



Contents

Editorial

Papers

- 357 What makes a board director better connected? Evidence from graph theory
Laleh Samarbakhsh, Boža Tasić
- 379 Missing Data Imputation in Cardiometabolic Risk Assessment:
A Solution Based on Artificial Neural Networks
Dunja Vrbaški, Aleksandar Kupusinac, Rade Doroslovački, Edita Stokić, Dragan Ivetić
- 403 Real-Time Tracking and Mining of Users' Actions over Social Media
Ejub Kajan, Noura Faci, Zakaria Maamar, Mohamed Sellami, Emir Ugljanin, Hamamache Kheddouci,
Dragan H. Stojanović, Djamel Benslimane
- 427 Land-use classification via ensemble dropout information discriminative extreme
learning machine based on deep convolution feature
Tianle Zhang, Muzhou Hou, Tao Zhou, Zhaode Liu, Weirong Cheng, Yangjin Cheng
- 445 Study of cardiac arrhythmia classification based on convolutional neural network
Yonghui Dai, Bo Xu, Siyu Yan, Jing Xu
- 459 Option predictive clustering trees for multi-target regression
Tomaž Stepišnik, Aljaž Osojnik, Sašo Džeroski, Dragi Kocev
- 487 A Robust Reputation System using Online Reviews
Hyun-Kyo Oh, Jongbin Jung, Sunju Park, Sang-Wook Kim
- 509 Fuzzy Logic and Image Compression Based Energy Efficient Application
Layer Algorithm for Wireless Multimedia Sensor Networks
Arafat SENTURK, Resul KARA, Ibrahim OZCELIK
- 537 Variational Neural Decoder for Abstractive Text Summarization
Huan Zhao, Jie Cao, Mingquan Xu, Jian Lu
- 553 Adaptive E-Business Continuity Management: Evidence from the Financial Sector
Milica Labus, Marijana Despotović-Zrakić, Zorica Bogdanović, Dušan Barać, Snežana Popović
- 581 Human Activities Recognition with a Single Wrist IMU via a Variational Autoencoder
and Android Deep Recurrent Neural Nets
Edwin Valarezo Añazco, Patricio Rivera Lopez, Hyemin Park, Nahyeon Park, Tae-Seong Kim
- 599 Production of linked government datasets using enhanced LIRE architecture
Nataša Veljković, Petar Milić, Leonid Stoimenov, Kristijan Kuk
- 619 Damaged Buildings Recognition of Post-Earthquake High-Resolution Remote Sensing
images based on Feature Space and Decision Tree Optimization
Chao Wang, Xing Qiu, Hui Liu, Dan Li, Kaiguang Zhao, Lili Wang
- 647 Multi-Agent Cooperation Q-Learning Algorithm Based on Constrained Markov Game
Yangyang Ge, Fei Zhu, Wei Huang, Peiyao Zhao, Quan Liu
- 665 A K-means algorithm based on characteristics of density applied to network intrusion detection
Jing Xu, Dezhi Han, Kuan-Ching Li, Hai Jiang
- 689 Cognitive Computation on Consumer's Decision Making of Internet
Financial Products Based on Neural Activity Data
Hongzhi Hu, Yunbing Tang, Yanqiang Xie, Yonghui Dai and Weihui Dai



Computer Science and Information Systems

Published by ComSIS Consortium

Volume 17, Number 2
June 2020

ComSIS is an international journal published by the ComSIS Consortium

ComSIS Consortium:

University of Belgrade:

Faculty of Organizational Science, Belgrade, Serbia
Faculty of Mathematics, Belgrade, Serbia
School of Electrical Engineering, Belgrade, Serbia

Serbian Academy of Science and Art:

Mathematical Institute, Belgrade, Serbia

Union University:

School of Computing, Belgrade, Serbia

University of Novi Sad:

Faculty of Sciences, Novi Sad, Serbia
Faculty of Technical Sciences, Novi Sad, Serbia
Faculty of Economics, Subotica, Serbia
Technical Faculty "Mihajlo Pupin", Zrenjanin, Serbia

University of Montenegro:

Faculty of Economics, Podgorica, Montenegro

EDITORIAL BOARD:

Editor-in-Chief: Mirjana Ivanović, University of Novi Sad

Vice Editor-in-Chief: Ivan Luković, University of Novi Sad

Managing Editor:

Miloš Radovanović, University of Novi Sad

Editorial Assistants:

Vladimir Kurbalija, University of Novi Sad

Jovana Vidaković, University of Novi Sad

Ivan Pribela, University of Novi Sad

Slavica Aleksić, University of Novi Sad

Srdan Škrbić, University of Novi Sad

Miloš Savić, University of Novi Sad

Editorial Board:

C. Badica, *University of Craiova, Romania*

M. Bajec, *University of Ljubljana, Slovenia*

L. Bellatreche, *ISAE-ENSMA, France*

I. Berković, *University of Novi Sad, Serbia*

M. Bohanec, *Jožef Stefan Institute Ljubljana, Slovenia*

D. Bojić, *University of Belgrade, Serbia*

Z. Bosnic, *University of Ljubljana, Slovenia*

S. Bošnjak, *University of Novi Sad, Serbia*

D. Brđanin, *University of Banja Luka, Bosnia and Hercegovina*

Z. Budimac, *University of Novi Sad, Serbia*

C. Chesñevar, *Universidad Nacional del Sur, Bahía Blanca, Argentina*

P. Delias, <https://pavlosdeliasite.wordpress.com>

B. Delibašić, *University of Belgrade, Serbia*

G. Devedžić, *University of Kragujevac, Serbia*

D. Đurić, *University of Belgrade, Serbia*

J. Eder, *Alpen-Adria-Universität Klagenfurt, Austria*

V. Filipović, *University of Belgrade, Serbia*

M. Gušev, *Ss. Cyril and Methodius University Skopje, North Macedonia*

M. Heričko, *University of Maribor, Slovenia*

L. Jain, *University of Canberra, Australia*

D. Janković, *University of Niš, Serbia*

J. Janousek, *Czech Technical University, Czech Republic*

Z. Jovanović, *University of Belgrade, Serbia*

Lj. Kaščelan, *University of Montenegro, Montenegro*

P. Kefalas, *City College, Thessaloniki, Greece*

S-W. Kim, *Hanyang University, Seoul, Korea*

J. Kratica, *Institute of Mathematics SANU, Serbia*

D. Letić, *University of Novi Sad, Serbia*

Y. Manolopoulos, *Aristotle University of Thessaloniki, Greece*

M. Mernik, *University of Maribor, Slovenia*

B. Milašinović, *University of Zagreb, Croatia*

A. Mishev, *Ss. Cyril and Methodius University Skopje, North Macedonia*

N. Mitić, *University of Belgrade, Serbia*

G. Nenadić, *University of Manchester, UK*

N-T. Nguyen, *Wroclaw University of Science and Technology, Poland*

P. Novais, *University of Minho, Portugal*

B. Novikov, *St Petersburg University, Russia*

S. Ossowski, *University Rey Juan Carlos, Madrid, Spain*

M. Paprzycki, *Polish Academy of Sciences, Poland*

P. Peris-Lopez, *University Carlos III of Madrid, Spain*

J. Protić, *University of Belgrade, Serbia*

M. Racković, *University of Novi Sad, Serbia*

B. Radulović, *University of Novi Sad, Serbia*

H. Shen, *Sun Yat-sen University/University of Adelaide, Australia*

J. Sierra, *Universidad Complutense de Madrid, Spain*

M. Stanković, *University of Niš, Serbia*

B. Stantic, *Griffith University, Australia*

L. Šereš, *University of Novi Sad, Serbia*

H. Tian, *Griffith University, Gold Coast, Australia*

N. Tomašev, *Google, London*

G. Trajčevski, *Northwestern University, Illinois, USA*

M. Tuba, *John Naisbitt University, Serbia*

K. Tuyls, *University of Liverpool, UK*

D. Urošević, *Serbian Academy of Science, Serbia*

G. Velinov, *Ss. Cyril and Methodius University Skopje, North Macedonia*

F. Xia, *Dalian University of Technology, China*

K. Zdravkova, *Ss. Cyril and Methodius University Skopje, North Macedonia*

J. Zdravković, *Stockholm University, Sweden*

ComSIS Editorial Office:

University of Novi Sad, Faculty of Sciences,

Department of Mathematics and Informatics

Trg Dositeja Obradovića 4, 21000 Novi Sad, Serbia

Phone: +381 21 458 888; **Fax:** +381 21 6350 458

www.comsis.org; Email: comsis@uns.ac.rs

Volume 17, Number 2, 2020
Novi Sad

Computer Science and Information Systems

ISSN: 1820-0214 (Print) 2406-1018 (Online)

The ComSIS journal is sponsored by:

Ministry of Education, Science and Technological Development of the Republic of Serbia
<http://www.mps.gov.rs/>



Computer Science and Information Systems

AIMS AND SCOPE

Computer Science and Information Systems (ComSIS) is an international refereed journal, published in Serbia. The objective of ComSIS is to communicate important research and development results in the areas of computer science, software engineering, and information systems.

We publish original papers of lasting value covering both theoretical foundations of computer science and commercial, industrial, or educational aspects that provide new insights into design and implementation of software and information systems. In addition to wide-scope regular issues, ComSIS also includes special issues covering specific topics in all areas of computer science and information systems.

ComSIS publishes invited and regular papers in English. Papers that pass a strict reviewing procedure are accepted for publishing. ComSIS is published semiannually.

Indexing Information

ComSIS is covered or selected for coverage in the following:

- Science Citation Index (also known as SciSearch) and Journal Citation Reports / Science Edition by Thomson Reuters, with 2019 two-year impact factor 0.927,
- Computer Science Bibliography, University of Trier (DBLP),
- EMBASE (Elsevier),
- Scopus (Elsevier),
- Summon (Serials Solutions),
- EBSCO bibliographic databases,
- IET bibliographic database Inspec,
- FIZ Karlsruhe bibliographic database io-port,
- Index of Information Systems Journals (Deakin University, Australia),
- Directory of Open Access Journals (DOAJ),
- Google Scholar,
- Journal Bibliometric Report of the Center for Evaluation in Education and Science (CEON/CEES) in cooperation with the National Library of Serbia, for the Serbian Ministry of Education and Science,
- Serbian Citation Index (SCIndeks),
- doiSerbia.

Information for Contributors

The Editors will be pleased to receive contributions from all parts of the world. An electronic version (MS Word or LaTeX), or three hard-copies of the manuscript written in English, intended for publication and prepared as described in "Manuscript Requirements" (which may be downloaded from <http://www.comsis.org>), along with a cover letter containing the corresponding author's details should be sent to official journal e-mail.

Criteria for Acceptance

Criteria for acceptance will be appropriateness to the field of Journal, as described in the Aims and Scope, taking into account the merit of the content and presentation. The number of pages of submitted articles is limited to 20 (using the appropriate Word or LaTeX template).

Manuscripts will be refereed in the manner customary with scientific journals before being accepted for publication.

Copyright and Use Agreement

All authors are requested to sign the "Transfer of Copyright" agreement before the paper may be published. The copyright transfer covers the exclusive rights to reproduce and distribute the paper, including reprints, photographic reproductions, microform, electronic form, or any other reproductions of similar nature and translations. Authors are responsible for obtaining from the copyright holder permission to reproduce the paper or any part of it, for which copyright exists.

Computer Science and Information Systems

Volume 17, Number 2, June 2020

CONTENTS

Editorial

Papers

- 357 What makes a board director better connected? Evidence from graph theory**
Laleh Samarbakhsh, Boža Tasić
- 379 Missing Data Imputation in Cardiometabolic Risk Assessment: A Solution Based on Artificial Neural Networks**
Dunja Vrbaški, Aleksandar Kupusinac, Rade Doroslovački, Edita Stokić, Dragan Ivetić
- 403 Real-Time Tracking and Mining of Users' Actions over Social Media**
Ejub Kajan, Noura Faci, Zakaria Maamar, Mohamed Sellami, Emir Ugljanin, Hamamache Kheddouci, Dragan H. Stojanović, Djamel Benslimane
- 427 Land-use classification via ensemble dropout information discriminative extreme learning machine based on deep convolution feature**
Tianle Zhang, Muzhou Hou, Tao Zhou, Zhaode Liu, Weirong Cheng, Yangjin Cheng
- 445 Study of cardiac arrhythmia classification based on convolutional neural network**
Yonghui Dai, Bo Xu, Siyu Yan, Jing Xu
- 459 Option predictive clustering trees for multi-target regression**
Tomaž Stepišnik, Aljaž Osojnik, Sašo Džeroski, Dragi Kocev
- 487 A Robust Reputation System using Online Reviews**
Hyun-Kyo Oh, Jongbin Jung, Sunju Park, Sang-Wook Kim
- 509 Fuzzy Logic and Image Compression Based Energy Efficient Application Layer Algorithm for Wireless Multimedia Sensor Networks**
Arafat SENTURK, Resul KARA, Ibrahim OZCELIK
- 537 Variational Neural Decoder for Abstractive Text Summarization**
Huan Zhao, Jie Cao, Mingquan Xu, Jian Lu
- 553 Adaptive E-Business Continuity Management: Evidence from the Financial Sector**
Milica Labus, Marijana Despotović-Zrakić, Zorica Bogdanović, Dušan Barać, Snežana Popović
- 581 Human Activities Recognition with a Single Writs IMU via a Variational Autoencoder and Android Deep Recurrent Neural Nets**
Edwin Valarezo Añazco, Patricio Rivera Lopez, Hyemin Park, Nahyeon Park, Tae-Seong Kim

- 599 Production of linked government datasets using enhanced LIRE architecture**
Nataša Veljković, Petar Milić, Leonid Stoimenov, Kristijan Kuk
- 619 Damaged Buildings Recognition of Post-Earthquake High-Resolution Remote Sensing images based on Feature Space and Decision Tree Optimization**
Chao Wang, Xing Qiu, Hui Liu, Dan Li, Kaiguang Zhao, Lili Wang
- 647 Multi-Agent Cooperation Q-Learning Algorithm Based on Constrained Markov Game**
Yangyang Ge, Fei Zhu, Wei Huang, Peiyao Zhao, Quan Liu
- 665 A K-means algorithm based on characteristics of density applied to network intrusion detection**
Jing Xu, Dezhi Han, Kuan-Ching Li, Hai Jiang
- 689 Cognitive Computation on Consumer's Decision Making of Internet Financial Products Based on Neural Activity Data**
Hongzhi Hu, Yunbing Tang, Yanqiang Xie, Yonghui Dai and Weihui Dai

Editorial

Mirjana Ivanović¹ and Miloš Radovanović¹

University of Novi Sad, Faculty of Sciences
Novi Sad, Serbia
{mira,radacha}@dmi.uns.ac.rs

In this second issue of *Computer Science and Information Systems* for 2020, we are happy to announce the impact factors of our journal, updated for 2019: the two-year IF has risen to 0.927, and the five-year IF to 0.775. We would like to thank all our productive authors, whose articles in challenging and exciting areas helped attract additional citations which led to the increased impact of our journal. We hope to continue in the same direction and that the current issue will offer our readers more interesting articles in contemporary and emerging research areas.

This issue consists of 16 regular articles. We are grateful for the hard work and enthusiasm of our authors and reviewers, without which the current issue, as well as the publication of the journal itself, would not have been possible.

In the first article, “What makes a board director better connected? Evidence from graph theory,” Laleh Samarbakhsh and Boža Tasić study the relationships between network centrality measures of board directors (degree, betweenness and closeness) and external variables depicting their age, gender, and other attributes. The study revealed strong correlations between age and connectedness, in addition to offering other interesting insights into the factors influencing centrality and their mutual relationships.

Dunja Vrbaški et al. in their article “Missing Data Imputation in Cardiometabolic Risk Assessment: A Solution Based on Artificial Neural Networks,” explore neural networks as a tool for imputing univariate missing laboratory data during cardiometabolic risk assessment. The proposed neural network approach is compared with other simple baseline methods. Experimental evaluation showed that neural networks outperform other algorithms for a diverse fraction of missing data and mechanisms causing their absence.

“Real-Time Tracking and Mining of Users’ Actions over Social Media,” by Ejub Kajan et al. presents a system called Social Miner that allows companies to make decisions about what, when, and how to respond to users’ actions over social media. The system allows real-time tracking of users’ actions such as raising concerns, commenting, and sharing recommendations and answering various questions, including what actions were frequently executed and why certain actions were executed more than others.

In “Land-Use Classification via Ensemble Dropout Information Discriminative Extreme Learning Machine Based on Deep Convolution Feature,” Tianle Zhang et al. consider combining convolutional neural networks (CNN) and information discriminating extreme learning machines (IELM) in order to produce a dropout-based ensemble learning method to be applied to remote sensor image classification in the land-use domain. The CNN is used to learn effective and robust features, which are then forwarded to the IELM classifier, with satisfactory results.

“Study of Cardiac Arrhythmia Classification Based on Convolutional Neural Network,” authored by Yonghui Dai et al. presents another application of CNNs, in this case studying feature classification of three kinds of electrocardiogram (ECG) signals, includ-

ing sinus rhythm (SR), ventricular tachycardia (VT) and ventricular fibrillation (VF). The article offers two main contributions: insight into the effects of setting different convolution layers of the CNN on arrhythmia data classification, and optimization of classification performance through use of different time periods according to the characteristics of the three kinds of ECG signals.

In “Option Predictive Clustering Trees for Multi-Target Regression,” Tomaž Stepišnik et al. propose to learn option trees (decision trees that allow alternative splits) for multi-target regression (MTR) based on the predictive clustering framework – option predictive clustering trees (OPCTs) – representing a kind of “internal” decision tree ensemble. Experimental evaluation show that OPCTs achieve statistically significantly better predictive performance than a single predictive clustering tree (PCT) and are competitive with bagging and random forests of PCTs.

“A Robust Reputation System using Online Reviews” by Hyun-Kyo Oh et al. addresses the issue of bias in buyer reviews caused by intentional attacks from malicious users and conflation between a buyer’s perception of seller performance and item satisfaction. The proposed approach decouple the measures of seller performance and item quality, while reducing the impact of malignant reviews.

Arafat Senturk et al., in “Fuzzy Logic and Image Compression Based Energy Efficient Application Layer Algorithm for Wireless Multimedia Sensor Networks,” propose an algorithm to minimize energy consumption during image data transmission between sensor nodes in a wireless multimedia sensor network (WMSN) by ensuring the nodes use their most important source. The approach, termed energy-aware application layer algorithm based on image compression (EALAIC), makes use of the top three image compression algorithms for WMSN and decides which one is the most efficient based on three parameters: the distance between the nodes, total node number, and data transmission frequency.

“Variational Neural Decoder for Abstractive Text Summarization,” authored by Huan Zhao et al. proposes a variational neural decoder text summarization model (VND) which steps away from a determined internal transformation structure of recurrent neural networks (RNNs). The model introduces a series of implicit variables by combining variational RNN and variational autoencoder, which is used to capture complex semantic representation at each step of decoding. Experimental results show that the proposed model offers significant improvement over the baseline.

In the article entitled “Adaptive E-Business Continuity Management: Evidence from the Financial Sector,” Milica Labus et al. focus on business continuity management in organizations that use modern e-business technologies: the Internet, mobile computing, e-services, and virtual infrastructure. The authors define a comprehensive framework for the implementation of an adaptive e-BCM adjustable to changes in the business environment, and evaluate it within three financial organizations.

“Human Activities Recognition with a Single Wrist IMU via a Variational Autoencoder and Android Deep Recurrent Neural Nets” by Edwin Valarezo Añazco et al. proposes a human activity recognition (HAR) system that is based on an autoencoder for denoising and recurrent neural network (RNN) for classification with a single inertial measurement unit (IMU) located on a dominant wrist. Experimental results demonstrate significant improvement of accuracy, achieving reliable HAR while relying only on one smart device.

The article “Production of Linked Government Datasets Using Enhanced LIRE Architecture,” by Nataša Veljković et al. describes the enhanced LInked RElations (LIRE) architecture for creating relations between datasets available on open government portals. The architecture is improved to be applicable on different open government data platforms using minimal configuration at the data processing layer. Besides evaluating and describing the advantages and disadvantages of the enhanced LIRE system, the article also introduces a LINDAT indicator that reflects the percentage of linked data in the total possible number of linked data on open government data portals.

Chao Wang et al. in their article “Damaged Buildings Recognition of Post-Earthquake High-Resolution Remote Sensing images based on Feature Space and Decision Tree Optimization,” propose a method for solving the task described in the title by a combination of potential building object set extraction (only) from post-earthquake data, adaptive decision tree number extraction, selection of spectrum, texture and geometric morphology features, and classification based on the optimized random forest (RF) model.

In “Multi-Agent Cooperation Q-Learning Algorithm Based on Constrained Markov Game,” Yangyang Ge et al. tackle the problem of applying reinforcement learning to agents in a multi-agent environment, where agents may fall into unsafe states where it can experience difficulty in bypassing obstacles, receiving information from other agents, etc. This is achieved by adding safety constraints are added to the set of actions, and each agent, when interacting with the environment to search for optimal values, should be restricted by the safety rules, so as to obtain an optimal policy that satisfies the security requirements. A new solution is introduced for calculating the global optimum state-action function that satisfies the safety constraints.

The article “A K-means Algorithm Based on Characteristics of Density Applied to Network Intrusion Detection,” by Jing Xu et al. addresses the sensitivity of the K-means clustering algorithm to the initial cluster selection by using density to choose the initial cluster seed. This is achieved through the use of the Kd-tree index structure, coupled with improved Kd-tree nearest neighbor search to prune the search space and optimize the operation for speed.

Finally, “Cognitive Computation on Consumer’s Decision Making of Internet Financial Products Based on Neural Activity Data,” by Hongzhi Hu et al. addresses the inherently uncertain nature of Internet financial products (possibility of capital loss and liquidity restrictions, and profit as well), by considering consumer’s cognition in the decision making of Internet financial products, conducting an EEG-fNIRS experiment, and proposing an effective cognitive computation method based on neural activity data through the BP-GA algorithm.

What makes a board director better connected? Evidence from graph theory

Laleh Samarbakhsh¹ and Boža Tasić²

¹ Ted Rogers School of Business Management
Ryerson University
350 Victoria Street
Toronto, ON
Canada M5B 2K3
lsamarbakhsh@ryerson.ca

² Ted Rogers School of Business Management
Ryerson University
350 Victoria Street
Toronto, ON
Canada M5B 2K3
btasic@ryerson.ca

Abstract. We are interested in quantifying and uncovering the relationships that form between the board directors of companies. Using these relationships we compute three network centrality measures for each director in the network and employ them in the analysis of connectedness of directors. Our focus in this study is on the attributes that make a board member better connected. The biological, educational and experiential attributes are used as independent variables to develop a regression model measuring the impact on the three connectivity measures (degree, betweenness and closeness). Our results show that “Age” has a direct significant impact on all connectedness measures of a board member. We also find that female directors have a higher measure of degree centrality and betweenness centrality, but lower closeness. The number of foreign degrees increases the degree centrality and betweenness centrality but not closeness. The three identified characteristics of “Age”, “Gender”, and “Education” are supporting the idea that a high level of social connection can in part be expected by the characteristics of individual board members and can explain up to 25% of the board member’s connectivity.

Keywords: Board of director networks, Centrality Measures

1. Introduction

Abstractly speaking, a network is a set whose elements are linked in some way. Depending on the nature of elements and links we can often model social, economic, and political activities and phenomena using networks. Various networks have been objects of study of many researchers in the last half a century. From early studies in sociology to the latest economic analysis, it is evident that networks play an important role in our understanding of behaviour, influence, performance and information exchange of subjects or entities. The interest in networks is even more pronounced these days because of the tremendous amounts of data that are becoming available to researchers, and powerful computers that are being employed in the analysis of the data.

In business and finance research, especially in regards to corporate finance and corporate governance, we have seen a growing interest in the analysis of networks between company directors or between companies (Cohen et al., [6], El-Khatib et al. [8], Fogel et al. [13], Larcker and Wang [18], Larcker and Tayan [17], Renneboog and Zhao [21]). Networks of directors or companies can be equipped with numerous links. The most common linkage for directors is derived from their relationship through the same company (Fogel et al. [13], Larcker and Wang [18], Larcker and Tayan [17], Renneboog and Zhao [21]). Two directors will be connected in the network if and only if they sit on the board of the same company. Using affiliation as a starting relationship between directors one can see how interlocking directors³ will immediately become links between companies. Beside the affiliation ties one may consider several other links, for example, two directors in a network are linked if and only if they have mutual alma mater, regional background or belong to the same social circle. These links are then used to measure director's or company's level of connectedness in the network.

To evaluate the level and specific notions of connectedness we use centrality measures. Network theorists have defined several distinct centrality measures that are correlated (Li et al. [19], Valente et al. [27].) An interlocking director will certainly be connected to more directors in the network than a director sitting on only one board. The number of connections of a director is measured by degree centrality. A director that is in the proximity of other directors can instantly communicate and exchange information with them. This type of connectedness is measured by closeness centrality. Lastly, a director that lies on the shortest path between two other directors will have power to control the flow of information or resources. This type of connectedness is measured by betweenness centrality. From the corporate governance point of view, a Director's Network plays an important role in the identification of relationships among board members and the overall effectiveness of the board.

Besides their regular responsibilities, directors on corporate boards are also expected to counsel and guide CEOs on major corporate strategic decisions. Directors' networks including educational, social, and professional relationships can have a positive impact on a director's advisory role. This is particularly important as firms with well-connected directors are known to outperform their peers and earn superior risk-adjusted stock returns as observed by Larcker, So and Wang [18]. We believe that with the increased impact of social media and accessibility of social networks, the flow of information can assist boards of directors in performing more effectively when their directors are well connected.

The study of connectedness through directors network seems to have significant implications for our understanding of the corporate governance as one of the most important areas of finance. This is particularly important at the present time as firms see broad demand and pressure from the public on understanding board members roles and responsibilities and the need for directors to be well connected to ensure global information exchange. In this regard, increasing the level of board's connectedness through the directors network is considered to be of great importance, and we would like to shed further light on this aspect through expected data-supported findings.

The paper is structured in the following manner. We start with the comprehensive review of the related literature in Section 2. Our main research question, "What attributes can be assigned to a director based on their centrality measures in the directors network",

³ Directors that are members of more than one company board

and hypotheses are introduced in Section 3. Data used in the study and sample selections are described in Section 4. Details about modeling of networks using graph theory and centrality measures' statistics are discussed in Section 5. Finally, our empirical findings and conclusions are presented in Sections 6 and 7.

2. Literature Background

Various networks have been studied in the last few decades using mathematical tools. One of the most frequently used tools in network analysis is certainly graph theory (Proctor and Loomis [20], Sabidussi [22], Freeman [11], Borgatti and Everett [5], Schoch and Brandes [23]). It is now evident that a position of an individual or an entity in a network matters. Simply, individuals or entities that are at the center of the network, as opposed to the periphery, have more access to information and greater power to control the flow of information. Researchers have been predominantly using four network centrality measures (degree, closeness, betweenness and eigenvector centrality) to capture the position of an individual or an entity in a network. Although the four centrality measures are different in nature they are all functions of "nodal statistics" (a vector whose coordinates describe position of a node in a network) as noted recently by Bloch, Jackson and Tebaldi in [3]. They further show that all four centrality measures can be characterized by three axioms (monotonicity, symmetry and additivity), and argue that the centrality measures differ according to which nodal statistics they use, not the way in which the measure processes that information.

Taking advantage of the accessible data from various aspects of day to day life, researchers are trying to understand and explain behaviour of subjects in different social networks.

Fracassi and Tate [10] study the impact that CEOs may have on appointments of directors. They found that firms who have powerful CEOs are likely to appoint directors who have ties to the CEO. Furthermore, they suggest that network ties with the CEO lessens the efficacy of board governance and consumes corporate value.

Renneboog and Zhao investigated the relationship between director networks and takeovers [21]. They found that companies that are better connected are more effective bidders, and the presence of interlocking directorate between the bidder and target impacts the negotiations. They have also found that probability of a successful takeover transaction is higher and the negotiation time is shorter when two firms are directly connected through their directors.

El-Khatib, Fogel and Jandik evaluated the impact of the CEOs network centrality measures to merger and acquisition outcomes in [8]. Their analysis shows that the CEOs with high network centrality easily access and control private information, influence other parties in the network and use these advantages for more frequent acquisition decisions. In spite of all these advantages, the merger and acquisition deals initiated by the CEOs with high network centrality are more likely to incur higher value losses to both the acquirer and the merged entity than the deals initiated by the low network centrality.

Cohen, Frazzini and Malloy investigated the connections between sell-side analysts and management of publicly listed companies [6]. Their focus is on connections through alumni networks. They argue that analysts with educational ties to the senior management

of companies benefit from exclusive information available to them and exceed on their stock recommendations.

Hwang and Kim add a social ties dimension to the conventional independence of boards [14] and find that social ties do matter. Directors are considered independent if they have no financial and familial ties to the CEO or to the company. They look into alumni links, same regional background, military service, academic discipline and type of industry as informal ties between the directors and CEOs. While 87% of boards from the dataset analyzed are categorized as conventionally independent, only 62% remain independent when a social component is added. They also argue that companies whose boards are both conventionally and socially independent compensate their directors at a lower rate and are more sensitive to performance-related pay.

Engelberg, Gao and Parsons consider CEOs personal connections (“rolodex”⁴) with eminent individuals (executives, directors, senior management) of other companies in [7]. They found that CEOs with a hefty rolodex earn more than those with a narrow circle of personal connections. Moreover, they computed that on average an additional file on the CEOs rolodex is worth at least \$17000.

Kramarz and Thesmar study in [15] impact of social networks on board composition and governance of French corporations listed on Paris stock exchange between 1992 and 2003. They found a powerful relationship between the CEO’s social network and the one of directors sitting on the board. They also argued that the governance of these companies is affected in a bad way by these social networks.

While most research related to social networks of directors focuses on the negative side of the inter-board-connections, Larcker and Tayan in [17] stress that these connections can contribute to the value of a company and its shareholders. Larcker, So, and Wang consider the network of boards of directors of U.S. corporations and define board’s centrality as a measure of each board’s well-connectedness in [18]. A position of a firm in the network is determined by using four standard centrality measures and the derived score from these measures called “N-score”. By ranking the firms according to their board’s centrality, they find that the firms whose boards are best-connected on average bring in significantly higher future surplus and stock price returns than the least-connected companies. They also observed that their findings are most evident among newly formed companies with a high growth potential. Overall, their conclusion is that the networks of boards of directors do have an “important and positive impact on the economic performance of a firm.”

The relationship between specific attributes of directors on the boards of companies and companies performance, board independence, decision making, social responsibility etc., has been a theme of many research papers across various academic disciplines. Our focus in this study are “Biological attributes” such as age and gender; “Educational background” such as total number of degrees, number of degrees obtain in North America and number of foreign degrees; and “Corporate experience background” like number of years on boards of either quoted and/or private companies.

Sharma has looked into gender, age, ethnic and cultural diversity of corporate boards of U.S firms from 2000-2006 and its impact on innovation in [24]. While these diversity attributes can lead to conflicts between decision makers in general, they also can contribute to higher level of innovation.

⁴ A Rolodex is a rotating file used to store business cards

Studying gender diversity in the corporate world has attracted attention of many researchers lately. There have been research showing how corporate world is male dominated and that female directors are not only under-represented but also often underpaid. Fortin, Bell and Bohm look into these observations in [9]. Based on the data from Canada, Sweden and United Kingdom, they found that women are in fact under-represented among top earners in these countries, and that explains a considerable part of a difference in gender pay. Lalanne and Seabright investigate the impact of the social network of directors on their salaries [16]. They found that in case of male directors the size of their networks (knowing influential individuals) is positively correlated to their earnings while that is not true for female directors. Their findings also apply to non-salaried compensation.

An interesting real world example of intervention to fix gender gap in corporate governance is the case of Norway. In 2003 the government passed the law requiring Norwegian firms to have a minimum of 40% women directors on their boards. At that time only 9% of directors were women. This particular law that regulated diversification of the boards of companies has been a topic of research in previous literature [28] and [1]. While Ahern and Dittmar [1] focus mostly on the negative impact of government's intervention on firms performance, Wang and Kelan argue [28] that mandatory gender quota has had a positive impact on the appointment of the female board chairs and CEOs.

3. Research Questions and Hypotheses

Our comprehensive review of the related literature shows that most studies that employ centrality measures, focus on one of the following scenarios:

- I) The impact of centrality measures on the CEO's (or director's) decision making and governance skills. There is evidence for stronger centrality measures resulting in higher power to appoint new directors from their own networks [10], influencing merger and acquisitions decisions [8, 21], sharing information with analysts [6], receiving higher compensation [7, 14, 16],
- II) The impact of centrality measures on individual company's performance and its corporate governance with no consensus: evidence for positive effects [17, 18], evidence for negative effects [15], the corporation's position in the corporate networks [18], impact on innovation [24],
- III) The impact of various attributes (biological, educational, experiential) of BODs⁵ on company's performance [1], board's diversity [1, 9, 24], company's innovation [24], board's decision making [24], company's hiring practices [14, 15], salaries of CEO's and directors based on gender [16].

To the best of our knowledge no study has been done that looks into the relationship between the centrality measures and attributes of directors. A simple question that we ask is what attributes can be assigned to a director based on their centrality measures in the directors network. Are directors' centrality measures related to certain age, gender, education, foreign descent, previous service on many boards? Do directors with high centrality come from certain age, specific gender, level of education, foreign education, or have they previously served on significant count of boards?

⁵ Board of Directors

In our study, we look into three centrality measures individually (degree, closeness, and betweenness)⁶ and have tested hypotheses that each centrality measure is related to, on the following attributes: Age, Gender, Number of boards to date (quoted and private), Average number of years on quoted boards, Number of domestic (North American) degrees, Number of foreign degrees. The definitions of these attributes can be found in Table 3 on page 364.

Our goal is to determine what drives the connectedness of the board members. What attributes make a board member better connected? For each director-year, we design three measures of connectedness. We investigate directors' network by estimating a panel regression with firm and year fixed effects and account for several attributes of the board members. In light of this observation, the findings of our study will shed light on these questions.

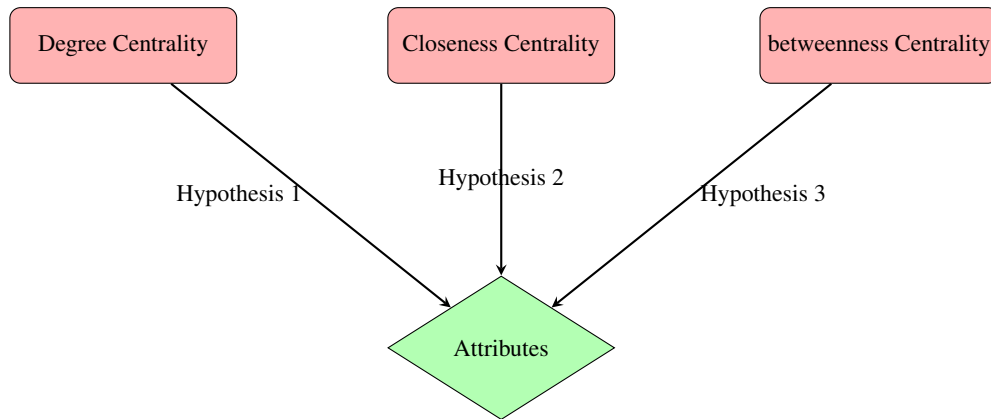


Fig. 1. Hypotheses

Null Hypothesis 1: *The directors with higher degree centrality*

- a) *have identifiable characteristics, and*
- b) *will be an asset for the firm because of their direct access to many sources of communication and information exchange.*

Null Hypothesis 2: *The directors with higher closeness centrality*

- a) *have identifiable characteristics, and*
- b) *will be an asset for the firm because of their proximity to other directors and opportunity to communicate and influence them without having to go through many intermediaries.*

⁶ The three centrality measures were selected as historically the oldest and most commonly used centrality measures in the analysis of networks ([4], [12]). They also have very established interpretation in the directors network context

Null Hypothesis 3: *The directors with higher betweenness centrality*

- a) *have identifiable characteristics, and*
- b) *will be an asset for the firm because they lie on numerous shortest paths between other directors in the network and often are vital brokers of information and resources exchange.*

Our hypotheses are formulated in a compact way to include several attributes rather than listing them individually. For example, by testing age and degree centrality we expect to find that older directors will know more people in the network; female directors will have greater degree centrality, and directors with more education degrees will have greater degree centrality, as well as other measures.

4. Data Description and Sample Statistics

In order to conduct the empirical analysis, we used the data from the North American (United States and Canada) market for the period between 2005 and 2015. Our sample has been compiled using the BoardEx database⁷ provided by Management Diagnostics Ltd. This database collects biographical information on corporate directors and top executives of publicly listed firms and large private firms in North America and around the world. It includes more than 400,000 individuals and over 14,500 companies.

Using the BoardEX database for North America we construct the graphs that describe the networks between directors on monthly basis and we compute metrics (centrality measures) on these graphs. The graphs are constructed and the measures are computed using Mathematica. Using the “Director Networks” files, that have information on connections between directors who were on the boards of the same company during a period of time, we have included the entries from Table 1 using the filters described in Table 2.

We considered only the new connections between directors that started in January 2005 or later. The connections that started prior to January 2005 and were still existing after January 2005 are not included in our analysis. We found that the total of 52352 directors were on the boards of 8270 companies during the decade 2005-2015. Out of 52352 directors approximately about 84% were NEDs and 16% were EDs.

Additionally, we have also included the information from the “Director Profile” files related to biological, educational and experiential attributes given in Table 3. The common characteristics such as age, gender, and nationality are defined in this table. In addition, the two variables “Quoted” and “Private” count the number of public and private company boards the director is sitting on respectively.⁸

Our data shows that boards of directors of public companies consist mostly of male directors. Only 12% of all directors are female. The information on nationality of directors is unknown for approximately 46.5% of all directors considered in our study. Out of remaining 53.5% of directors with known nationality, majority of them are American, almost 46.5%. Following are the directors with Canadian and British nationality who together comprise about 5%. The remaining 1.89% of directors are of 89 different nationalities.

⁷ <http://corp.boardex.com>

⁸ It is noteworthy to mention that total sum of Quoted and Private can be a measure of “Busyness” of each director.

Table 1. Director Profiles Characteristics

Characteristic	Definition
DirectorID*	A unique identifier assigned to each Director
Linked DirectorID*	Director ID of the individual linked to the starting Director
Connected CompanyID*	Company ID of the company through which the starting individual and the linked individual are connected
Connected Company Type	Type ^a of company through which the starting individual and the linked individual are connected
Date of overlap	Starting and ending dates through which the two individuals overlap at the given company or organization
Overlapping Person's Role Title	Role or title of the individual connected to the starting individual at the time of the overlap
Individual's Role Title	Role or title of the starting individual at the time of the overlap
ED/NED/SM	Executive Director, Non-Executive Director, or Senior Management Indicator

^a Quoted - Publicly listed company, Private - Private company, University, Club etc.

Table 2. Filters applied to the Director Profiles Characteristics

Connected Company Type	Quoted
Date of overlap	January 2005 - December 2015
Overlapping Person's Role Title	ED, NED
Individual's Role Title	ED, NED

Table 3. Director Profiles Characteristics

Characteristic	Definition
Age	Director's Age either calculated from DOB ^a or known from disclosure
Gender	The Director's gender
Nationality	The Director's nationality
Education Country	The country in which director obtained education degree
Number of Degrees	Total number of director's education degrees
Number of Foreign Degrees	Total number of degrees obtained in the countries outside of North America
Number of Quoted Boards	Number of boards of publicly listed companies the director has sat on over his career
Number of Private Boards	Number of boards of private companies the director has sat on over his career

^a Date of birth

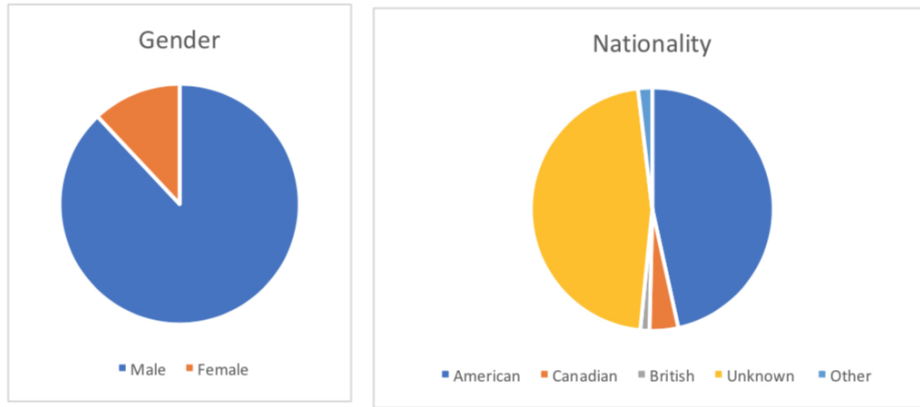


Fig. 2. Pie Charts of Director Attributes

5. Methodology

We consider the network of directors a graph. The individual directors are the vertices of the graphs and the edges represent relationships of connectivity. We then apply graph theory to detect, using a quantitative method, some of the characteristics of the board of directors' connectedness.

The directors attributes (age, gender, nationality, education, number of quoted boards to date) will be exogenous variables to the model of centrality measures. In order to measure the strength of the relationship, we use the three centrality measures (degree, closeness and betweenness) as the observed dependent variables. In other words, we expect to find that the following relationships hold

$$\text{Degree Centrality} = f_1(\text{Attributes})$$

$$\text{Closeness Centrality} = f_2(\text{Attributes})$$

$$\text{betweenness Centrality} = f_3(\text{Attributes})$$

5.1. Modeling Director Network Using Graphs

A graph is an ordered pair $G = (V(G), E(G))$, where $V(G)$ is a nonempty set of vertices and $E(G)$ is the set of edges. Vertices are elements of a nonempty set. A set can be finite or infinite. Edges are links between 2 vertices and can be oriented if needed.

Given a set of directors D_1, D_2, \dots, D_n from a network we define the set of vertices to be $V = \{D_1, D_2, \dots, D_n\}$. Edges represent relationships that may exist between directors. We studied the edges defined by the following relation

1. there is an edge between two directors if and only if they were on a board of a company at the same time,

but various other types of edges can be defined on the set of directors. For example,

2. there is an edge between two directors if and only if they attended the same university,
3. there is an edge between two directors if and only if they share the same regional background.

For each month between January 2005 and December 2015 we constructed a graph whose vertices are directors (only NED and ED), and there is an edge between two directors if and only if they were on a board of a company during a particular month period. We used three centrality metrics (degree, betweenness, and closeness) to explain how director's connections are related to their professional and personal characteristics.

5.2. Director and Edge Count

The chart in Figure 3 presents the average number of directors and edges per year from 2005-2015. As mentioned earlier, this only includes directors who started serving on boards of companies from January 2005. We see that 22684 directors (identified by their unique IDs) joined boards of various companies during the year 2005. This number grew over a decade to 50392. The growth has plateaued after 2008 to approximately 49000 directors per year. In spite of this, the number of edges (ED and NED connections through boards of companies) has continued to grow steadily. In 2005 there were 39213 edges, but by the end of 2015 this number almost quintupled to 191911 edges. Considering that the number of directors became more or less constant after 2008, this suggests a very significant increase of interlocking directors.

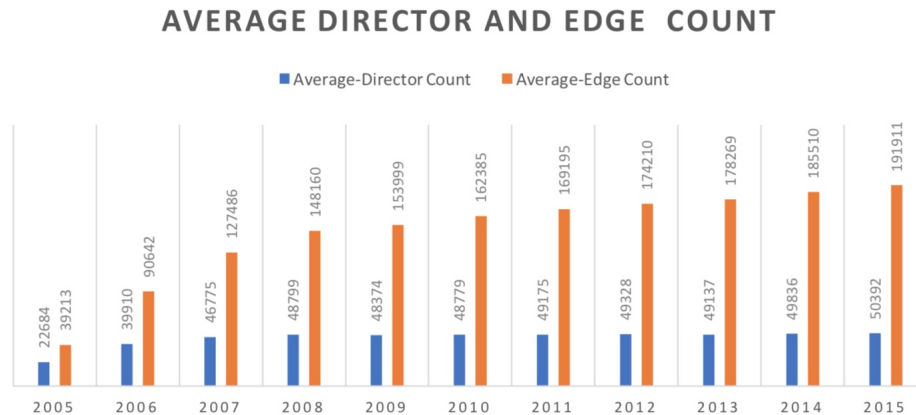


Fig. 3. Average Director and Edge Count

The average number of directors per year indicates the approximate size of an average monthly graph. These monthly graphs were not connected. A usual monthly scenario would be a few big connected subgraphs and many small isolated connected subgraphs that represent companies with no interlocking directorate. The centrality measures were computed on the connected components.

