilog. In: IJSCE, vol. 2(2), ISSN:2231-2307 (2012)
5. Liu, T., Wang, Y.: IP Design of Universal Multiple Devices SPI Interface. Department of Electronic Engineering, Xiamen University, IEEE (2011)
6. Sandya, M., Rajasekhar, K.: Design and verification of serial peripheral interface. Int. J. Eng. Trends Technol. **3**(4), 522–524 (2012)
7. SPI Block Guide V03.06, Document number S12SPIV3/D, Original Release Date: 21 JAN (2000), Revised: 04 FEB (2003), MOTOROLA INC
8. Oudjida, A.K., Berrandjia, M.L., Liacha, A., Tiar, R., Tahraoui, K., Alhoumays, Y.N.: Design and test of general purpose SPI Master/Slave IPs on OPB bus.In: 2010 7th International Multi-conference on Systems Signals and Devices (SSD), 27–30 June. IEEE Press, Amman (2010)
9. Oudjida, K., Berrandjia, M.L., Tiar, R., Liacha, A., Tahraoui, K.: FPGA implementation of I$^2$C and SPI protocols: a comparative study. In: Proceedings 16th IEEE International Conference on Electronics, Circuits, and Systems, pp. 507–510 (2009)

# Retraction Note to: Real-Life Facial Expression Recognition Systems: A Review

**Samta Jain Goyal, Arvind K. Upadhyay, R. S. Jadon and Rajeev Goyal**

**Retraction Note to:**
**Chapter "Real-Life Facial Expression Recognition Systems: A Review" in: S. C. Satapathy et al. (eds.),** *Smart Computing and Informatics*, **Smart Innovation, Systems and Technologies 77,**
https://doi.org/10.1007/978-981-10-5544-7_31

The editor has retracted this chapter [1] because of significant overlap with an earlier published article by Benta and Vaida [2]. Authors Samta Jain Goyal, Arvind K. Upadhyay and Rajeev Goyal agree with this retraction. Author R.S. Jadon has not responded to any correspondence from the publisher about this retraction.

[1] Goyal S.J., Upadhyay A.K., Jadon R.S., Goyal R. (2018) Real-Life Facial Expression Recognition Systems: A Review. In: Satapathy S., Bhateja V., Das S. (eds) Smart Computing and Informatics. Smart Innovation, Systems and Technologies, vol 77. Springer, Singapore
https://doi.org/10.1007/978-981-10-5544-7_31

[2] K.-I. Benta, M.-F. Vaida, "Towards Real-Life Facial Expression Recognition Systems," Advances in Electrical and Computer Engineering vol. 15, no. 2, pp. 93–102, 2015, https://doi.org/l0.4316/AECE.2015.02012

---

The retracted online version of this chapter can be found at
https://doi.org/10.1007/978-981-10-5544-7_31

# Author Index