5.3. Centrality Measures

Centrality as a concept was first introduced in social network analysis as a tool to identify influential individuals in the network and to explain how networks of individuals or organizations behave. The most commonly used centrality measures include degree, closeness, betweenness, eigenvector, clustering coefficient, page rank centrality etc. As previously mentioned, we considered only the first three measures. In this section we first define these measures using graph theory language and afterwards we interpret their meaning in relation to the networks of directors constructed using BoardEx database.

Degree Centrality An important vertex in a graph is involved in large number of interactions. The degree of a vertex in a graph measures the amount of direct links between that vertex and other vertices.

Let $G = (V(G), E(G))$ be a graph with n vertices and let $x \in V(G)$ be a vertex. Degree centrality of a vertex x , denoted by $d_G(x)$, is defined by

$$d_G(x) = \text{the number of edges connecting } x \text{ with other vertices.}$$

In order to compare networks of different sizes we normalize the degree centrality of a vertex. The normalised degree centrality of a vertex x , denoted by $Nd_G(x)$, is defined by

$$Nd_G(x) = \frac{d_G(x)}{n - 1}$$

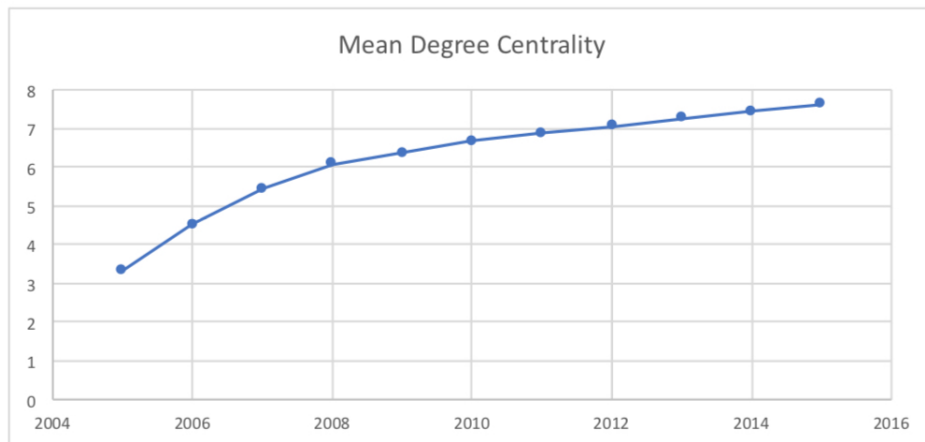
($n - 1$ is the maximum possible degree of a vertex in an undirected graph). It can be expressed either as a proportion or percentage.

In relation to the graphs constructed using BoardEx database, where vertices are directors, the higher degree of a vertex (director) implies better connectedness of that director in the network (graph). A director with high degree will have immediate access to more information, resources, and communication channels through his direct connections. As we can see from the summary Table 4 of degree centrality statistics, the minimal degree was 1. This may look like an anomaly on the first sight, implying that we have a company whose board has only 2 directors. The reason for this is the fact that some North American directors also serve on the boards of international companies and not all details of those boards were accessible to us. For example, if North American directors A, B in a given month were only sitting on a board of a British company, their connection will be listed in BoardEx data but there will be no additional information of the other board members of the British company unless they were also North American. This allows for a minimum degree to be 1. We were not aware of this fact until recently.

Maximum degrees range between 44 and 110 approximately. High degree like these implies that a director was sitting on several large boards of companies. Mean degree is more realistic number of direct connections of a typical director. The column Mean DC of Table 4 shows a 125% increase of the number of direct connections of a typical director over the period 2005-2015. The polygon line given in Figure 4 indicates that although the mean degree grew continuously it did not go above 8. This suggests that an average director was likely sitting on the board of only one company over the decade 2005-2015.

Table 4. Degree Centrality Statistics

Degree Centrality: DC	Min DC	Max DC	Mean DC
2005	1.00	43.75	3.33
2006	1.00	87.25	4.51
2007	1.00	104.00	5.45
2008	1.00	110.00	6.07
2009	1.00	106.64	6.37
2010	1.00	109.27	6.66
2011	1.00	104.36	6.88
2012	1.00	98.27	7.06
2013	1.00	90.91	7.26
2014	1.00	89.82	7.44
2015	1.00	85.73	7.62

**Fig. 4.** Polygon line of Mean Degree Centrality

Closeness Centrality Closeness measures proximity of a vertex in a graph to other vertices. Let $G = (V(G), E(G))$ be a connected graph with n vertices $V = \{x_1, x_2, \dots, x_n\}$ and let $x_i \in V(G)$ be a vertex. Geodesic distance (also known as the shortest path) between vertices x_i and x_j is denoted by $d(x_i, x_j)$. Closeness centrality of a vertex x_i , denoted by $C_G(x_i)$, is defined by

$$C_G(x_i) = \frac{n-1}{\sum_{j=1}^n d(x_i, x_j)}$$

$C_G(x_i)$ is the inverse of the average of all shortest paths between the vertex x_i and any other vertex x_j .

Closeness of a vertex is always between 0 and 1. A vertex will have closeness centrality 1 if it has direct access (the path of length 1) to all other vertices. The smaller the average of all shortest paths between the vertex x_i and any other vertex x_j , the larger the

closeness centrality of the vertex x_i . In relation to the graphs constructed using BoardEx database, where vertices are directors, the larger closeness centrality means that a director can quickly reach other directors in a network without having to go through several intermediaries. Closeness is often referred to as a measure of influence rather than information flow.

Table 5. Closeness Centrality Statistics

Closeness Centrality: CC	Min CC	Max CC	Mean CC
2005	0.04	1.00	0.31
2006	0.06	1.00	0.26
2007	0.07	1.00	0.25
2008	0.07	1.00	0.26
2009	0.06	1.00	0.26
2010	0.06	1.00	0.26
2011	0.06	1.00	0.26
2012	0.08	1.00	0.26
2013	0.07	1.00	0.26
2014	0.07	1.00	0.26
2015	0.06	1.00	0.26

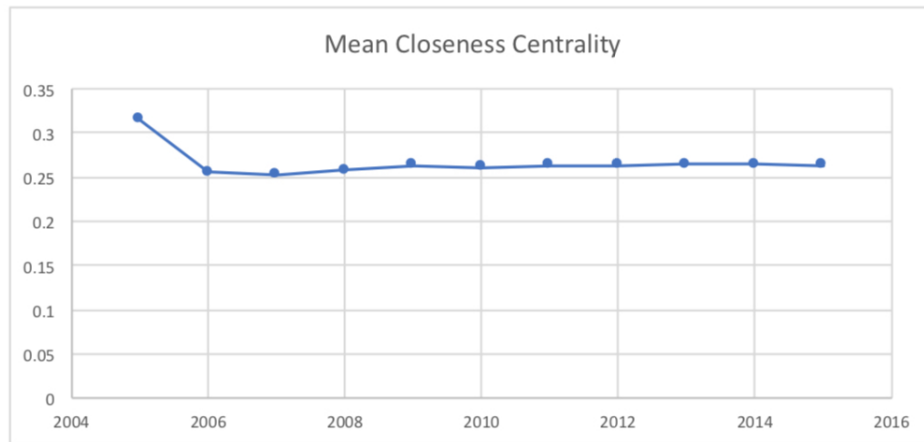


Fig. 5. Polygon line of Mean Closeness Centrality

Table 5 shows the closeness centrality statistics. We see that the closeness centrality ranges between 0.04 and 1. The maximum closeness of 1 is attributed to the above-mentioned directors with degree of 1, or to directors of an isolated company with no interlocking directorate. We can ignore these because no director who is a part of a bigger

subnetwork (connected subgraph) will have centrality score of 1. The minimum closeness is between 0.04 and 0.08 approximately. These would be the scores of the most remote directors in a bigger subnetwork. The polygon line of the mean closeness, given in Figure 5, was constantly around 0.26 except for the year 2005. This implies that during the decade 2005-2015, the average of all shortest paths between a typical director and any other director in the connected subnetwork was around 4.

Betweenness Centrality A vertex that is often in the shortest path (geodesic) between two other vertices has a central role in the network. Such a vertex is essential in the information transfer between other vertices. Let $G = (V(G), E(G))$ be a connected graph with n vertices $V = \{x_1, x_2, \dots, x_n\}$ and let $x_i \in V(G)$ be a vertex. $\sigma(x_i, x_j)$ denotes the total number of shortest paths from the vertex x_i to the vertex x_j and $\sigma_k(x_i, x_j)$ is the number of those paths that pass through the vertex x_k . Betweenness centrality of a vertex x_k , denoted by $B_G(x_k)$, is defined by

$$B_G(x_k) = \sum_{i \neq j \neq k} \frac{\sigma_k(x_i, x_j)}{\sigma(x_i, x_j)}$$

$B_G(x_k)$ is the sum of the shortest paths between all vertices x_i and x_j that pass through the vertex x_k , scaled by the total number of shortest paths between the vertices x_i and x_j . Normalizing the betweenness centrality is done by rescaling $B_G(x_k)$ by $\frac{(n-1)(n-2)}{2}$, which is the largest possible betweenness metric in a connected graph with n vertices. Our betweenness centrality measure is not normalized.

Table 6. betweenness Centrality Statistics

betweenness Centrality: BC	Min BC	Max BC	Mean BC
2005	0.00	4388166.40	46558.70
2006	0.00	9466666.67	90428.00
2007	0.00	8503636.36	103280.84
2008	0.00	9489090.91	103055.71
2009	0.00	8703636.36	99524.35
2010	0.00	9955454.55	100807.40
2011	0.00	9423636.36	100757.29
2012	0.00	9482727.27	101169.78
2013	0.00	7861818.18	100558.75
2014	0.00	8582727.27	101971.89
2015	0.00	8409090.91	103332.10

Table 6 shows the betweenness centrality statistics. Peripheral directors in a network, and directors of boards of isolated companies with no interlocking directorate will have betweenness 0. As the size and the number of edges of the monthly graphs grew, some directors took up a between position on many more geodesics connecting other directors in a network. The maximal values come from directors who were part of subnetworks that had between 4000 and 4500 connected directors. Directors with high betweenness

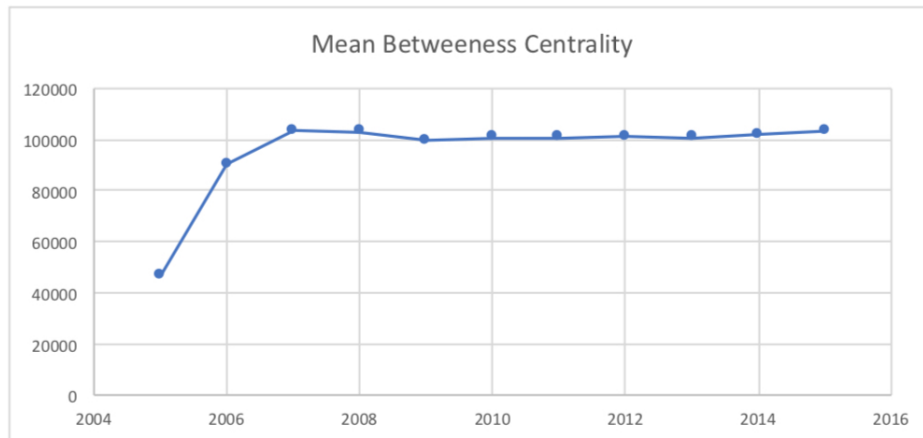


Fig. 6. Polygon line of Mean Betweenness Centrality

are exposed to a lot of information being exchanged among other directors and as such have power to influence other directors and collect information faster. As indicated on the polygon line given in Figure 6, the mean betweenness has plateaued after 2007. This suggests that the number of shortest communication channels between directors, that a typical director was a part of, stayed relatively constant for the most part of the decade 2005-2015.

6. Empirical Results

The calculated measures of connectivity provide an important input into the empirical analysis of our research. In this section we explain the results from the statistical and econometrics findings.

By employing 2,407,752 firm-year-board observations of 52,352 unique board members from 8270 companies in North America from 2005 to 2015, we create a pooled panel dataset that includes the three connectivity measures as the key dependent variables of the study. The sample shows some interesting findings which supports our hypotheses.

We conduct a series of multiple regression analyses that reveal a significantly positive relationship between board connectivity and key characteristics of the board members. The empirical results confirm that such individual characteristics increase the network connectivity of each board member and hence facilitate corporate governance and consequently promote a better connected board. Our results can also be interpreted that a good institutional environment may benefit from the effect of board members' social connectedness facilitated by diverse characteristics of the board members.

6.1. Correlation analysis

Table 7 shows the correlation between the variables under study. Based on the correlation results, we find that Age, Number of degrees, Number of foreign degrees, and gender,

are correlated with our connectivity measures. This supports conducting the regression analysis later. We can then test our hypotheses by running a linear regression on these characteristics against the key dependent variables.

Table 7. Correlation Analysis.

Table shows the correlation between the variables under study. Based on the correlation results, we find that Age, Number of degrees, Number of foreign degrees, and gender, are correlated with our connectivity measures. This supports conducting the regression analysis later.

	DC	CC	BC	Age	Qtd	Pvt	Fem	AvgYr	NumDeg	NumFDeg
DC	1									
CC	-0.1389	1								
BC	0.7522	-0.1328	1							
Age	0.0659	-0.0047	0.071	1						
Quoted	0.4651	-0.1798	0.4683	0.1581	1					
Private	0.195	-0.1205	0.1865	0.0689	0.4629	1				
Female	0.0634	-0.0278	0.0407	-0.0545	-0.0213	-0.0706	1			
Average years on board	-0.1467	0.0059	-0.0621	0.2955	-0.0088	0.0224	-0.0544	1		
Number of degrees	0.053	-0.0567	0.0566	0.0789	0.0709	0.0696	0.0585	-0.0063	1	
Number of Foreign Degrees	0.0901	-0.0457	0.0876	-0.0419	0.0823	0.088	-0.0248	-0.0586	0.2181	1

6.2. Regression analysis

In this section we run three key regressions to test our hypotheses. We test whether directors with higher number of connections have any identifiable characteristics. This as a result will be an asset for the firm because of the facilitation and speed of information exchange in the global market.

For each director-year, we have designed and calculated three measures of connectedness. These three measures act as dependent variables in each regression model. Tables 8, 9, and 10 display the results of the regression analysis. We also control the analytic tests by estimating a panel regression with firm and year fixed effects. We account for several characteristics of the board members.

The results show that “Age” has a direct significant impact on all connectedness measures of the board member. This is intuitive, given the longer the experience of the board member, the more individuals they would know. Our next finding is on role of gender: Female directors have a higher measure of degree centrality and betweenness centrality, but lower closeness.

Expanding on this result, there are implications for directors as well as board composition: nominating committee will be more likely to support that additional female directors will be added to the board. Having a higher degree of betweenness for female directors can also motivate policy makers to encourage inclusion of female directors which will impact board effectiveness as well as more diverse board decisions. See Bernilea, Bhagwal, and Yonker [2] and Sila, Gonzales, and Hagendorff [26].

Finally, the number of degrees and foreign degrees both increase the connectedness degree and betweenness centrality but not closeness. The results support the value of global education and the resulting network of such experience for a board member.

The three identified characteristics of “Age”, “Gender”, and “Education” are supporting the idea that a high level of social connection can in part be expected by the describing characteristics of individual board members. These characteristics can explain up to 25% of the board member’s connectivity (based on the average R-squared values).⁹

Table 8. Degree Centrality Regression.

Regression results showing the relationship between independent variables (described in Table 3) Age, Quoted, Private, Gender, Average Years on Board, Number of Degrees, Number of Foreign Degrees and the dependent variable Degree Centrality

Dependent Variable: DC	Coefficient	Standard Error	t
Age	0.0247	0.0004	69.2300
Quoted	1.1764	0.0017	712.2200
Private	-0.0219	0.0009	-23.9000
Female	1.7203	0.0138	124.4300
Average Years on Board	-0.1535	0.0006	-241.8400
Number of Degrees	0.0234	0.0046	5.0600
Number of Foreign Degrees	0.7132	0.0078	91.6100
Intercept	5.1521	0.0217	237.2500
R-squared	24.10%		

Table 9. Closeness Centrality Regression. Regression results showing the relationship between independent variables (described in Table 3) Age, Quoted, Private, Gender, Average Years on Board, Number of Degrees, Number of Foreign Degrees and the dependent variable Closeness Centrality

Dependent Variable: CC	Coefficient	Standard Error	t
Age	0.0006	0.0000	37.3200
Quoted	-0.0162	0.0001	-232.3000
Private	-0.0023	0.0000	-58.9500
Female	-0.0201	0.0006	-34.4500
Average Years on Board	-0.0006	0.0000	-23.2100
Number of Degrees	-0.0122	0.0002	-62.4200
Number of Foreign Degrees	-0.0106	0.0003	-32.3100
Intercept	0.3134	0.0009	341.4900
R-squared	4.04%		

⁹ In panel data analysis, it is customary to rely more on individual significance and overall significance of the model instead of R^2 or adjusted R^2 . In our panel data due to large number of directors and heterogeneity of cross sections, R^2 is not too high.

Table 10. Betweenness Centrality Regression. Regression results showing the relationship between independent variables (described in Table 3) Age, Quoted, Private, Gender, Average Years on Board, Number of Degrees, Number of Foreign Degrees and the dependent variable Betweenness Centrality

Dependent Variable: BC	Coefficient	Standard Error	<i>t</i>
Age	683.9481	23.33007	29.32
Quoted	78699.05	107.9423	729.08
Private	-3118.249	59.85688	-52.1
Female	79559.59	903.524	88.05
Average Years on Board	-4032.804	41.48484	-97.21
Number of Degrees	5636.232	302.2457	18.65
Number of Foreign Degrees	46183.41	508.7785	90.77
Intercept	-71865.91	1419.166	-50.64
R-squared	22.64%		

Table 11. Log of Betweenness Centrality Regression. Regression results showing the relationship between independent variables (described in Table 3) Age, Quoted, Private, Gender, Average Years on Board, Number of Degrees, Number of Foreign Degrees and the Log of dependent variable Betweenness Centrality. Since the betweenness centrality scores used are inherently large values, this table provides a robustness test for Table 10 results.

Dependent Variable: LogBC	Coefficient	Standard Error	<i>t</i>
Age	0.0112	0.0004	29.0900
Quoted	0.6228	0.0015	425.6900
Private	-0.0081	0.0008	-9.7400
Female	0.6008	0.0134	44.8700
Average Years on Board	-0.0286	0.0007	-41.0200
Number of Degrees	0.2027	0.0046	44.0600
Number of foreign degrees	0.3150	0.0072	43.6500
Intercept	6.2436	0.0233	267.7600
R-squared	17.22%		

7. Conclusion

In the context of the directors network and the impact of the directors attributes (age, gender, education) to the centrality measures we found that “Age” has a direct significant impact on all connectedness measures of a board member. This suggests that older directors are likely to get to know or meet many directors in their long tenure and career lifetime, and use these connections to channel swiftly the information exchange either directly or through liaisons.

Another important finding is the multi-dimensional role of gender: Female directors have a higher measure of degree and betweenness centrality, but lower closeness. This indicates that female directors despite having more direct connections and being on the

shortest paths between many directors in the network, are less influential than male directors. Our study is unique in quantifying such important aspect.

In addition, the number of degrees and foreign degrees increase the degree and betweenness centrality measures but not closeness. The directors who obtained one or more education degrees in North America (or internationally), have potentially met some directors outside of business world. So, they are likely to have more connections in general and to be intermediaries for the information exchange, but they are not necessarily more influential individuals in their current business network.

While this research explores multiple attributes of board members and makes key contributions on the role of connectivity, our future research is directed towards measuring the numerical impact of these attributes on firms' performance. Analyzing firm's return and boards turnover is another interesting expansion that will be further studied.

Acknowledgments. The authors would like to thank the editor and five anonymous referees for their suggestions and helpful comments. Dr. Samarbakhsh would like to thank Royal Bank of Canada for funding support as part of RBC IDI grant.

References

1. Ahern, K. R., Dittmar, A., K., 2012. The changing of the boards: the impact on firm valuation of mandated female board representation. *The Quarterly Journal of Economics* 127, 137–197. doi:10.1093/qje/qjr049.
2. Bernilea, G., Bhagwat, V., Yonker, S., 2018. Board diversity, firm risk, and corporate policies. *Journal of Financial Economics* Vol. 127 Issue 3, 588–612.
3. Bloch, F., Jackson, M. O., Tebaldi, P., 2017. Centrality Measures in Networks. Available at SSRN: <http://dx.doi.org/10.2139/ssrn.2749124>
4. Boldi, P., Vigna S., 2014. Axioms for centrality. *Internet Mathematics* 10, 222–262.
5. Borgatti, P., Everett, G., 2006. A graph-theoretic perspective on centrality. *Social Networks*, 28, 466–484.
6. Cohen, L., Frazzini, A., Malloy C., 2010. Sell-Side School Ties. *The Journal of Finance* 65.4, 1409–1437.
7. Engelberg, J., Gao, P., Parsons, C., 2013. The price of a CEO's rolodex. *Review of Financial Studies* 26, 79–114.
8. El-Khatib, R., Fogel, K., Jandik, T., 2015. CEO network centrality and merger performance. *Journal of Financial Economics* 116, 349–382.
9. Fortin, M. N., Bell, B., Bohm B., 2017. Top earnings inequality and the gender pay gap: Canada, Sweden, and the United Kingdom. *Labour Economics* Vol. 47, 107–123.
10. Fracassi, C., Tate, G., 2012. External Networking and Internal Firm Governance. *The Journal of Finance* Vol LXVII, No 1, 153–194.
11. Freeman, L.C., 1977. A set of measures of centrality based on betweenness. *Sociometry* 40, 35–41.
12. Freeman L., 1979. Centrality in social networks: Conceptual clarification. *Social Networks* 1, 212–239.
13. Fogel, K., Ma. L, and Morck, R., 2014. Powerful Independent Directors. European Corporate Governance Institute (ECGI) - Finance Working Paper No. 404/2014. Available at <http://dx.doi.org/10.2139/ssrn.2377106>
14. Hwang, B., Kim, S., 2011. It pays to have friends. *Journal of Financial Economics* 100, 382–401.

15. Kramarz, F., Thesmar, D., 2013. Social networks in the boardroom. *J. Eur. Econ. Assoc.* 11, 780–807.
16. Lalanne, M., Seabright, P., 2011. The Old Boy Network: Gender Differences in the Impact of Social Networks on Remuneration in Top Executive Jobs. CEPR Discussion Paper No. DP8623. Available at <https://ssrn.com/abstract=1952484>
17. Larcker, D.F., Tayan B., 2010. Director Networks: Good for the Director, Good for Shareholders. Rock Center for Corporate Governance at Stanford University Closer Look Series: Topics, Issues and Controversies in Corporate Governance No. CGRP-08
18. Larcker, D.F., So E.C., Wang C.C., 2013. Boardroom centrality and firm performance. *Journal of Accounting and Economics* 55.2, 225–250.
19. Li, C., Li, Q., Van Mieghem, P. et al., 2015. Correlation between centrality metrics and their application to the opinion model. *Eur. Phys. J. B* 88: 65. <https://doi.org/10.1140/epjb/e2015-50671-y>
20. Proctor, C., Loomis, C., 1951. Analysis of sociometric data. In: Holland, P., Leinhardt, S. (Eds.), *Research Methods in Social Relations*. Dryden Press, New York, NY.
21. Renneboog, L., Zhao, Y., 2014. Director networks and takeovers. *Journal of Corporate Finance* 28, 218–234.
22. Sabidussi, G., 1966. The centrality index of a graph. *Psychometrika* 31, 581–603.
23. Schoch, D., Brandes, U. 2016. Re-conceptualizing centrality in social networks. *European Journal of Applied Mathematics*, 27(6), 971–985. doi:10.1017/S0956792516000401
24. Sharma, Z. 2016. Board Composition and Innovation. *Applied Finance Letters*, Vol 05, Issue 02.
25. Shahgholian, A., Theodoulidis, B., Diaz, D., 2015. Social Network Metrics: The BoardEx case study, <http://ssrn.com/abstract=2613156>
26. Sila, V., Gonzalez, A., Hagendorff, J., 2016. Women on board: Does boardroom gender diversity affect firm risk? *Journal of Corporate Finance* 36, 26–53.
27. Valente T.W. et al., 2008. How Correlated Are Network Centrality Measures? *Connect (Tor)*, 28(1), 16–26. Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2875682/>.
28. Wang, M., Kelan, E., 2013. The Gender Quota and Female Leadership: Effects of the Norwegian Gender Quota on Board Chairs and CEOs. *Journal of Business Ethics* Vol. 117, Issue 3, 449–466.

Dr. Laleh Samarbakhsh is an assistant professor of Finance at the Ted Rogers School of Management at Ryerson University. Her research interests include empirical corporate finance, investments, fixed income and derivatives. Dr. Samarbakhsh has also performed and published research in financial education and pedagogy. She has presented her research in various international academic conferences, including Financial Management Association, European Financial Management Association, Northern Finance Association, Asian Financial Management Association, Eastern Finance Association, Midwest Finance Association, and the Administrative Sciences Association of Canada. She currently coordinates and teaches Managerial Finance at TRSM and has taught courses in Personal Finance, Financial Modeling, and Corporate Finance. She is the recipient of the 2019 Dean's Faculty Award of Excellence at TRSM.

Dr. Boža Tasić is currently an Assistant Professor in Global Management Studies at the Ted Rogers School of Management. He joined Ryerson University in 2006 as an assistant professor in the mathematics department. He taught calculus and linear algebra courses for engineering and science students as well as mathematics for economics. During the

period 2010-2013 he was the coordinator of the Math Centre at Ryerson. One part of his role as the coordinator was developing online teaching tools and systems for assessment and data capturing. He designed interactive basic math modules for all five faculties. Another part of his role was introducing an academic peer support model known as FA-ST (Facilitated study) to the first year engineering math courses. He joined the Ted Rogers School of Management in 2013. He teaches QMS (quantitative methods) courses at the Ted Rogers School of Management. His research interest is in universal algebra and mathematical logic. He has been a member of the Yeats School of Graduate Studies since 2007. He was nominated for TVO's Best Lecturer in 2010.

Received: June 28, 2019; Accepted: November 5, 2019.

Missing Data Imputation in Cardiometabolic Risk Assessment: A Solution Based on Artificial Neural Networks

Dunja Vrbaški¹, Aleksandar Kupusinac^{1*}, Rade Doroslovački¹, Edita Stokić², and Dragan Ivetić¹

¹ University of Novi Sad, Faculty of Technical Sciences
Trg Dositeja Obradovića 6, 21000 Novi Sad, Republic of Serbia
{dunja.vrbaski, sasak, rade.doroslovacki, ivetic}@uns.ac.rs

² University of Novi Sad, Faculty of Medicine
Hajduk Veljkova 3, 21000 Novi Sad, Republic of Serbia
edith@sezampro.rs

Abstract. A common problem when working with medical records is that some measurements are missing. The simplest and the most common solution, especially in machine learning domain, is to exclude records with incomplete data. This approach produces datasets with reduced statistical power and can even lead to biased or erroneous final results. There are, however, many proposed imputing methods for missing data. Although some of them, such as multiple imputation, are mature and well researched, they can be prone to misuse and are not always suitable for building complex frameworks. This paper explores neural networks as a potential tool for imputing univariate missing laboratory data during cardiometabolic risk assessment, comparing it to other simple methods that could be easily set up and used further in building predictive models. We have found that neural networks outperform other algorithms for diverse fraction of missing data and different mechanisms causing their missingness.

Keywords: missing data, cardiometabolic risk, artificial neural networks

1. Introduction

Missing data is a well known and commonly present problem in both research and industry. Many datasets contain information that is incomplete, due to a variety of reasons: data could have been unavailable, recorded incorrectly or not collected at all, damaged or lost. Health data is not an exception. Actually, those datasets are very prone to having different types of missing data according to its structure, volume and relation with observed data. Not dealing with missing data properly could represent a significant problem for further analysis or building predictive models [23] [33] [28] [12] [1] [41].

The subject of this study is to explore whether a machine learning method, such as an artificial neural network (ANN), could serve as an imputation tool for missing data in cardiometabolic risk assessment (CMR), comparing it to several other single imputation statistics methods. More specifically, we are interested in imputation of missing values

* Corresponding Author: Aleksandar Kupusinac (sasak@uns.ac.rs)

that constitute outcome CMR values and which would not be further used in building final predictive models, rather than missing values for predictors or outcomes on their own.

1.1. Cardio-Metabolic Risk Assessment

There is a number of factors that participate in development of cardiovascular diseases and which are related to CMR [37] [32] [16]. In a clinical routine for obese patients, an evaluation of CMR starts with an anamnesis, an estimation of nutritional status and an adipose tissue distribution. These steps are usually performed by simple anthropometric measurements. The following steps include laboratory analyses of lipids and lipoproteins levels, glycemia, insulinemia and other indicators of obesity comorbidities. Since those procedures have a higher level of invasiveness and they induce additional costs, there is an interest in building predictive models for risk scores that rely only upon inexpensive, commonly available and non-laboratory data.

But, in order to build that cost-effective, prognostic machine learning model, such as the one based on simple parameters and neural network algorithm [19], we need this laboratory data, since it is used for computing the outcome (CMR) values in the data pre-processing phase. If all such subjects, for which only some laboratory data is missing, are omitted, then the number of outcomes, and accordingly, the dataset used for training, could be significantly reduced, making learning more challenging and the results potentially erroneous.

1.2. Missing Data and Machine Learning

In statistics, analysis of the missing data itself is an important task. Data can be missing as an univariate or a multivariate, missingness can follow some pattern and an origin of the missingness could be explained by the observed or the missing data itself. We are using standard notion when differentiating origin of missingness, in literature usually called the missingness mechanism [34] [24]:

- When missing data does not depend neither on observed nor on unobserved values (missing completely at random; MCAR).
- When missing data does not depend on the unobserved, but may depend on the observed values (missing at random; MAR).
- When missing data depends on the unobserved values themselves (missing not at random; MNAR).

Based on the analysis of missing data itself, an appropriate statistic method could be performed for dealing with missing data. Removing the data that contains the missing information is the easiest and, very often, a method of choice in machine learning application domain. But, this approach, where only the complete cases are used, could lead to flawed or unreliable results. Even in the case of MCAR data, where deletion produces unbiased results, we could end up reducing analysis power and weakening some of our tests [10]. In the case of MAR or MNAR data, when there is an underlying reason for missingness, and especially when proportion of missing data is larger, we must exercise

caution. In contrast to ignoring data with missing values, one could choose some imputation method if assumptions for the selected method are met [26]. There are many proposed imputation methods. Properly chosen and used method can significantly reduce the impact which missing data has on the research result. That means that an approach to a missing data problem should be careful, in order to make an educated decision whether the data could safely be deleted, or how and why some imputation method is chosen.

On the other hand, in an engineering and machine learning domain, such analysis is usually not performed, mostly because engineers are not very interested in explaining the data, but more in building and validating their models using that data [4] [40]. If an imputation method that is not sensitive to the nature or volume of missing data and that does not require previous analysis could be developed, that would enable automation of the data preprocessing and feature engineering task, which many recognize as one of the holy grails of machine learning [5].

Machine learning methods themselves are good candidates for such a task, since in their essence is to learn complex relations and learn them from the data without additional instructions. Therefore, we choose to explore how one machine learning algorithm handles imputation in different scenarios with missing data in CMR assessment. We have chosen ANN, amongst other possible algorithms, since it is already shown that it can successfully predict CRM from non-laboratory values, and we presume that there will also exist a dependency function between CMR and laboratory data which some ANNs can approximate. Accordingly, in this paper, we have considered different amount with different missingness mechanism of univariate missing laboratory data, hypothesizing that there are ANNs which could successfully deal with it, regardless of the nature and volume of the data that is missing.

There are some previous results that explore comparison between methods for imputation [25] [11] [14] and even study ANNs within missing data problem [3] [22] [21] [29] [36]. However, all of the research which has been pointed out is somewhat different than our goal. These papers either: explore imputation for predictor values; handle missing data to explore and describe the data, observe the estimates such as regression coefficients and standard errors; use multiple imputation; use a large number of predictors and big data to train machine learning models. On the other side, our subject is imputation through ANN that explores imputing values that are used to calculate outcome values, further leading to building predictive, machine learning, models and one that uses small, simple structured data and explore all three missingness mechanisms. Moreover, although ANNs are used to predict CMR, they are not researched as imputation tool in this domain and in this manner.

2. Methodology

Through simulation and analysis, we have compared neural networks with other single imputation methods in univariate missing data for laboratory values through different settings that reflect different missingness scenarios.

Firstly, we have established structures of ANNs that will be used in comparison. Then we have simulated different scenarios of missing data occurrence and compared performance of ANNs with other methods using several measurements.

2.1. Data

Dataset was produced as a result of the study at the Department of Endocrinology, Diabetes and Metabolic Disorders of the Clinical Center of Vojvodina in Novi Sad, Serbia. The inquired group consisted of 2985 individual respondents, 1980 women and 1005 men, aged 18 to 69 years. The study was conducted in accordance with the Declaration of Helsinki and approved by Ethical Committee of the Clinical Center of Vojvodina (No. 0020/649).

Dataset contains the following CMR risk factors:

- non-laboratory: gender (GEN), age (AGE), body mass index (BMI), waist-to-height ratio (WHtR)
- laboratory: triglycerides (TG), total cholesterol (TCH), low-density lipoprotein (LDL), high-density lipoprotein (HDL) and glycemia (GLY).

Descriptive statistics for the data is shown in the Table 1.

Table 1. Descriptive statistics for risk factors in the CMR dataset

	Mean	St.Dev.	Min	Max
AGE	43.413	10.615	18	69
BMI	29.732	6.472	16.600	50.440
WHtR	0.565	0.091	0.338	0.899
LDL	3.762	0.950	2.030	10.140
HDL	1.124	0.262	0.460	2.090
TCH	5.952	1.376	2.770	13.240
TG	2.057	1.819	0.350	27.320
GLY	5.145	1.321	2.800	13.800

Since adipose tissue for men shows a tendency towards central or abdominal accumulation, male gender bears higher potential risk of cardiovascular diseases. With women, the risk increases with aging, due to the adipose tissue centralization. It is known that acceleration of atherosclerosis increases with age and that the cardiometabolic risk increases with age. Beside the gender and genetic predisposition, this is an additional risk factor that cannot be controlled.

BMI is an indication of nutritional state and is used to quantify the level of obesity. Despite of a lot of controversy about its reliability in fat mass prediction, it shows high efficiency in cardiovascular risk prediction. Values of BMI over 25 kg/m^2 correspond to being overweight, and values over 30 kg/m^2 correspond to obesity [27]. BMI is calculated as a ratio of body mass and body height squared. Body weight is measured with a balance beam scale. Body height is measured with Harpenden anthropometer (Holtain Ltd, Croswell, UK) with precision of 0.1 cm .

Waist circumference is correlated with the amount of visceral abdominal adipose tissue, but also with the level of lipids, lipoproteins and insulin and it is a significant predictor of the obesity comorbidity. An index calculated as a ratio of waist circumference and body height (WHtR) has been shown to be a better risk indicator and the values $\text{WHtR} \geq 0.5$

are considered to indicate increased risk [2] [18]. Waist circumference is measured with a measurement tape with precision of 0.1 *cm*. It is measured at half the distance between the lowest point of the costal arch and the highest point of the iliac crest.

Disturbances of lipid and lipoprotein metabolism are present in 30% of obese persons. They are manifested as one or more of the following disruptions: hypercholesterolemia, hypertriglyceridemia, protective HDL-cholesterol level drop off, raised level of LDL-cholesterol and increased fraction of small, dense, atherogenic LDL-particles. In our study, cholesterol and triglycerides levels are determined by the standard enzyme procedure. The values of HDL-cholesterol were determined by precipitation procedure with sodium-phosphor-wolframate. The values of LDL-cholesterol were calculated using Friedewald's formula [8].

Hyperglycemia is also a risk factor for cardiovascular diseases. Increased level of glucose accelerates the process of atherosclerosis by increasing the oxidative stress and protein glycolization [20]. In this research, the glycemia values were determined using Dialab glucose GOD-PAP method. All inquiries were taken during the morning hours (after fasting overnight).

In this research, we have observed missing data imputation for: high-density cholesterol (HDL), low-density cholesterol (LDL), total cholesterol (TCH), triglycerides (TG) and glycemia (GLY) using following cut off values as indication of cardiometabolic risk [7] [13]:

- HDL < 1.29 (woman) and < 1.03 (man)
- LDL \geq 3.3
- TCH \geq 5.2
- TG \geq 1.71
- GLY \geq 6.1

Distribution of HDL, LDL, TCH, TG and GLY is shown in Figure 1 with highlighted CMR threshold values. For HDL, since different values are used for man and woman, plot has two lines annotated with *M* (male) and *F* (female) respectively.

2.2. Structure of the Neural Networks

For each variable of interest, we have tested several networks in order to find their optimal structure. Input vectors for ANNs are all variables from dataset except variable of interest which represent output value. Since the data is simple, and not massive, we have opted for single layered feedforward artificial network. For each variable, ANNs with 2 to 10 neurons in the hidden layer are tested. Data split for training and testing is performed through bootstrap resampling [6]. For each number of hidden neurons, 100 simulations of ANN testing are performed.

Experiments were performed using *R* software v3.5.0 [30] with packages *neuralnet* v1.33 [9] and *caret* v6.0 [17] with following settings for ANN:

- training algorithm: resilient backpropagation
- starting weights: random initialization
- activation function: tanh (tangent hyperbolicus transfer function)
- output neuron: linear function

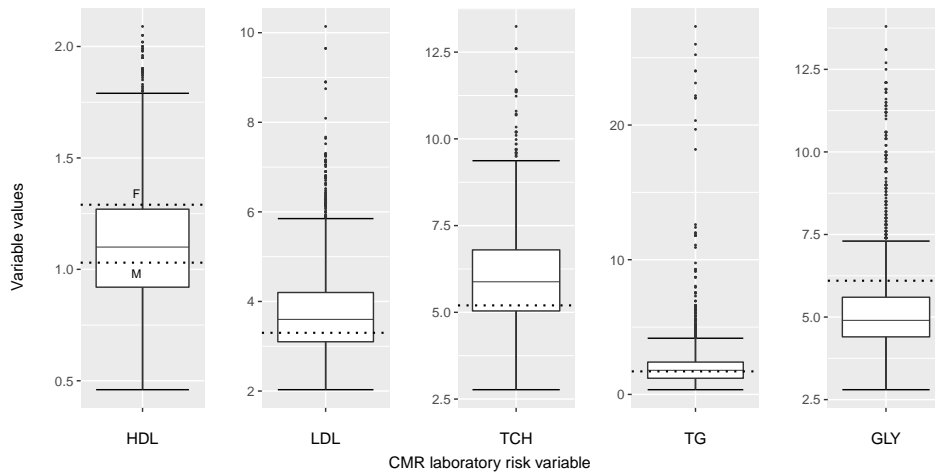


Fig. 1. Boxplots representing distribution of laboratory CMR risk values in the dataset. Cutoff values, used as indication of CMR, are highlighted with dashed lines. For HDL, two distinctive cutoff lines are shown corresponding to different genders

- stopping criteria: threshold for the partial derivatives of the error function or maximum number of steps, whichever is reached before.

Instead of the regular back propagation algorithm, resilient backpropagation (rprop+) is used [15] [31], as a faster method that does not require learning rate as a parameter for its training. Since the study covers a number of simulations, improving speed and reducing grid search for hyperparameter tuning in this phase of work was beneficial.

An optimal number of neurons for each variable was selected according to the number of neurons which, on average, produced the smallest root mean squared error (RMSE). According to this criteria, networks trained in the second part of the research for HDL, LDL, TCH, TG and GLY were networks with single hidden layer and number of neurons: 6, 9, 5, 2 and 2, respectively. Results are shown in Table 2.

2.3. Comparison of Imputation Methods

In the second part of the study, missing data is artificially introduced in the dataset. Accuracy of imputation for neural networks and several other methods is addressed. For each variable, the missing data is introduced by varying percentage (10%, 20%, 30% and 50%) and according to different missingness mechanism (MCAR, MAR and MNAR). For each combination of settings (60 in total), comparison of 5 imputation methods was performed through 100 simulations.

Depending on the selected mechanism, percentage of the original data is removed from the dataset. Data that is missing under MCAR assumption is randomly selected based solely on the percentage of missing data. For both MAR and MNAR mechanism, logistic function for probability that data is missing was used. To simulate MAR mechanism, values for AGE and BMI were used to introduce missingness. Observations with

Table 2. Average **RMSE** errors produced through 100 simulations during testing of neural networks with different number of hidden neurons. Smallest values are highlighted. For each variable minimum and maximum values are displayed

Variable (min - max)	LDL (2.03-10.14)	HDL (0.46 - 2.09)	TCH (2.77 - 13.24)	TG (0.35 - 27.32)	GLY (2.80 - 13.80)
Hidden neurons					
2	0.74019	0.24994	0.93934	1.78078	1.25208
3	0.72843	0.24584	0.92250	1.84684	1.26054
4	0.71176	0.24282	0.91913	1.86890	1.27322
5	0.70994	0.24334	0.91842	1.93557	1.28708
6	0.70764	0.24189	0.92108	1.94169	1.28131
7	0.70391	0.24341	0.92424	1.82648	1.28666
8	0.70358	0.24198	0.92970	1.84911	1.27962
9	0.70343	0.24529	0.92365	1.83521	1.27707
10	0.70500	0.24631	0.92339	1.84555	1.28290

lower AGE and BMI are considered to have bigger probability of missing laboratory data. For MNAR mechanism, the same variable which was investigated was used as a cause for missingness. Observations with lower values, except for HDL, are considered to have greater probability of missing data. For HDL, lower values had lower probability of missing data. Function *ampute* from *R* library *mice* v3.3.0 was used to carry out those deletions (amputations) [38] [39] [35].

For all data sets produced in the described manner, 5 different imputation methods are performed: neural network (NN), predictive mean matching (PMM), stochastic linear regression (SLR), random forest (RF) and mean imputation (MEAN). Neural network training and imputation is performed using the settings and architectures obtained in first set of experiments. Other methods are implemented using *R* package *mice*. SLR and MEAN methods did not require special parameters. For RF, training number of trees was set to 10. For PMM, all variables except the one of interest, are used for finding five possible donors and distance between predicted and drawn values was used as a matching distance.

Comparison between methods was based on imputation performance which is evaluated by: root mean squared error (RMSE), mean absolute percentage error (MAPE) and classification accuracy (CA) between imputed and original values. RMSE measures deviation of this difference and is used as common metric for model comparison. MAPE explain accuracy as percentage error and is given due to its intuitive interpretation. CA measures what portion of imputed values will be correctly classified as risk factor for CRM according to their threshold values. Lower values for RMSE and MAPE denote better performance, and higher CA values indicate better results.

3. Results

Comparison results for imputation of HDL, LDL, TCH, GLY and TG are respectively shown in tables 3, 4, 5, 6 and 7. For each variable, missingness mechanism and percent of missing data average results for RMSE, MAPE and CA are shown where best performance values are highlighted. Also, due to space constraints, graphical illustration of obtained results is given in Appendix (figures: 2, 3, 4, 5, 6).

Table 3. Algorithm comparison for variable **HDL** considering three types of missing data mechanisms (MAR, MCAR, MNAR) and different levels of missingness (10, 20, 30 and 50%) based on imputation accuracy (root mean squared error, mean absolute percentage error and classification accuracy)

HDL	RMSE	MAPE	CA	RMSE	MAPE	CA
<i>MCAR</i>						
	10%			20%		
NN	0.2802	20.6955	0.7316	0.2785	20.5654	0.7202
PMM	0.3608	26.1986	0.6050	0.3617	26.3017	0.5969
SLR	0.3667	27.1279	0.5854	0.3660	27.1464	0.5831
RF	0.2946	18.8455	0.7010	0.2953	19.2711	0.6936
MEAN	0.2610	19.2869	0.6953	0.2612	19.3118	0.6903
	30%			50%		
NN	0.2794	20.7115	0.6382	0.2797	20.6359	0.6398
PMM	0.3617	26.4378	0.5953	0.3628	26.4785	0.5941
SLR	0.3649	27.0686	0.5857	0.3659	27.1615	0.5856
RF	0.3035	20.0461	0.6826	0.3143	21.3192	0.6623
MEAN	0.2617	19.4449	0.6902	0.2615	19.3626	0.6924
<i>MAR</i>						
	10%			20%		
NN	0.2389	16.9590	0.7391	0.2406	17.2297	0.7290
PMM	0.3598	26.2777	0.6071	0.3587	26.2944	0.6063
SLR	0.3694	27.7332	0.5881	0.3646	27.4427	0.5913
RF	0.2933	19.3031	0.7019	0.2982	20.0395	0.6914
MEAN	0.2601	19.7928	0.6932	0.2592	19.8083	0.6914
	30%			50%		
NN	0.2437	17.4571	0.7261	0.2453	17.5117	0.7200
PMM	0.3597	26.2950	0.6077	0.3616	26.3412	0.5996
SLR	0.3651	27.3646	0.5928	0.3663	27.2623	0.5881
RF	0.3083	21.0124	0.6747	0.3168	21.7676	0.6611
MEAN	0.2597	19.8737	0.6903	0.2609	19.6912	0.6901
<i>MNAR</i>						
	10%			20%		
NN	0.2394	23.3216	0.6742	0.2581	25.5564	0.6472
PMM	0.3582	32.0925	0.5974	0.3724	33.7851	0.5734
SLR	0.3635	33.6794	0.5693	0.3772	35.2030	0.5442
RF	0.2856	23.5159	0.6986	0.3080	26.2282	0.6604
MEAN	0.2575	26.5496	0.6231	0.2773	28.7606	0.6227
	30%			50%		
NN	0.2818	28.3264	0.6240	0.2675	23.1169	0.6560
PMM	0.3889	35.7334	0.5439	0.3762	30.7928	0.5674
SLR	0.3928	37.0190	0.5188	0.3783	31.4889	0.5472
RF	0.3305	29.1086	0.6188	0.3238	25.3834	0.6331
MEAN	0.3025	31.5800	0.6235	0.2835	25.5359	0.6444

Table 4. Algorithm comparison for variable **LDL** considering three types of missing data mechanisms (MAR, MCAR, MNAR) and different levels of missingness (10, 20, 30 and 50%) based on imputation accuracy (root mean squared error, mean absolute percentage error and classification accuracy)

LDL	RMSE	MAPE	CA	RMSE	MAPE	CA
<i>MCAR</i>						
	10%			20%		
NN	0.6517	13.5230	0.7495	0.6649	13.6613	0.7549
PMM	0.9998	19.5170	0.6591	1.0105	19.5738	0.6638
SLR	1.0302	22.1732	0.6555	1.0307	22.1106	0.6581
RF	0.8660	14.4301	0.7272	0.8870	15.0236	0.7171
MEAN	0.9451	19.1504	0.6505	0.9451	19.1504	0.6552
	30%			50%		
NN	0.6788	13.8307	0.7500	0.7112	14.3998	0.7425
PMM	1.0119	19.6088	0.6607	1.0171	19.6594	0.6610
SLR	1.0301	22.0974	0.6589	1.0318	22.1538	0.6560
RF	0.9087	15.5668	0.7101	0.9508	16.8752	0.6884
MEAN	0.9495	19.2054	0.6522	0.9521	19.2136	0.6507
<i>MAR</i>						
	10%			20%		
NN	0.6387	13.5412	0.7508	0.6704	14.0882	0.7409
PMM	0.9611	19.2735	0.6520	0.9655	19.2723	0.6549
SLR	1.0151	22.4694	0.6440	1.0253	22.7197	0.6429
RF	0.8646	15.0221	0.7086	0.8904	15.7338	0.6998
MEAN	0.9612	20.7286	0.6022	0.9659	20.9315	0.6040
	30%			50%		
NN	0.6903	14.3724	0.7374	0.7274	14.7458	0.7349
PMM	0.9738	19.4508	0.6519	1.0069	19.6371	0.6562
SLR	1.0262	22.6180	0.6442	1.0343	22.3459	0.6507
RF	0.9206	16.6055	0.6841	0.9551	17.1912	0.6780
MEAN	0.9696	21.1739	0.6037	0.9659	20.1857	0.6341
<i>MNAR</i>						
	10%			20%		
NN	0.6149	15.7534	0.6277	0.6518	16.7199	0.6057
PMM	0.8784	20.9437	0.5788	0.9152	22.0134	0.5657
SLR	0.9903	25.8651	0.5446	1.0144	26.4461	0.5435
RF	0.7268	15.8575	0.6670	0.7731	17.3816	0.6369
MEAN	0.8324	25.4867	0.3621	0.9006	27.7790	0.3618
	30%			50%		
NN	0.6947	18.0946	0.5807	0.7176	16.3713	0.6638
PMM	0.9458	22.9883	0.5546	0.9969	21.4488	0.6240
SLR	1.0457	27.2633	0.5358	1.0464	24.5275	0.6051
RF	0.8347	19.4666	0.6006	0.9259	18.9541	0.6466
MEAN	0.9902	30.8539	0.3622	0.9869	25.4807	0.5285

Table 5. Algorithm comparison for variable **TCH** considering three types of missing data mechanisms (MAR, MCAR, MNAR) and different levels of missingness (10, 20, 30 and 50%) based on imputation accuracy (root mean squared error, mean absolute percentage error and classification accuracy)

TCH	RMSE	MAPE	CA	RMSE	MAPE	CA
<i>MCAR</i>						
	10%			20%		
NN	0.9044	11.2777	0.8465	0.9096	11.3255	0.8447
PMM	1.3849	17.9032	0.7646	1.3858	17.8030	0.7667
SLR	1.4125	19.3143	0.7096	1.4071	19.3211	0.7167
RF	1.1452	13.2419	0.8244	1.1760	13.6968	0.8172
MEAN	1.3765	19.6452	0.7120	1.3768	19.5699	0.7130
	30%			50%		
NN	0.9159	11.3713	0.8477	0.9254	11.5042	0.8448
PMM	1.3831	17.7930	0.7669	1.3798	17.7741	0.7672
SLR	1.4101	19.3680	0.7142	1.4124	19.3726	0.7138
RF	1.2203	14.5303	0.8060	1.2504	15.1500	0.7976
MEAN	1.3772	19.6186	0.7127	1.3788	19.5867	0.7135
<i>MAR</i>						
	10%			20%		
NN	0.8916	11.6872	0.8259	0.9012	11.8159	0.8257
PMM	1.3595	18.3046	0.7359	1.3640	18.2839	0.7381
SLR	1.4035	20.3694	0.6782	1.4078	20.4852	0.6761
RF	1.1616	14.3625	0.7973	1.2017	15.1341	0.7838
MEAN	1.4254	22.5705	0.6222	1.4333	22.8869	0.6230
	30%			50%		
NN	0.9152	11.9812	0.8235	0.9405	11.8581	0.8342
PMM	1.3671	18.3581	0.7371	1.3727	17.9005	0.7556
SLR	1.4080	20.4700	0.6749	1.4083	19.7614	0.6981
RF	1.2422	16.0273	0.7751	1.2683	15.9046	0.7849
MEAN	1.4480	23.3124	0.6225	1.4142	21.3220	0.6764
<i>MNAR</i>						
	10%			20%		
NN	0.8253	13.4357	0.7565	0.8647	14.2280	0.7442
PMM	1.3150	20.6076	0.6800	1.3530	21.4554	0.6689
SLR	1.3586	23.1945	0.6075	1.3939	23.7496	0.6021
RF	1.1121	16.0774	0.7388	1.1956	17.8133	0.7139
MEAN	1.3765	26.6161	0.4466	1.4726	28.6865	0.4482
	30%			50%		
NN	0.9302	15.5389	0.7213	0.9445	13.6014	0.7896
PMM	1.4125	22.4742	0.6621	1.4210	20.3798	0.7268
SLR	1.4522	25.0296	0.5859	1.4406	21.9095	0.6678
RF	1.2946	19.8007	0.6870	1.3073	18.0855	0.7481
MEAN	1.6098	31.7311	0.4478	1.4896	25.4515	0.6143

Table 6. Algorithm comparison for variable **TG** considering three types of missing data mechanisms (MAR, MCAR, MNAR) and different levels of missingness (10, 20, 30 and 50%) based on imputation accuracy (root mean squared error, mean absolute percentage error and classification accuracy)

TG	RMSE	MAPE	CA	RMSE	MAPE	CA
<i>MCAR</i>						
	10%			20%		
NN	1.6781	45.8713	0.6947	1.7419	47.0029	0.6794
PMM	2.4287	63.7588	0.6149	2.4450	63.1532	0.6062
SLR	2.4359	103.3195	0.5578	2.4596	102.7914	0.5577
RF	2.0806	44.0431	0.7074	2.0547	44.7400	0.7023
MEAN	1.7646	57.4358	0.5290	1.8152	57.5840	0.5279
	30%			50%		
NN	1.7228	47.7940	0.6762	1.7834	48.2514	0.6755
PMM	2.3925	63.0147	0.6056	2.4565	64.0044	0.6048
SLR	2.4405	104.0732	0.5564	2.4431	103.1582	0.5599
RF	2.1222	47.7587	0.6877	2.1983	51.0667	0.6650
MEAN	1.7804	57.7058	0.5268	1.8126	57.6760	0.5264
<i>MAR</i>						
	10%			20%		
NN	1.7303	49.2283	0.6911	1.8401	49.4026	0.6982
PMM	2.3680	65.6342	0.6064	2.4370	65.0048	0.6053
SLR	2.4674	109.4875	0.5521	2.5070	109.1327	0.5512
RF	2.0619	46.5159	0.7054	2.2181	49.7677	0.6896
MEAN	1.8133	65.2662	0.4868	1.8857	65.7658	0.4868
	30%			50%		
NN	1.8198	51.1575	0.6783	1.8659	50.0006	0.6637
PMM	2.4042	65.0420	0.6074	2.4305	63.5184	0.6034
SLR	2.4843	108.4079	0.5566	2.4544	102.3056	0.5595
RF	2.2407	53.1501	0.6720	2.2509	52.7012	0.6629
MEAN	1.8649	66.6036	0.4857	1.8821	61.2814	0.5108
<i>MNAR</i>						
	10%			20%		
NN	0.8652	68.9414	0.5467	1.0042	80.4539	0.4532
PMM	1.7352	82.8811	0.5752	1.7816	88.5791	0.5513
SLR	1.9790	145.7663	0.4997	2.0904	153.9056	0.4926
RF	1.3241	59.3654	0.6693	1.5014	67.5492	0.6345
MEAN	0.9872	95.4777	0.1981	1.0842	104.9737	0.1980
	30%			50%		
NN	1.1993	96.6311	0.3699	1.7054	74.3427	0.4908
PMM	1.9182	98.8680	0.5078	2.3500	81.8861	0.5489
SLR	2.2228	164.3346	0.4758	2.4714	131.9337	0.5173
RF	1.6640	79.4840	0.5810	2.1402	68.1128	0.6036
MEAN	1.2115	117.6751	0.1988	1.6757	87.7802	0.3676

Table 7. Algorithm comparison for variable **GLY** considering three types of missing data mechanisms (MAR, MCAR, MNAR) and different levels of missingness (10, 20, 30 and 50%) based on imputation accuracy (root mean squared error, mean absolute percentage error and classification accuracy)

GLY	RMSE	MAPE	CA	RMSE	MAPE	CA
<i>MCAR</i>						
	10%			20%		
NN	1.2561	16.1328	0.8387	1.2590	16.2744	0.8405
PMM	1.7727	23.3703	0.7566	1.7698	23.4156	0.7550
SLR	1.7831	26.4223	0.6979	1.7701	26.3160	0.7001
RF	1.4650	16.5648	0.8279	1.4737	17.1581	0.8244
MEAN	1.3308	17.3669	0.8480	1.3223	17.3917	0.8499
	30%			50%		
NN	1.2636	16.2883	0.8384	1.2759	16.4178	0.8375
PMM	1.7739	23.4765	0.7552	1.7703	23.4657	0.7549
SLR	1.7712	26.3031	0.7000	1.7812	26.4517	0.6970
RF	1.4894	17.5696	0.8186	1.5417	18.7614	0.8089
MEAN	1.3230	17.3702	0.8474	1.3253	17.4029	0.8484
<i>MAR</i>						
	10%			20%		
NN	1.1772	16.4996	0.8693	1.1963	16.7634	0.8712
PMM	1.6304	23.0727	0.8011	1.6436	23.0982	0.8001
SLR	1.7147	26.8912	0.7427	1.7331	27.2023	0.7401
RF	1.3543	17.0068	0.8581	1.3818	17.7424	0.8502
MEAN	1.2193	18.5149	0.8802	1.2283	18.8077	0.8822
	30%			50%		
NN	1.2200	17.2983	0.8735	1.2587	17.0065	0.8539
PMM	1.6478	23.1569	0.8019	1.7262	23.3995	0.7719
SLR	1.7483	27.4624	0.7383	1.7811	27.0636	0.7138
RF	1.4170	18.6302	0.8426	1.5060	19.1082	0.8203
MEAN	1.2356	19.1823	0.8826	1.2910	18.3311	0.8644
<i>MNAR</i>						
	10%			20%		
NN	1.0146	20.9776	0.9733	1.0910	22.8274	0.9642
PMM	1.5630	26.4797	0.8574	1.6378	27.9884	0.8395
SLR	1.6236	31.2847	0.7685	1.6987	32.8102	0.7441
RF	1.1721	18.8606	0.9312	1.2621	20.9540	0.9162
MEAN	1.0738	23.5457	1.0000	1.1629	25.7221	1.0000
	30%			50%		
NN	1.2012	25.6262	0.9562	1.2858	21.3637	0.8708
PMM	1.7052	29.6050	0.8210	1.7721	26.6352	0.7625
SLR	1.7859	34.6600	0.7103	1.8149	30.4577	0.6779
RF	1.3510	23.2115	0.8989	1.4924	21.4904	0.8303
MEAN	1.2823	28.7493	1.0000	1.3336	23.2268	0.8966

Performances from tables 3, 4, 5, 6 and 7 are summarized in Table 8. The table shows the number of cases where the selected method has the best performance values. Only NN, MEAN and RF methods are shown as competing algorithms with best results overall. One case represents combination of: performance measure, variable, missingness mechanism and volume of missing data. Since there are: 3 metrics, 5 variables, 3 mechanisms and 4 percentages, there are 180 cases in total. For MEAN and RF there are also, in parentheses, values which indicate how many cases have NN as next best performance. In addition to the total winning scores, also given in the table are scores grouped by different context: by variables that are missing, by missingness mechanism and by percentage of missing data.

Table 8. Overview of the number of different tested cases where selected algorithm shows best performance, in total and grouped by: variable, mechanism and volume of missing data. Number of cases where NN has next best performance is shown in parentheses

Method	Total number of winning cases				
NN	133				
MEAN	21 (18)				
RF	26 (20)				
by variable	HDL	LDL	TCH	TG	GLY
NN	24	33	36	19	21
MEAN	8 (5)	0	0	1 (1)	12 (12)
RF	4 (1)	3 (3)	0	16 (13)	3 (3)
by mechanism	MCAR	MAR	MNAR		
NN	40	54	39		
MEAN	12 (9)	4 (4)	5 (5)		
RF	8 (6)	2 (2)	16 (12)		
by volume	10%	20%	30%	50%	
NN	31	33	33	36	
MEAN	4 (4)	4 (4)	6 (4)	7 (6)	
RF	10 (9)	8 (5)	6 (5)	2 (1)	

3.1. Discussion

Observing all obtained performance values, ANNs prevail as the best method for imputation, considering different missing mechanisms and proportion of missing data.

For HDL, MEAN imputation shows the best results for all missing frequencies but only for MCAR data. Also, all those MEAN results are closely followed with ANN and RF algorithm. As missing mechanism changes to mechanisms that enclose dependency in missing data, ANNs emerge as a method with the best performance.

For LDL, there are only three cases for CA metric where RF shows better results than ANN and only in case where data is missing according to MNAR. Even in those cases, ANNs are closely behind.

For TCH, ANNs display the best performance results, considering all three measures across all missing pattern cases.

For TG, mixed results can be observed. MAPE errors are very large for each algorithm in all scenarios, which gives fairly inaccurate imputation overall. In MCAR and MAR case, ANN and RF are competing with similar imputation accuracy (both MAPE and CA) while the other methods perform notably worse. In an MNAR setting, the disparity between performance of ANN and RF results is bigger. Still, ANNs demonstrate the smallest RMSE errors through all settings.

For GLY, it can be observed that MEAN shows best results according to classification, especially in the MNAR case where it shows 100% accuracy. From data distribution (Figure 1), it can be observed that GLY has CMR cutoff value higher than the upper quartile with lots of outliers. MNAR mechanism is simulated by removing lower GLY values with higher probability. When the volume of missing data is smaller (10, 20, 30%) mean of sample set with complete cases stays lower than the cutoff value, therefore imputation gives values that are, as original data, lower than the cutoff, which explains very high classification accuracy. Similarly as for previous variables, ANNs have the smallest RMSE error through all simulation settings.

Regarding the overview of performance given in the table 8, it is noticeable that ANNs are winning in most scenarios. Even for cases where other method is the winner, for a large number of them, ANN is the next best imputation method.

Here should be noted that the obtained performance results (RMSE, MAPE, CA), considering the distribution of values for each variable (Table 1, Figure 1), except maybe for TCH, indicate that neither explored imputation method should be used as a final prediction model for variables separately. Nevertheless, in this research, we are not interested in prediction models for those values as separate models, but rather in imputation in preprocessing phase of building CMR models. That is why we are solely interested in relative, comparison values. If a goal of some future research would be to build prediction models for HDL, LDL, TCH, TG and GLY by and of itself, one can use this research as a supporting ground and explore other sets of predictors, as well as other models for each variable independently.

Limitations and Further Research What should be examined is how ANN and RF with different architectures compare solely, especially for TG, since RFs exhibit good performance for some scenarios. In that case, fine tuning of RFs hyperparameters should be performed since RFs with fixed number of trees is used in this research. Also, although this work provides promising results, it should be explored how distribution around CMR cutoff values is correlated with imputation accuracy, especially for GLY, and should be determined if performance is a result of this specific clinical dataset. Additionally, the used MAR and MNAR settings demonstrate just some examples of these mechanisms. Simulations with different, more sophisticated MAR and MNAR mechanisms could be performed and ascertain if the results are agreeable. Lastly, it could be tested whether introducing additional hidden layers or tuning process and parameters of ANNs could further improve accuracy and performance of neural networks as an imputation method. At the end, final networks could be ensembled in order to broaden the proposed methodology for multivariate imputation, along with exploring how those imputations affect final results in development of new and enhanced CMR prediction models.

4. Conclusions

Prognostic CMR models require simple anthropometric measures as well as some laboratory values. In order to build machine learning models that solely use small, low-cost set of predictors - laboratory values are necessary, but only in the preprocessing phase when outcome CMR values are produced. If some of those laboratory values are missing, a dataset used for learning could be fairly reduced and even produce flawed end results, which can then make process of building final CMR model more difficult. In statistics, missingness is often analyzed separately. In engineering, during machine learning research, this step is sometimes skipped or overlooked. Therefore, we have explored how one machine learning model (ANN) behaves as an imputation method in CMR risk assessment to enable fully independent algorithmic building process for CMR model, diminishing the need for separate analysis of missingness mechanisms.

To explore how neural networks perform as an imputation method for laboratory data used for calculating outcome CMR values which could further be used in building predictive model for CMR, we have explored a number of ANN structures and compared their imputation performance with other simple single imputation methods. First, we have built and tested single layered neural networks with different numbers of hidden neurons to propose optimal settings for univariate imputation of missing values for each variable. Those settings have afterwards been used for comparison of ANNs with other imputation methods. We have simulated three missingness mechanisms (MCAR, MAR and MNAR) and performed simulations for different volume of missing data (10, 20, 30 and 50%). Through all scenarios, ANNs showed strong performance according to different measures of imputation accuracy. They outperformed or were closely behind other methods in almost all the cases, considering both proportion of missing data and missingness mechanism.

Considering the results, we propose that an ANN should be considered and used in imputation of laboratory values, in preprocessing phase, as a step in pipeline framework which could lead to development of more robust and precise CMR prediction models.

Although this work calls for next steps in the future research such as: ensembling obtained networks or development of new types of ANNs which will deal with multivariate missing data and analysis of impact of those imputations in the final model development, this step was necessary in order to explore versatile settings of missing data, systemize results and conclusions and prepare the basis for final CMR assessment.

It is worth noting that this research does not exclude other ML algorithms as potential imputation tools. On the contrary, we also propose that ML algorithms, in general, should be considered, researched and used as imputation methods for both predictors and outcomes values, since they could enable automatic integration of imputation in model development process without a need for separate analysis of data distribution and missingness mechanisms, leaving datasets used for learning complete and final results less prone to error due to missingness mechanisms and volume. It is self-evident that any method for dealing with missing data cannot substitute real data, but machine learning could provide us tools for imputation that can be automatic, self-sufficient and domain independent. In that context, the proposed comparison methodology for this specific problem justifies the effort and could be used as a guideline for further research.

Acknowledgments. This work was partially supported by the Ministry of Education, Science and Technological Development of the Republic of Serbia within the project TR-32044.

References

1. Altman, D.G., Bland, J.M.: Missing data. *Bmj* 334(7590), 424–424 (2007)
2. Ashwell, M., Gunn, P., Gibson, S.: Waist-to-height ratio is a better screening tool than waist circumference and BMI for adult cardiometabolic risk factors: systematic review and meta-analysis. *Obes Res* 13(3), 275–286 (Mar 2012)
3. Beaulieu-Jones, B.K., Moore, J.H.: Missing Data Imputation in the Electronic Health Records Using Deeply Learned Autoencoders. *Pacific Symposium on Biocomputing* 22, 207–218 (2016)
4. Breiman, L.: Statistical modeling: The two cultures. *Statistical science* 16(3), 199–215 (2001)
5. Domingos, P.M.: A few useful things to know about machine learning. *Commun. acm* 55(10), 78–87 (2012)
6. Efron, B.: Estimating the error rate of a prediction rule: improvement on cross-validation. *Journal of the American statistical association* 78(382), 316–331 (1983)
7. Expert Panel on Detection Evaluation and Treatment of High Blood Cholesterol in Adults: Executive Summary of the Third Report of the National Cholesterol Education Program (NCEP) (Adult Treatment Panel III). *JAMA: The Journal of the American Medical Association* 285, 2486–2497 (2001)
8. Friedewald, W.T., Levy, R.I., Fredrickson, D.S.: Estimation of the concentration of low-density lipoprotein cholesterol in plasma, without use of the preparative ultracentrifuge. *Clinical chemistry* 18(6), 499–502 (1972)
9. Fritsch, S., Guenther, F.: *neuralnet: Training of Neural Networks* (2016), <https://CRAN.R-project.org/package=neuralnet>, r package version 1.33
10. Graham, J.W., Cumsille, P.E., Elek-Fisk, E.: Methods for Handling Missing Data. In: *Handbook of Psychology*, pp. 87–114. John Wiley & Sons, Inc. (2003)
11. Greenland, S., Finkle, W.D.: A critical look at methods for handling missing covariates in epidemiologic regression analyses. *American journal of epidemiology* 142(12), 1255–1264 (1995)
12. Groenwold, R.H., White, I.R., Donders, A.R.T., Carpenter, J.R., Altman, D.G., Moons, K.G.: Missing covariate data in clinical research: when and when not to use the missing-indicator method for analysis. *CMAJ* 184(11), 1265–1269 (2012)
13. Grundy, S.M., Cleeman, J.I., Daniels, S.R., Donato, K.A., Eckel, R.H., Franklin, B.A., Gordon, D.J., Krauss, R.M., Savage, P.J., Smith, S.C., Spertus, J.A., Costa, F.: Diagnosis and Management of the Metabolic Syndrome. *Circulation* 112, 2735–2752 (2005)
14. Hughes, R.A., Heron, J., Sterne, J.A., Tilling, K.: Accounting for missing data in statistical analyses: multiple imputation is not always the answer. *International journal of epidemiology* 1, 11 (2019)
15. Igel, C., Toussaint, M., Weishui, W.: Rprop using the natural gradient. In: *Trends and applications in constructive approximation*, pp. 259–272. Springer (2005)
16. Kannel, W.B., McGee, D., Gordon, T.: A general cardiovascular risk profile: The framingham study. *The American Journal of Cardiology* 38(1), 46–51 (1976)
17. Kuhn, M.: Building predictive models in r using the caret package. *Journal of Statistical Software, Articles* 28(5), 1–26 (2008)
18. Kupusinac, A., Stokić, E., Srdić, B.: Determination of WHtR Limit for Predicting Hyperglycemia in Obese Persons by Using Artificial Neural Networks. *TEM J* 1(4), 270–272 (2012)
19. Kupusinac, A., Doroslovački, R., Malbaški, D., Srdić, B., Stokić, E.E.: A primary estimation of the cardiometabolic risk by using artificial neural networks. *Computers in Biology and Medicine* 43(6), 751–757 (2013)
20. Laakso, M.: Hyperglycemia and cardiovascular disease in type 2 diabetes. *Diabetes* 48(5), 937–942 (1999)
21. Leke, C., Marwala, T.: Missing data estimation in high-dimensional datasets: A swarm intelligence-deep neural network approach. In: *Advances in Swarm Intelligence*. pp. 259–270. Springer International Publishing (2016)

22. Leke, C., Marwala, T., Paul, S.: Proposition of a Theoretical Model for Missing Data Imputation using Deep Learning and Evolutionary Algorithms. arXiv (2015)
23. Little, R.J., D'Agostino, R., Cohen, M.L., Dickersin, K., Emerson, S.S., Farrar, J.T., Frangakis, C., Hogan, J.W., Molenberghs, G., Murphy, S.A., Neaton, J.D., Rotnitzky, A., Scharfstein, D., Shih, W.J., Siegel, J.P., Stern, H.: The Prevention and Treatment of Missing Data in Clinical Trials. *New England Journal of Medicine* 367(14), 1355–1360 (oct 2012)
24. Little, R.J.A., Rubin, D.B.: *Statistical analysis with missing data*. Wiley (2019)
25. Marshall, A., Altman, D.G., Royston, P., Holder, R.L.: Comparison of techniques for handling missing covariate data within prognostic modelling studies: a simulation study. *BMC medical research methodology* 10(1), 7 (2010)
26. Masconi, K.L., Matsha, T.E., Echouffo-Tcheugui, J.B., Erasmus, R.T., Kengne, A.P.: Reporting and handling of missing data in predictive research for prevalent undiagnosed type 2 diabetes mellitus: a systematic review. *The EPMA Journal* (2015)
27. Organization, W.H.: Obesity: preventing and managing the global epidemic. report of a world health organization consultation. Tech. Rep. 894, WHO Obesity Technical Report Series (2000)
28. Papageorgiou, G., Grant, S.W., Takkenberg, J.J.M., Mokhles, M.M.: Statistical primer: how to deal with missing data in scientific research? *Interactive CardioVascular and Thoracic Surgery* 27(2), 153–158 (05 2018)
29. Pesonen, E., Eskelinen, M., Juhola, M.: Treatment of missing data values in a neural network based decision support system for acute abdominal pain. *Artificial Intelligence in Medicine* 13(3), 139 – 146 (1998)
30. R Core Team: *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria (2018), <https://www.R-project.org/>
31. Riedmiller, M.: Advanced supervised learning in multi-layer perceptrons — from backpropagation to adaptive learning algorithms. *Computer Standards and Interfaces* 16(3), 265 – 278 (1994)
32. Rosolova, H., Nussbaumerova, B.: Cardio-metabolic risk prediction should be superior to cardiovascular risk assessment in primary prevention of cardiovascular diseases. *The EPMA journal* 2, 15–26 (2011)
33. Salgado, C.M., Azevedo, C., Proença, H., Vieira, S.M.: *Missing Data*, pp. 143–162. Springer International Publishing, Cham (2016)
34. Schafer, J.L., Graham, J.W.: Missing data: our view of the state of the art. *Psychological methods* 7(2), 147–77 (2002)
35. Schouten, R.M., Lugtig, P., Vink, G.: Generating missing values for simulation purposes: a multivariate amputation procedure. *Journal of Statistical Computation and Simulation* (2018)
36. Silva-Ramírez, E.L., Pino-Mejías, R., López-Coello, M., de-la Vega, M.D.C.: Missing value imputation on missing completely at random data using multilayer perceptrons. *Neural Networks* 24(1), 121 – 129 (2011)
37. Stokić, E., Tomić-Naglić, D., Derić, M., Jorga, J.: Therapeutic options for treatment of cardiometabolic risk. *Medicinski preglad* 62 Suppl 3, 54–8 (2009)
38. Van Buuren, S.: *Flexible imputation of missing data*, vol. 20125245. Chapman and Hall/CRC (2012)
39. Van Buuren, S., Groothuis-Oudshoorn, K.: mice: Multivariate imputation by chained equations in r. *Journal of Statistical Software* 45(3), 1–67 (2011)
40. Waljee, A.K., Higgins, P.D., Singal, A.G.: A primer on predictive models. *Clinical and translational gastroenterology* 5(1), e44 (2014)
41. Ware, J.H., Harrington, D., Hunter, D.J., D'Agostino, R.B.: Missing data. *New England Journal of Medicine* 367(14), 1353–1354 (2012)

Dunja Vrbaški is a PhD student at the Faculty of Technical Sciences at the University of Novi Sad, Serbia. She received her BSc and MSc degrees in Computer Science from

the Faculty of Sciences at the University of Novi Sad in 2003 and 2013, respectively. Prior to her enrolment to doctoral studies she has been working as a software developer and project manager in the IT industry. Her research interests are algorithms, machine learning and data visualization.

Aleksandar Kupusinac is an Associate Professor at the Department of Computing and Control Engineering at the Faculty of Technical Sciences, University of Novi Sad, Serbia. He received the Dipl. Ing. degree in Electrical Engineering in 2005, and the MSc and PhD degrees in Computer Science, in 2008 and 2010 from the Faculty of Technical Science, University of Novi Sad. His research interests include data science and artificial intelligence.

Rade Doroslovački is a Professor of Mathematics at the Faculty of Technical Sciences, University of Novi Sad. He received his BSc, MSc and PhD in Mathematics from the Faculty of Sciences at the University of Novi Sad in 1976, 1984 and 1989, respectively. His research interests include combinatorics, graph theory and discrete mathematics. He is currently serving as the Dean of the Faculty of Technical Sciences, University of Novi Sad.

Edita Stokić holds a PhD in Internal Medicine and is a full professor at University of Novi Sad, Faculty of Medicine. She is also employed in the Clinic of Endocrinology, Diabetes and Metabolic Disorders of the Clinical Centre of Vojvodina.

Dragan Ivetić was born in Subotica and received the Dipl. Ing. degree in Electrical Engineering in 1990, and the MSc and PhD degrees in computer science, in 1994 and 1999 respectively, from the Faculty of Technical Science, University of Novi Sad. Since 2000 he has been Professor of Computer Science at the Department for Computing and Automatics, Faculty of Technical Sciences, University of Novi Sad. His research interests include HCI, medical informatics, and computer graphics. Professor Ivetić received DAAD scholarship in 1997 and ACM scholarship in 1998.

Received: July 10, 2019; Accepted: January 20, 2020.

5. Appendix

Graphical representations of algorithm comparisons for different percent and mechanism of missingness are given on the following figures.

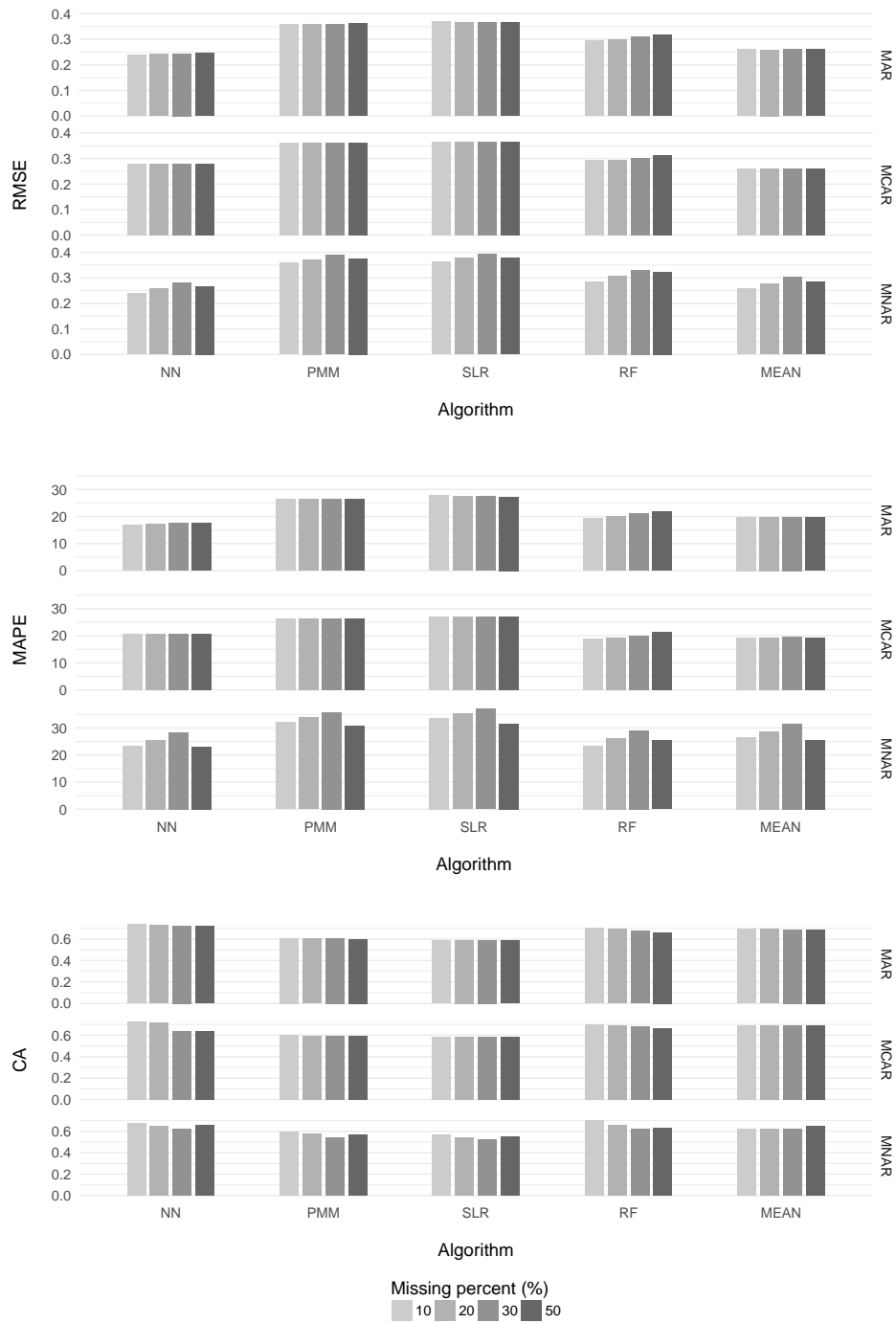


Fig. 2. Results of algorithms comparison for variable HDL. Results are grouped by accuracy performance (RMSE, MAPE, CA) shown on the right and missing data mechanisms (MAR, MCAR, MNAR) shown on the left. Every volume of missingness (10, 20, 30 and 50%) is displayed by separate color

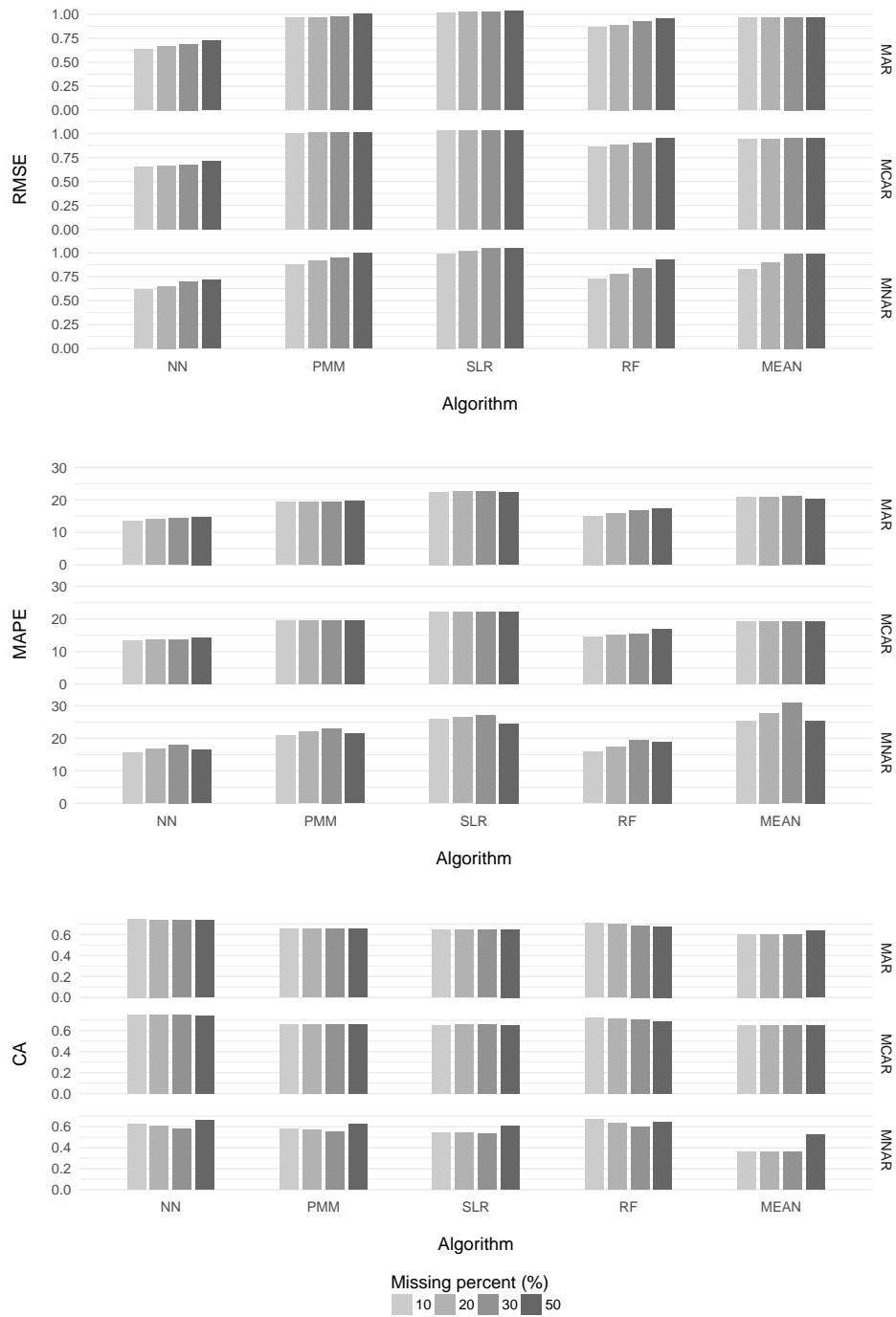


Fig. 3. Results of algorithms comparison for variable LDL. Results are grouped by accuracy performance (RMSE, MAPE, CA) shown on the right and missing data mechanisms (MAR, MCAR, MNAR) shown on the left. Every volume of missingness (10, 20, 30 and 50%) is displayed by separate color

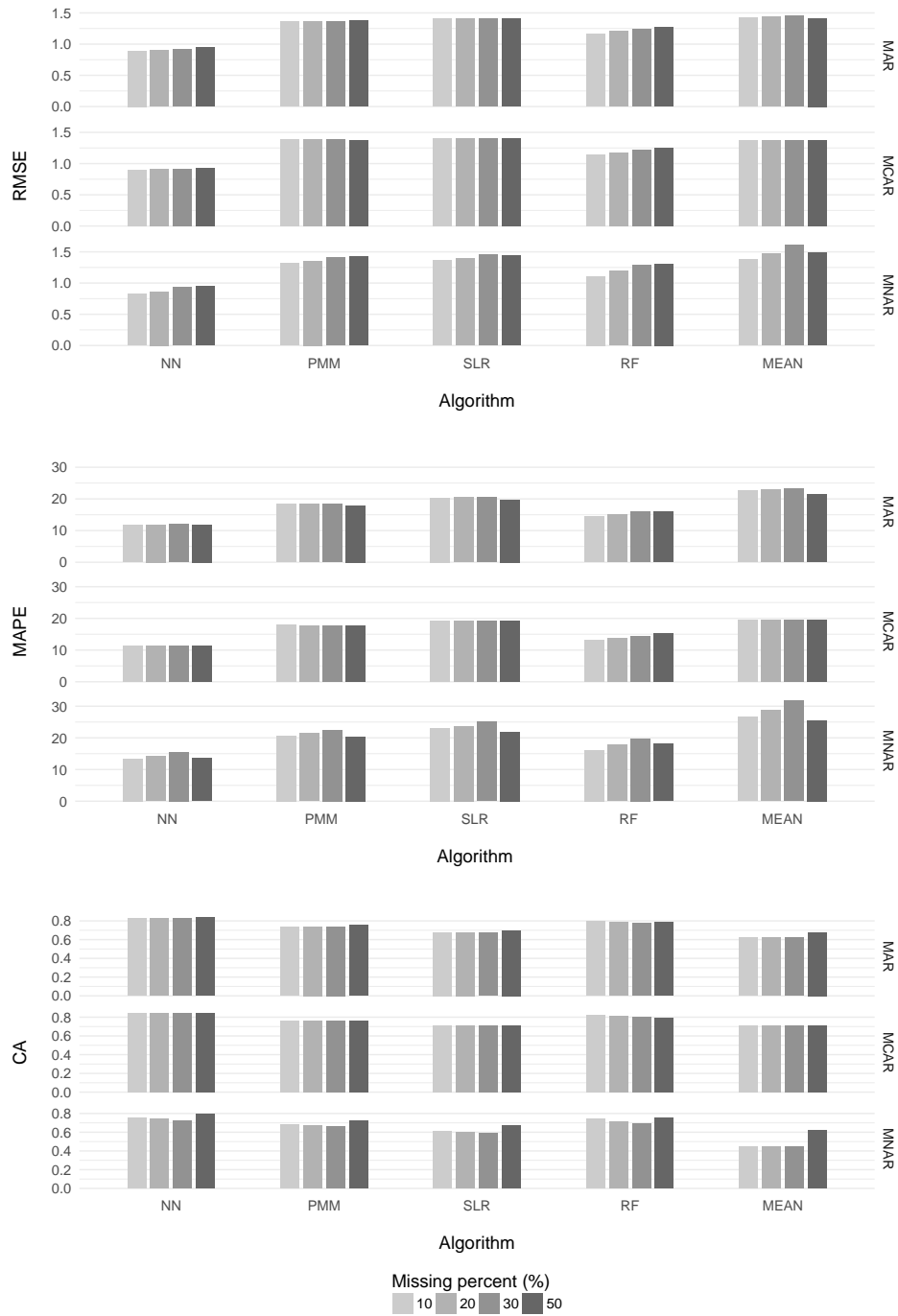


Fig. 4. Results of algorithms comparison for variable TCH. Results are grouped by accuracy performance (RMSE, MAPE, CA) shown on the right and missing data mechanisms (MAR, MCAR, MNAR) shown on the left. Every volume of missingness (10, 20, 30 and 50%) is displayed by separate color

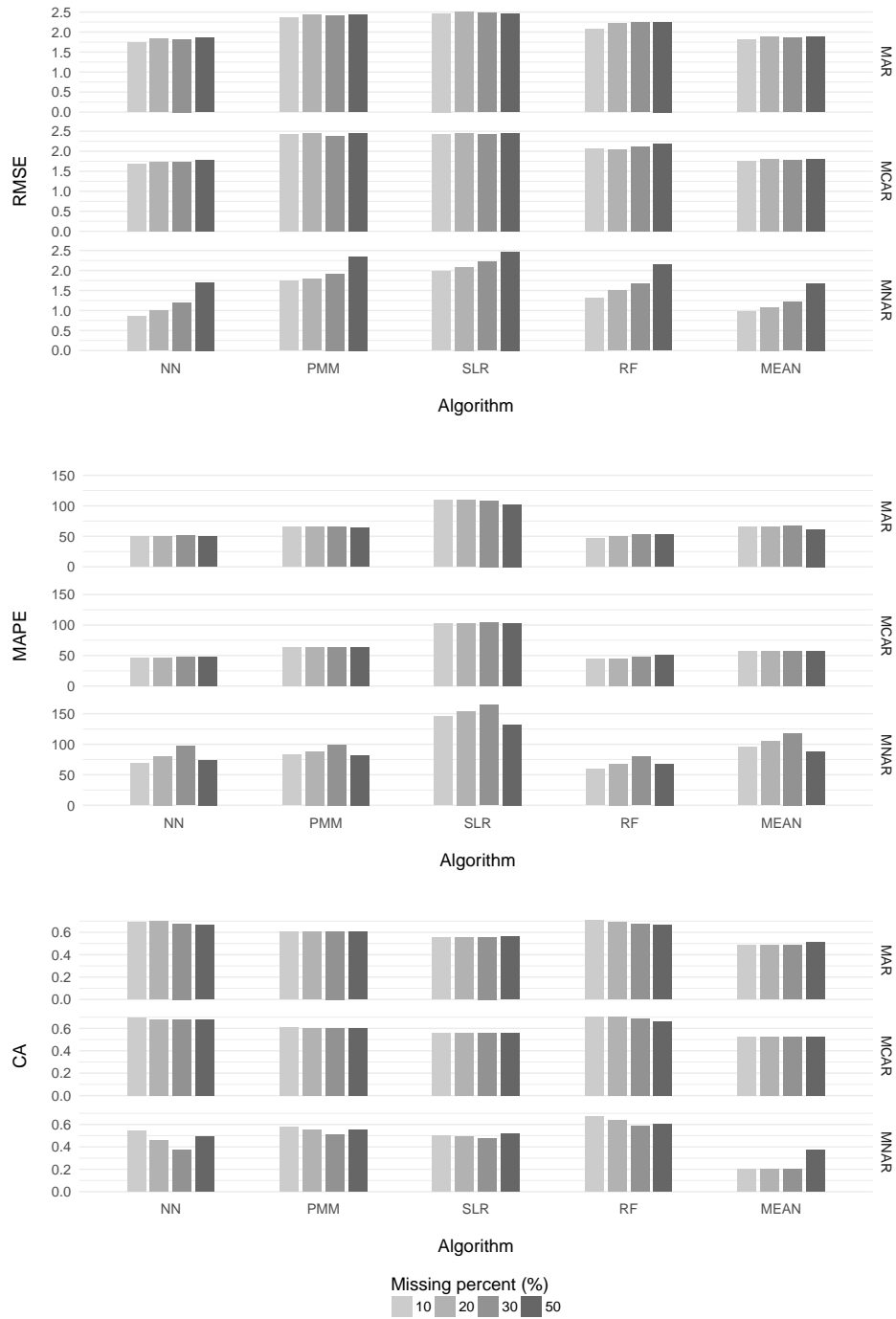


Fig. 5. Results of algorithms comparison for variable TG. Results are grouped by accuracy performance (RMSE, MAPE, CA) shown on the right and missing data mechanisms (MAR, MCAR, MNAR) shown on the left. Every volume of missingness (10, 20, 30 and 50%) is displayed by separate color

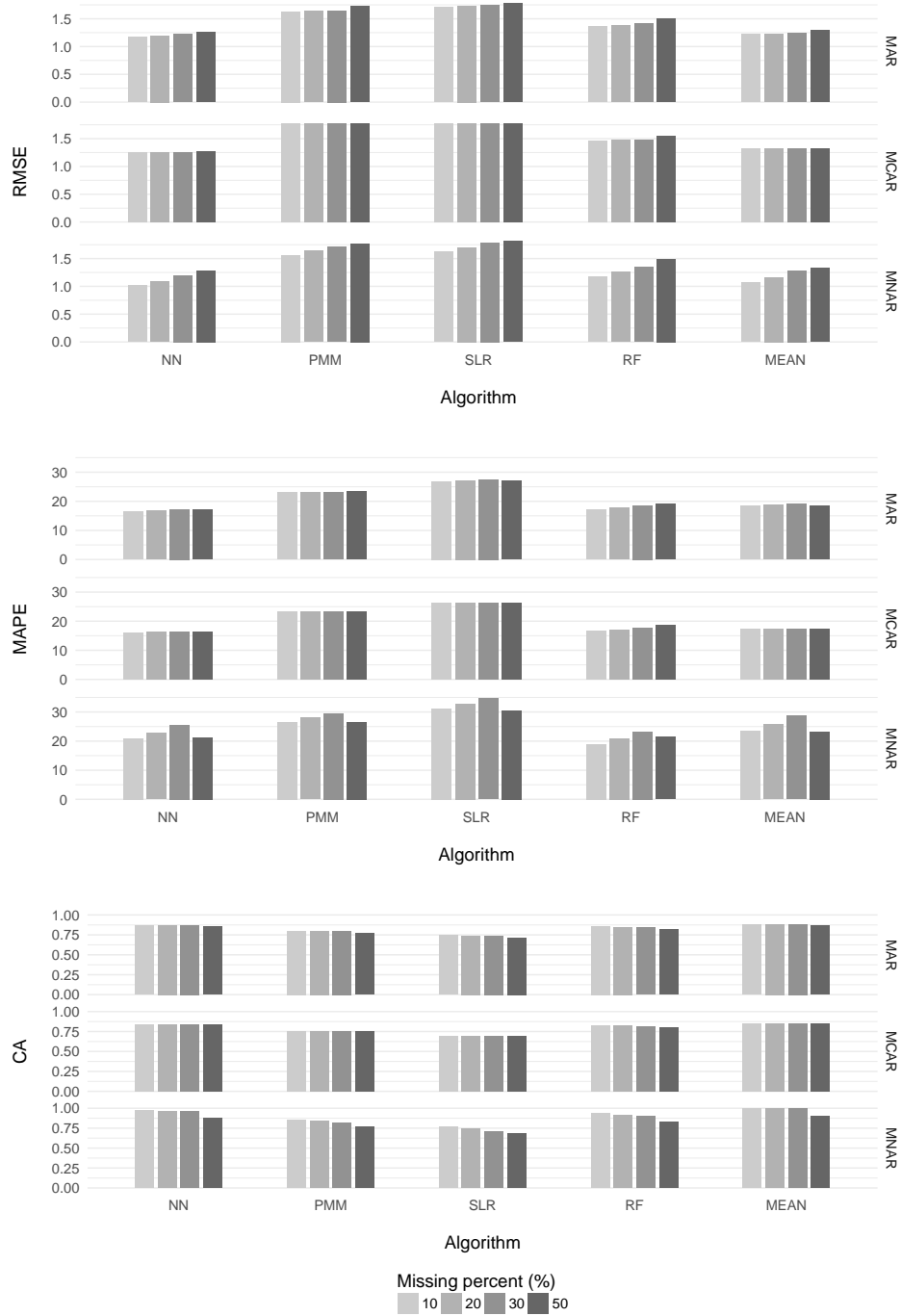


Fig. 6. Results of algorithms comparison for variable **GLY**. Results are grouped by accuracy performance (RMSE, MAPE, CA) shown on the right and missing data mechanisms (MAR, MCAR, MNAR) shown on the left. Every volume of missingness (10, 20, 30 and 50%) is displayed by separate color

Real-Time Tracking and Mining of Users' Actions over Social Media *

Ejub Kajan¹, Noura Faci², Zakaria Maamar³, Mohamed Sellami⁴, Emir Ugljanin⁵,
Hamamache Kheddouci², Dragan H. Stojanović⁵, and Djamel Benslimane²

¹ State University of Novi Pazar
Vuka Karadžica bb, 36300 Novi Pazar, Serbia
dr.ejubkajan@gmail.com

² Univ Lyon, Université Claude Bernard Lyon 1, LIRIS
69622, Villeurbanne Cedex, France
firstname.lastname@univ-lyon1.fr

³ Zayed University
Po Box 19282, Dubai, U.A.E
zakaria.maamar@zu.ac.ae

⁴ Télécom SudParis, SAMOVAR, Institut Polytechnique de Paris
91011, Evry Cedex, France
mohamed.sellami@telecom-sudparis.eu

⁵ University of Niš
Aleksandra Medvedeva 14, 18106 Niš, Serbia
emirugljanin@gmail.com, dragan.stojanovic@elfak.ni.ac.rs

Abstract. With the advent of Web 2.0 technologies and social media, companies are actively looking for ways to know and understand what users think and say about their products and services. Indeed, it has become the practice that users go online using social media like Facebook to raise concerns, make comments, and share recommendations. All these actions can be tracked in real-time and then mined using advanced techniques like data analytics and sentiment analysis. This paper discusses such tracking and mining through a system called *Social Miner* that allows companies to make decisions about what, when, and how to respond to users' actions over social media. Questions that *Social Miner* allows to answer include what actions were frequently executed and why certain actions were executed more than others.

Keywords: Data analytics, Facebook, Sentiment analysis, Social media.

1. Introduction

Business Processes (BP) are a cornerstone to the success of any company that wishes to sustain its growth and remain competitive. According to [25], "*a process is nothing more than the coding of a lesson learned in the past, transformed into a standard by a group of experts and established as a mandatory flow for those who must effectively carry out the work*". More precisely, a BP consists of tasks (t) connected to each other according to

* This is an extended version of a 2-page WETICE2016 demo paper [36].

a process model defining, at design-time, who does what, when, and where. At run-time, tasks are assigned to persons (p) and/or machines (m) so for execution.

Since the advent of social media, the traditional view of how companies operate has completely changed. Contrary to top-down commands and bottom-up feedback that limit innovation and creativity in companies, interactions that social media allows to happen in companies are crossing all levels and occurring in all directions. This organizational shift reduces cost, improves efficiency, facilitates innovation, among other benefits [7],[35]. Social media is also impacting the design of BPs. Earlier, we looked into this design to shed light on social interactions between tasks ($t2t$), between persons ($p2p$), and between machines ($m2m$) in BPs [14],[21]. These interactions reveal for instance, which task is “easy” to replace with other tasks, which person is mostly solicited for partnership with other persons, and which machine works well with other machines.

As a follow-up to our work on BP social-design we also looked into bridging the gap between the business world (hosting BPs) and social world (hosting Facebook as an illustrative application of social media) [20]. While the business world continuously attracts the R&D community’s attention [13],[23], the social world’s surface is barely “scratched” and hence, several opportunities are untapped. To reverse this trend we raise many questions that need responses such as who are the social world’s stakeholders, what actions can the stakeholders perform, and how to track the interactions in the social world. While we detail in [10] the actions that stakeholders perform in the context of social media, we focus in this paper on the interactions that arise in the social world in response to both events triggered and actions taken in the business world and then, to what extent these interactions would impact BPs. For instance, increasing the delivery fees of goods in a process could raise concerns over social media that decision makers in the business world would like to be aware of, should corrective actions need to be taken to address these concerns. For this purpose, we associate the business world with control flow and the social world with social flow, define the constituents of each flow, and manage the cross-flow interactions. We present the design and development of Social Miner, a **real-time** tracking and mining tool on top of Facebook.

The rest of this paper is organized as follows. Section 2 gives a short overview of social media mining in favor of business. Section 3 presents a motivating scenario to stress out the gap in examining the interactions between the business and social worlds. Section 4 briefly discusses data mining (with focus on sentiment analysis) and business process modeling. Section 5 details the real-time tracking and mining of users’ actions over social media where Facebook is used for illustration purposes. Section 6 implements this tracking and mining through a **real** marketing campaign on Facebook. Finally, Section 7 discusses the importance of connecting the business and social worlds together. Finally, Section 8 concludes the paper and identifies some future work.

2. Related work

Mining social networks in favor of business application is not new. Bonchi et al. present an overview of key problems in this domain and the techniques in social network analysis in an infant stage and emphasize, among others, potential business benefits and technical challenges [5]. These challenges include data preparation, network dynamics, propaga-

tion, evaluation, etc, whilst value-added benefits may be summarized in a short sentence: *"Social networking may allow increased revenue"*.

In essence, 2 types of social network mining strategies are in use. The first strategy mimics social networks by tracking activities of employees and their relationships inside an organization during business process execution in order to make them more flexible and productive. Examples include, but not limited to, MiSoN (Mining Social Networks) and SUPER (Social-based bUusiness Process managEmEnt fRamework). In MisoN, Aalst et al. use events logs for social networks mining [1]. These event logs are made by employees of an enterprise when they use some of enterprise information systems (e.g., ERP and CRM) and transform into sociograms by MisoN, which are later use for workflow analysis. SUPER is based on social relations between employees who are in charge to execute particular BPs. These relations are delegation, substitution and peering and results of their use are reported in [14,21]

The second strategy goes beyond an enterprise and use public social networks (e.g., Facebook and Tweeter) to mine customers' opinions about companies' products and services to facilitate CRM via customer needs anticipation and reputation monitoring, identify their churn and reasons for it, find experts, etc. [5]. These findings are obtained using sentiment analysis, opinion mining, and at large scale social influence mining. For the later, Tang et al. present how this mining may help expert finding [33]. A general approach that allows to identify the node that influences others is presented in [4].

3. Motivating scenario

GreenUtility is a utility company that is going to launch an awareness campaign about renewable energies on Facebook. To achieve the campaign's targets like increasing the number of green advocates, the marketing team must look after this campaign's business aspects (e.g., develop the campaign's content and layout, secure the necessary approvals, and assess the results of the campaign) and social aspects (e.g., announce the campaign on Facebook page, engage in discussions with this Facebook page's subscribers, post, and refresh the marketing content if necessary). Both business and social aspects become intertwined when the campaign is in a full swing. For instance, posting content on Facebook needs the marketing director's approval. And, collecting subscribers' comments from Facebook permits to decide on extending the campaign.

After securing the necessary approvals to launch the campaign, GreenUtility's Facebook page is updated with details like tips for saving energy and actions for contributing to a green world. Afterwards, the subscribers to this Facebook page **could** (in fact, the subscribers are not obliged) react to the campaign by posting responses, initiating new communication threads, and expressing feelings over some ecological incidents (e.g., Gulf of Mexico oil spill).

Putting all the social actions (e.g., post comments and invite others) that subscribers perform on Facebook page together should permit to develop social flows that would give GreenUtility better insights into what is being discussed over its Facebook page instead of relying on some quantitative performance indicators (e.g., number of visitors), only. Mining the social flows would help GreenUtility answer many questions like when is it appropriate to post a campaign so that a good response rate is achieved, who are the main

supporters/opponents of/to a campaign so that their feedback/concerns are shared/dealt with, and who should respond to supporters/opponents so that fruitful discussions occur.

4. Some preliminaries

This section is an overview of data mining and business process modeling.

Data mining. The phenomenal growth of online social media (e.g., discussion fora and blogs) is backed by the abundance and richness of user content such as comments, reviews, and feeds that could correspond to opinions on events like visited places, tried restaurants, and consulted books. Opinions can also be associated with sentiments reflecting users' attitudes and feelings towards events like anger and happiness. For a better understanding of opinions, many Sentiment Analysis (SA) techniques and tools are reported in the literature (e.g., Batrinca and Treleaven [3] and Tang et al. [32]). Generally speaking, Pang and Lee suggest 3 stages that SA goes through: *opinion retrieval*, *sentiment classification*, and *opinion summarization* [26]. The first determines which textual sources (e.g., documents, posts, blogs, and news) should be considered when looking for opinionated material with respect to a certain granularity level (e.g., entire documents, phrases, and separate words). The second identifies the overall sentiment that each textual source conveys. Finally, the third provides an integrated view of sentiments expressed by multiple textual sources. Fig. 1 depicts a comprehensive view of SA's outcomes according to the type of source to analyze (single *versus* multiple).

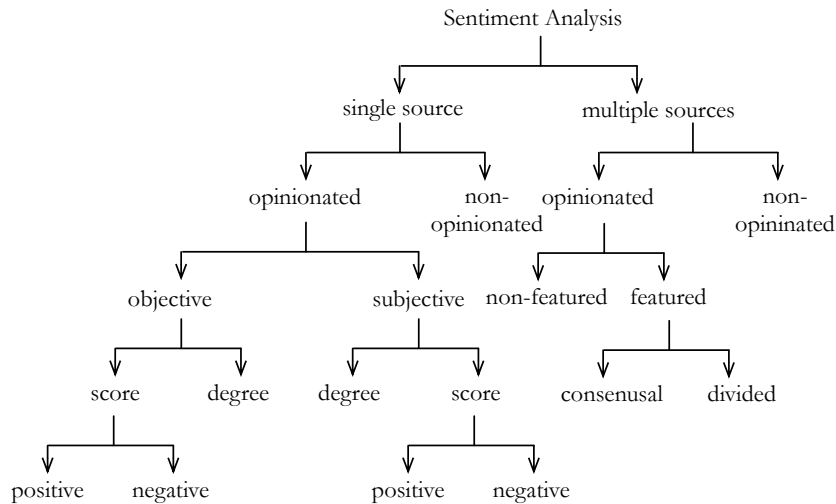


Fig. 1. Sentiment analysis decomposition

The aforementioned 3 stages can require different techniques and tools that are at the crossroad of Natural Language Processing (NLP), Information Retrieval (IR),

and Machine Learning (ML) [26]. Since textual sources could be opinion-free, the *opinion retrieval* stage relies on IR techniques to rank content in terms of *relevancy* (i.e., whether the content is topic-related) and *opinionatedness* (i.e., whether the content contains opinions) (e.g., Luo et al. [18] and Soboroff et al. [30]). The *sentiment classification* stage adopts NLP-based ML techniques to agree on content's *subjectivity/objectivity* and *polarity* (e.g., Liu [16] and Ting et al. [34]). Polarity that can be either qualitative or quantitative, permits to tag opinions with positive/negative sentiment scores (e.g., good/bad and like/dislike) or degrees (e.g., good/excellent and bad/worst). Last but not least, the *opinion summarization* stage uses IR techniques related to a content's *featurability* and NLP-based ML techniques related to *consensus* and *division* (e.g., Archak et al. [2,17]). The former identifies important features that would characterize some events. And, the latter decides whether multiple sources contain similar and/or contradictory opinions related to features so that opinions are differentiated.

Despite the benefits of existing SA techniques, many shortcomings can be reported. Indeed, they consider textual sources as one block and, thus, overlook nested exchanges in blocks. Plus, they do not capture users' attitudes towards received content (should the user respond, delay to respond, or ignore). In this work, we move one step-forward by analyzing a content's structure and growth along with users' sentiments so that we understand what happened, might happen as well as when opinion changes happened.

Business process modeling in brief. BP modeling is about documenting and displaying BPs graphically to help stakeholders analyze process models and find possible ways of improvement. A modeling language consists of 3 parts: (i) syntax that provides constructs and rules to combine constructs, (ii) semantics that gives meaning to constructs, and (iii) notation that includes graphical symbols to visualize constructs. In [27], Pourshahid et al. state that all 3 parts together should allow to model various aspects of a BP such as tasks, events, resources, roles, constituents, functions, organization, and hierarchy. Over the years, Business Process Model and Notation (BPMN) has rapidly become a standard for BP modeling as per the Object Management Group (OMG) [24]. BPMN provides graphical elements to develop multiple flows like control flow to show the partial order (i.e., conditional and concurrent) between tasks in a BP and communication flow to show the exchange of messages between a BP's stakeholders. In addition to BPMN flows, other types of flows are reported in the literature. Sadiq et al. use data flow for process specification following a data/artifact-centric perspective and process verification according to 3 properties: *correctness*, *soundness*, and *variability* [29]. Reichert and Dadam use control flow and data flow to specify BPs following a process-centric and data-driven perspective, and verify a BP's correctness properties like *reachability* and *termination* [28]. Finally, Maamar et al. develop and synchronize control, communication, and navigation flows to monitor BP execution [22].

Completely different from flows for BP modeling that are structured and known in advance, social flows are built on-the-fly and capture social actions over social media. We expect tapping into social flows to understand why users execute certain social actions, what business tasks triggered certain social actions, which social actions are

triggered because of some social actions, what is the execution chronology of social actions, and what content like feelings and opinions does social actions convey.

5. Tracking and mining approach

This section details the approach for real-time tracking and mining users' actions over social media. It starts with some definitions and examples and then, discusses how putting social actions together lead into the development of social flows.

5.1. Foundations

To formalize our approach, we first, refer readers to [19] where a definition for social action (e.g., send, co-author, and tweet) is proposed. It is an operation that a Web 2.0 application allows users to execute whether online or offline. Also, as per the same reference, a social action falls into one of the following categories (Table 1): *communication* (e.g., chat and poke), *sharing* (e.g., publish and upload), and *enrichment* (e.g., comment and tag). In this paper A is the set of all social actions available for execution over all Web 2.0 applications ($A = \{\text{poke, chat, send, co-author, tweet, post, comment, reply, tag, upload, } \dots\}$) and $A_{app2.0}$ ($A_{app2.0} \subseteq A$) is the set of all social actions available for execution over a particular Web 2.0 application ($app_{2.0}$), Facebook in our case ($A_{Facebook} = \{\text{poke, chat, send, post, upload, comment, tag, } \dots\}$).

Table 1. Representative categories of social actions ([19])

Category	Description	Examples of social actions
Communication	Includes actions that establish back-and-forth interactions between users, which should engage them in joint operations	Chat with a user or group of users, Poke someone, Send direct messages to a user's inbox
Sharing	Includes actions that establish one-way interactions and allow to create and edit shared content and to facilitate this content's consumption	Co-author a text/media on a Wiki, Publish a post on a Blog Web site, Upload a photo/video on a public repository, or any other data (e.g., sensor reading)
Enrichment	Includes actions that provide additional [meta] data on shared content by providing opinions and/or ranking	Comment a post, Tag users' photos, videos, activities, etc.

Fig. 2 depicts how the control and social flows are anchored to the business and social worlds, respectively, with focus on the interactions from the business to social worlds that trigger forming social flows. Interactions from the social to business worlds are the result of mining social flows. For clarity purposes, these interactions are not represented in Fig. 2 but are handled through metrics defined in Section 6.2. It happens that the successful execution of some business tasks in the control flow makes users execute some social actions over Web 2.0 applications (e.g., after approving the marketing content, GreenUtility posts the content on Facebook). The outcomes of executing these social actions could motivate (same and/or other) users to execute additional social actions and so

on. This execution chain of social actions expands until some termination conditions are met (e.g., deadline has passed and response-rate has become low). By putting all the social actions together with respect to when they occurred, social flows are obtained. Fig. 2 shows 1 control flow (F_1^c) and 3 social flows (F_1^s , F_2^s , and F_3^s). Later it will be shown that social flow could branch into sub social-flows (*aka* nested flows).

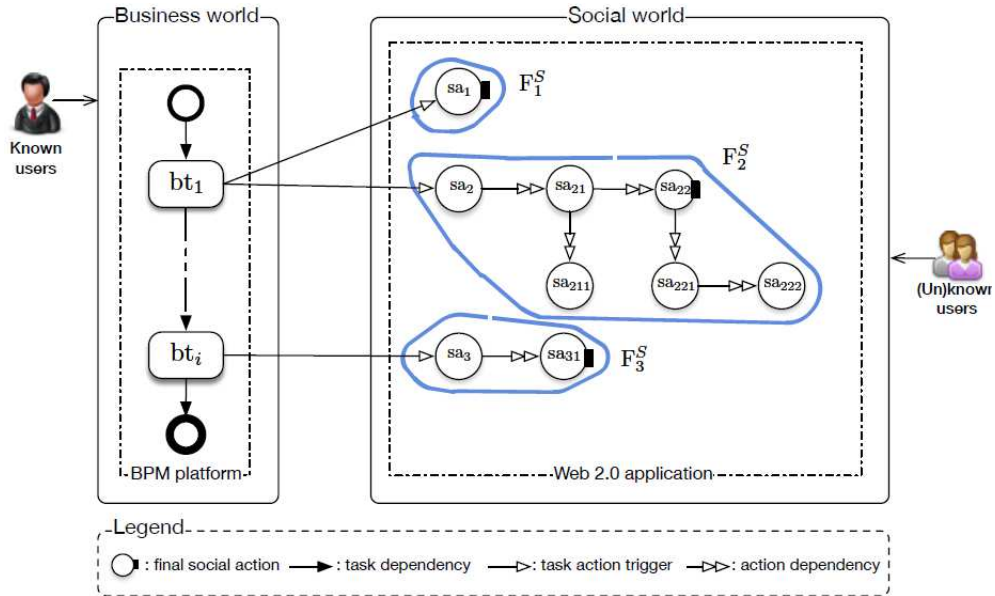


Fig. 2. Business and social worlds from a flow perspective

5.2. Definitions and examples

In the following all examples are drawn from the motivating scenario and Fig. 2. The definitions given in this section are integrated into the automatic building of social flows as demonstrated in Section 6.

Definition 1. Control Flow (F^c). It represents the process model of a BP and consists of business tasks (bt) and dependencies between business tasks. Formally, F^c is a 4-tuple $\langle T^c, D^c, IT^c, FT^c \rangle$ where: T^c contains all business tasks in a BP; $D^c \subseteq T^c \times T^c$ is the set of all dependencies between business tasks; $IT^c \subseteq T^c$ is the set of all initial business tasks; and, $FT^c \subseteq T^c$ is a set of all final business tasks.

Example: $F^c = \langle \{bt_1, bt_2, \dots, bt_i\}, \{(bt_1, bt_2), \dots, (bt_{i-1}, bt_i)\}, \{bt_1\}, \{bt_i\} \rangle$ where $bt_1 = \text{secure-necessary-approvals}$ and $bt_2 = \text{develop-campaign-design}$.

Definition 2. Business-to-Social Link ($L(b2s)$). It captures the statement that “upon the successful execution of a business task, a user **may** execute social actions in response to this execution”. We aim at conditionally mapping each business task in a control flow onto a set of initial social actions that will form the roots of the future social flows. Formally, $L(b2s) : T^c \times C \rightarrow 2^{A^{s_0}}$ is a function where: T^c contains all business tasks in a BP, C is a set of conditions; and, A^{s_0} is a set of all initial social actions in the social flows.

Example:

- $L(b2s) : (bt_1, [cond_{bt_1}]) \rightarrow \{sa_1, sa_2\}$ where $bt_1 =$ secure-necessary-approvals, $cond_{bt_1} =$ is-campaign-approved?, $sa_1 =$ post-benefits-of-the-campaign-on-Facebook, $sa_2 =$ post-benefits-of-the-campaign-on-Twitter; sa_1 and sa_2 are executed if the campaign is approved.
- $L(b2s) : (bt_2, \phi) \rightarrow \phi$ where $bt_2 =$ develop-campaign-design; no social action is associated with bt_2 .
- $L(b2s) : (bt_3, [cond_{bt_3}]) \rightarrow \{sa_3\}$ where $bt_3 =$ analyze-customer-application, $cond_{bt_3} =$ is customer application rejected?, and $sa_3 =$ post-a-note-on-customer-wall; sa_3 is executed if the customer application is rejected.

Definition 3. Social Flow (F^s). It is a set of social actions put together on-the-fly. One of these social actions is initial (i.e., linked to a business task as per Definition 2) and the rest are either intermediaries or finals. First, the connection between social actions is dependent on (i) the authorized relations that Web 2.0 applications allow to have between their social actions (Section 5.3) and (ii) the nested levels of exchange that a Web 2.0 application allows to happen⁶. Second, the selection of the next social actions to execute is based on contextual elements that do not fall into this paper’s scope. Formally, F^s is a 4-tuple $\langle A_{app2.0}^s, STR_{app2.0}^s, sa_0^s, FA^s \rangle$ where $A_{app2.0}^s \subseteq A_{app2.0}$ contains those social actions in a Web 2.0 application that end-users have **voluntarily** decided to execute; $STR_{app2.0}^s : A_{app2.0}^s \times L_{app2.0} \rightarrow A_{app2.0}^s$ is a function that corresponds to a time-stamped authorized relation connecting a social action, that occurred at a certain level of exchange ($l \in L_{app2.0}$), to another social action; $sa_0^s \in A_{app2.0}^{s_0}$ is the initial social action; and, $FA^s \subseteq A_{app2.0}^s$ is a set of final social actions.

Example:

- $F_1^s = \langle A_{app2.0_1}^s, STR_{app2.0_1}^s, sa_{0_1}^s, FA_1^s \rangle$ where $A_{app2.0_1}^s = \{sa_1\}$, $STR_{app2.0_1}^s$ not applicable, $sa_{0_1}^s = sa_1$, and $FA_1^s = \{sa_1\}$.
- $F_2^s = \langle A_{app2.0_2}^s, STR_{A_{app2.0_2}^s}^s, sa_{0_2}^s, FA_2^s \rangle$ where $A_{app2.0_2}^s = \{sa_2, sa_{21}, sa_{221}, \dots\}$, $STR_{A_{app2.0_2}^s}^s = [(sa_2, sa_{21}), (sa_{21}, sa_{22}), [(sa_{21}, sa_{211})], [(sa_{22}, sa_{221}), (sa_{221}, sa_{222})]]$, $sa_{0_2}^s = sa_2$, and $FA_2^s = \{sa_{21}, \dots\}$. In this flow, 2 levels of exchange represented by [] and [[]], respectively, exist. Note that F_2^s contains 1 primary sub-flow referring to social actions between [] and 2 secondary sub-flows referring to social actions between [[]].

⁶ E.g., 1 nested-level of exchange would support 2 types of social flows: primary and secondary. In Facebook *comment* and *reply* can trigger *post*. Thus, *comment* and *reply* are part of the primary social-flow and *post* is part of the secondary social-flow.

- $F_3^s = \langle A_{app2.0_3}^s, STR_{app2.0_3}^s, sa_{0_3}^s, FA_3^s \rangle$ where $A_{app2.0_3}^s = \{sa_3, sa_{31}\}$,
 $STR_{app2.0_3}^s = \{(sa_3, sa_{31})\}$, $sa_{0_3}^s = sa_3$, and $FA_3^s = \{sa_{31}\}$
- $A_{app2.0}^{s_0} = \{sa_1, sa_2, sa_3\}$ for $F_{1,2,3}^s$.

In addition to control flow, business-to-social link, and social flow definitions, extra concepts and definitions are deemed necessary to allow the mining of social flows. Among these concepts we cite *scores* of social actions (*nodes* for short) that are calculated while the different social flows are under-development. As per Fig 1, each content to analyze is related to a single user, and is both opinionated and subjective. So, a *score* is either *local* that is about user's feedback, *global* that aggregates local scores using direct neighbors' scores, or cumulative that aggregates global scores using direct neighbors' scores, as well.

Definition 4. Local Score Function (LS). It quantifies a user's feedback on a social action using for instance, sentiment analysis techniques such as CoreNLP ([31]). CoreNLP's scores are -2 for very negative, -1 for negative, 0 for neutral, 1 for positive, and 2 for very positive. With respect to social actions that do not have content such as "like" and "wow", their sentiment values are assigned using "common sense", for example. Formally, $LS : A_{app2.0}^s \rightarrow \mathbb{Z}$ is defined as per Equation 1:

$$LS(sa) = \begin{cases} \text{sentimentAnalysis}(sa(\text{feedback})), & \text{sa has content} \\ \text{selfAssignment}(sa()), & \text{sa has no content} \end{cases} \quad (1)$$

where *selfAssignment* assigns a sentiment value to *sa* based on "common sense".

Definition 5. Global Score Function (GS). It represents the cumulative feedback of a social action's free-of-content and secondary neighbors at time *t*. The number of secondary neighbors depends on the Web 2.0 application's nested levels (1 for illustration purposes). Formally, $GS : A_{app2.0}^s \times T \rightarrow \mathbb{Z}$ is defined as per Equation 2:

$$GS(sa, t) = \begin{cases} \text{sign}(sa_{parent}) \times LS(sa) \\ \text{sign}(sa_{parent}) \times LS(sa) + \sum_{i=1}^k GS(sa_i, t), & \text{sa is secondary} \\ LS(sa) + \sum_{i=1}^k GS(sa_i, t) + \sum_{j=1}^m GS(sa_j, t), & \text{sa is primary and } \neq sa_0^s \text{ (Definition 3)} \end{cases} \quad (2)$$

where sa_{parent} is a social action's parent; *sign* is a function that returns +1 if $LS(sa_{parent})$ is positive or neutral, otherwise -1; sa_i is a free-of-content neighbor of *sa*; and, sa_j refers to all secondary neighbors of *sa*. Note that $GS(sa_0^s, t) = \sum_{j=1}^m GS(sa_j, t)$

where sa_j refers to all primary neighbors of sa_0^s .

From an implementation perspective, Algorithm 1 illustrates how Equations 1 and 2 for local and global score calculations, were programmed. At the end of each execution round, each social action gets its respective local and global score, that is used later for analysis and presentation. After a certain time of launching a campaign, collecting necessary details about social actions begins allowing to proceed as follows:

- Lines 4-8 list all reactions on the initial post. The self assignment value for each reaction is multiplied by the sign of the post to get the global score of the reaction, and added to the global score of the post.
- Lines 9-28 list all the comments on the initial post including reactions, replies, and reactions on replies, as well. All of these affect the global score of posts and comments via the nested loops. A comment's sentiment value is checked using CoreNLP tool, multiplied by the post's sign to get its global score and then added to the post's global score.
- Lines 12-26 proceed with the same analysis targeting this time reactions on comments, replies to comments, and reactions to replies, respectively.

Algorithm 1 Local/Global score calculation

```

1: for all posts as p do
2:   p.LS = 0
3:   p.GS = 0
4:   for all reactionsOnPost as rp do
5:     rp.LS = selfAssignment(rp)
6:     rp.GS = sign(p) * rp.LS
7:     p.GS = p.GS + rp.GS
8:   end for
9:   for all commentsOnPost as cp do
10:    cp.LS = sentimentAnalysis(cp.feedback)
11:    cp.GS = sign(p) * cp.LS
12:    for all reactionsOnComment as rc do
13:      rc.LS = selfAssignment(rc)
14:      rc.GS = sign(cp) * rc.LS
15:      cp.GS = cp.GS + rc.GS
16:    end for
17:    for all replyOnComment as rpc do
18:      rpc.LS = sentimentAnalysis(rpc.feedback)
19:      rpc.GS = sign(cp) * rpc.LS
20:      for all reactionsOnReply as rr do
21:        rr.LS = selfAssignment(rr)
22:        rr.GS = sign(rpc) * rr.LS
23:        rpc.GS = rpc.GS + rr.GS
24:      end for
25:      cp.GS = cp.GS + rpc.GS
26:    end for
27:    p.GS = p.GS + cp.GS
28:  end for
29: end for

```

Definition 6. Cumulative Score function (CS). It represents the cumulative score of a social action's direct neighbors at time t . The number of direct neighbors depends on the

Web 2.0 application's nested levels (1 for illustration purposes). Formally, CS is defined as per Equation 3:

$$CS(sa) = \begin{cases} GS(sa) + CS(sa_i, t), & sa \text{ is secondary} \\ GS(sa) + CS(sa_i, t) + CS(sa_j, t), & sa \text{ is primary} \end{cases} \quad (3)$$

where sa_i is both a secondary node and a certain direct neighbor of sa ; and, sa_j is both a primary node and a certain direct neighbor of sa . Note that Equation 3 provides different cumulative scores whether the social action is primary or secondary, and, therefore, considers the social flow's structure. Regarding the initial social actions, their CS is computed recursively using all nodes' cumulative scores. When a change happens in a social flow (e.g., new like, new comment, and new share), CS automatically changes.

When the development of a social flow is in progress, some scores automatically change (e.g., if like is connected to comment at time $t+1$, like's score at $t+1$ will be different from time t). This change would impact other nodes in the social flow through score propagation. We rely on asynchronous self-stabilization principle to propagate impacted global scores after each update (i.e., a newly-added social action to the social flow) [9]. This principle consists of re-computing the scores of initial social actions' neighbors. Each node checks its direct neighbors and detects any change of scores among their neighbors. If there is a change, the node computes its score again and again. Thanks to this domino effect, all nodes update their scores until reaching all initial social actions.

5.3. Authorized relations between social actions

The ongoing expansion of social flows is dependent on the authorized relations that a Web 2.0 application supports in order to connect social actions together ($STR_{app2.0}^s$ in Definition 3). Each Web 2.0 application allows a limited number of (next) social actions from which users can select for execution. Although these relations are not explicitly shown in Web 2.0 applications, we expose them for 2 reasons: enumerate the next possible social actions and recommend some next possible social actions to users with respect to what has been executed earlier. Enumerating the next possible social actions is relevant when building social flows; it permits to track exchanges online and to connect social actions on the fly.

Table 2 suggests examples of next possible social actions in some representative Web 2.0 applications. In this table, $[0/1.. * (resp.1)]sa$ means zero/one to many (resp. only one) social action(s) will be executed, and $(||)$ and (\oplus) are *or* and *xor* logical operators, respectively. To define some authorized relations in Web 2.0 applications, we analyzed Decker and Lesser's coordination relations between tasks namely *facilitates*, *cancel*, *inhibits*, *constrains*, *enables*, and *causes* [8]. Due to the inappropriateness of the first 3 relations for our work, we discuss the remaining ones:

1. $enables(sa_i, \{sa_j\})$: upon the successful execution of a social action sa_i , the Web 2.0 application activates other social actions $\{sa_j\}$ from which users can execute some (i.e., zero to many) and many times. Examples are $enables(share, \{like\})$ and $enables(post, \{share, like, comment\})$ in Facebook and $enables(tweet, \{reply, retweet, post - to - Facebook\})$ in Twitter.

2. $constrains(sa_i, \{sa_j\})$: upon the successful execution of a social action sa_i , the Web 2.0 application activates other social actions $\{sa_j\}$ from which users can execute one social action sa_j , only. Examples are $constrains(request, \{confirm, delete\})$ in Facebook and $constrains(follow, \{accept, deny\})$ in Instagram.
3. $causes(sa_i, \{sa_j\})$: upon the successful execution of a social action sa_i , another social action sa_j is automatically executed. Example is $causes(add, \{follow\})$ in Facebook.

Table 2. Illustration of some authorized relations between social actions

Web 2.0 application	Social action	Next possible social actions
Facebook	post	[0..*]like [0..*]comment [0..*]share [0..1]delete
	share	[0..*]like [0..1]delete
	like	[0..1]unlike
	follow	[0..1]unfollow
	comment	[0..*]reply [0..1]edit [0..1]delete
	friend request	[0..1]confirm \oplus [0..1]delete
Twitter	tweet	[0..*]reply [0..*]re-tweet [0..1]like [0..1]delete
	reply	[0..*]reply [0..*]like [0..*]re-tweet [0..1]delete
	quote tweet	[0..*]reply [0..*]re-tweet [0..1]delete
	like	[0..1]unlike
Instagram	post	[0..*]send to [0..*]comment [0..*]like [0..*]share-to
	comment	[0..*]reply [0..*]like [0..1]delete
	follow	[0..1]accept \oplus [0..1]deny
	send to	[0..1]like \oplus [0..*]comment

Let us consider GreenUtility and Facebook's social actions defined in Table 2. When no social action is executed, GreenUtility administrator executes post so that texts, images, or videos are displayed on the company's Facebook page. This post enables the administrator and other (un)known Facebook members to like that post (like), comment that post (comment), and/or share it (share). This happens because of $enables(post, \{like, comment, share\})$ authorized relation. In the same way, executing one of the recently enabled social actions will allow executing other social actions in a chain reaction. For instance, executing comment after post enables to like that comment (like) and/or to reply to that comment (reply) in compliance with $enables(comment, \{like, reply\})$ authorized relation.

5.4. Illustration

Let us apply the different definitions to GreenUtility. On the one hand, the campaign's business aspect refers to a control flow that includes many business tasks (e.g., prepare-campaign-material and approve-campaign-material) and dependencies between these

tasks (e.g., approve-campaign-material requires finalize-material before). On the other hand, the campaign's social aspect refers to multiple social flows initiated depending on the outcomes of executing certain business tasks. Let us assume that approve-campaign-material is successfully executed. Next is to share this material with the community on Facebook. The administrator logs into GreenUtility's Facebook page and executes post as a social action. Now that the campaign's material is online, subscribers of GreenUtility's Facebook page can share the material with others, comment the material, or like the material as per the authorized relations associated with post (Table 2). If a person makes a comment, then comment as a social action is executed and will be connected to the first social action that is post. At this stage, the under-development social flow consists of two social actions: post then comment.

In Fig. 3a, we show GreenUtility's post on Facebook at time t . This post has resulted into executing additional actions by people like Alison, Bob, and the administrator of GreenUtility's Facebook page. In Fig. 3a, we map this execution onto an under-development social flow that will grow over time. Since Facebook supports 1 nested level of exchange, the social flow is represented as 1 primary (level 0) *caterpillar* (i.e., a tree such that its internal vertices constitute a path and the other vertices are the "hairs" of the tree and 2 secondary (level 1) *caterpillars* [13]. The primary *caterpillar* has GreenUtility's post as a root with 2 like and 2 subsequent comment while the first secondary *caterpillar* has Bob's reply as a hairless root. The nodes are labeled with 3 values (sentiment score, global score, and cumulative score). All under-development social flows are acyclic and temporal.

6. System development

This section consists of 2 parts. The first part describes the architecture supporting real-time monitoring and mining of users' social actions over Facebook. The second part describes the experiments that were carried out along with the results of these experiments.

6.1. Architecture

We developed a tool, named Social Miner (SM), for tracking and mining users' actions over social media with Facebook as a targeted Web 2.0 application. SM's architecture is given in Fig. 4 (a demo video is available at <https://youtu.be/crBsEk2pSzo>) and consists of 4 modules: *dashboard*, *social-action manager*, *social-action tracker*, and *social-flow analyzer*.

The *dashboard* is the interface provided to employees and BP engineers to manage campaigns on GreenUtility's Facebook page like launching a new campaign with the assistance of the *social-action manager* (1.1) and to perform the necessary analysis (1.2). The *social action tracker* "keeps-an-eye" on any change over this Facebook page while the *social-flow analyzer* obtains insights into the activities over the Facebook page such as, which subscribers are (un)supportive of a campaign and which campaign is most attractive according to our mining analysis. To this end, the *social action tracker* uses Webhooks⁷ to subscribe to changes in the Facebook page. These changes are stored in the *social actions*

⁷ developers.facebook.com/docs/graph-api/webhooks..

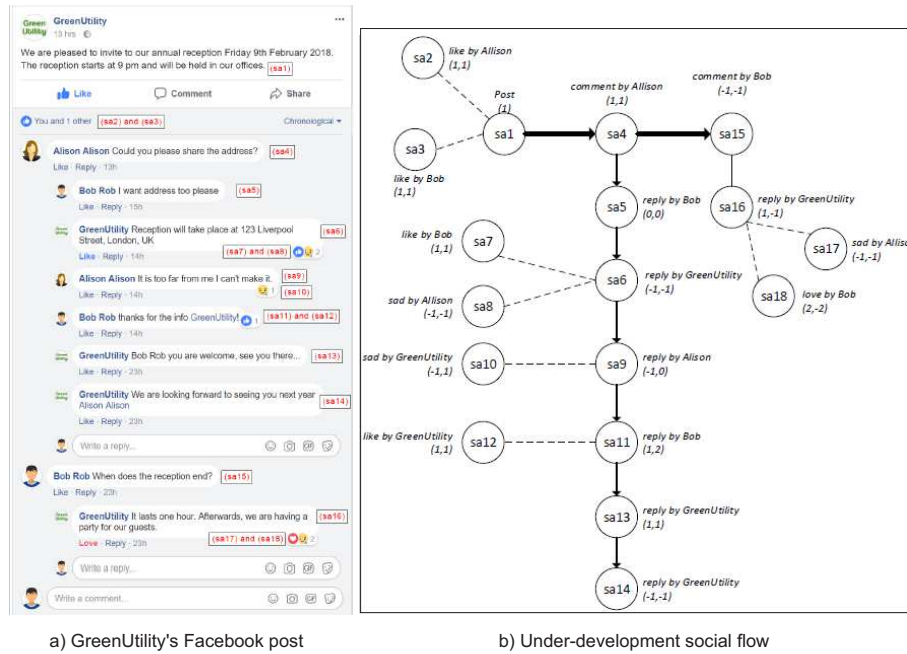


Fig. 3. GreenUtility's Facebook page and its associated social flows

repository and then made available to the *social-flow analyzer* for building the necessary *social flows* so they are mined at a later stage.

1. The *social-action manager* uses Facebook SDK library for PHP so that requests (e.g., publish campaign and reply to some comment message) are submitted to Facebook Graph API⁸ and published on a Facebook page.
2. The *social-actions repository* is a MySQL database that stores details like time about the social actions executed over GreenUtility's Facebook page.
3. The *subscribers repository* is another MySQL database that stores details about the subscribers (e.g., user id, user name, and weight) who take part in the discussions over GreenUtility's Facebook page.
4. The *social action tracker* includes 2 modules: *CoreNLP* and *Score calculator*. On top of time-stamped details about the executed social actions (e.g., user id, page id, post id, and parent id) collected via Facebook Webhooks' notifications, CoreNLP as a sentiment-analysis tool annotates these actions with sentiment scores. The result of CoreNLP analysis is formatted as JavaScript Object Notation (JSON) and then, translated into a relational format (through an in-house script) prior to storing it in the *social-actions repository*. The relational format has eased the storage of different details in multiple tables and running queries over these tables when building and analyzing the social flows. Any notification from Webhooks (e.g., newly added/updated social actions) triggers the *score calculator* that computes new scores (Equa-

⁸ developers.facebook.com/docs/graph-api.

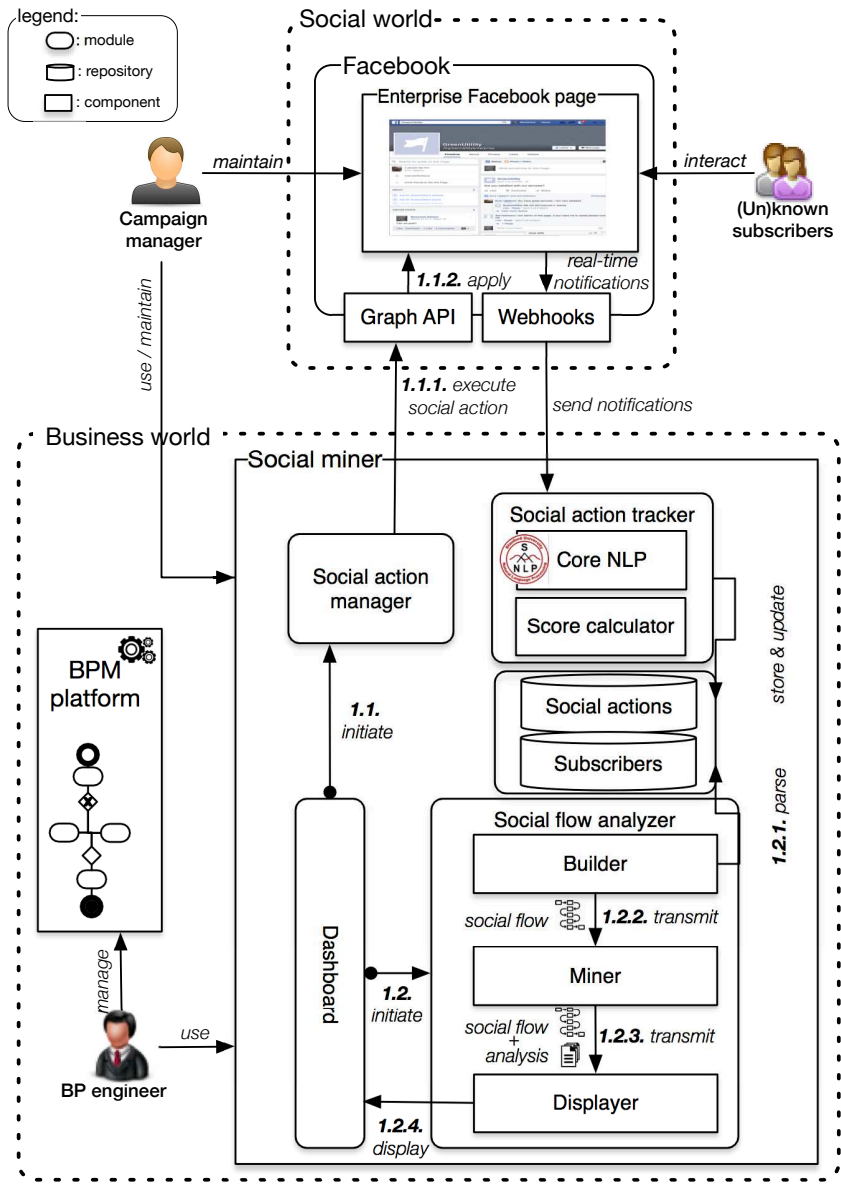


Fig. 4. Architecture of Social Miner

tions 4 and 5) and/or revises some existing scores as per the score propagation algorithm (Definition 6). For local scores, the *score calculator* considers subscribers' weights (e.g., reputation) using the *subscribers repository* (Definition 4).

- The *social flow analyzer* includes 3 sub-modules: *builder*, *miner*, and *displayer*. The *builder* parses the content of the *social actions repository* (1.2.1) to generate the necessary social flows enriched with scores and transmit the enriched social flows to the *miner* for further analysis (1.2.2). The *miner* performs 3 types of metrics and one analysis discussed in Section 6.2. The *displayer* visualizes real-time social flows along with the obtained analysis transmitted by the *miner* (1.2.3) on a browser showing how GreenUtility promotes its services to customers and seeks their feedback through Facebook (Fig. 2). The *displayer* uses Cytoscape graph-theory library for analysis and visualization [11]. It also highlights with assistance of the *miner* the relevant social actions that form the social flows. Different shapes are used to differentiate the social actions: star for post, rectangle for comment, hexagon for reply and triangle for reaction. The *displayer* also uses colors to emphasize whether a social action is positive, negative, or neutral so that a manager can easily identify the points of interest. GreenUtility uses the different flows to identify what social actions that (un)known subscribers have executed over its Facebook page. This could lead into reviewing BPs if their feedback were deemed relevant. Finally, the *displayer* is developed in HTML 5 and JavaScript while the *builder* and *miner* are developed as PHP programs and deployed on an Apache Web server.

Besides installing Facebook Graph API, companies interested in using Social Miner do not require any additional installation or configuration to track their campaigns on Facebook.

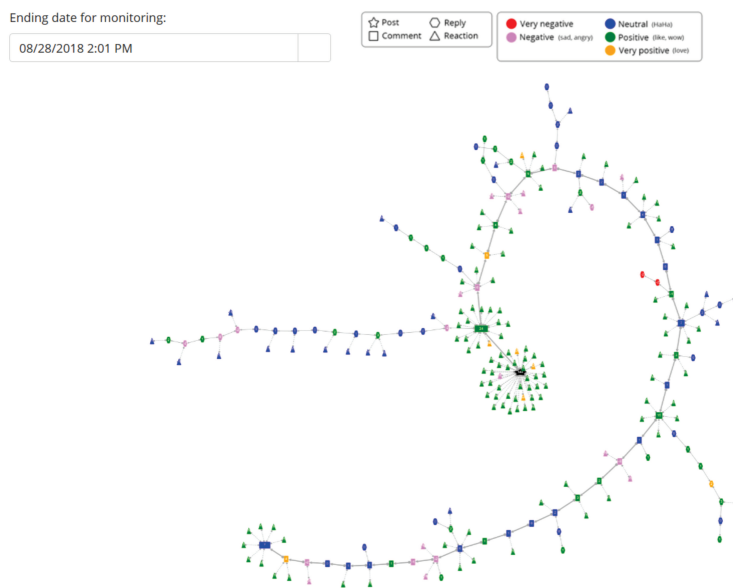


Fig. 5. Example of social flow

6.2. Experiments

To evaluate the benefits of using SM to decision makers, we carried out many experiments associated with a real campaign known as “*we are announcing Universe 11 plus, the greatest phone ever*” that was active from March 11, 2018 to March 16, 2018 on Facebook. The metrics that result out of these experiments are discussed below and assessed over one-day long time intervals. These metrics are implemented in the *miner* sub-module, part of the *social flow analyzer* (Fig. 4).

1. *Campaign attractiveness* metric (M_1 , Fig. 6) defines how appealing a campaign was to respondents by tracking their positive, neutral, and negative responses over different time intervals. Formally, M_1 is defined by Equation 4. By considering attractiveness, managers could extend campaigns, for example.

$$M_1(t_i) = \frac{new(t_i)}{returning(t_i) + new(t_i)} \quad (4)$$

Where t_i is a certain time interval [from, to] that could be days, weeks, months, etc., $new(t_i)$ is the number of new respondents who executed some social actions during t_i , and $returning(t_i)$ is the number of returning respondents who first, executed some social actions during t_i and second, were included in the previous time interval (t_{i-1}) that was used for defining the attractiveness metric. At t_0 , all respondents are treated as new. We rely on the time-stamped authorization relation $STR_{Facebook}^s$ (Definition 3) to compute $new(t_i)$ and $returning(t_i)$.

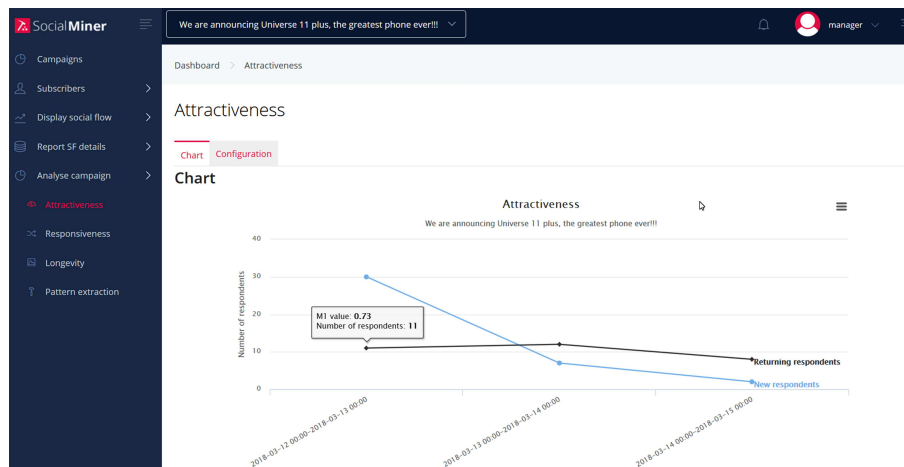


Fig. 6. Chart associated with campaign attractiveness

Since M_1 enables a campaign’s manager to discuss attractiveness from a global perspective, the focus is on the number of (new and returning) respondents. It would be useful for the manager to study attractiveness from a local perspective by identifying

the social actions that led for instance, to a major increase/decrease in the number of new/returning respondents both compared to previous interval times. To this end, we define *local-attractiveness* metric (M'_1) with focus on returning respondents' responsiveness levels (Equation 5):

$$M'_1(sa, t_i, t_{i+1}) = \frac{\text{returning}(sa, t_i, t_{i+1})}{\text{returning}(sa, t_{i-1}, t_i)} \quad (5)$$

Where sa is a certain social action that is subject to analysis, t_{i-1}, t_i, t_{i+1} are 3 homogeneous time intervals such that t_i happened earlier than t_{i+1} , $\text{returning}(sa, t_i, t_{i+1})$ is the number of returning respondents who were "new" at t_i and executed any action that came after sa during t_{i+1} , and finally, $\text{returning}(sa, t_{i-1}, t_i)$ is the number of returning respondents who were "new" at t_{i-1} and executed any action that came after sa during t_i . Similar to $\text{new}(t_i)$ and $\text{returning}(t_i)$, $\text{returning}(sa, t_i, t_{i+1})$ and $\text{returning}(sa, t_{i-1}, t_i)$ are computed based on $STR_{Facebook}^s$. Since a social action may appear many times in the time interval t_{i+1} , the manager points the desired social action that he wishes to analyze using first occurrence, for example.

2. *Campaign responsiveness* metric (M_2 , Fig. 7) indicates how a campaign is perceived by the community of respondents based on their feedback, whether positive (supportive), negative (opponent), or neutral. We rely on the local score functions (Definition 4) to formally define M_2 with focus on positive feedback in Equation 6:

$$M_2(t_i) = \frac{|sa|_{positive}}{|sa|_{positive} + |sa|_{neutral} + |sa|_{negative}} \quad (6)$$

Where t_i is a certain time interval [from, to], $|sa|_{positive}$ is the number of social actions executed during t_i such that $\text{sign}(GS(sa, t))$ is positive ($t \in t_i$), $|sa|_{negative}$ is the number of social actions executed during t_i such that $\text{sign}(GS(sa, t))$ is negative, and $|sa|_{neutral}$ is the number of social actions executed during t_i such that $LS(sa)$ is zero.

3. *Campaign longevity* metric (M_3 , Fig. 8) indicates how long a campaign remained "alive/active", i.e., respondents have continuously (without "big" gaps) provided feedback on the campaign so that a certain activity level over Facebook for example, is maintained. The longevity metric refers to the longitudinal dispersion of the provided feedback over a certain time interval t_k ([from, to]) that shall fall into a certain accepted activity level set by the campaign's manager (e.g., minimum number of actions in a day). We define this activity level with respect to a standard deviation (σ) upon which the decision of putting the campaign either on hold (suspend) or offline (stop). We rely on the time-stamped authorization relation $STR_{Facebook}^s$ to formally define this standard deviation as per Equation 7:

$$\sigma = M_3(t_k) = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}} \quad (7)$$

where t_k is a time interval [from, to] sliced into n equal time intervals (t_i) (e.g., in days and in weeks), x_i is the number of social actions executed during the slice

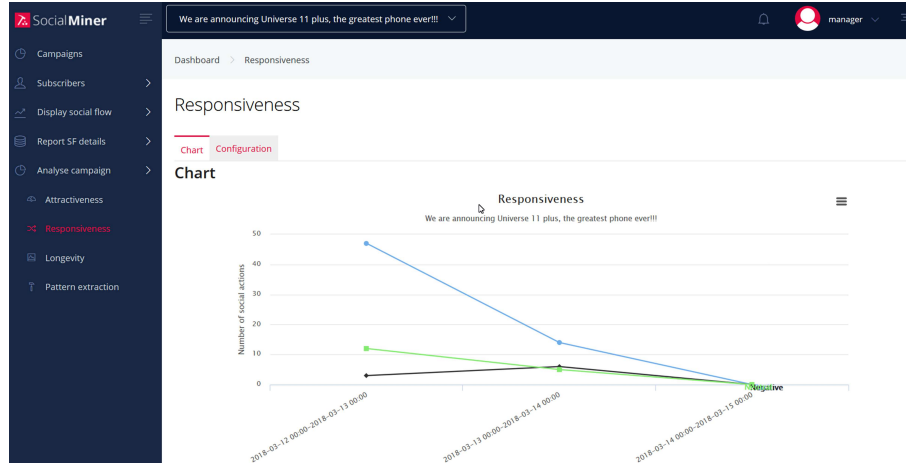


Fig. 7. Chart associated with campaign responsiveness

$t_i \subset t_k$ where $i \in [1, n]$, and \bar{x} is the average number of all social actions executed during the time interval $t_k (\sum_{i=1}^n x_i/n)$. The number of social actions over a certain time interval remains “acceptable” if this number falls into $[\bar{x} - \sigma, \bar{x} + \sigma]$. “Acceptable” could be defined over time and by benchmarking different time intervals together. We, thus, analyze the longevity metric according to \bar{x} and σ so that this metric’s time space corresponds to the set of time intervals where the number of social actions is declared “acceptable”.

In addition to the different metrics that SM produces, a *reversal trend analysis* of a campaign is implemented to help a manager identify the reasons behind a change in a campaign’s perception, for example. This perception could be based on a series of positive/negative and then negative/positive feedback. SM relies on the social flows’ secondary caterpillars to look for potential patterns such as 2 consecutive positive feedback followed by 3 consecutive negative feedback, and so on. To achieve the reversal analysis, we adopted gSpan algorithm for mining labeled graphs [37]. This algorithm uses a set of graphs \mathbf{D} and the minimum frequency (i.e., number of subgraphs before claiming that these subgraphs are repetitive) as inputs. In our case, \mathbf{D} could be a set or portions of secondary caterpillars. Because many social actions can be executed over time, we “cleaned” the secondary caterpillars from irrelevant time intervals (e.g., those where the campaign is inactive) for quality purposes. In Fig. 9, the minimum support threshold is 3 as an example, and all the social actions until August 28, 2018 are considered when extracting the patterns.

7. Discussions

From a business point of view, according to Zion Market Research, “*Global enterprise 2.0 technologies market expected to reach around USD 14,955 million by 2024, at a*

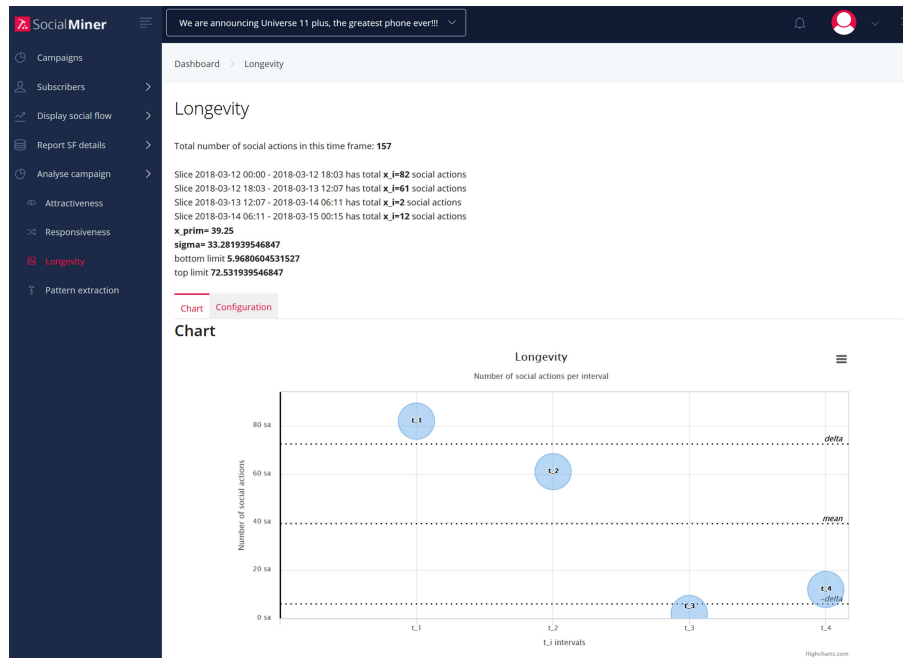


Fig. 8. Chart associated with campaign longevity

*CAGR approximately 46.87 for the forecast period from 2018 to 2024*⁹. Despite this heavy investment and the social fever that has caught every single activity of people's daily lives, many are still reluctant to embracing social media whether for personal use or for business use. Different concerns are continuously raised, including whether social media is bringing any value-added to companies. Gartner clearly states that "... many large companies are embracing internal social networks, but for the most part, they are not getting much from them" [15]. And, social media is also seen as the source of new forms of security threats, privacy breaches, and distraction to employees [10]. Contrarily to these "skeptical" views, a London-based think tank, Demos, encourages companies to allow their employees to embrace social network applications in order to establish and foster contacts with stakeholders such as colleagues, customers, and suppliers [12]. Striking the right balance between social media's pros and cons requires strict guidelines that social media users should comply with. Burégio et al. present such guidelines from 3 perspectives known as technology, organization, and management [6]. The technology perspective identifies the appropriate type of social media that should sustain a company growth and fall into its mission. The organization perspective puts in place the necessary procedures that should ensure an efficient use of social media to avoid misuses, for example. Finally, the management perspective identifies the metrics (or key performance

⁹ www.zionmarketresearch.com/report/enterprise-technologies-market.., last visited November, 14, 2019

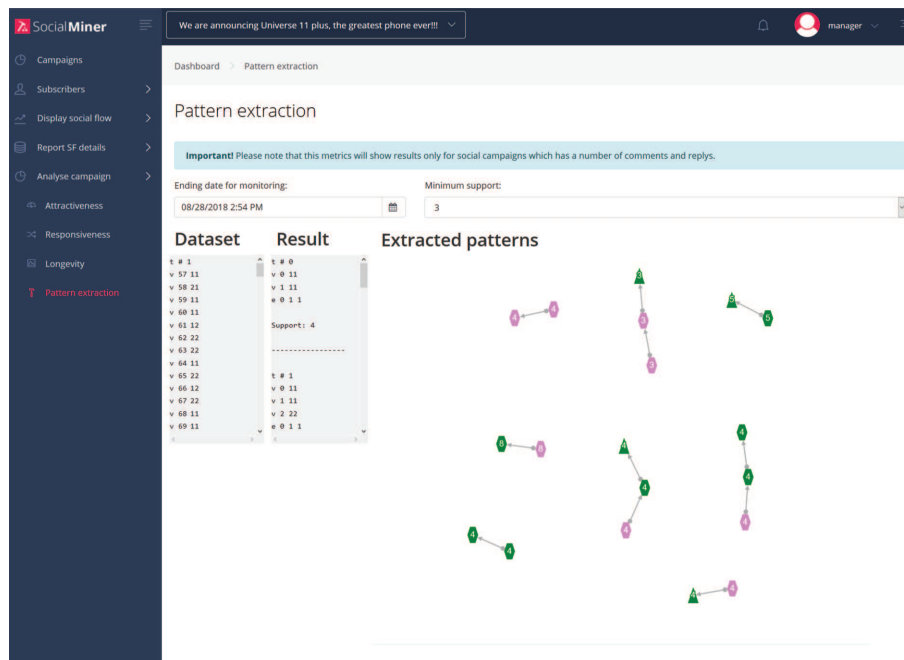


Fig. 9. Pattern extraction during a campaign

indicators) that should permit to evaluate the efficient use of social media based on some tangible benefits.

8. Conclusion

We presented an approach for developing flows in the context of companies that wish to tap into social media's opportunities. The flows are specialized into control and social. The former consists of tasks that form business processes. The latter consists of social actions that are executed over social media in response to specific events. Social flows are enriched with scores based on sentiment analysis so that companies would secure a better understanding of what-happened and what-might-happen in their ecosystems. For validation purposes, we developed a tool, Social Miner, on Facebook allowing to track and mine users' posts, comments, responses, etc. The system permits to answer questions like what actions were frequently executed, why certain actions were executed more than others, and when such actions were executed. In term of future work we would like to examine the impact of contextual factors on the next social actions to execute and the deployment of Social Miner on another social media such as Twitter.

References

1. van der Aalst, W., Song, M.: Mining Social Networks: Uncovering Interaction Patterns in Business Processes. In: Proceedings of the Second International Conference on Business Process

- Management (BPM'2004). Potsdam, Germany (2004)
2. Archak, N., Ghose, A., Ipeirotis, P.: Show me the Money!: Deriving the Pricing Power of Product Features by Mining Consumer Reviews. In: Proceedings of the 13th ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD'2007). San Jose, CA, USA (2007)
 3. Batrinca, B., Treleven, P.: Social Media Analytics: A Survey of Techniques, Tools, and Platforms. *AI Soc.* 30(1) (2015)
 4. Bonchi, F.: Influence Propagation in Social Networks: A Data Mining Perspective. *IEEE Intelligent Informatics Bulletin* 12(1) (2011)
 5. Bonchi, F., Castillo, C., Gionis, A., Jaimes, A.: Social Network Analysis and Mining for Business Applications. *ACM TIST* 2(3) (2011)
 6. Burégio, V.A., Maamar, Z., Meira, S.L.: An Architecture and Guiding Framework for the Social Enterprise. *IEEE Internet Computing* 19(1) (2015)
 7. Cross, R., Gray, P., Cunningham, S., Showers, M., Thomas, R.: The Collaborative Organization: How to Make Employee Networks Really Work. *IEEE Engineering Management Review* 39(1) (2011)
 8. Decker, K., Lesser, V.: Generalizing the Partial Global Planning Algorithm. *Int. J. Cooperative Inf. Syst.* 1(2) (1992)
 9. Dolev, S.: Self-stabilization. MIT press (2000)
 10. Faci, N., Maamar, Z., Burégio, V., Ugljanin, E., Benslimane, D.: Web 2.0 Applications in the Workplace: How to Ensure their Proper Use? *Computers in Industry* 88 (2017)
 11. Franz, M., Lopes, C., Huck, G., Dong, Y., Sümer, S., Bader, G.: Cytoscape.js: A Graph Theory Library for Visualisation and Analysis. *Bioinformatics* 32(2) (2016)
 12. Howe, J.: Crowdsourcing: Why the Power of the Crowd is Driving the Future of Business. *Crown Business* (2009)
 13. Jalali, A.: Aspect Mining in Business Process Management. In: Proceedings of the 13th International Conference on Perspectives in Business Informatics Research (BIR'2014). Lund, Sweden (2014)
 14. Kajan, E., Faci, N., Maamar, Z., Loo, A., Pljaskovic, A., Sheng, Q.Z.: The Network-based Business Process. *IEEE Internet Computing* 18(2) (2014)
 15. Kanaracus, C.: Gartner: Social Business Software Efforts Largely Unsuccessful, For Now (2013), <https://tinyurl.com/y5lc9eg4>
 16. Liu, B.: Sentiment Analysis and Subjectivity. Taylor and Francis Group, Boca (2010)
 17. Lu, Y., Zhai, C.: Opinion Integration Through Semi-supervised Topic Modeling. In: Proceedings of the 17th International Conference on World Wide Web (WWW'2008). Beijing, China (2008)
 18. Luo, Z., Osborne, M., Wang, T.: Opinion Retrieval in Twitter. In: Proceedings of the 6th International AAAI Conference on Weblogs and Social Media (ICWSM'2012). Dublin, Ireland (2012)
 19. Maamar, Z., Burégio, V.A., Faci, N., Benslimane, D., Sheng, Q.Z.: "Controlling" Web 2.0 Applications in the Workplace. In: Proceedings of the 19th IEEE International Conference on Enterprise Distributed Object Computing (EDOC'2015). Adelaide, Australia (2015)
 20. Maamar, Z., Burégio, V., Sellami, M., Rosa, N.S., Peng, Z., Subin, Z., Prakash, N., Benslimane, D., Silva, R.: Bridging the Gap Between the Business and Social Worlds: A Data Artifact-Driven Approach. *T. Large-Scale Data- and Knowledge-Centered Systems* 35 (2017)
 21. Maamar, Z., Faci, N., Kajan, E., Sakr, S., Boukhebouze, M., Barnawi, A.: How to Make Business Processes "Socialize"? *EAI Endorsed Trans. Indust. Netw. & Intellig. Syst.* 2(5) (2015)
 22. Maamar, Z., Faci, N., Sellami, M., Boukadi, K., Yahya, F., Barnawi, A., Sakr, S.: On Business Process Monitoring using Cross-Flow Coordination. *Service Oriented Computing and Applications* 11(2) (2017)

23. Martinho, D., Rito Silva, A.: Non-intrusive Capture of Business Processes Using Social Software. In: Daniel, F., Barkaoui, K., Dustdar, S. (eds.) Business Process Management Workshops help in conjunction with BPM'2012. Tallinn, Estonia (2012)
24. OMG: Business Process Model and Notation (BPMN) version 2.0. Tech. rep., Object Management Group (2011)
25. OpenKnowledge: Social Business Process Reengineering. Tech. rep., OpenKnowledge (2012), <http://socialbusinessmanifesto.com/social-business-process-reengineering>
26. Pang, B., Lee, L.: Opinion Mining and Sentiment Analysis. *Foundations and Trends in Information Retrieval* 2(1-2) (2008)
27. Pourshahid, A., Amyot, D., Peyton, L., Ghanavati, S., Chen, P., Weiss, M., Forster, A.J.: Business Process Management with the User Requirements Notation. *Electronic Commerce Research* 9(4) (2009)
28. Reichert, M., Dadam, P.: ADEPT_{flex}-Supporting Dynamic Changes of Workflows without Losing Control. *J. Intell. Inf. Syst.* 10(2) (1998)
29. Sadiq, S.W., Orlowska, M.E., Sadiq, W., Foulger, C.: Data Flow and Validation in Workflow Modelling. In: Proceedings of the Fifteenth Australasian Database Conference (ADC'2004). Dunedin, New Zealand (2004)
30. Soboroff, I., Ounis, I., Macdonald, C.: Overview of the TREC-2008 Blog Track (Jan 2010)
31. Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning, C.D., Ng, A.Y., Potts, C.: Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank. In: Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing (EMNLP'2013). Seattle, Washington, USA (2013)
32. Tang, J., Chang, Y., Liu, H.: Mining social media with social theories: A survey. *SIGKDD Explor. Newsl.* 15(2) (2014)
33. Tang, J., Sun, J., Wang, C., Yang, Z.: Social influence analysis in large-scale networks. In: Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'2009). Paris, France (2009)
34. Ting, S., Ip, W., Tsang, A.: Is Naive Bayes a Good Classifier for Document Classification? *International Journal of Software Engineering and Its Applications* 5(3) (2011)
35. Turban, E., Bolloju, N., Liang, T.: Enterprise Social Networking: Opportunities, Adoption, and Risk Mitigation. *J. Org. Computing and E. Commerce* 21(3) (2011)
36. Ugljanin, E., Faci, N., Sellami, M., Maamar, Z.: Tracking Users' Actions over Social Media: Application to Facebook. In: Proceedings of the 25th IEEE International Conference on Enabling Technologies: Infrastructure for Collaborative Enterprises (WETICE'2016). Paris, France (2016)
37. Yan, X., Han, J.: gSpan: Graph-Based Substructure Pattern Mining. In: Proceedings of the 2002 IEEE International Conference on Data Mining!(ICDM'2002). Maebashi City, Japan (2002)

Ejub Kajan is an Associate Professor of Computer Science at State University of Novi Pazar, Serbia. His research interests include social computing, interoperability, service computing, business process management and IoT. He has a PhD in computer science from University of Nis, Serbia. He has published with almost all major publishers like ACM, Elsevier, IEEE, Springer, Wiley, etc. He is a Senior Member of ACM.

Noura Faci is an Associate Professor at the University Lyon 1, France since October 2008. Her current research interests are service computing, social computing, business process management, and IoT. She has published several papers in high-quality journals

and conferences and actively contributes to the IEEE TSC, IEEE IC, and Computer journal's review process, as well.

Zakaria Maamar is a Professor in the College of Information Technology at Zayed University, Dubai, United Arab Emirates. His research interests include Web services, social networks, and context-aware computing. He has a PhD in computer science from Laval University, Quebec City, Canada.

Mohamed Sellami (http://www-public.imtbs-tsp.eu/sellam_m/) is an Associate Professor of Computer Science at Telecom SudParis member of Institut Polytechnique de Paris and the SAMOVAR Laboratory (Evry, France). He received his PhD in computer science from Telecom SudParis in 2011. His main research interests are related to service computing and cloud computing. His publication list includes international journals and conferences and he serves on the program and organizing committees of numerous international conferences and workshops.

Emir Ugljanin is a PhD student at the University of Nis. He has MA degree from Technical faculty at the University of Novi Sad and engineer's degree from the State University of Novi Pazar. His current research interests are Internet of Things, social web and business process management.

Hamamache Kheddouci is full Professor in Computer Science at Lyon 1 University since 2004. He received his PhD degree in Computer Science from Paris XI University in 1999. In 2003, he obtained his research supervision habilitation in Computer Science from the Burgundy University, Dijon. His research interest includes combinatorial and algorithmic aspects of graphs and their applications, in particular, in big data and social networks. For more details, see: <http://perso.univ-lyon1.fr/hamamache.kheddouci/>

Dragan Stojanović is a Professor in Computer Science, at the Faculty of Electronic Engineering, University of Niš, Serbia. His research interests include Big Data processing and analytics, spatio-temporal and mobility data management, as well as mobile and ubiquitous systems and IoT in Smart Cities. He has published widely in those and related topics.

Djamal Benslimane (<http://www710.univ-lyon1.fr/dbenslim/>) is currently a Professor at Lyon 1 University in France. He is the head of the computer science department at the IUT Lyon 1. His research interests lie in the areas of Services computing and database interoperability. His recent works were published in IEEE Transaction on Knowledge and Data Engineering, IEEE Internet Computing, IEEE Transactions on Services Computing, IEEE Transactions on Systems, Man, and Cybernetics, Communications of the ACM, ACM Transactions on Internet Computing, ACM Transactions on Software Engineering and Methodology, and WWW Journal.

Received: August 20, 2019; Accepted: February 20, 2020.

Land-use classification via ensemble dropout information discriminative extreme learning machine based on deep convolution feature

Tianle Zhang¹, Muzhou Hou¹, Tao Zhou^{2*}, Zhaode Liu², Weirong Cheng³, and Yangjin Cheng⁴

¹ School of Mathematics and Statistics, Central South University
410083 Changsha, China
houmuzhou@sina.com, csuztl@csu.edu.cn

² School of Economics, Guangdong University of Finance and Economics
510320, Guangzhou, China
zhoutaoscut@hotmail.com*, lzhaode@163.com
*Corresponding author

³ School of Mathematics and Science, Peking University
100871 Beijing, China
chenweirong@pku.edu.cn

⁴ School of Mathematics and Computational Science Xiangtan University
411105, Xiangtan, Hunan, China
yjcheng@xtu.edu.cn

Abstract. Classifying land-use scenes with high quality and accuracy is an important research direction in current hyperspectral remote sensing images, which is conducive to scientific management and utilization of land. An effective classifier and feature extractor can improve classification stability and accuracy. Therefore, based on deep learning technique, a dropout-based ensemble learning method is proposed in this paper, which combines convolutional neural network (CNN) and information discriminating extreme learning machine (IELM). Pre-trained CNN is used to learn effective and robust features, and deep convolution features are fed to the IELM classifier. Then the adoption of dropout technique and ensemble method can improve generalization capabilities and stability. The effectiveness of the proposed algorithm is tested by hyperspectral remote sensing image classification experiments. The experimental results show that the proposed E-CNN-dropIELM has achieved satisfactory results compared to state-of-the-art methods in terms of classification accuracy and stability.

Keywords: Extreme learning machine (ELM), Convolutional neural network (CNN), Dropout, Ensemble, Land-use classification.

1. Introduction

Nowadays remote sensing imaging technology and deep learning have developed rapidly, and now high-resolution remote sensing (HRRS) images have become a hot research field. Its most important research goal is to classify high-resolution remote sensing images, for instance, agricultural, airplane, beach, etc. It is of great importance in many remote sensing applications, such as land resource management, urban development and nature

conservation. However, there are lots of challenges in the land-use scenes classification such as insufficient training sample data, extracting effective classification features, and improving classifier accuracy.

In recent decades, in order to overcome these problems, a large number of previous work use machine learning methods to improve the accuracy of land use classification, such as support vector machine (SVM) [28], manifold learning (ML) [2], random forests (RF) [24], artificial neural network (ANN) [11], unsupervised learning (UL), etc. Recently, many methods have made very good progress in the field of remote sensing image classification. Most of the previous works mainly focused on classifying by low-level features. Lowe extracted distinctive invariant features from images [16]. Oliva, et. used the gist descriptor to achieve probable semantic category of scenes [18]. Ren, et. proposed a novel scheme for local binary patterns (LBP) to tackle pixel correlation [23]. Although these low-level features have proven their validity for scene classification, these features are not ideal due to the inability to reflect semantic information [18,23,37]. Many researchers have developed statistical methods to process shallow features for mid-level features. Bag-of-visual-words (BOVWs) model is the most famous approach, which has been great successful in scene classification [33,38]. Of course, there are some other ways to get the semantic information of the images, such as locality-constrained linear coding (LLC) [15], vector of locally aggregated descriptors (VLAD) [10], improved fisher kernel (IFK) [22], et. However, the semantic information of hand-crafted features is insufficient. In the field of image recognition, convolutional neural network (CNN) as one of deep learning models has achieved successful achievements and rapid development [26,35,12,13].

Recently, CNN which can extract discriminative and robust features, has been applied to land-use image classification [8,19]. Features extracted from the pretrained CNN show outstanding performance in scene classification[21,1]. So instead of using these CNN architectures directly as the final classifiers, many researchers use pre-trained CNNs as feature extractors and combined it with effective and simple classifiers. The Extreme Learning Machine (ELM) is a powerful learning algorithm based on single hidden layer feedforward neural networks (SLFNs). It obtains the output weight by solving the least squares solution of the SLFNs while the input weight of the hidden node is randomly generated. At the same time the ELM can avoid the possibility of slow convergence and local optimal solution of back propagation algorithm (BP). Because of its very efficient and impressive approximation and generalization capabilities, ELM has been applied to many areas [6,34,36,17,27].ELM is applied to land cover classification as a classifier [20,14,29]. ELM achieves better classification results compared with BP and SVM, and its computational time is much smaller than BP and SVM. Yan, et. proposed information discriminative extreme learning machine (IELM) to overcome the shortcoming that ELM is insufficient with limited hyperspectral remote sensing images [31]. Qian, et. proposed a model for land-use scene classification by integrating CNN and constrained extreme learning machine (CELM) [30]. Zhu, et. integrated CNN and KELM, and KELM is used to enhance the discriminative ability of classifier by kernel function [39]. While the effectiveness of such features and classifiers have been verified, some papers consider ensemble methods. Actually, it is able to provide a more effective way to solve the land-use classification problem. Although, the ensemble algorithm is very effective in the general pattern recognition tasks[5,3,4], it is still a challenging task to land use classification in

practical applications. In order to make full use of effective deep convolution features and overcome the limitations of traditional classifiers and fully connected layer (FCL), a new framework for land-use scene classification is proposed in this paper. First, CNN is trained with a large land-use scene image dataset. The pre-trained CNN removing FCL is used as a feature extractor to learn deep and discriminative features. Then, the IELM-based basic classifiers with dropout technology are built, and excellent classification result is obtained via majority voting. To sum up, the following contributions have been made in this paper:

- i A new framework for land-use scene classification is present combining deep convolutional features by CNN with ensemble classifier via majority voting.
- ii By introducing the dropout method, IELM-based basic classifiers with different subsets of discriminative features are constructed. Such a classification framework has better generalization ability and can be effectively for extensive applications.
- iii The experimental results on the remote scene image land-use classification show that E-CNN-dropIELM can improve the classification performance compared to the several the-state-of-art techniques.

The rest of this paper is organized as follows. Section 2-5 introduces the related background knowledge of IELM and CNN, et. Section 6 describes in detail the proposed E-CNN-dropIELM approach. Experiments and the results are presented in Section 7, while some concluding remarks are drawn in Section 8.

2. The pretrained CNN

CNN is considered to be one of the most successful structures for deep learning due to its good performance. A typical CNN consists of two parts. One is the feature extraction part, which is usually alternately connected to the pooling layers by alternating convolution layers. The output of the upper layer is the input of the next layer. The shallow layer extracts relatively specific features, and the deep layer often extracts relatively abstract features by features learned in previous layers. The other part is the classifier, fully-connected layer. And a typical feed-forward neural network uses a Softmax layer as a classifier to calculate the likelihood of each class. Parameters of CNN are usually optimized through a classical stochastic gradient descent algorithm based on the backpropagation. A typical CNN architecture for land-use scene classification in Fig. 1.

The deep convolution features learned by the pre-training CNN are effective for the classification of land use scenarios, referring to the previous papers [7]. Therefore, the pre-training architecture of CNN provided by MatConvNet toolbox is used, which is pre-trained on Imagenet [25] and has similar structure with AlexNet. The architecture used in this paper is shown in Fig. 2, and it consists of five convolution layers and three fully-connected layer. One pooling layer is connected behind each convolution layer. We remove full-connection layers and use the remaining CNNs as a feature extractor.

In order to achieve better performance in a small sample dataset, this paper uses a series of random transforms to augment images for obtaining better features to increasing accuracy. The formulas are as follows,

$$\mathbf{I}' = s(\lambda)\mathbf{R}(\theta)\mathbf{I} + \mathbf{T}(\delta) \quad (1)$$

where \mathbf{I} is the input image matrix. $s(\lambda)$ is the scaling factor. $\mathbf{R}(\theta)$ is the rotation matrix. $\mathbf{T}(\delta)$ is the shift matrix. Then bilinear interpolation is used to resize the image.

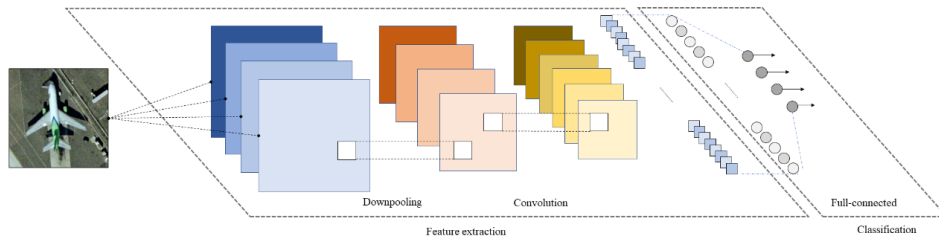


Fig. 1. A typical CNN architecture for land-use scene classification.

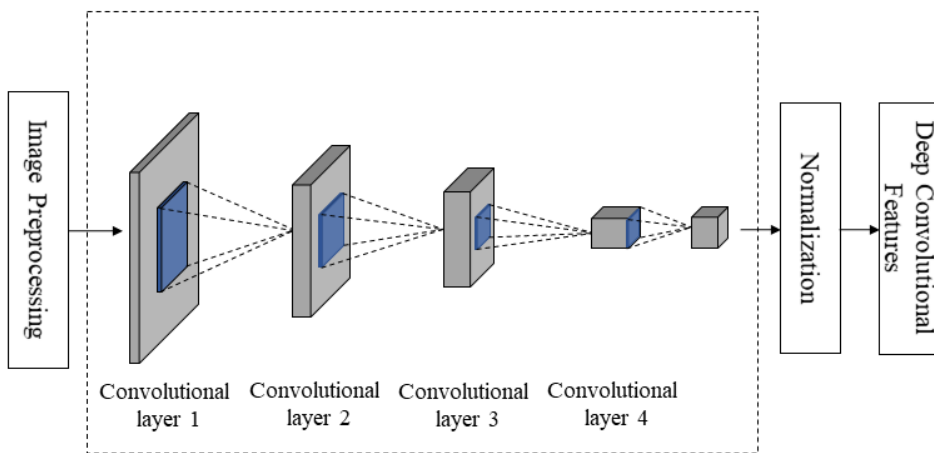


Fig. 2. The CNN architecture used in this paper, and the max-pool layers are not displayed.

3. Dropout method

Dropout means that neurons are randomly removed during network training, and it is an effective way to improve network's generalization capabilities. As shown in Fig. 3, removing a neuron equal to drop it temporarily from the network with all its input and output connections. Choosing which neuron to be removed is random, and usually each neuron is preserved with a fixed probability independent of other neurons. We introduce the basic idea of dropout by linear neurons. Given an input vector $\mathbf{I} = (I_1, I_2, I_3, \dots, I_n)$

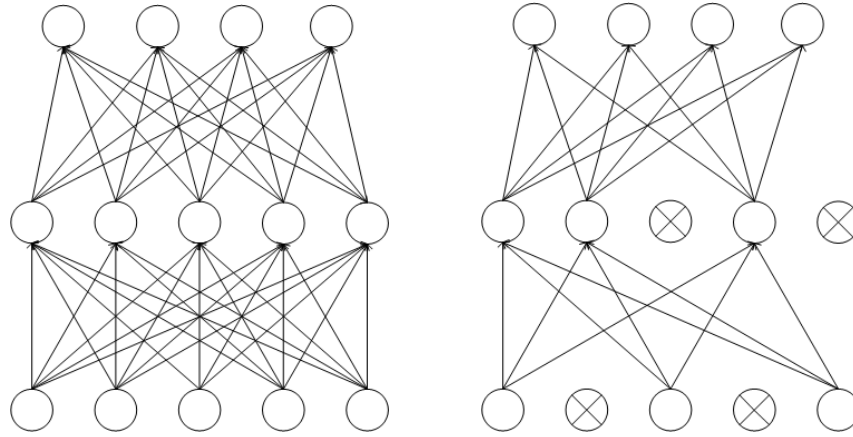


Fig. 3. The diagram of basic idea of the proposed algorithm.

, the output of the linear neuron is $S(\mathbf{I}) = \sum_{i=1}^n \omega_i I_i$. If I_i is deleted with a uniform distribution, or deleted with an equal probability of 0.5, 2^n sub-networks will be generated including the empty network. The average of outputs of all sub-networks is:

$$E(S) = \frac{1}{2^n} \sum S(\omega_i, I) \tag{2}$$

where $S(\omega_i, I)$ means a set all sub-networks.

Suppose $\delta_i (1 \leq i \leq n)$ is a stochastic variable satisfying Bernoulli distribution and they are independent with each other. $p_i = P(\delta_i = 1), q_i = P(\delta_i = 0) = 1 - p_i$. If remove $I_i (1 \leq i \leq n)$ or $\omega_i (1 \leq i \leq n)$ with probability q_i , then the output of the linear element can be expressed as follow:

$$S(\mathbf{I}) = \sum_{i=1}^n \omega_i \delta_i I_i \tag{3}$$

According to the linear property by mathematical expectation,

$$E(S) = \sum_{i=1}^n \omega_i E(\delta_i) I_i = \sum_{i=1}^n \omega_i p_i I_i \tag{4}$$

If $p_1 = p_2 = \dots = p_n = p$, Eq. (4) can be written as

$$E(S) = \sum_{i=1}^n \omega_i E(\delta) I_i = \sum_{i=1}^n \omega_i p I_i \tag{5}$$

If there is a bias in the linear neuron, output of the linear neuron can be written as:

$$S(\mathbf{I}) = \sum_{i=1}^n \omega_i \delta_i I_i + b \delta_b \tag{6}$$

Similarly,

$$E(S) = \sum_{i=1}^n \omega_i p_i I_i + b p_b \tag{7}$$

where $p_b = P(\delta_b = 1)$.

4. ELM algorithm

The extreme learning machine is a novel learning algorithm based on the structure of single-hidden layer feedforward networks (SLFNs) that was proposed by G B Huang et al. in 2006 and then extended to “generalized” SLFNs where the hidden layer neurons need not be neuron alike [9]. The hidden neuron’s parameters are randomly generated without iterative tuning and are independent of the data. The network topology of ELM is shown in Fig. 4.

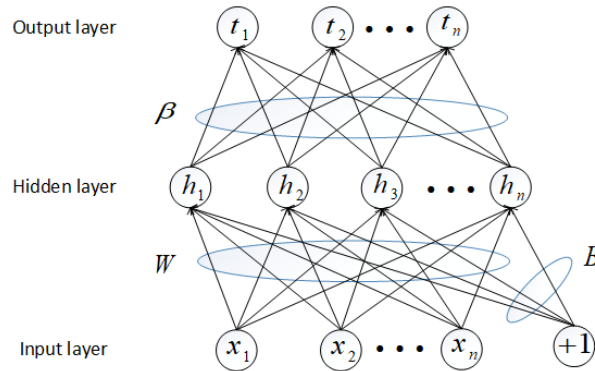


Fig. 4. The network topology of ELM.

Given a training set with N distinct samples $\Gamma = \{(\mathbf{x}_j, \mathbf{t}_j) | \mathbf{x}_j \in \mathbb{R}^n, \mathbf{t}_j \in \mathbb{R}^m, \mathbf{j} = 1, \dots, N\}$, where X is the predictive variable with dimension n and t is the objective variable with dimension m . The mathematical model of SLFNs with M hidden neurons and an activation function $G(\cdot)$ are described as

$$\sum_{i=1}^n \beta_i G(\mathbf{w}_i, b_i, \mathbf{x}_j), j = 1, 2, \dots, N \tag{8}$$

where $b_i \in R$ is the randomly assigned bias of the i th hidden node and $\mathbf{w}_i \in R$ is the randomly assigned input weight vector connecting the i th hidden node to the output node. $G(\mathbf{w}_i, b_i, \mathbf{x}_j)$ is the output of the i th hidden node with respect to the input sample \mathbf{x}_j .

When the SLFN is completely approximating the data, that is, when the error between the output \hat{t}_i and the actual t_i is zero, the relationship is:

$$\sum_{i=1}^n \beta_i G(\mathbf{w}_i, b_i, \mathbf{x}_j) = t_j, j = 1, 2, \dots, N \quad (9)$$

Eq. (9) can be written as

$$\mathbf{H}\beta = \mathbf{T} \quad (10)$$

where

$$\mathbf{H} = [\mathbf{h}_1^T \ \mathbf{h}_2^T \ \dots \ \mathbf{h}_N^T] = \begin{bmatrix} G(\mathbf{w}_1, b_1, \mathbf{x}_1) & \dots & G(\mathbf{w}_n, b_n, \mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ G(\mathbf{w}_1, b_1, \mathbf{x}_N) & \dots & G(\mathbf{w}_n, b_n, \mathbf{x}_N) \end{bmatrix}_{N \times n} \quad (11)$$

$$\beta = \begin{bmatrix} \beta_1^T \\ \beta_2^T \\ \vdots \\ \beta_N^T \end{bmatrix}_{n \times m}, T = \begin{bmatrix} t_1^T \\ t_2^T \\ \vdots \\ t_N^T \end{bmatrix}_{N \times m} \quad (12)$$

The output weights can be computed by finding the least-square solutions of the SLFNs in Eq. (10), which is given as

$$\beta = \mathbf{H}^\dagger \mathbf{T} \quad (13)$$

where \mathbf{H}^\dagger is the Moore-Penrose generalized inverse of the matrix \mathbf{H} . If $\mathbf{H}^\dagger \mathbf{H}$ is not a singular matrix, then, Eq. (13) can be written as

$$\beta = \mathbf{H}^\dagger \mathbf{T} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{T} \quad (14)$$

In addition, we add an l -2 regularization constraint to the loss function of ELM which can improve the generalization and robustness. Eq. (14) becomes the following formula:

$$\beta = \begin{cases} \mathbf{H}^T (\frac{I}{C} + \mathbf{H}^T \mathbf{H})^{-1} \mathbf{T} & \text{if } N \leq L \\ (\frac{I}{C} + \mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{T} & \text{if } N > L \end{cases} \quad (15)$$

where C assigns the regularization, I assigns the unit matrix, N assigns the number of samples and L assigns the number of hidden layer neurons.

5. IELM algorithm

The paper [32] proposed a kind of regularized extreme learning machine based on discriminative information (IELM). For the classification problem, IELM considers the geometric features and the discriminant information of data samples, and optimizes the output

weight of the extreme learning machine by maximizing the heterogeneous dispersion and minimizing the same kind of dispersion. It can improve accuracy and generalization of the traditional ELM.

Given N distinct samples $\Gamma = \{(\mathbf{x}_j, \mathbf{t}_j) | \mathbf{x}_j \in \mathbb{R}^n, \mathbf{t}_j \in \mathbb{R}^m, j = 1, \dots, N\}$, S_B is the inter-class scatter matrix, S_W is the within-class scatter matrix. They are described as:

$$S_B = \sum_{i=1}^C \sum_{j=1}^{N_c} (x_j - u^i) (x_j - u^i)^T \quad (16)$$

$$S_W = \sum_{i=1}^C (u^i - u) (u^i - u)^T \quad (17)$$

where C is the number of categories of data samples, u_i is the mean of i th class, and u is the mean of data samples. Formation difference matrix can be described as:

$$S = S_B - (1 - \eta)S_W \quad (18)$$

where η is a parameter from 0 to 1. The parameter η plays the role of adjusting the intra-class and the within-class discriminant information. So IELM makes good use of the discriminant information in the data and improves the classification ability of ELM.

The loss function of IELM is described as:

$$\begin{aligned} \min_{\beta} \quad & \frac{1}{2} \beta^T S \beta + \frac{1}{2} C \sum_{i=1}^N \varepsilon^2 \\ \text{s.t.} \quad & \sum_{i=1}^n \beta_i G(\mathbf{w}_i, b_i, \mathbf{x}_j) - \mathbf{t}_j = \varepsilon_i, j = 1, 2, \dots, N \end{aligned} \quad (19)$$

where \mathcal{E} assigns training error.

The Lagrange function of Eq. (19) is:

$$L = \frac{1}{2} \beta^T S \beta + \frac{1}{2} C \sum_{i=1}^N \varepsilon^2 + \sum_{i=1}^N \sum_{j=1}^m (\alpha_{ij} \beta_j G(\mathbf{x}_i) - \mathbf{t}_{ij} + \varepsilon_{ij}) \quad (20)$$

where β_j is the output weight of j th output node, and the KKT condition is:

$$\frac{\partial L}{\partial \beta_j} = 0 \rightarrow \beta_j = \sum \alpha_{ij} G(x_i)^T \rightarrow WS = H^T \alpha \quad (21)$$

$$\frac{\partial L}{\partial \varepsilon_j} = 0 \rightarrow \alpha_i = C \varepsilon_i, i = 1, \dots, N \quad (22)$$

$$\frac{\partial L}{\partial \alpha_j} = 0 \rightarrow G(x_i) \beta - \mathbf{t}_i^T + \varepsilon_i^T = 0 \quad (23)$$

where $\alpha_i = [\alpha_{i1}, \dots, \alpha_{im}]^T$, $\alpha = [\alpha_1, \dots, \alpha_n]$

According to the above, Eq. (15) becomes the following formula:

$$\beta = \begin{cases} \mathbf{h}^r \left(\frac{S}{C} + \mathbf{h}^r \mathbf{h} \right)^{-1} S \mathbf{T} & \text{if } N \leq L \\ \left(\frac{S}{C} + \mathbf{H}^T \mathbf{h} \right)^{-1} \mathbf{h}^T \mathbf{T} & \text{if } N > L \end{cases}$$

And the output of IELM is:

$$g(x) = G(x)\beta = \begin{cases} G(x) \left(\mathbf{H}^T \left(\frac{S}{C} + \mathbf{H}^T \mathbf{H} \right)^{-1} S \mathbf{T} \right) & \text{if } N \leq L \\ G(x) \left(\frac{S}{C} + \mathbf{H}^T \mathbf{H} \right)^{-1} \mathbf{H}^T \mathbf{T} & \text{if } N > L \end{cases}$$

6. The proposed E-CNN-dropIELM algorithm

In this section, we introduce the basic idea of the proposed E-CNN-dropIELM. The accuracy of image classification is mainly affected by two parts, feature extractor and classifier. Firstly, this method uses CNN as a feature extractor, which is necessary because effective features can improve the performance of the classification. And then add dropout method to IELM, which considers the inter-class distribution and discriminant information of the data. And neurons are randomly removed, so the IELM is changing for each training. These trained IELM networks are used as basic classifiers by repeatedly “dropout” the neurons. Therefore, all basic classifiers work on different feature subspaces. Finally, these basic classifiers used for data classification are ensembled by simple majority voting method, and it can effectively avoid over-fitting of a local feature by the network.

The diversity and complementarity of basic classifiers have a great influence on the performance of the ensembled classification. Diversity is the error diversity of basic classifiers, which means the samples misclassified by basic classifiers should be different. If the samples classified by each basic classifier are different, the diversity of basic classifiers is best. Obviously, such basic classifiers are also the best complementarity. In other words, classifiers can learn from each other when classifying the sample. The typical method of constructing basic classifiers is to divide the training set into several subsets of data, and then train basic classifiers with these subsets of data. Because these basic classifiers are trained with different subsets of data, basic classifiers trained in this way have better diversity. In this paper, the method of constructing basic classifiers is different from the typical method, and use the Dropout technique to construct the basic classifier, IELM. These basic classifiers are trained with the same dataset, but the network topology and feature subspace of basic classifiers are different. Naturally, when classifying the samples, they have good diversity and can better learn from each other.

The proposed E-CNN-dropIELM algorithm can be described as follow:

7. Remote sensing image land-use classification experiments

7.1. The UC-Merced land use Dataset

In this paper, we investigate the performance of the proposed E-CNN-dropIELM on the UC-Merced (UCM) data set [33]. The data set contains high-resolution remote sensing images of 21 different land-use classes, such as airplanes, beaches, etc. Each class has

Algorithm 1 Ensemble dropout information discriminative extreme learning machine based on deep convolutional features for image classification.

Input:

$T = \{(\mathbf{x}_j, \mathbf{t}_j) | \mathbf{x}_j \in \mathbb{R}^n, \mathbf{t}_j \in \mathbb{R}^m, j = 1, \dots, N\}$, training set; T , testing set; l , the number of basic classifiers; C , regularization parameter; η , discriminant information parameter

Output:

j^* , the class label of $\mathbf{x} \in T$

- 1: Extract feature via a fix CNN removed the last layer softmax;
 - 2: Initialize a larger IELM with m hidden nodes;
 - 3: // Train l basic classifiers: $L = L_1, L_1, \dots, L_l$;
 - 4: **for all** ($i = 1; i \leq l; i = i + 1$) **do**
 - 5: Generate a m -dimensional binary random vector whose elements are drown independently from a Bernoulli distribution;
 - 6: Dropout some hidden nodes with the generated m -dimensional binary random vector, obtain a IELM _{i} denoted by L_i ;
 - 7: Train L_i with IELM and soft-maximize its outputs, obtain a probability distribution $(p_{i1}(\mathbf{x}), p_{i2}(\mathbf{x}), \dots, p_{ik}(\mathbf{x}))$;
 - 8: **end for**
 - 9: // Test $\mathbf{x} \in T$ via basic classifiers: $L = L_1, L_1, \dots, L_l$;
 - 10: **for all** ($\mathbf{x} \in T$) **do**
 - 11: Extract feature via a fix CNN removed the last layer softmax;
 - 12: **for all** ($i = 1; j \leq l; j = j + 1$) **do**
 - 13: Calculate label vector $(p_{i1}(\mathbf{x}), p_{i2}(\mathbf{x}), \dots, p_{ik}(\mathbf{x}))$ via L_i with IELM;
 - 14: **end for**
 - 15: Calculate $p_{j^*}(\mathbf{x}) = \max_{1 \leq j \leq k} \sum_{i=1}^l p_{ij}(\mathbf{x})$ via majority voting;
 - 16: **end for**
 - 17: **return** j^* ;
-

100 images and their size is 256 x 256 pixels. The sample images of 12 classes in UCM is shown in Fig. 5.



Fig. 5. Image examples from UCM land use data set.

All samples were resized to 227 x 227 pixels because input image of the pre-trained CNN requires. 4096-dimensional vector is fed into basic classifiers IELMs as feature.

7.2. Experiment and results

In order to evaluate the performance of the proposed algorithm, we compare it with the classical algorithms such as SVM, ELM and some successful algorithms in mentioned papers in UCM dataset [37,35,21,29,7]. In this paper, all the experiments were performed 30 times using MATLAB R2017b on a computer with 2.60 GHz CPU and 16.0 GB RAM.

In this experiment, 70% of each class is used as training set and the rest is used as testing set. We average accuracies of 30 times experiments on testing set and use grid search method to obtain the optimal values. The optimal regularization parameter C is 2^{11} by employing $C = 2^x, x = -5, \dots, 30$, and the optimal discrimination information pa-

parameter η and the optimal dropout probability p are 0.5 and 0.9 from a set $\{0, 0.1, \dots, 1\}$. Impacts of the proposed E-CNN-dropIELM by different factors are shown in Fig. 6.

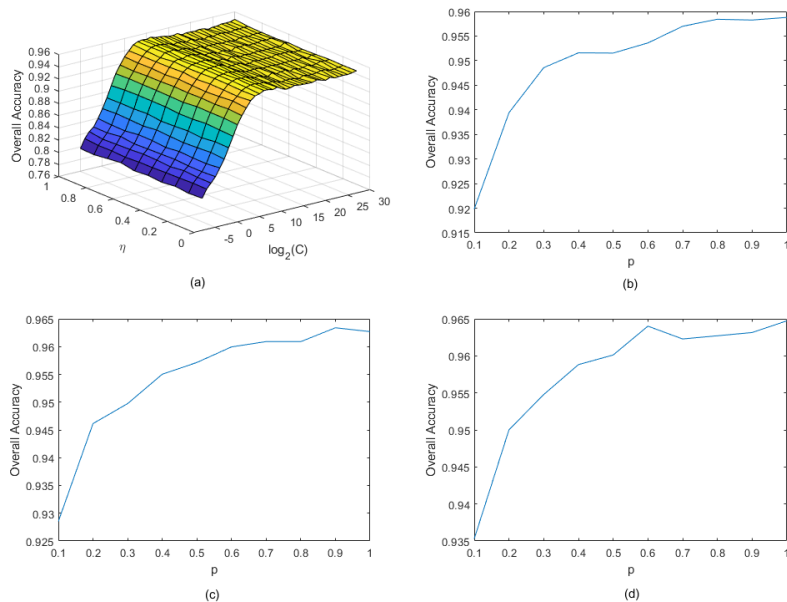


Fig. 6. Impacts on OA of E-CNN-dropIELM by different factors. (a) Overall Accuracy by discrimination information parameter η and regularization parameter C (b) Overall Accuracy by dropout probability p when the number of basic classifiers is 3. (c) Overall Accuracy by dropout probability p when the number of basic classifiers is 5. (d) Overall Accuracy by dropout probability p when the number of basic classifiers is 7.

In order to test the performance of the proposed classifier, compare the performance of the proposed algorithm, IELM, ELM, SVM and KNN with different numbers of training sets on CNN features. Fig. 7 shows the proposed E-CNN-dropIELM has better performance on the CNN features, in contrast to these traditional algorithms and the accuracy is higher with the increase of training data.

The confusion matrix of the classification results shows that the most error-prone categories are buildings and river, and the most confused categories are buildings and beach, medium residential and overpass, as shown in Fig. 8. But actually, we can find the proposed method has resulted in more than 96% accuracy for most classes. We compare the best results with the reported classification accuracy of some existing methods on the UCM Land Use dataset. The results are shown in Table 1. As expected, the high classification accuracy of the proposed algorithm indicates that our method performs primarily, and the reason is twofold. First, the features extracted by the pretrained CNN are sufficiently discriminative. On the other hand, IELMs with the random feature subspace constructed via the dropout method are used as basic classifiers, and then effectively ensembles these basic classifiers to achieve satisfactory result. In addition, to verify the classification accuracy of the proposed algorithm, experimental results show that the proposed method outperform some existing method proposed in [37,35,21,29,7].

Table 1. Classification accuracy (%) of existing methods and the proposed method on the UCM Land Use data set

Method	Accuracy(%)
MCMI-LBP [37]	88.20
CNN-SVM [21]	93.42
CaffeNet [7]	94.43
AlexNet [7]	94.37
VGG-S [7]	94.60
RCNet [35]	94.53
CNN-ELM [29]	95.62
proposed	96.90

8. Conclusion

In this paper, a novel land scene classification framework is proposed. This framework uses the pre-trained CNN with the fully connected layer removed as the feature extractor, and then it allows IELM with dropout method to be the basic classifier. Finally, the excellent classification result is obtained via majority voting, which effectively utilizes these basic classifiers based on different sub-feature spaces. The performance of the proposed E-CNN-dropIELM for land-use scene classification is proved on the UCM Land Use database. The experimental results show that the proposed algorithm can achieve more competitive results compared with the algorithm mentioned above paper. The properties of the proposed hybrid model is twofold: One is the proposed model uses pre-trained CNN to automatically extract high-level features, while most other traditional methods rely heavily on hand-designed low-level and mid-level features. The other is that the proposed model employs IELM with dropout as the basic classifiers, and then effectively ensembles these basic classifiers via majority voting. This is able to improve generalization and accuracy, and also can avoid sample shortages and overfitting problems.

Acknowledgments. This study was funded by the Major project of National Social Science Foundation (15zdc012), and by the National Natural Science Foundation of China (Grant number 61375063, 61271355, 11301549 and 11271378) and by the Fundamental Research Funds for the Central Universities of Central South University (2018zzts322).

References

1. Castelluccio, M., Poggi, G., Sansone, C., Verdoliva, L.: Land use classification in remote sensing images by convolutional neural networks. arXiv preprint arXiv:1508.00092 (2015)
2. Chen, Y., Crawford, M.M., Ghosh, J.: Applying nonlinear manifold learning to hyperspectral data for land cover classification. In: Proceedings. 2005 IEEE International Geoscience and Remote Sensing Symposium, 2005. IGARSS'05. vol. 6, pp. 4311–4314. IEEE (2005)
3. Dietterich, T.G.: Ensemble methods in machine learning. In: International workshop on multiple classifier systems. pp. 1–15. Springer (2000)
4. Dietterich, T.G.: An experimental comparison of three methods for constructing ensembles of decision trees: Bagging, boosting, and randomization. *Machine learning* 40(2), 139–157 (2000)
5. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
6. Hou, M., Zhang, T., Weng, F., Ali, M., Al-Ansari, N., Yaseen, Z.M.: Global solar radiation prediction using hybrid online sequential extreme learning machine model. *Energies* 11(12), 3415 (2018)
7. Hu, F., Xia, G.S., Hu, J., Zhang, L.: Transferring deep convolutional neural networks for the scene classification of high-resolution remote sensing imagery. *Remote Sensing* 7(11), 14680–14707 (2015)
8. Hu, F., Xia, G.S., Zhang, L.: Deep sparse representations for land-use scene classification in remote sensing images. In: 2016 IEEE 13th International Conference on Signal Processing (ICSP). pp. 192–197. IEEE (2016)
9. Huang, G.B., Zhu, Q.Y., Siew, C.K.: Extreme learning machine: theory and applications. *Neurocomputing* 70(1-3), 489–501 (2006)
10. Jegou, H., Perronnin, F., Douze, M., Sánchez, J., Pérez, P., Schmid, C.: Aggregating local image descriptors into compact codes. *IEEE transactions on pattern analysis and machine intelligence* 34(9), 1704–1716 (2011)
11. Kadavi, P.R., Lee, C.W.: Land cover classification analysis of volcanic island in aleutian arc using an artificial neural network (ann) and a support vector machine (svm) from landsat imagery. *Geosciences Journal* 22(4), 653–665 (2018)
12. Karpathy, A., Fei-Fei, L.: Deep visual-semantic alignments for generating image descriptions. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3128–3137 (2015)
13. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: Advances in neural information processing systems. pp. 1097–1105 (2012)
14. Li, Y., Xie, W., Li, H.: Hyperspectral image reconstruction by deep convolutional neural network for classification. *Pattern Recognition* 63, 371–383 (2017)
15. Liu, G., Liu, Y., Guo, M., Liu, P., Wang, C.: Non-negative locality-constrained linear coding for image classification. In: International Conference on Intelligent Science and Big Data Engineering. pp. 462–471. Springer (2015)
16. Lowe, D.G.: Distinctive image features from scale-invariant keypoints. *International journal of computer vision* 60(2), 91–110 (2004)
17. Muzhou, H., Ming, C., Yangchun, Z.: A self-organizing mixture extreme leaning machine for time series forecasting. In: Proceedings of ELM-2014 Volume 1, pp. 225–236. Springer (2015)

18. Oliva, A., Torralba, A.: Modeling the shape of the scene: A holistic representation of the spatial envelope. *International journal of computer vision* 42(3), 145–175 (2001)
19. Othman, E., Bazi, Y., Alajlan, N., Alhichri, H., Melgani, F.: Using convolutional features and a sparse autoencoder for land-use scene classification. *International Journal of Remote Sensing* 37(10), 2149–2167 (2016)
20. Pal, M.: Extreme-learning-machine-based land cover classification. *International Journal of Remote Sensing* 30(14), 3835–3841 (2009)
21. Penatti, O.A., Nogueira, K., Dos Santos, J.A.: Do deep features generalize from everyday objects to remote sensing and aerial scenes domains? In: *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*. pp. 44–51 (2015)
22. Perronnin, F., Sánchez, J., Mensink, T.: Improving the fisher kernel for large-scale image classification. In: *European conference on computer vision*. pp. 143–156. Springer (2010)
23. Ren, J., Jiang, X., Yuan, J.: Learning lbp structure by maximizing the conditional mutual information. *Pattern Recognition* 48(10), 3180–3190 (2015)
24. Ruiz Hernandez, I.E., Shi, W.: A random forests classification method for urban land-use mapping integrating spatial metrics and texture analysis. *International journal of remote sensing* 39(4), 1175–1198 (2018)
25. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. *International journal of computer vision* 115(3), 211–252 (2015)
26. Sharif Razavian, A., Azizpour, H., Sullivan, J., Carlsson, S.: Cnn features off-the-shelf: an astounding baseline for recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*. pp. 806–813 (2014)
27. Sun, H., Hou, M., Yang, Y., Zhang, T., Weng, F., Han, F.: Solving partial differential equation based on bernstein neural network and extreme learning machine algorithm. *Neural Processing Letters* 50(2), 1153–1172 (2019)
28. Ustuner, M., Sanli, F.B., Dixon, B.: Application of support vector machines for landuse classification using high-resolution rapideye images: a sensitivity analysis. *European Journal of Remote Sensing* 48(1), 403–422 (2015)
29. Weng, Q., Mao, Z., Lin, J., Guo, W.: Land-use classification via extreme learning classifier based on deep convolutional features. *IEEE Geoscience and Remote Sensing Letters* 14(5), 704–708 (2017)
30. Weng, Q., Mao, Z., Lin, J., Liao, X.: Land-use scene classification based on a cnn using a constrained extreme learning machine. *International journal of remote sensing* 39(19), 6281–6299 (2018)
31. Yan, D., Chu, Y., Li, L., Liu, D.: Hyperspectral remote sensing image classification with information discriminative extreme learning machine. *Multimedia Tools and Applications* 77(5), 5803–5818 (2018)
32. Yan, D., Chu, Y., Zhang, H., Liu, D.: Information discriminative extreme learning machine. *Soft Computing* 22(2), 677–689 (2018)
33. Yang, Y., Newsam, S.: Bag-of-visual-words and spatial extensions for land-use classification. In: *Proceedings of the 18th SIGSPATIAL international conference on advances in geographic information systems*. pp. 270–279 (2010)
34. Yang, Y., Hou, M., Luo, J., Liu, T.: Neural network method for lossless two-conductor transmission line equations based on the ielm algorithm. *AIP Advances* 8(6), 065010 (2018)
35. Zhang, L., Zhang, L., Du, B.: Deep learning for remote sensing data: A technical tutorial on the state of the art. *IEEE Geoscience and Remote Sensing Magazine* 4(2), 22–40 (2016)
36. Zhang, T., Hou, M., Weng, F., Yang, Y., Sun, H., Wang, Z., Gao, Z., Luo, J.: An online learning algorithm for voice activation detection based on a pretrained online extreme learning machine. In: *Proceedings of the 2nd International Conference on Computer Science and Application Engineering*. pp. 1–8 (2018)

37. Zhao, F., Sun, H., Liu, S., Zhou, S.: Combining low level features and visual attributes for vhr remote sensing image classification. In: MIPPR 2015: Remote Sensing Image Processing, Geographic Information Systems, and Other Applications. vol. 9815, p. 98150C. International Society for Optics and Photonics (2015)
38. Zhao, L.J., Tang, P., Huo, L.Z.: Land-use scene classification using a concentric circle-structured multiscale bag-of-visual-words model. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* 7(12), 4620–4631 (2014)
39. Zhu, X., Li, Z., Zhang, X.Y., Li, P., Xue, Z., Wang, L.: Deep convolutional representations and kernel extreme learning machines for image classification. *Multimedia Tools and Applications* 78(20), 29271–29290 (2019)

Tianle Zhang is a graduate student in the School of Mathematics and Statistics of Central South University. His research interests lie in applied statistics and deep learning.

Muzhou Hou received the B.S. degrees from Hunan Normal University, Changsha, China, M.S. and Ph.D. degree from Central South University, Changsha, China in 1985, 2002, and 2009, respectively. He is currently a Professor and doctoral advisor in School of Mathematics and Statistics, Central South University, PR China. His current research interests include machine learning, computational intelligence, neural networks, data mining and numerical approximation.

Tao Zhou is received the B.S. degrees from Hunan Normal University, Changsha, China, M.S. degree from Central South University, Changsha, China in 2002, 2010 respectively. and a Ph.D.degree from South China University of Technology, Guangzhou, Guangdong, China, in 2018. He is currently a teacher in School of Economics, Guangdong University of Finance & Economics, Guangzhou, Guangdong, China. His current research interests include machine learning, computational intelligence, neural networks, numerical approximation and Management science and Real options.

Zhaode Liu received a M.S. degree from Chongqing Normal University, Chongqing, China, in 2002 and a Ph.D.degree from South China University of Technology, Guangzhou, Guangdong, China, in 2013.candidate at Guangdong University of Finance & Economics, Guangzhou, Guangdong, China, in 2002. he is an associate professor in the Faculty of Economics, Guangdong University of Finance & Economics, Guangzhou, Guangdong, China. His research in interests lie in Comprehensive evaluation, Management science and Real options.

Weirong Chen is a graduate student in the School of Mathematics Sciences, Peking University. Her research interests include machine learning and natural language processing.

Yangjin Cheng is currently a Professor in School of Mathematics and Computational Science Xiangtan University, PR China. His current research include interests computational intelligence, neural networks and Uncertainty optimization theory and its application in Finance.

Received: December 22, 2019; Accepted: March 20, 2020.

Study of cardiac arrhythmia classification based on convolutional neural network

Yonghui Dai¹, Bo Xu¹, Siyu Yan^{2,*}, and Jing Xu²

¹ Management School, Shanghai University of International Business and Economics, Shanghai 201620, China

daiyonghui@suike.edu.cn, brianxubo@163.com

² School of Information Management and Engineering, Shanghai University of Finance and Economics, Shanghai 200433, China

ysy18956@163.com, 1053517221@qq.com

Abstract. Cardiovascular disease is one of the diseases threatening the human health, and its diagnosis has always been a research hotspot in the medical field. In particular, the diagnosis technology based on ECG (electrocardiogram) signal as an effective method for studying cardiovascular diseases has attracted many scholars' attention. In this paper, Convolutional Neural Network (CNN) is used to study the feature classification of three kinds of ECG signals, which including sinus rhythm (SR), Ventricular Tachycardia (VT) and Ventricular Fibrillation (VF). Specifically, different convolution layer structures and different time intervals are used for ECG signal classification, such as the division of 2-layer and 4-layer convolution layers, the setting of four time periods (1s, 2s, 3s, 10s), etc. by performing the above classification conditions, the best classification results are obtained. The contribution of this paper is mainly in two aspects. On the one hand, the convolution neural network is used to classify the arrhythmia data, and different classification effects are obtained by setting different convolution layers. On the other hand, according to the data characteristics of three kinds of ECG signals, different time periods are designed to optimize the classification performance. The research results provide a reference for the classification of ECG signals and contribute to the research of cardiovascular diseases.

Keywords: Cardiac arrhythmia classification, convolutional neural network, ECG signals, segment size.

1. Introduction

In recent years, with the improvement of information technology and medical technology, the diagnosis technology of some diseases has been improved, which provides a better guarantee for people's health. Among the major diseases that threaten human health, cardiovascular disease has always been at the forefront of easily-occurring diseases and it is extremely threatening. SCD (Sudden Cardiac Death) as a typical cardiovascular disease, it causes many people to face life threatening or even lose their lives. Its cause is closely related to VT (Ventricular Tachycardia) and VF (Ventricular Fibrillation). For example, it is about 390,000 people die of malignant arrhythmias each year in the United States, and it is about 540,000 people die of sudden cardiac death in China each year [19]. When VT signal appears in heart, doctors should observe and diagnose it in time. Otherwise,

VT is easy to induce VF. If the patient does not defibrillate in time when VF occurs, it may cause damage to cardiac tissue and function and even death of myocardial tissue [1]. Therefore, the diagnosis of arrhythmia has always been a research hotspot in the medical community. Because the electrocardiogram (ECG) can clearly reflect the health of the heart. When the waveform or value of the ECG signal changes abnormally, it indicates that the body is not very healthy [3]. Therefore, Electrocardiogram-based automatic analysis and diagnosis technology has been adopted by many hospitals as an effective method for determining heart disease. At present, ECG arrhythmia signal classification and recognition are regarded as an important basis for the diagnosis of cardiovascular disease, and this method greatly reduces the workload of doctors, which has important practical value for the timely diagnosis and treatment of heart disease.

Because of the advantages of computers in data calculation, processing, and analysis, many computer applications and analysis tools are used in the medical field. Computer-aided diagnosis is used as an important way of medical diagnosis. For example, medical image recognition based on computer image processing, disease prediction based on computer algorithm, etc. ECG signal reflects objective indicators of people's heart function, and it is an important reference basis for the judgment and diagnosis of cardiovascular and cerebrovascular diseases. The ECG signal has a huge amount of data, and these data are non-linear data, which makes it difficult to use. In order to find rules from these data, it needs a series of data processing to realize. Considering the good adaptability of convolution neural network in nonlinear aspect, this paper uses convolution neural network to study the classification of ECG signal, and studies three kinds of arrhythmia classification of sinus rhythm, ventricular tachycardia and ventricular fibrillation.

This paper is organized as follows. In section 2, related work of cardiac arrhythmia classification is introduced. In section 3, the methodology of convolution neural network technology is shown, which including convolution neural network structure and model for cardiac arrhythmia classification. In section 4, experiment and analysis are illustrated. In section 5, the conclusion of this paper is described.

2. Literature review

In the research of arrhythmia, scholars mainly focus on the preprocessing and classification of ECG data. For example, one-dimensional convolutional neural network was used to study arrhythmia detection, and a six-layer convolutional neural network structure was adopted, which including the number of convolutional layers, pooling layers, and forward fully connected layers are two [8]. The overall accuracy rate of this detection method reached 97%, but the sensitivity of recognition of supraventricular ectopic heart-beat was 60.3%, and the positive predictive value was 63.5%. Guo (2017) et al studied simple convolution neural network for image classification, different parameters were adjusted to achieve good image classification, and they analyzed for different optimization algorithms for the influence of learning rate and optimal parameters on image classification [5]. Thereafter, Some scholars studied five types of the arrhythmias and proposed a new method [17]. In their study, wavelet transform and mean filter were utilized to remove noise in the data preprocessing stage, and the label output by the convolutional neural network was combined with the db6 wavelet coefficient in the feature extraction phase. The functions were normalized before entering libsvm. ECG signals were applied

to represent the heartbeat, and The experimental results show that the new method they mentioned is better than the existing method, which is a suitable approach to increase the accuracy of classification. Qiu (2018) et al studied five types of arrhythmia, and they put forward a new method for data pre-processing. In order to decrease the noise, mean filter and wavelet transform were applied, and 250-point signal of ECG were selected to appear the heartbeat, and an accuracy of 98.79% was achieved after the data preprocessing [16].

3. Methodology

3.1. Convolutional Neural Network

Nowadays, deep learning has been successfully applied in various fields. Especially the convolutional neural network, it has been widely used in image recognition and speech recognition, such as image classification, face recognition, speech emotion recognition, video analysis as an important content of deep learning [6][18]. Convolutional neural networks are made up of many neurons connected to each other. After each neuron accepts a linear combination of inputs, it begins with a simple linear weighting, and then adds a nonlinear activation function to each neuron for nonlinear transformation and output, it has a very good performance in terms of picture recognition and classification. [11]. Therefore, if there is a picture of a car, the thing that CNN needs to do is to extract the features on the picture, perform feature recognition, and determine which category and which specific thing the object in the picture belongs to through the established model. The connection between every two neurons represents a weight value, called a weight. Different weights and activation functions result in different outputs of the neural network [21].

Compared with convolution neural network, there are still some problems in full connection neural network [2] [7]. For example, If the image is expanded into vector information, it will lose some spatial information. It has too much parameters, and its efficiency will be affected, which increased the difficulty of training. At the same time, a large number of parameters will also lead to network over-fitting [15][22]. Relatively speaking, convolutional neural networks can solve the above problems well.

The structure of CNN includes three kind of layers, which including input layer, hidden layer and output layer [10]. Every different kinds of neural network have different hidden layers. In this structure, convolutional layer, pooling layer, relu layer and fully connected layer are included by the hidden layers of CNN. The structure is shown as in Fig. 1.

In the CNN structure, there are four main contents, the first content is padding [?]. Because the original image will become smaller after being convolved by the filter, the edge count of the image in the convolution is less. At this time, the edge information is easy to lose. In order to solve this problem, we can use the padding method. Before each convolution, we fill in a blank space around the image, so that the image is as large as the original after convolution, and the original edge is also calculated more times [9]. The second content is the step size. In general, the default step size is 1, but in fact, step size can be set to other values. The third content is pooling, which is to extract the main features of a certain area, and reduce the number of parameters to prevent the model from overfitting. When the merge is executed, there is both the maximum pooling and the average value

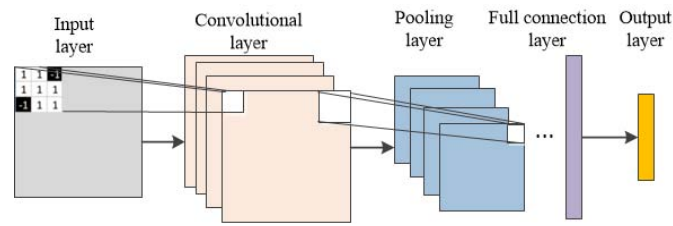


Fig. 1. The structure of CNN

of the area, which is the average pooling. The fourth content is the convolution of multi-channel pictures, color images, usually three channels of RGB, so the input data has three dimensions: (length, width, channel). For example, an image with 28×28 RGB, and the dimension is $(28, 28, 3)$. If the input image is three-dimensional such as $(8, 8, 3)$, and the dimension of the filter is $(3, 3, 3)$ and the last dimension will keep the same dimension as the dimension of input channel. At this time, the convolution is that the elements of the three channels are multiplied and summed.

Therefore, for example, if there is a picture of a car, the thing that CNN needs to do is to extract the features on the picture, perform feature recognition, and determine which category and which specific thing the object in the picture belongs to through the established model.

3.2. The CNN structure for cardiac arrhythmia classification

Because there are many types of arrhythmia, and each type has its own characteristics, which makes the diagnosis of arrhythmia is difficult. For example, the arrhythmia caused by a single irregular heartbeat is called morphological arrhythmia. The arrhythmia caused by a group of irregular heartbeat is called rhythmic arrhythmia. If the data of arrhythmia is detected manually, the whole detection process is not only very tedious and time-consuming, but also human errors may occur in the analysis and detection process. In order to solve the above problems, support vector machine, genetic algorithm, artificial neural network have been used in the detection and recognition of arrhythmia and have achieved good results. In our study, the CNN method was adopted to the classification method of cardiac arrhythmia. In order to get a satisfactory model, the parameters of convolutional neural network model are tried many times. An example of a CNN model's parameters is shown as in Fig. 2.

It can be seen from Fig. 2 that after defining the sequential model, the convolution layer `conv1d` is added, which is the layer convoluting one-dimensional data in the keras framework. The number of convolution cores is 64, and the number of filters is 1. In the input data, the length of one-dimensional array is 500, and the activation function is `relu()`. The use of `relu()` function can accelerate convergence, increase the sparsity of the network, and make the extracted data more representative. The pooling layer uses `maxpooling1d`, which reduces the dimension of data and prevents over coupling in training. In order to prevent the gradient of neural network from disappearing in the process of back propagation, the batch normalization layer is added to the BN layer to force each input value back to the standard normal distribution of 0 to 1 by some standardized means. The

Layer (type)	Output Shape	Param #
conv1d_27 (Conv1D)	(None, 1, 64)	32064
max_pooling1d_21 (MaxPooling1D)	(None, 1, 4)	0
batch_normalization_18 (Batch Normalization)	(None, 1, 4)	16
flatten_14 (Flatten)	(None, 4)	0
dense_36 (Dense)	(None, 16)	80
dropout_18 (Dropout)	(None, 16)	0
dense_37 (Dense)	(None, 3)	51
Total params: 32,211		
Trainable params: 32,203		
Non-trainable params: 8		

Fig. 2. Example of the structure parameters of CNN

flatten layer can process the input multidimensional data in one dimension. In the above model, two full connection layers and one dropout layer are added to prevent over fitting during training. The confidence rate of this model is about 92.5%. Its graph converges quickly. However, due to the existence of BN layer and flatten layer, the accuracy will be virtual high. It has the phenomenon of over fitting.

In order to get the best results, during the construction of the convolutional neural network, the subsequent layers are tested one by one. In addition, the number of convolution pooling layers is also changed, its value changed from two to four to verify which is the better result. The structure parameters of the modified convolution neural network are shown in Fig. 3 [20].

Layer (type)	Output Shape	Param #
conv1d_1 (Conv1D)	(None, 291, 100)	1100
max_pooling1d_1 (MaxPooling1D)	(None, 145, 100)	0
conv1d_2 (Conv1D)	(None, 136, 160)	160160
global_average_pooling1d_1 (Global Average Pooling1D)	(None, 160)	0
dropout_1 (Dropout)	(None, 160)	0
dense_1 (Dense)	(None, 3)	483
Total params: 161,743		
Trainable params: 161,743		
Non-trainable params: 0		

Fig. 3. The structure parameters of the modified CNN

Convolutional layer: The layer is defined by function Conv1D(), as a one-dimensional convolutional layer, Conv1D has a good recognition performance on sequence data. The number of convolution kernels were defined as 64, the step sizes are 10, and the padding variable is default values of 'same'. In this paper, the activation function of the convolution layer uses the relu() function, because the activation function can make the relationship between input and output close to any function, rather than a simple linear relationship. Moreover, the activation function can accelerate the convergence of neural network, increase the sparsity of neural network, and make the data extracted by deep neural network more representative. The relu function is shown as follows.

$$f(x) = x^+ = \max(0, x), \quad (1)$$

Relu function is a piecewise linear function, which changes all negative values to 0, while positive values remain unchanged. Pooling layer: The first layer is defined by function maxpooling1D(), the second pooling layer is defined by global-average-pooling1d. The maximum pooling is used in the paper, which is the most common method. It means that the maximum value will be the selected in a field. The pooling kernel size is 2. The step size is also default. Pooling is actually a process of subsample. Pooling can also improve the fault tolerance of neural networks and prevent overfitting.

Dropout layer: This layer is defined by function dropout. In order to prevent overfitting, this article refers to the research of previous scholars, and set up a dropout layer to discard 30% of the training data, and then put the remaining 70% of the training data into the second dense.

Connected layer and dropout layer: This layer is defined by function Dense(). In this paper, there is only one dense layer. The final output is three categories, which are SR, VT and VF. The overall network structure is shown as in Fig. 4.

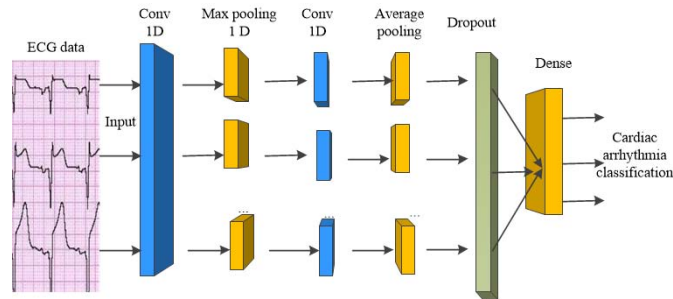


Fig. 4. The CNN structure for Cardiac arrhythmia classification

3.3. Confusion Matrix

Confusion matrix is defined by function confusion_matrix(). After the neural network training is completed, the confusion matrix was added to evaluate the accuracy of the neural network [4]. In this paper, the confusion matrix is a matrix of 3 rows and 3 columns.

Confusion matrix can reflect the accuracy of data classification from different sides. In addition, the confusion matrix can be used to compare the classification results with the average value of the trained neural network, and to compare the calculation to obtain the sensitivity and reliability.

4. Experiment and Analysis

Data pre-processing The data were taken from Physiobank and the extended American Heart Association Database (AHADB). The databases from Physiobank include the European ST-T Database (EDB), the Creighton University Ventricular Tachyarrhythmia Database (CUDB), the MIT-BIH Arrhythmia Database (MITDB), and the MIT-BIH Malignant Ventricular Arrhythmia Database (VFDB). From these databases only records containing examples of VT or VF were used, which is 98 records in total. All records were resampled to 100 Hz sampling rate and processed by a 0.5 Hz high pass FIR filter for removal of baseline wander. The ECG records were further normalized to unit sample variance [13].

Each piece of data that is not originally the same length, segment the data as the same length, and the excess data after division is discarded. In order to explore the influence of the division of different seconds on the result, data is divided into different lengths, which included 100, 200, 300, 1000 (i.e. 1s, 2s, 3s, 10s).

At the same time, in order to ensure the same training and testing for each type of data, we separately segment the segmentation data again, 80% of them are used as training data and 20% of them are used as testing data. And in order to ensure the sample balance, VF is the smallest number of samples, therefore, the large number of samples, SR and VT should be down-sampling, and the same amount of data as the VF is selected from SR and VT [12]. Finally, based on data segmentation, the data should be normalized. In this paper, gradient descent iteration method is used when training the model, the feature scaling function can ensure that our different data features are in a similar range, so that the model can converge more quickly during the gradient descent.

In the method of data normalization, the unit variance method is applied [14], and the specific formula is as follows.

$$x[n] = \frac{x[n]}{\sqrt{\frac{\sum_0^n |x[n]|^2}{N}}} \quad (2)$$

After normalizing the data using the unit variance method, the absolute values of all the data are in the interval [0, 1], and put the processed data into the input layer for training, which means the data preprocessing phase ends.

4.1. Data balance method

Class imbalances in data are often encountered when using machine learning algorithms, also known as class skew. For example, if you need to predict whether a patient has a rare disease, however, the percentage of positives in historical data may be only 1%. In this case, it is difficult to learn to get a good classifier, and in this case, the conclusion is

usually a big mistake. Some scholars came up with some data balancing method to deal with this situation.

The main methods to deal with data imbalance include that expanding the data set, sampling the data set, generate artificial data samples, different classification algorithms, etc. This paper uses the method of sampling data set to achieve data balance. Sampling includes under-sampling and over-sampling. Under-sampling is to reduce the number of groups with more data, and over-sampling is to increase the number of groups with less data to achieve the balance results. However, the former is usually less used in practical applications, in case of avoiding the data lose. For this reason, new data generation based on existing data is widely adopted, such as synthetic sampling, which is implemented by modifying the data distribution method. In this paper, according to the one-dimensional signal given by the data set and after drawing the time-domain image, the data of SR is much more than that of VT and VF. Moreover, because there is huge difference in data length, data preprocessing is also needed, which is processing samples with unbalanced data. In the data preprocessing stage, the down-sampling was used to balance data, and confusion matrix was utilized for detecting result. Down-sampling is a specific method to resolve sample imbalance, which is to extract the same amount of data from three types of data to ensure sample balance. Confusion matrices are a way to reflect sample results than single precision more closely.

4.2. Experiments and results

In order to classify the cardiac arrhythmia data, this paper tries to use two kinds of CNN architecture to find the best structure and get the best precision results. The first structure includes four convolutional and pooling layers, dropout layer and dense layer. The other includes two convolutional and pooling layers, dropout layer and dense layer. When building the neural network, it is necessary to adjust the neural network itself to obtain the best training effect, which is mainly achieved by controlling different training parameters. For example, Table 1 is the sample data training result of adjusting the different matrix parameter.

Table 1. The result of different matrix parameter

Time	Matrix	SR acc (%)	VF acc (%)	VT acc (%)	Total acc (%)
1s	5000	74.9	28.5	70.0	57.8
1s	10000	83.4	38.7	72.1	64.7
1s	15000	74.3	39.3	68.7	60.7
1s	30000	75.4	41.1	72.1	62.9
Remarks: Batch Size:157, Number of hidden unit:20, Learning rate:0.01					

It can be seen from Table 1 that the learning effect of the whole model is the best when the Matrix is 10,000. While keeping the Matrix, changing the learning Rate to explore the impact of different learning rates on the model. For the adjustment scheme, the increase 5% each time. If the accuracy is lower than the previous one, we will change the range of change to 50%. The results of test are shown as Table 2.

Table 2. The result of different Learning Rate

Time	Learning Rate	SR acc (%)	VF acc (%)	VT acc (%)	Total acc (%)
1s	0.0125	72.3	48.3	82.4	67.6
1s	0.0100	83.4	38.4	72.1	64.6
1s	0.0075	82.6	37.8	71.4	63.9
1s	0.015	\	\	\	\
1s	0.005	\	\	\	\
Remarks: Batch Size:157, Number of hidden unit:20, Matrix:10000					

From the previous parameter adjustment test, the parameters most suitable for the current CNN model are obtained. Next, we try to change the size of the input data to get the best results. The parameters are set to keep the Matrix at 10000 and the learning rate at 0.01. The test is carried out by changing the size of input data to 100, 200, 300, 400. The result is shown as Table 3.

Table 3. The result of different Learning Rate

Time	SR acc (%)	VF acc (%)	VT acc (%)	Total acc (%)
1s	83.4	38.7	72.1	64.7
2s	75.5	36.1	67.6	59.7
3s	62.9	38.6	65.6	55.7
10s	45.5	35.6	71.2	50.8
Average	66.8	37.3	69.1	57.7
Remarks: Batch Size:157, Number of hidden unit:20, Matrix:10000				

At the first time, the first structure includes four convolutional pooling layers, dropout layer and dense layer. Another one includes two convolutional pooling layers, dropout layer and dense layer. Table 4 is the accuracy of first structure including stacking four convolutional max pooling layers.

Table 4. The result of epoch 50 with four convolutional pooling layers

Time	SR acc (%)	VF acc (%)	VT acc (%)	Total acc (%)
1s	70.8	65.7	32.4	56.3
2s	67.7	63.2	30.6	53.8
3s	65.4	56.8	29.7	50.6
10s	30.9	28.3	10.8	23.3
Average	58.7	53.5	25.9	46.1
Remarks: Batch Size:128, epoch: 50				

Obviously, it shows that the average accuracy is very low, the classification performance is not good enough. The accuracy is around from 20% to 60%, which is not a perfect result. However, when the segment size is 1s, the sensitivity is much higher than the segment size is 10s. For example, the SR sensitivity dropped from 70.8% to 30.9%.

Therefore, it is obvious that the segment size will affect the accuracy of the final result. When the segment size is small, the accuracy will perform better. When the neural network was constructed, it is necessary to adjust the neural network itself to obtain the best training effect. By controlling different training parameters, we can obtain the accuracy of each type under different training levels, and then choose the optimal one. The final result is show as Table 5.

Table 5. The result of different epoch

Time	Epoch	SR acc (%)	VF acc (%)	VT acc (%)	Total acc (%)
1s	25	85.3	60.4	82.2	75.9
1s	50	90.7	81.7	92.0	88.1
1s	100	91.7	78.1	92.4	87.4
1s	300	86.6	71.4	79.3	79.1
Average		88.6	72.9	86.5	82.6
Remarks: Batch Size:128					

By adjusting the comparison of the number of epochs. When the epoch is 25, the model didn't converge apparently. As the increase of the epoch, the model gradually converged, the total accuracy also grows up, until the accuracy has reached up the best result, which is 88.1%. When epoch is 300, the model is overfitting, therefore, the accuracy is 79.1% at this time, which is lower than the former one. Next, keep the batch size is 128, epoch is 50, and change the segment size. Then change the segment size from 1s to 3s, and 10s. The result of epoch 50 is show as the Table 6.

Table 6. The result of epoch 50 with two convolutional pooling layers

Time	SR acc (%)	VF acc (%)	VT acc (%)	Total acc (%)
1s	90.7	81.7	92.0	88.1
2s	92.2	77.3	93.0	87.5
3s	90.1	82.0	89.7	87.3
10s	45.5	35.6	71.2	50.8
Remarks: Batch Size:128, epoch: 50				

When the segment size was changed, and the data was put into the neural network for training. According to the different data length from 1s to 3s, the accuracy changed from 88.1% to 87.3%, and it doesn't have a dramatic drop. So, the segment size was set to 10s, the sensitivity had a great decrease. So, it can be concluded that the segment size will affect the model training effect. The model performance is better when the segment size is smaller. As far as the results of training model are concerned, the accuracy of identifying SR and VT are higher than those of VF. It can be seen from the ECG signal of these three in time domain, the signal of SR and VT are more regular and more representative, while VF are more chaotic. This is the reason that the sensitivity of VF is not as good as those two, because it is not easy to extract features and learn for neural networks. In addition, in this paper, two different kind of CNN architectures were established, one of the structures

is two convolutional-pooling layer model, another is four convolutional-pooling layers. Compared to Table 3 and Table 5, the average accuracy in table 3 is approximately 50%, which is apparently lower than Table 5. Therefore, the second structure performed better than first structure. It is because the amount of data is limited, if I put the limited data into a complex model, it may easier lead to over-fitting. So, the model was changed into the present model, using the convolutional neural network to implement it, and the accuracy and sensitivity of the result is more stable than former model. The result was shown as in Fig. 5.

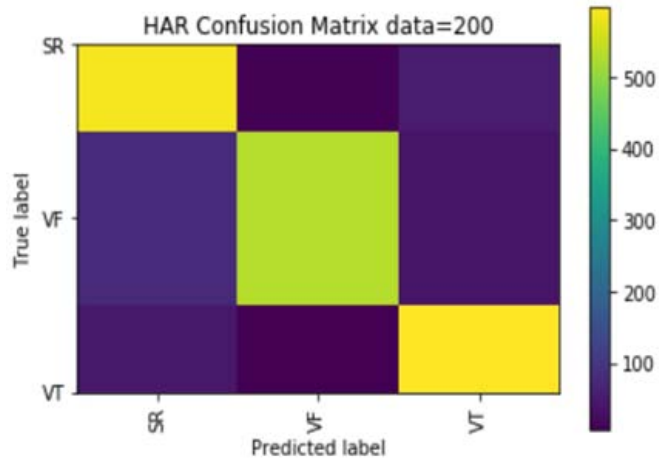


Fig. 5. The confusion matrix

5. Conclusions

In order to study cardiac arrhythmia classification, two CNN structures are used in our study, one is the four-convolutional pooling layer architecture, the other is the two-convolutional pooling layer architecture. From the results of Table 3 and Table 4, CNN with two convolutional pooling layers is better, its average recognition rate is 57.7%, which is higher than 46.1% of the four convolutional pooling layers. In addition, different parameters of the CNN are tested, including the parameters of matrix, learning rate, epoch and segment size. According to the test results in Table 1, when the matrix is 10000, the sensitivity of SR reaches 83.4%, and the total accuracy is 64.7%, which is much higher than the sensitivity when the matrix is 5000 or 15000. Therefore, when the matrix is 10000, the performance of the model is better. According to the test results in Table 2, when the learning rate is 0.0125, the sensitivity of SR is the highest, 83.4%, and the total accuracy is the highest, 67.6%. However, when the learning rate is 0.015 or 0.005, the model does not converge and cannot be classified.

The classification of arrhythmias has always been a research hotspot in the medical field. Before the application of CNN, the classification of heart rate data is a difficult issue. There are two main reasons, one of them is the amount of data to be processed is too large, resulting in high cost and low efficiency. In addition, it is difficult to retain the original features in the process of digitization, resulting in low accuracy of data processing. In order to deal with the issues, two convolutional pooling layers are used to classify the ECG signals, and four time periods (1s, 2s, 3s, 10s) are set and compared to get the best classification results. The research of this paper provides a reference for the classification of ECG signals and is helpful for the research of cardiovascular diseases. In the future research, more sample parameters and continuous optimization of convolutional layer are worth further exploration.

Acknowledgements. This work is supported by the project of Shanghai Philosophy and Social Sciences Plan (No. 2018BGL023).

References

1. Acharya, U.R., Fujita, H., Lih, O.S., Hagiwara, Y., Tan, J.H., Adam, M.: Automated detection of arrhythmias using different intervals of tachycardia ecg segments with convolutional neural network. *Information sciences* 405, 81–90 (2017)
2. Aljarah, I., Faris, H., Mirjalili, S.: Optimizing connection weights in neural networks using the whale optimization algorithm. *Soft Computing* 22(1), 1–15 (2016)
3. Benitez, D., Gaydecki, P.A., Zaidi, A., Fitzpatrick, A.P.: The use of the hilbert transform in ecg signal analysis. *Computers in biology and medicine* 31(5), 399–406 (2001)
4. Deng, S., Tian, Y., Hu, X., Wei, P., Qin, M.: Application of new advanced cnn structure with adaptive thresholds to color edge detection. *Communications in Nonlinear Science and Numerical Simulation* 17(5), 1637–1648 (2012)
5. Guo, T.M., Dong, J.W., Li, H.J., Gao, Y.X.: Simple convolutional neural network on image classification. In: *IEEE 2nd International Conference on Big Data Analysis (ICBDA)*. pp. 721–724 (2017)
6. Hong, S., Kwak, S., Han, B.: Weakly supervised learning with deep convolutional neural networks for semantic segmentation: Understanding semantic layout of images with minimum human supervision. *IEEE Signal Processing Magazine* 34(6), 39–49 (2017)
7. Kapanova, K.G., Dimov, I., Sellier, J.: A genetic approach to automatic neural network architecture optimization. *Neural Computing and Applications* 39, 1481–1492 (2016)
8. Kiranyaz, S., Ince, T., Gabbouj, M.: Real-time patient-specific ecg classification by 1-d convolutional neural networks. *IEEE Transactions on Biomedical Engineering* 63(3), 664–675 (2016)
9. Limonova, E., Sheshkus, A., Ivanova, A., Nikolaev, D.: Convolutional neural network structure transformations for complexity reduction and speed improvement. *Pattern Recognition & Image Analysis* 28(1), 24–33 (2018)
10. Liu, H., Li, C., Xu, L.: Dimension reduction and classification of hyperspectral images based on neural network sensitivity analysis and multi-instance learning. *Computer Science and Information Systems* 16(2), 443–468 (2019)
11. Meng, H.H., Zhang, Y.: Classification of electrocardiogram signals with deep belief networks. In: *IEEE International Conference on Computational Science and Engineering (CSE)*. pp. 7–12 (2015)
12. Mesin, L.: Heartbeat monitoring from adaptively down-sampled electrocardiogram. *Computers in Biology and Medicine* 84, 217–225 (2017)

13. Mohan, N., Sachin, K.S., Poornachandran, P., Soman, K.: Modified variational mode decomposition for power line interference removal in ecg signals. *International Journal of Electrical and Computer Engineering* 6(1), 151–159 (2016)
14. Montminy, D.P., Baldwin, R.O., Temple, M.A., Laspe, E.D.: Improving cross-device attacks using zero-mean unit-variance normalization. *Journal of cryptographic engineering* 3(2), 99–110 (2013)
15. Piotrowski, A.P., Napiorkowski, J.J.: A comparison of methods to avoid overfitting in neural networks training in the case of catchment runoff modelling. *Journal of Hydrology Amsterdam* 476(1), 97–111 (2013)
16. Qiu, L., Li, W., Cai, W., Zhang, M., Zhu, W., Wang, L.: Heartbeat classification using convolution neural network and wavelet transform to extract features. In: *IEEE SmartWorld, Ubiquitous Intelligence & Computing, Advanced & Trusted Computing, Scalable Computing & Communications, Cloud & Big Data Computing, Internet of People and Smart City Innovation*. pp. 139–143 (2018)
17. Salem, M., Taheri, S., Yuan, J.: Ecg arrhythmia classification using transfer. learning from 2-dimensional deep cnn features. In: *IEEE Biomedical Circuits and Systems Conference (BioCAS)*. pp. 1–4 (2018)
18. Vikramjit, M., Ganesh, S., Hosung, N., Carol, E.W., Elliot, S., Mark, T.: Hybrid convolutional neural networks for articulatory and acoustic information based speech recognition. *Speech Communication* 89, 103–112 (2017)
19. Wei, C., Jia, Z.H., Yuan, G.Q., Wu, Y.L.: Integration and regulation- new strategy of arrhythmia drug intervention from "anti;"heart rhythm" to "regulate heart rhythm". *Chinese Journal of Cardiac Arrhythmias* 18(1), 76–78 (2014)
20. Yan, S.Y. (ed.): *Cardiac arrhythmia classification on convolution neural network*. King's College London (2019)
21. Yu, S.N., Chen, Y.H.: Noise-tolerant electrocardiogram beat classification based on higher order statistics of subband components. *Artificial intelligence in medicine* 46(2), 165–178 (2009)
22. Zhou, T., Hou, M.Z., Liu, C.H.: Forecasting stock index with multi-objective optimization model based on optimized neural network architecture avoiding overfitting. *Computer Science and Information Systems* 15(1), 211–236 (2018)

Yonghui Dai is currently a lecturer at the Management School, Shanghai University of International Business and Economics, China. He received his Ph.D. in Management Science and Engineering from Shanghai University of Finance and Economics, China in 2016. His current research interests include management information systems, affective computing and artificial intelligence. His works have appeared in international journals more than thirty papers. Contact him at daiyonghui@suibe.edu.cn.

Bo Xu is a professor at the Management School, Shanghai University of International Business and Economics, China. He received his Ph.D. in Management Science from Fudan University, China in 1997. His interests include strategic management and educational technology. Contact him at brianxubo@163.com.

Siyu Yan is the corresponding author of this paper, she is a Ph.D. candidate at School of Information Management and Engineering at Shanghai University of Finance and Economics. She received her MS in King's College London. Her research interests include data mining and big data analytics. Contact her at ysy18956@163.com.

Jing Xu is a Ph.D candidate at School of Information Management and Engineering at Shanghai University of Finance and Economics. She received her MA. in School of Business at Shanghai University of Finance and Economics, China, in 2013. Her research interests include personalized recommendation. E-mail: 1053517221@qq.com.

Received: December 29, 2019; Accepted: May 5, 2020.

Option predictive clustering trees for multi-target regression^{*}

Tomaž Stepišnik^{1,2}, Aljaž Osojnik^{1,2}, Sašo Džeroski^{1,2}, and Dragi Kocev^{1,2}

¹ Department of Knowledge Technologies, Jožef Stefan Institute, Ljubljana, Slovenia

² Jožef Stefan International Postgraduate School, Ljubljana, Slovenia
{tomaz.stepisnik,aljaz.osojnik,saso.dzeroski,dragi.kocev}@ijs.si

Abstract. Decision trees are one of the most widely used predictive modelling methods primarily because they are readily interpretable and fast to learn. These nice properties come at the price of predictive performance. Moreover, the standard induction of decision trees suffers from myopia: a single split is chosen in each internal node which is selected in a greedy manner; hence, the resulting tree may be sub-optimal. To address these issues, option trees have been proposed which can include several alternative splits in a new type of internal nodes called option nodes. Considering all of this, an option tree can be also regarded as a condensed representation of an ensemble. In this work, we propose to learn option trees for multi-target regression (MTR) based on the predictive clustering framework. The resulting models are thus called option predictive clustering trees (OPCTs). Multi-target regression is concerned with learning predictive models for tasks with multiple numeric target variables. We evaluate the proposed OPCTs on 11 benchmark MTR data sets. The results reveal that OPCTs achieve statistically significantly better predictive performance than a single predictive clustering tree (PCT) and are competitive with bagging and random forests of PCTs. By limiting the number of option nodes, we can achieve a good trade-off between predictive power and efficiency (model size and learning time). We also perform parameter sensitivity analysis and bias-variance decomposition of the mean squared error. Our analysis shows that OPCTs can reduce the variance of PCTs nearly as much as ensemble methods do. In terms of bias, OPCTs occasionally outperform other methods. Finally, we demonstrate the potential of OPCTs for multifaceted interpretability and illustrate the potential for inclusion of domain knowledge in the tree learning process.

Keywords: multi-target regression, option trees, interpretable models, predictive clustering trees, bias-variance decomposition of error.

1. Introduction

Supervised learning is one of the most widely researched and investigated areas of machine learning. The goal in supervised learning is to learn, from a set of examples with known class, a function that outputs a prediction for the (scalar-valued) class of a previously unseen example. However, in many real life problems of predictive modelling the output (target) is structured, e.g. it is a vector of class values or a tuple of target variables. There can be dependencies between the class values/targets (e.g., they can be organized

* The authors wish it to be known that the first two authors should be regarded as joint first authors.

into a tree-shaped hierarchy or a directed acyclic graph) or some internal relations between the class values may exist (e.g., as in sequences).

In this work, we concentrate on the task of predicting multiple numeric variables. Examples thus take the form $(\mathbf{x}_i, \mathbf{y}_i)$, where $\mathbf{x}_i = (x_{i1}, \dots, x_{ik})$ is a vector of k input variables and $\mathbf{y}_i = (y_{i1}, \dots, y_{it})$ is a vector of t target variables. This task is known under the name of *multi-target regression* (MTR) [4, 23, 24] (also known as multi-output or multivariate regression). MTR is a type of structured output prediction task which has applications in many real life problems, where we are interested in simultaneously predicting multiple numeric variables. Prominent examples come from ecology and include predicting the abundance of different species living in the same habitat [12] and predicting properties of forests [21, 29]. Due to its applicability to a wide range of domains, this task is recently gaining increasing interest in the research community.

Several methods for addressing the task of MTR have been proposed [24, 31]. These methods can be categorized into two groups of methods [2]: (1) local methods, that predict each of the target variable separately and then combine the individual model predictions to get the overall model prediction and (2) global methods, that predict all of the variables simultaneously (also known as ‘big-bang’ approaches). In the case of local models, for a domain with t target variables one needs to construct t predictive models – each predicting a single target. The prediction vector (that consists of t components) of an unseen example is then obtained by concatenating the predictions of the multiple single-target predictive models. Conversely, in the case of global models, for the same problem one needs to construct only one model. In this case, the prediction vector of an unseen example is obtained by passing the example through the model and getting its (complete) prediction.

In the past, several researchers proposed methods for solving the task of MTR directly and demonstrated their effectiveness [1, 8, 10, 21, 24, 20, 30]. The global methods have several advantages over the local methods. First, they exploit and use the dependencies that exist between the components of the structured output in the model learning phase, which can result in better predictive performance. Next, they are typically more efficient: it can happen that the number of components in the output is very large (e.g., predicting the bioactivity profiles of compounds described with their quantitative structure-activity relationships on a large set of proteins), in which case executing a basic method for each component is not feasible. Furthermore, they produce models that are typically smaller than the sum of the sizes of the models built for each of the components.

The state-of-the-art methods for MTR are based on tree and ensemble learning [8, 24, 27, 31]. Trees for MTR (from the predictive clustering framework) inherit the properties of regression trees: they are interpretable models, but their construction is greedy. The performance of trees is significantly improved when they are used in an ensemble setting [24, 20]. However, the myopia, i.e., greediness, of the tree construction process can lead to learning sub-optimal models. One way to alleviate this is to use a beam-search algorithm for tree induction [22], while another approach is to introduce option splits in the nodes [9, 25].

Beside the many appealing properties of trees, their predictive performance is often limited. In a variety of machine learning tasks including MTR, learning ensembles of trees typically significantly improve the predictive power of the single models [24]. However, learning ensembles has a significant drawback – lack of a possibility to interpret the obtained predictive model. It is not feasible to analyze each tree in an ensemble due to both

the number of trees and the randomized procedure that is used to learn them. Typically, this means that there is a trade-off between interpretability and predictive performance.

To address this issue, we propose to extend predictive clustering trees (PCTs) for MTR towards option trees, i.e., we introduce option predictive clustering trees (OPCTs). An option tree can be seen as a condensed representation of an ensemble of trees which shares a common substructure. More specifically, the heuristic function for split selection can return multiple values that are close to each other within a predefined range. These splits are then used to construct an option node. For illustration, see Figure 1.

The contributions of this work can be summarized as follows:

- We introduce a new method for addressing the task of MTR that provides a balance between interpretability and predictive performance: An option PCT can be treated as an ensemble, or the best subtree can be extracted and analyzed as a single tree.
- We analyze the computational complexity of the proposed method and compare it to the complexity of the competing methods.
- We perform an empirical evaluation of the performance of the proposed method along three dimensions: predictive power, time efficiency and size of the models.
- We further analyse the predictive performance by examining the bias-variance decomposition of the models' errors. We perform the analysis for the proposed method as well as the competing methods. As far as we are aware, such an analysis has not been performed before for MTR methods.

This work extends our earlier work presented in [26] along two major directions. First, we perform a theoretical analysis of the computational complexity of the proposed method and compared it to the computational complexity of learning a single PCT and ensembles of PCTs. Second, we calculate the bias-variance decompositions of mean squared errors to determine which components of the error are addressed by different parameter selections and competing methods.

The remainder of this paper is organized as follows. Section 2 proposes the algorithm for learning option PCTs for MTR. Next, Section 3 outlines the design of the experimental evaluation and the details on the bias-variance decomposition of the error. Section 4 continues with a discussion of the results. Finally, Section 5 concludes and provides directions for further work.

2. Option predictive clustering trees

The predictive clustering trees framework views a decision tree as a hierarchy of clusters. The top-node corresponds to one cluster containing all data, which is recursively partitioned into smaller clusters while moving down the tree. The PCT framework is implemented in the CLUS system [3], which is available for download at <http://clus.sourceforge.net>.

PCTs are induced following the *top-down induction of decision trees* (TDIDT) algorithm [7]. There are however two major differences: the heuristic score and the prototype function are instantiated by the PCT algorithm based on the task that is being addressed (classification, regression or multi-target prediction etc.). The heuristic score used for selecting the tests in the internal nodes of a regular PCT is the maximization of the variance reduction resulting from the partitioning of the instances into subtrees corresponding to

the outcome of the tests. By doing so, we also maximize cluster homogeneity, and, consequently, improve the predictive performance. The PCT algorithm also defines the prototype functions used in each tree leaf to give predictions for new examples based on the task at hand (e.g., uses averaging for MTR).

Option predictive clustering trees (OPCT) extend the usual PCT framework, by introducing option nodes into the tree building procedure outlined in Algorithm 1. Option decision trees were first introduced as classification trees by [9] and then analyzed in more detail by [25]. [17] analyzed regression option trees in the context of data streams.

The major motivation for the introduction of option trees is to address the myopia of the *top-down induction of decision trees* (TDIDT) algorithm [7]. Viewed through the lens of the predictive clustering framework, a PCT is a non-overlapping hierarchical clustering of the whole input space. Each node/subtree corresponds to a clustering of a subspace and prediction functions are placed in the leaves, i.e., lowest clusters in the hierarchy. An OPCT, however, allows the construction of an overlapping hierarchical clustering. This means that, at each node of the tree several alternative hierarchical clusterings of the subspace can appear instead of a single one.

When using an OPCT for prediction on a new example, we produce the prediction by aggregating over the predictions of the alternative subtrees (overlapping clusters) the example may encounter. However, as not all parts of the tree (hierarchical clustering) are necessarily overlapping, the example may encounter only nonoverlapping (sub)clusters. In that case, we produce the prediction as we would with a regular PCT.

When using TDIDT to construct a predictive clustering tree, and in particular when partitioning the data, all possible splits are evaluated by using a heuristic and the best one is selected. However, other splits may have very similar heuristic values. The best partition could be obtained with another split as a consequence of noise or of the sampling that generated the data. In this case, selecting a different split could be optimal. To address this concern, the use of option nodes was proposed [25].

An option node is introduced into the tree when it would be hard to determine the best split, that is when the best splits have similar heuristic values. When this occurs, instead of selecting only the best split, we select several of them. Specifically, we select up to 5 splits s , called *options*, that satisfy the following

$$\frac{\text{Heur}(s)}{\text{Heur}(s_{best})} \geq 1 - \varepsilon \cdot d^{level},$$

where s_{best} is the best split, ε determines how similar the heuristics must be, $d \in [0, 1]$ is a decay factor and $level$ is the level in the tree of the node we are attempting to split. This equation assumes that heuristic scores are to be maximized (higher value is better). After we have determined the candidate splits, we introduce an option node whose children are split nodes obtained by using the selected splits, i.e., an option node contains options as its children. Selecting more than 5 options is possible, but, uses more resources and is not advised [25]. A drawback of this method of split selection is that if some attributes are highly correlated to the attribute that produces the best split, those attributes will likely be selected for other splits in the option node. In this case the option node will only increase the learning time.

As usual, we define the level of a node to be the number of its ancestor nodes, however, we do not count option nodes. This is motivated by the predictive clustering viewpoint,

i.e., the option nodes only mark that there are overlapping clusters and not that there is an additional level of clustering.

The use of a decay factor makes the selection criterion more stringent in the lower nodes of the tree. The intuition behind this is that higher up, the split selection is more important and a larger error would be inferred by introducing a non-optimal split. However, as we get deeper into the tree, the use of a non-optimal split makes decreasing impact. This intuition also allows us to prohibit the use of option nodes on levels 3 and greater, which severely mitigates the problem of combinatorial explosion.

Outline of the entire tree building process is presented in algorithm 1. The variable *candidates* is a list of at most 5 tests with the highest heuristic values. Every test is represented as a triplet (t, h, P) , where t is the actual test function, h its heuristic value and P the partitions produced by the test.

When using a small ε , e.g., $\varepsilon = 0.1$, we are selecting only options whose heuristics are within 10% of the best split. However, the use of larger ε , in the extreme case even $\varepsilon = 1$, can also be motivated through the success of methods such as random forests and ensembles of extremely randomized trees. Allowing the selection of splits whose heuristic values are considerably worse than the heuristic value of the best split might not necessarily reduce the performance of the tree, but actually increase it.

Algorithm 1 The top-down induction algorithm for option PCTs.

Procedure OptionPCT

Input: A data set E , parameter ε , decay factor d , current tree level l

Output: An option predictive clustering tree

```

candidates = FindBestTests( $E$ , 5)
if |candidates| > 0 then
  if |candidates| = 1 or  $l > 2$  then
     $(t^*, h^*, \mathcal{P}^*) = \text{candidates}[0]$ 
    for each  $E_i \in \mathcal{P}^*$  do
       $tree_i = \text{OptionPCT}(E_i, \varepsilon, d, l + 1)$ 
    return  $\text{node}(t^*, \bigcup_i \{tree_i\})$ 
  else
     $(t_0^*, h_0^*, \mathcal{P}_0^*) = \text{candidates}[0]$ 
     $nodes = \{\}$ 
    for each  $(t_i^*, h_i^*, \mathcal{P}_i^*) \in \text{candidates}$  do
      if  $\frac{h_i^*}{h_0^*} \geq 1 - \varepsilon \cdot d^l$  then
        for each  $E_j \in \mathcal{P}_i^*$  do
           $tree_j = \text{OptionPCT}(E_j, \varepsilon, d, l + 1)$ 
           $nodes = nodes \cup \{\text{node}(t^*, \bigcup_j \{tree_j\})\}$ 
    if |nodes| > 1 then
      return  $\text{option\_node}(nodes)$ 
    else
      return  $nodes[0]$ 
  else
    return  $\text{leaf}(\text{Prototype}(E))$ 

```

Once an OPCT is built, we want to use it for prediction. In a regular PCT, it is simple to produce a prediction for a new example. It is sorted into a leaf (reached according to the splits of the tree) where a prediction is made by using a prototype function. When traversing an example through an OPCT, we behave the same when we encounter a split or leaf node. If we traverse an example to an option node, however, we clone the example for each of the options and traverse one of the copies down each of the options. This means that in an option node an example is (by proxy of its copies) traversed to multiple leaves, where multiple predictions are produced. To obtain a single prediction in an option node, we aggregate the obtained predictions. When addressing multi-target regression this is generally done by averaging all the predictions per target.

An option tree is usually seen as a single tree, however, it can also be interpreted as a compact representation of an ensemble. To generate the ensemble of the *embedded trees*, we start recursively from the root node and move in a top-down fashion. Each time we encounter an option node we copy the tree above (and in "parallel") for each of the options and replace the option node with only the option, i.e., single split. This produces one tree for each option while removing the option node in question. We repeat this procedure on all the generated trees until we are left with no option nodes. This is illustrated in Figure 1. For this reason, we will sometimes refer to option trees as pseudo-ensembles.

A given OPCT is also an extension of the PCT that would be learned on the same data. By definition, whenever we introduce an option node, we include the best split (in terms of the heuristic)³. In regular construction of PCTs, this is the only split we consider. Consequently, the PCT is embedded in the OPCT. We can extract it if we, in a top-down fashion, select the best option in each option node.

Let's consider a full option node, that is one where we have the full 5 options. Let's assume that there are no option nodes further down in the tree. Since we have 5 options and no options lower in the tree, we have a total of 5 embedded trees. Now, let's consider the size of the embedded ensemble, when we add two such nodes under a split node, i.e., each leaf of a split node was extended into a full option node. To construct an embedded tree we can now choose one of the 5 options when the test is satisfied and one of 5 options when it is not. This results in 25 different embedded trees. It can be inferred, that to calculate the number of embedded trees for an option node we need to sum up the number of embedded trees for each of its options. In a (binary) split node, however, we multiply the numbers of embedded trees of the subtree that satisfies the split and of the subtree that does not, to obtain the total number of embedded trees.

Given the construction constraints described above, we know that option nodes with up to 5 options can appear only on the first three levels, i.e., levels 0, 1 and 2. If we now consider a full option tree, we calculate a maximum of $(5^2 \cdot 5)^2 \cdot 5 = 5^7 = 78125$ embedded trees. However, many of these trees overlap to a large extent.

Note that a given example will not traverse the entire tree. For example, in Figure 1, if an example reaches S_1 and is traversed into the left child L_1 , the same result would happen in both the first and second embedded tree. The example is "agnostic" of any option nodes in the right child of S_1 . Therefore, in a general option tree, a given example will visit only up to $5^3 = 125$ leaves, as it will only traverse down one side of the tree in each split node.

³ Not only is the best split included, other splits are compared to it to determine their inclusion in the option tree.

In other words, there are a maximum of 125 different predictions that would be aggregated in order to obtain the final prediction of an option tree constructed this way. This can be compared to a single tree, where only 1 prediction will be made for each example, or to a tree ensemble, where 1 prediction would be made for each member of the ensemble and then aggregated to produce the final prediction.

[24] show that the computational complexity of building a multi-target PCT is $\mathcal{O}((S + \log N)DN \log N)$, where N is the number of examples in the data set, D is the number of descriptive variables and S is the number of target variables. Let A be the maximum number of options in an option node (in our case, we set A to 5) and B the maximum depth at which an option node can occur (in our case, we set B to 3). When constructing an OPCT, in the worst case scenario, computationally-wise, we have option nodes at all allowed levels with the maximum number of options. From that point onward we essentially construct A^B regular PCTs. Assuming that B is only a fraction of the final tree depth (i.e., $B \ll \log N$), we conclude that the computational complexity of building an OPCT is A^B times greater than the computational complexity of a regular PCT.

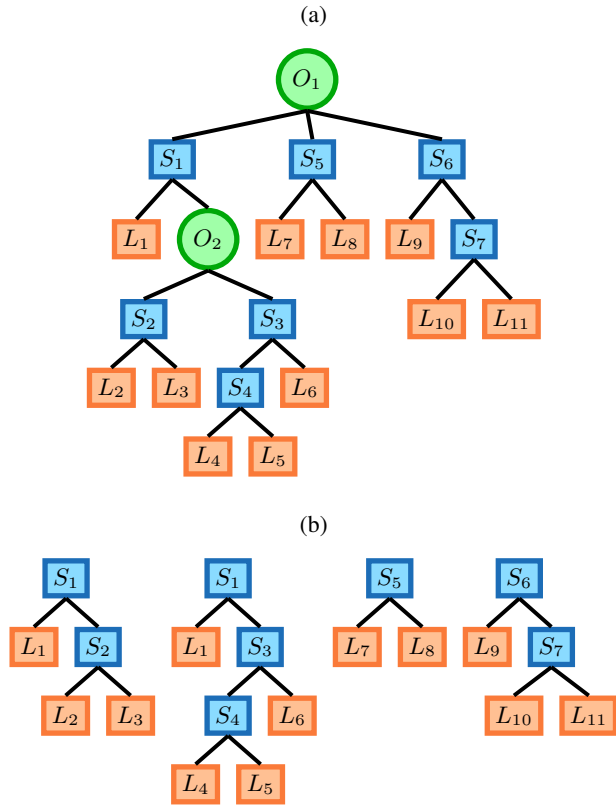


Fig. 1. An option tree (a) and the ensemble of its embedded trees (b). O_i are option nodes, S_j split nodes and L_k leaf nodes.

In comparison, a bagging ensemble constructs M independent PCTs on bootstrapped training sets, making their computational complexity M times that of a single PCT. Random forests, in addition to bootstrapping, only consider $f(D)$ descriptive variables at each node. In the regular case, i.e., when $f(D) = \sqrt{D}$, this yields $\mathcal{O}((S+\log N)\sqrt{DMN} \log N)$ complexity. To summarize, OPCTs and bagging ensembles are constant factor-times slower (A^B and M respectively) than single PCTs, while random forests can be faster than a single PCT for data sets with large numbers of descriptive variables.

3. Experimental design

To evaluate the performance and efficiency of the OPCT method, we construct OPCTs with two parameter configurations, as well as standard PCTs and ensembles of PCTs. We first present the benchmark data sets used for evaluation of the methods and then give the specific experimental setup: parameter selections and evaluation measures.

3.1. Data description

The 11 data sets with multiple numeric targets used in this study come mainly from the domain of ecological modelling. Table 1 outlines the properties of the data sets. The selection contains data sets with various numbers of examples described with different numbers of attributes. For more details on the data sets, we refer the reader to the referenced literature, the repositories available at: <http://mulan.sourceforge.net/datasets-mlc.html> and <http://kt.ijs.si/DragiKocev/PhD/resources/>, as well as [24] and the references therein.

3.2. Evaluation of predictive performance and efficiency

We parameterize OPCTs by selecting values for the parameters ε and d . We consider two parameter configurations: $(\varepsilon = 1, d = 1)$ and $(\varepsilon = 0.2, d = 1)$. When $\varepsilon = 1$, there are no constraints on the heuristic value of the selected splits with regards to the best

Table 1. Properties of the data sets with multiple numeric targets (regression data sets): M is the number of instances, $|D|$ the number of descriptive attributes, and T the number of target attributes.

Name of data set	M	$ D $	T
Collembola [18]	393	47	3
EDM [19]	154	16	2
Forestry-Kras [29]	60607	160	11
Forestry-Slivnica-LandSat [28]	6218	150	2
Forestry-Slivnica-IRS [28]	2731	29	2
Forestry-Slivnica-SPOT [28]	2731	49	2
Sigma real [11]	817	4	2
Soil quality [12]	1944	142	3
Vegetation Clustering [16]	29679	65	11
Vegetation Condition [21]	16967	40	7
Water quality [14]	1060	16	14

test. However, since only the 5 best splits are selected, the risk that a split which would decrease the predictive performance would be selected is relatively low. This setting most resembles the ensemble setting, where greater variation is desired. The other parameter configuration provides a balance between predictive performance and efficiency. In our previous work [26], we used $\varepsilon = 0.5$ for the efficient version, but it turned out to be too close to the ensembles still. So in this paper, we selected $\varepsilon = 0.2$.

Next, we define the parameter values used in the algorithms for constructing single PCTs and ensembles of PCTs. The multi-target PCTs are obtained using F-test pruning. This pruning procedure uses the exact Fisher test to check whether a given split/test in an internal node of the tree results in a reduction in variance that is statistically significant at a given significance level. If there is no split/test that can satisfy this, then the node is converted to a leaf. An optimal significance level was selected by using internal 3-fold cross validation, from the following values: 0.125, 0.1, 0.05, 0.01, 0.005 and 0.001.

We consider two ensemble learning techniques: bagging [5] and random forests [6]. These are the most widely used tree-base ensemble learning methods. The construction of both ensemble methods takes as an input parameter the size of the ensemble, i.e., number of base predictive models to be constructed. We constructed ensembles with 100 base predictive models [24]. Furthermore, the random forests algorithm takes as input the size of the feature subset that is randomly selected at each node. For this purpose, we use the square root of the number of descriptive attributes $\lceil \sqrt{|D|} \rceil$.

We use 10-fold cross-validation to estimate the predictive performance of the used methods. We assess the predictive performance of the algorithms using several evaluation measures. In particular, since the task we consider is MTR, we employed three well known measures: the correlation coefficient (CC), mean squared error (MSE) and relative root mean squared error (RRMSE). We present here the bias-variance analysis of MSE and statistical comparison of methods according to RRMSE, where similar conclusions hold also for the other two measures. Finally, the efficiency of the proposed methods is measured with the time needed to learn a model and the size of the models (in terms of total number of leaf nodes). While OPCTs can be learned in parallel on node splitting level, our implementation does not support it yet. For this reason we used no parallelization in our experiments, to provide a fairer comparison.

In order to assess the statistical significance of the differences in performance of the studied algorithms, we adopt the recommendations by [13] for the statistical evaluation of the results. In particular, we use the Friedman test for statistical significance. Afterwards, to detect where statistically significant differences occur (i.e., between which algorithms), we use the Nemenyi post-hoc test.

We present the results of the statistical analysis with *average ranks diagrams*. The diagrams plot the average ranks of the algorithms and connect those whose average ranks differ by less than a given value, called *critical distance*. The critical distance depends on the level of the statistical significance (set in our case to 0.05). The difference in performance of the algorithms connected with a line is not statistically significant at the given significance level.

3.3. Parameter sensitivity analysis

This set of experiments is designed to provide more insight into how *varepsilon* and *d* parameters affect the performance and efficiency of OPCTs. For ε we consider the values

$\{0.1, 0.2, 0.5, 1.0\}$, corresponding in order from the most stringent to the least stringent construction criterion. If we were to select $\varepsilon = 0$, the resulting OPCT would almost always directly coincide with a regular PCT, as no split would likely reach exactly the same heuristic value as the best split. Hence, the only way an option node would be induced is if two splits had the exact same heuristic value, therefore this configuration is not of interest.

As discussed above, the higher in the tree an option node is induced, the higher the variation in the learned subtrees. Induction of option nodes in lower levels of the tree not only contributes to the combinatorial explosion of the number of trees (and consequently the use of resources), but also generates less variation in the predictions, since the subtrees affected cover a smaller number of examples. Hence, we wish to curtail the number of options induced in option nodes lower in the tree. If we select a decay factor of $d = 1$, the depth of the option node induction will have no impact, while selecting a decay factor of 0.5 will effectively double the effective heuristic requirement $\varepsilon \cdot d^l$ at each level. For example, for $\varepsilon = 0.5$ and $d = 0.5$, the requirement would be 0.5 at the root level, 0.25 at the first level and 0.125 at the third level. We use the following values of the decay factor: $\{0.5, 0.9, 1\}$, these are the most to the least constrictive in terms of the construction of the OPCT.

Note that, on a given data set, all values of ε and d could produce the same OPCT, if the splits have very similar heuristic values, e.g., there could always be 5 splits that are within $10\% \cdot d^l$ of the heuristic value of the best split. Therefore, the evaluation of how ε and d affect both the predictive performance and efficiency must by design be evaluated on multiple data sets. The parameterized version of the OPCT method for a given ε and d is denoted $\text{OPCTe}\varepsilon d d$, e.g., $\text{OPCTe}0.5d0.9$.

In order to facilitate replication of all the experiments, we implemented the method proposed in this paper in the latest version of CLUS, already available in the public CLUS repository at <http://clus.sourceforge.net>.

3.4. Bias-variance decomposition

We analyze the predictive performance of OPCTs as well as the performance of the competing methods by using bias-variance decomposition of the mean squared error [15]. The analysis allows us to investigate the source of errors of the methods. We perform it for each target separately. To calculate the decomposition, we need to predict the targets of a single data sample several times using models trained on different training sets.

Suppose we have a set of m predictive models learned using the same method on different, but related, training sets. They give m predictions for a target variable for each of the n data samples. Let y_i be the real target value of the i -th data sample and let f_{ij} be the models' prediction for y_i obtained on the j -th training set. Then the mean squared error (MSE), bias and variance can be calculated as follows:

$$\bar{f}_i = \frac{1}{m} \sum_{j=1}^m f_{ij}$$

$$\begin{aligned}\text{Bias} &= \frac{1}{n} \sum_{i=1}^n (y_i - \bar{f}_i)^2 \\ \text{Variance} &= \frac{1}{mn} \sum_{i=1}^n \sum_{j=1}^m (f_{ij} - \bar{f}_i)^2 \\ \text{MSE} &= \frac{1}{mn} \sum_{i=1}^n \sum_{j=1}^m (f_{ij} - y_i)^2.\end{aligned}$$

It is easy to show that $\text{MSE} = \text{Bias} + \text{Variance}$.

More complex models tend to have lower bias because they can better accommodate individual variations in the data, but this also increases their sensitivity to changes in the training set, leading to increased variance. Therefore, to minimize the error of the model a balance between bias and variance should be attained.

As mentioned earlier, the bias-variance decomposition procedure requires several predictions for each example obtained from models trained on different training sets using the same method. Standard 10-fold cross validation will therefore not suffice. To this end, after we split a data set into 10 folds, we select each of the folds for test set once and use the remaining non-test folds as a source of training sets. From this source we generate 20 training sets using bootstrapping and we learn predictive models on each bootstrap sample. The learned models are then applied to the examples from the test folds to obtain the multiple predictions for each data example.

4. Results and discussion

We discuss the results from the experimental evaluation along four major dimensions. First, we compare the performance of OPCTs with a single PCT and ensembles of PCTs. We compare them based on their predictive performance, learning time and size of the produced models. Next, present the effect of different parameter values on the construction of OPCTs. Third, we inspect the biases and variances of the algorithms to gain a better understanding of the comparisons of the predictive performance. Finally, we discuss the interpretability of OPCTs in juxtaposition to PCTs.

When analyzing the predictive power of algorithms in sections 4.2 and 4.1, ranking was done separately for every target of every data set, giving us a total of 59 samples. For efficiency we compared times needed to induce a model on entire data sets and the total sizes of these models. The Friedman test showed significant differences at $p < 1 \cdot 10^{-8}$ in all cases and the results of Nemeny post-hoc tests are presented as average ranking diagrams.

4.1. Predictive performance and efficiency

Figure 2 shows the results of the statistical evaluation of the predictive performance of the proposed OPCT method in comparison to a single PCT and ensembles of PCTs. First of all, we note that there is no statistically significant difference in the performance of bagging of PCTs, random forests of PCTs and OPCTe1d1. Notwithstanding this, random forests of PCTs have a slightly better average rank compared to the other two methods

that are very close in rank. Second, while OPCTe0.2d1 is significantly worse than the three aforementioned methods, it is still significantly better than a single PCT.

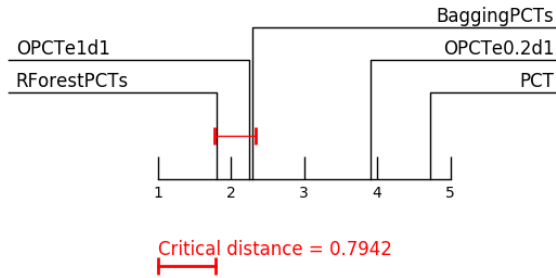


Fig. 2. Comparison of OPCTs predictive performance of OPCTs to the competing methods: Average rank diagram in terms of predictive performance.

In figure 3 we compare the methods in terms of their efficiency. It shows that learning PCTs is the most efficient method both in terms of time needed for model construction and model size, according to the average rank. However, learning OPCTe0.2d1 is not statistically significantly worse in terms of both time and size efficiency than learning a single PCT. Random forests of PCTs are close to single PCTs in terms of time, but not in terms of size, as they are significantly larger. Bagging of PCTs and OPCTe1d1 are the two slowest and largest methods.

Focusing on the efficiency of the OPCT models, we see that OPCTe1d1 is the least efficient: it produces models with the largest numbers of leaves and takes the most time for model construction. Recall that this is consistent with the analysis of computational complexity, as OPCTe1d1 can be seen as a condensed representation of what amounts to 125 PCTs, since the the algorithm selects the 5 best options at the first three levels. Furthermore, reducing the value of the ϵ parameter to 0.2 offers a good trade-off between predictive performance and efficiency. These models give significantly better predictive performance than a single PCT and are not significantly larger or slower to learn.

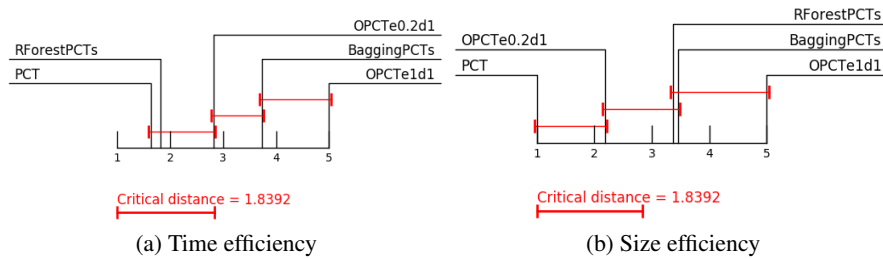


Fig. 3. Comparison of the efficiency of OPCTs to that of the competing methods as measured by the time needed to learn a model and the size of the models (number of leaf nodes).

4.2. Parametrization of OPCTs

Figures 4, 5 and 6 show the average rank diagrams for predictive performance, learning time and size of different OPCTs, respectively, obtained using the experimental design outlined above. Notably, low values of both parameters lead to degradation in predictive performance. This is to be expected, since lower values of parameters force the algorithm to introduce fewer option nodes, resulting in smaller OPCTs. On the opposite side of the spectrum are the OPCTs obtained with large values of the parameters. These achieve the best predictive performance and the corresponding OPCTs are the largest in terms of the number of leaves.

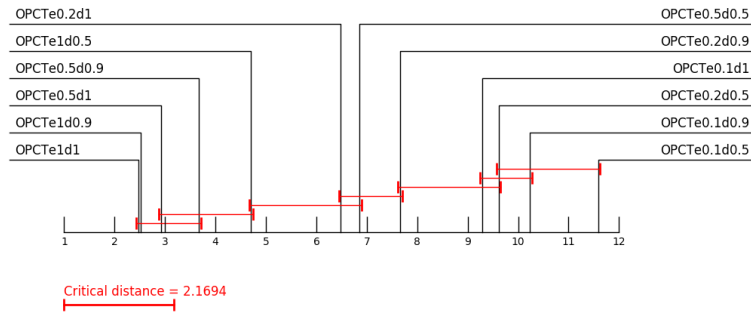


Fig. 4. Average rank diagram in terms of predictive performance for OPCTs obtained with different values of the parameters ε (e) and d .

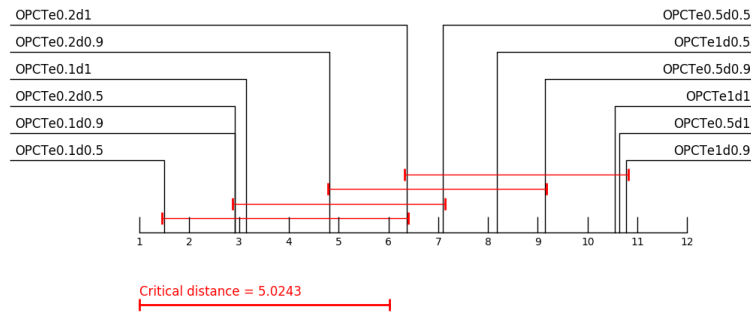


Fig. 5. Average rank diagram in terms of learning time for OPCTs obtained with different values of the parameters ε (e) and d .

Additionally, we observe that the ε parameter has a stronger influence on the performance than d . Namely, the OPCTs constructed using ε values of 0.1 and 0.2 (with the exception of OPCTe0.2d1) are the ones with the weakest predictive power. Furthermore, we also note that larger values for d also lead to better predictive performance. The best

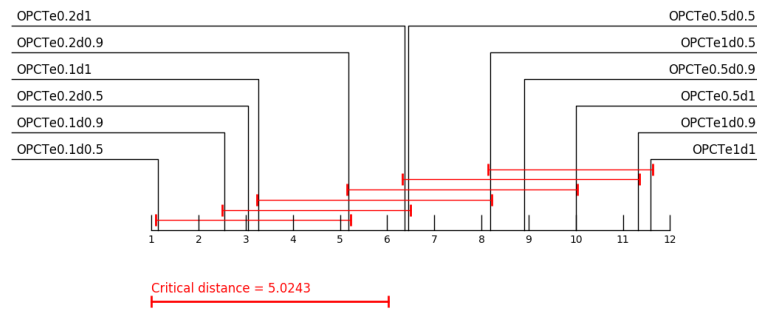


Fig. 6. Average rank diagram in terms of model size for OPCTs obtained with different values of the parameters ε (e) and d .

predictive performance is obtained when using 1 as the value of both parameters. This setting produces the largest OPCTs and requires the most learning time. Tables 2-4 in the Appendix present all values for learning times, model sizes and option node counts for different parameter values.

4.3. Bias-variance decomposition of the errors

Figures 7, 8 and 9 present the bias-variance decomposition of MSE for every target of the *Vegetation Condition*, *Sigma real* and *Forestry-Slivnica-LandSat* data sets, respectively. The graphs for the remaining data sets are presented in Appendix.

Algorithms are presented in the following order: PCT, bagging of PCTs, random forest of PCTs, then OPCTs with $e0.1d0.5$, $e0.1d0.9$, $e0.1d1$, $e0.2d0.5$, $e0.2d0.9$, $e0.2d1$, $e0.5d0.5$, $e0.5d0.9$, $e0.5d1$, $e1d0.5$, $e1d0.9$, $e1d1$. For every target of every data set values were scaled to $[0, 1]$ by the largest MSE of all algorithms for that target, so that they can be presented on one figure. While the biases (and variances) of the different targets are presented on the same scale, only the values of different algorithms on the same target variable can be reasonably compared among each other.

The decompositions of the errors are presented on the same graph for all targets of a data set to show that they produce very similar results. That is, if a method has smaller bias (or variance) than another method on one target of a data set, it tends to have smaller bias (or variance) on the other targets of that data set as well.

The results for the *Vegetation Condition* data set given in Figure 7 represent the most common behavior. The bias of the methods is fairly similar, with less restrictive OPCTs (larger parameter values) being slightly better in this regard. PCTs have the largest variance (and consequently total error), while random forests of PCTs have the smallest. Increasing the parameter values of OPCTs also reduces the variance, bringing it on par with that of bagging and random forests.

The error decomposition for the *Sigma real* data set given in Figure 8 stands out, since OPCTs have larger bias and variance than other algorithms for all parameter values. Increasing parameter values again reduces the variance although to a lesser extent, but bias for the second target is increased at the largest values.

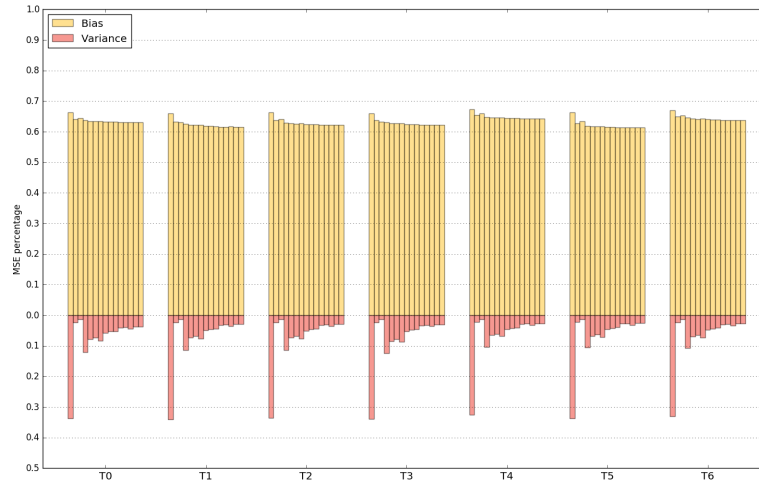


Fig. 7. Bias-variance decomposition of MSE for every target of the *Vegetation Condition* data set. Algorithms in order: PCT, bagging of PCTs, random Forest of PCTs and OPCTs with e0.1d0.5, e0.1d0.9, e0.1d1, e0.2d0.5, e0.2d0.9, e0.2d1, e0.5d0.5, e0.5d0.9, e0.5d1, e1d0.5, e1d0.9, e1d1.

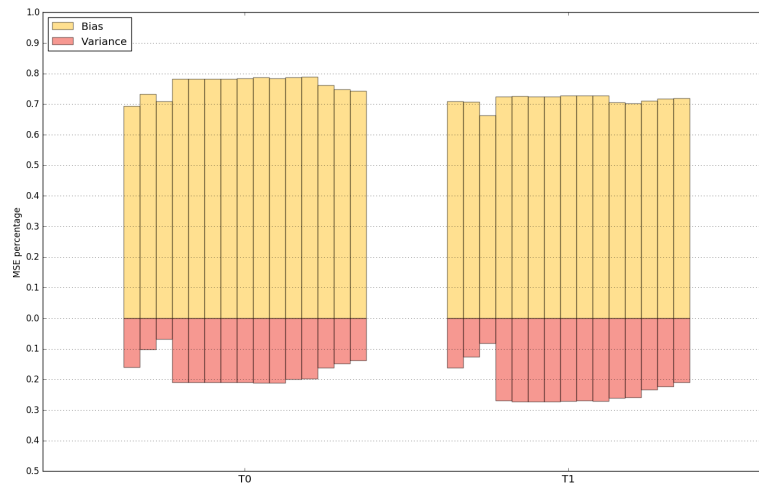


Fig. 8. Bias-variance decomposition of MSE for every target of the *Sigma real* data set. Algorithms in order: PCT, bagging of PCTs, random Forest of PCTs and OPCTs with e0.1d0.5, e0.1d0.9, e0.1d1, e0.2d0.5, e0.2d0.9, e0.2d1, e0.5d0.5, e0.5d0.9, e0.5d1, e1d0.5, e1d0.9, e1d1.

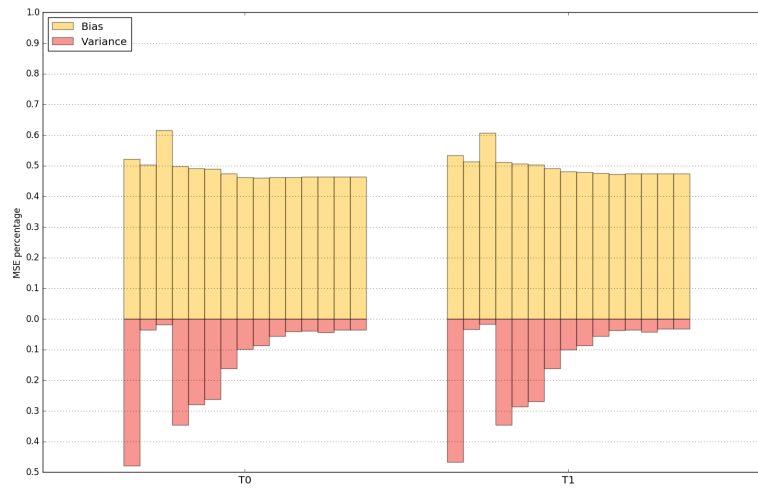


Fig. 9. Bias-variance decomposition of MSE for every target of the *Forestry-Slivnica-LandSat* data set. Algorithms in order: PCT, bagging of PCTs, random Forest of PCTs and OPCTs with $e0.1d0.5$, $e0.1d0.9$, $e0.1d1$, $e0.2d0.5$, $e0.2d0.9$, $e0.2d1$, $e0.5d0.5$, $e0.5d0.9$, $e0.5d1$, $e1d0.5$, $e1d0.9$, $e1d1$.

The variance results for the *Forestry-Slivnica-LandSat* data set in Figure 9 are similar to the results for the *Vegetation condition* data set, with the notable difference that the bias of random forest of PCTs is substantially larger than the bias of the other algorithms. This occurred also in a few data sets which are shown in Appendix 5 and is not surprising, since individual trees in a random forest only use a fraction of all attributes to learn.

The results show that the difference in the predictive performance mainly arises from the difference in their variances. With a few exceptions, PCTs have the largest variance, followed by OPCTs with small parameter values. Increasing the parameter values greatly reduces the variance of OPCTs, while bagging and especially random forests have the smallest variances. The biases of different methods are generally very similar. Single PCTs and random forests have larger values on some data sets, while biases of OPCTs produced with large parameter values are occasionally smaller.

4.4. Interpretability of OPCTs

OPCTs, like option trees in general, offer a much higher degree of interpretability than ensemble methods. This is expressed through both the fact that the "ensemble" of an option tree is represented in a compact form, i.e., a single tree, as well as the fact that many of the embedded trees overlap. Additionally, the regular PCT that would be learned from the same data is always present in the OPCT, as described in Section 2. A PCT and an OPCT learned on the *EDM* data set, and their relationship are illustrated in Figure 10.

Providing a domain expert with an option tree gives them a lot of choices with regards to the model. They can observe the selected options and attempt to determine which of the selected options were selected due to their actual importance as opposed to the sampling of the data set or other artifacts of the data. If they are able to discard all but one of the options in each of the option nodes, we can collapse the OPCT into a single tree,

specifically a PCT, which corresponds not only to the data, but also to the knowledge of the domain expert. In terms of the predictive clustering framework, this means selecting only one of the overlapping hierarchical clusters when multiple clusters are presented. This approach also has the advantage that the domain expert need not be available for interaction when the model is learned, but can assess the OPCT and chose the preferred options later on.

This process can also be looked at through a different lens. As we have introduced option trees (in part) to address myopia, by considering more options and later deciding on only one of them in each option node, we are essentially "looking ahead" of just the one split and utilizing, in the case of the domain expert, additional domain knowledge.

However, instead of using domain knowledge by proxy of interaction with a domain expert, we could also collapse the learned OPCT to a single PCT by using additional unseen data, i.e., by calculating an unbiased estimate of the predictive performance of the different options present in the OPCT. Since the collection and preparation of additional data examples could be expensive to the point of infeasibility, we could introduce a modified experimental setup. Part of the training data could be separated into a validation set which would not be used for the initial learning of the OPCT, but would be utilized to determine which of the selected options, and consequently embedded trees, has the best predictive performance. We would then collapse the OPCT into a PCT according to this validation set, after which we would test the obtained PCT on the test set. In this scenario, we could not only study the effect of myopia by comparing the collapsed OPCT to a PCT learned on the entire training data set, but also observe the effect of averaging multiple predictions on the predictive performance by comparing the collapsed PCT and the original (pseudo-ensemble) OPCT.

5. Conclusions

In this work, we propose an algorithm for learning option predictive clustering trees (OPCTs) for the task of multi-target regression (MTR). In contrast to standard regression, where the output is a single scalar value, in MTR the output is a data structure – a tuple/vector of numeric variables. We consider learning of a global model, that is a single model that predicts all target variables simultaneously.

More specifically, we propose OPCTs to address the myopia of the standard greedy PCT learning algorithm. OPCTs have the possibility to construct option nodes, i.e., nodes with a set of alternative sub-nodes, each containing a different split. These option nodes are constructed in the cases when the heuristic scores of the candidate splits are close to each other. Furthermore, OPCTs, and option trees in general, can be regarded as a condensed representation of an ensemble of trees.

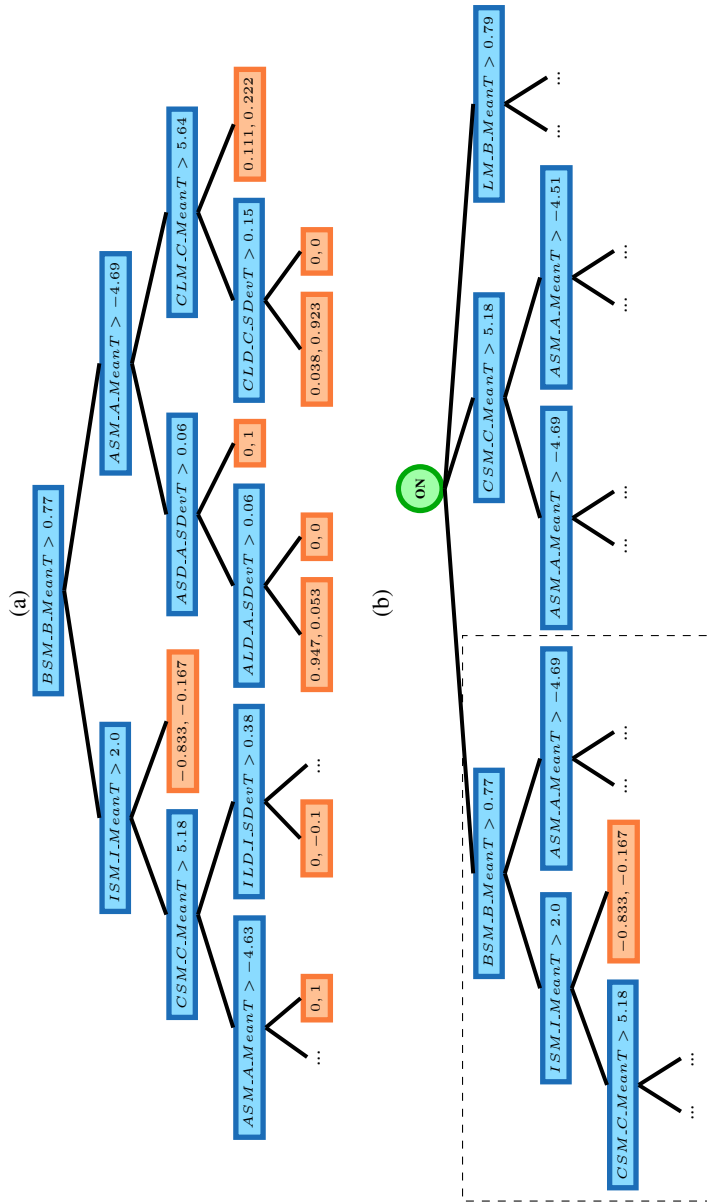


Fig. 10. A regular PCT (a) and an OPCT (b) learned on the EDM data set. The left child of a split node corresponds to the subtree where the test is satisfied. Note that the regular PCT is included in the OPCT as the subtree enclosed in the dashed rectangle.

The proposed method was experimentally evaluated on 11 benchmark MTR data sets. We selected two parameter configurations, OPCTe1d1 ($\varepsilon = 1, d = 1$), the largest OPCT which offers the best predictive performance, and OPCTe0.2d1 ($\varepsilon = 0.2, d = 1$), which provides a balance between predictive performance and efficiency. We compared this pair of OPCTs to regular PCTs and two ensemble learning methods – bagging and random forests of PCTs. The evaluation revealed that both OPCTe1d1 and OPCTe0.2d1 yield statistically significantly better predictive performance than single PCT. Next, the predictive performance of OPCTe1d1 is not statistically significantly different than that of the other two ensemble methods, however, it is slower to train and produces more leaves than the other compared methods. While the predictive performance of OPCTe0.2d1 was not on par with ensemble methods, it was not statistically significantly less efficient than PCTs, while offering significantly better performance.

We also performed parameter sensitivity analysis. The results show that both parameters that control the number of option nodes need large values to achieve good predictive performance. However, limiting the number of option nodes can lead to a more favorable trade-off between performance and efficiency.

We performed bias-variance decomposition of the mean squared error to determine the source of errors of the different methods. The results show the bias of the compared methods exhibits very similar behavior. OPCTs with larger parameter values were occasionally better in this regard, while random forests of PCTs sometimes had larger bias than other methods. On the other hand, the variance component often differed greatly between the algorithms. It was almost always the largest for single PCTs and smallest for random forests of PCTs. Increasing the parameter values for OPCTs reduced the variance and in some cases OPCTe1d1 achieved variance similar to that of the ensemble methods.

Finally, through an example, we illustrated the interpretability of the constructed OPCTs: they offer a multifaceted view on the data at hand.

We plan to extend this work along several directions. We will evaluate the OPCTs in the single tree context, i.e., we will use the induction of OPCTs as a beam-search algorithm for tree induction. Next, we will evaluate the influence of the two parameters at a more fine grained resolution. Finally, we will extend the algorithm towards other output types, i.e., machine learning tasks, such as multi-label classification, hierarchical multi-label classification and time series prediction.

Acknowledgments. We acknowledge the financial support of the European Commission through the grants ICT-2013-612944 MAESTRA and ICT-2013-604102 HBP, as well as the support of the Slovenian Research Agency through a young researcher grant (AO, TSP), the program Knowledge Technologies (P2-0103) and the project J2-9230.

References

1. Appice, A., Džeroski, S.: Stepwise induction of multi-target model trees. In: *Machine Learning: ECML 2007*, LNCS, vol. 4701, pp. 502–509. Springer (2007)
2. Bakir, G.H., Hofmann, T., Schölkopf, B., Smola, A.J., Taskar, B., Vishwanathan, S.V.N.: *Predicting structured data*. Neural Inf. Processing, The MIT Press (2007)
3. Blockeel, H., Struyf, J.: Efficient algorithms for decision tree cross-validation. *Journal of Machine Learning Research* 3, 621–650 (2002)

4. Borhani, H., Varando, G., Bielza, C., Larrañaga, P.: A survey on multi-output regression. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 5(5), 216–233 (2015)
5. Breiman, L.: Bagging predictors. *Machine Learning* 24(2), 123–140 (1996)
6. Breiman, L.: Random forests. *Machine Learning* 45(1), 5–32 (2001)
7. Breiman, L., Friedman, J., Olshen, R., Stone, C.J.: *Classification and Regression Trees*. Chapman & Hall/CRC (1984)
8. Breskvar, M., Kocev, D., Džeroski, S.: Ensembles for multi-target regression with random output selections. *Machine Learning* 107(11), 1673–1709 (2018)
9. Buntine, W.: Learning classification trees. *Statistics and Computing* 2(2), 63–73 (1992)
10. Corizzo, R., Pio, G., Ceci, M., Malerba, D.: DENCAST: distributed density-based clustering for multi-target regression. *Journal of Big Data* 6(1), 43 (2019)
11. Demšar, D., Debeljak, M., Džeroski, S., Lavigne, C.: Modelling pollen dispersal of genetically modified oilseed rape within the field. In: *The Annual Meeting of the Ecological Society of America*. p. 152 (2005)
12. Demšar, D., Džeroski, S., Larsen, T., Struyf, J., Axelsen, J., Bruns-Pedersen, M., Krogh, P.H.: Using multi-objective classification to model communities of soil. *Ecological Modelling* 191(1), 131–143 (2006)
13. Demšar, J.: Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research* 7, 1–30 (2006)
14. Džeroski, S., Demšar, D., Grbovič, J.: Predicting chemical parameters of river water quality from bioindicator data. *Applied Intelligence* 13(1), 7–17 (2000)
15. Geman, S., Bienenstock, E., Doursat, R.: Neural networks and the bias/variance dilemma. *Neural Computing* 4(1), 1–58 (Jan 1992)
16. Gjorgjioski, V., Džeroski, S., White, M.: Clustering analysis of vegetation data. Tech. Rep. 10065, Jožef Stefan Institute (2008)
17. Ikonovska, E., Gama, J., Zenko, B., Džeroski, S.: Speeding-up hoeffding-based regression trees with options. In: *Proceedings of the 28th International Conference on Machine Learning, ICML 2011*. pp. 537–544 (2011)
18. Kampichler, C., Džeroski, S., Wieland, R.: Application of machine learning techniques to the analysis of soil ecological data bases: relationships between habitat features and Collembolan community characteristics. *Soil Biology and Biochemistry* 32(2), 197–209 (2000)
19. Karalič, A.: First order regression. Ph.D. thesis, Faculty of Computer Science, University of Ljubljana, Ljubljana, Slovenia (1995)
20. Kocev, D., Ceci, M.: Ensembles of extremely randomized trees for multi-target regression. In: *Discovery Science: 18th International Conference (DS 2015)*. LNCS, vol. 9356, pp. 86–100 (2015)
21. Kocev, D., Džeroski, S., White, M., Newell, G., Griffioen, P.: Using single- and multi-target regression trees and ensembles to model a compound index of vegetation condition. *Ecological Modelling* 220(8), 1159–1168 (2009)
22. Kocev, D., Struyf, J., Džeroski, S.: Beam search induction and similarity constraints for predictive clustering trees. In: *Proc. of the 5th Intl Workshop on Knowledge Discovery in Inductive Databases KDID - LNCS 4747*. pp. 134–151 (2007)
23. Kocev, D., Vens, C., Struyf, J., Džeroski, S.: Ensembles of multi-objective decision trees. In: *ECML '07: Proceedings of the 18th European Conference on Machine Learning – LNCS 4701*. pp. 624–631. Springer (2007)
24. Kocev, D., Vens, C., Struyf, J., Džeroski, S.: Tree ensembles for predicting structured outputs. *Pattern Recognition* 46(3), 817–833 (2013)
25. Kohavi, R., Kunz, C.: Option decision trees with majority votes. In: *Proceedings of the 14th International Conference on Machine Learning*. pp. 161–169. ICML '97, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA (1997)

26. Osojnik, A., Džeroski, S., Kocev, D.: Option predictive clustering trees for multi-target regression. In: Discovery Science: 19th International Conference (DS 2016). LNCS, vol. 9956, pp. 118–133 (2016)
27. Spyromitros-Xioufis, E., Tsoumakas, G., Groves, W., Vlahavas, I.: Multi-target regression via input space expansion: Treating targets as inputs. *Machine Learning* 104(1), 55–98 (2016)
28. Stojanova, D.: Estimating forest properties from remotely sensed data by using machine learning. Master's thesis, Jožef Stefan IPS, Ljubljana, Slovenia (2009)
29. Stojanova, D., Panov, P., Gjorgjioski, V., Kobler, A., Džeroski, S.: Estimating vegetation height and canopy cover from remotely sensed data with machine learning. *Ecological Informatics* 5(4), 256–266 (2010)
30. Struyf, J., Džeroski, S.: Constraint based induction of multi-objective regression trees. In: Proc. of the 4th International Workshop on Knowledge Discovery in Inductive Databases KDID - LNCS 3933. pp. 222–233. Springer (2006)
31. Tsoumakas, G., Spyromitros-Xioufis, E., Vrekou, A., Vlahavas, I.: Multi-target regression via random linear target combinations. In: Machine Learning and Knowledge Discovery in Databases: ECML-PKDD 2014, LNCS 8726. pp. 225–240 (2014)

Appendix

Bias-Variance Graphs

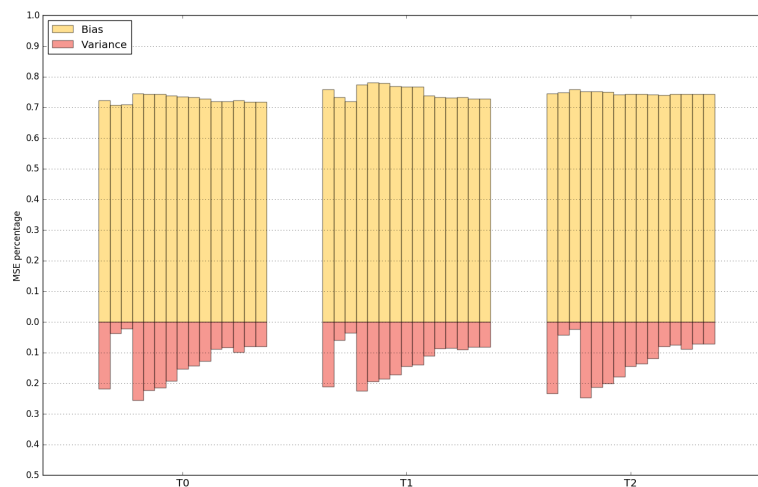


Fig. 11. Bias-variance decomposition of MSE for every target of the *Collembola* data set. Algorithms in order: PCT, bagging of PCTs, random Forest of PCTs and OPCTs with e0.1d0.5, e0.1d0.9, e0.1d1, e0.2d0.5, e0.2d0.9, e0.2d1, e0.5d0.5, e0.5d0.9, e0.5d1, e1d0.5, e1d0.9, e1d1.

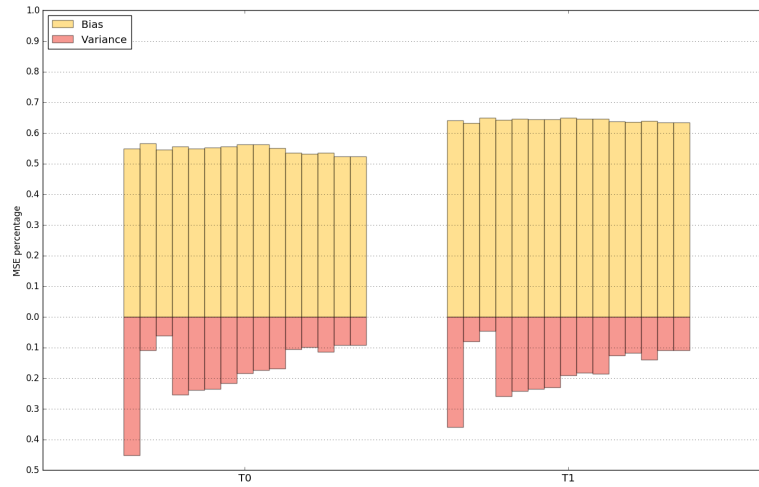


Fig. 12. Bias-variance decomposition of MSE for every target of the *EDM* data set. Algorithms in order: PCT, bagging of PCTs, random Forest of PCTs and OPCTs with e0.1d0.5, e0.1d0.9, e0.1d1, e0.2d0.5, e0.2d0.9, e0.2d1, e0.5d0.5, e0.5d0.9, e0.5d1, e1d0.5, e1d0.9, e1d1.

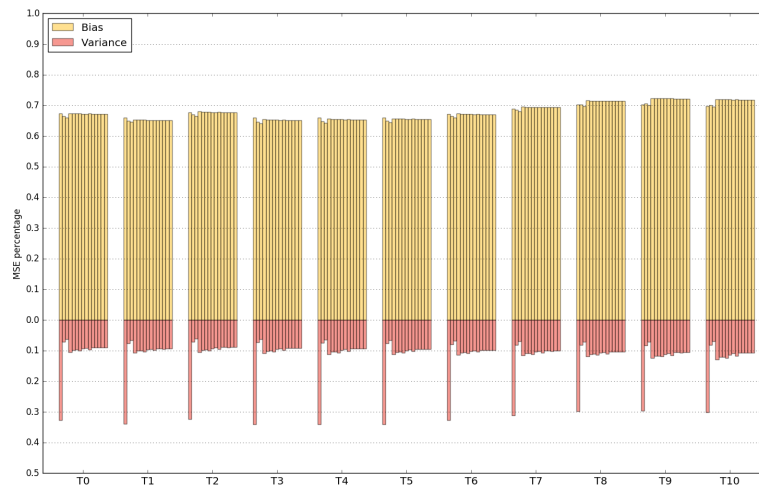


Fig. 13. Bias-variance decomposition of MSE for every target of the *Forestry-Kras* data set. Algorithms in order: PCT, bagging of PCTs, random Forest of PCTs and OPCTs with e0.1d0.5, e0.1d0.9, e0.1d1, e0.2d0.5, e0.2d0.9, e0.2d1, e0.5d0.5, e0.5d0.9, e0.5d1, e1d0.5, e1d0.9, e1d1.

Additional tables

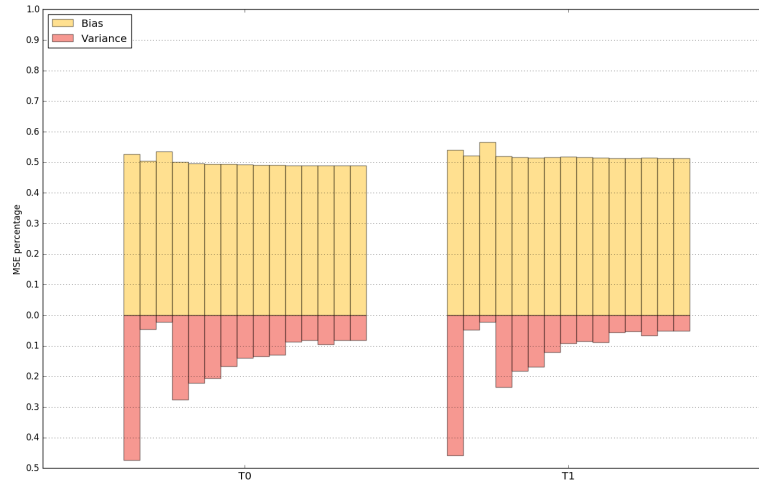


Fig. 14. Bias-variance decomposition of MSE for every target of the *Forestry-Slivnica-IRS* data set. Algorithms in order: PCT, bagging of PCTs, random Forest of PCTs and OPCTs with $\epsilon 0.1d0.5$, $\epsilon 0.1d0.9$, $\epsilon 0.1d1$, $\epsilon 0.2d0.5$, $\epsilon 0.2d0.9$, $\epsilon 0.2d1$, $\epsilon 0.5d0.5$, $\epsilon 0.5d0.9$, $\epsilon 0.5d1$, $\epsilon 1d0.5$, $\epsilon 1d0.9$, $\epsilon 1d1$.

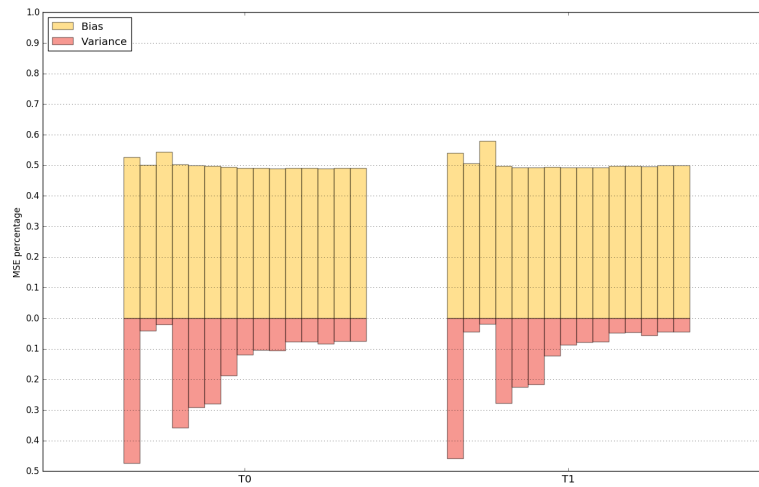


Fig. 15. Bias-variance decomposition of MSE for every target of the *Forestry-Slivnica-SPOT* data set. Algorithms in order: PCT, bagging of PCTs, random Forest of PCTs and OPCTs with $\epsilon 0.1d0.5$, $\epsilon 0.1d0.9$, $\epsilon 0.1d1$, $\epsilon 0.2d0.5$, $\epsilon 0.2d0.9$, $\epsilon 0.2d1$, $\epsilon 0.5d0.5$, $\epsilon 0.5d0.9$, $\epsilon 0.5d1$, $\epsilon 1d0.5$, $\epsilon 1d0.9$, $\epsilon 1d1$.

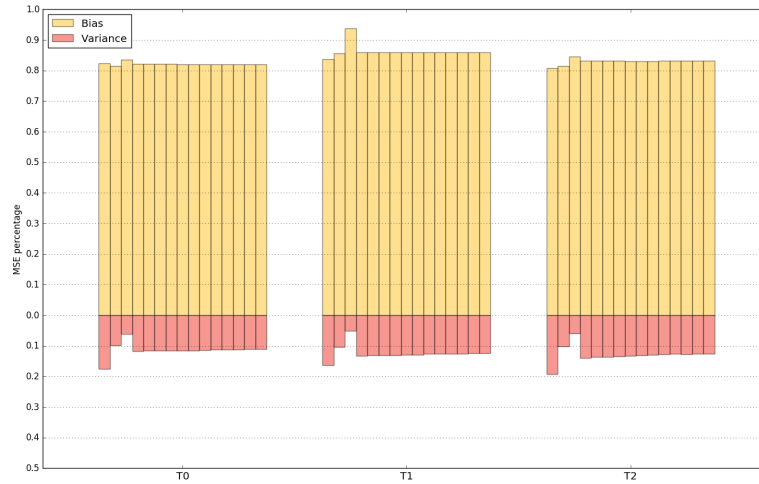


Fig. 16. Bias-variance decomposition of MSE for every target of the *Soil quality* data set. Algorithms in order: PCT, bagging of PCTs, random Forest of PCTs and OPCTs with e0.1d0.5, e0.1d0.9, e0.1d1, e0.2d0.5, e0.2d0.9, e0.2d1, e0.5d0.5, e0.5d0.9, e0.5d1, e1d0.5, e1d0.9, e1d1.

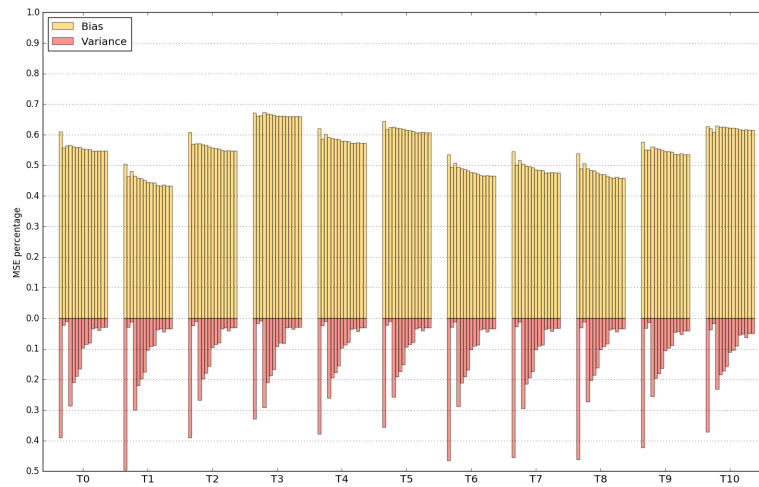


Fig. 17. Bias-variance decomposition of MSE for every target of the *Vegetation Clustering* data set. Algorithms in order: PCT, bagging of PCTs, random Forest of PCTs and OPCTs with e0.1d0.5, e0.1d0.9, e0.1d1, e0.2d0.5, e0.2d0.9, e0.2d1, e0.5d0.5, e0.5d0.9, e0.5d1, e1d0.5, e1d0.9, e1d1.

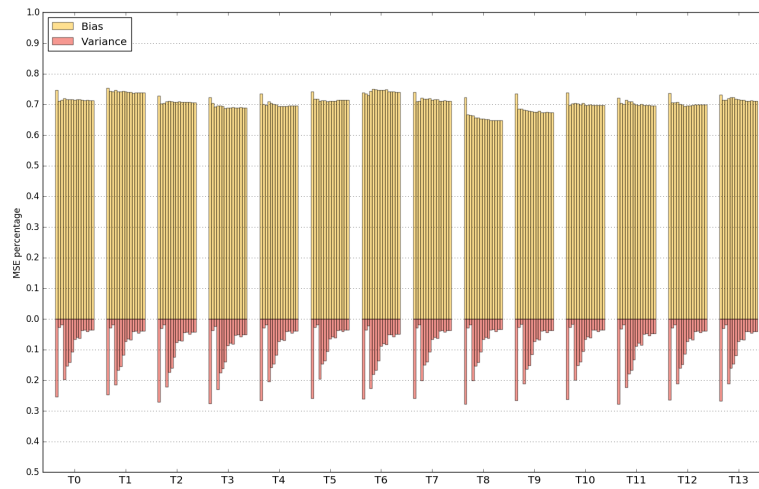


Fig. 18. Bias-variance decomposition of MSE for every target of the *Water quality* data set. Algorithms in order: PCT, bagging of PCTs, random Forest of PCTs and OPCTs with e0.1d0.5, e0.1d0.9, e0.1d1, e0.2d0.5, e0.2d0.9, e0.2d1, e0.5d0.5, e0.5d0.9, e0.5d1, e1d0.5, e1d0.9, e1d1.

	$\epsilon = 0.1$		$\epsilon = 0.2$		$\epsilon = 0.5$		$\epsilon = 1$	
	$d = 0.5$	$d = 0.9$	$d = 1$	$d = 0.5$	$d = 0.9$	$d = 1$	$d = 0.5$	$d = 0.9$
collembolaV2	1689	4137	4137	4260	7513	7559	14595	15836
edm1	202	239	239	283	464	601	623	3557
Forestry_Kras	447646	830108	1002953	661259	1498692	1890340	1157139	2273266
Forestry_LIDAR_IRS	10804	18389	18921	28213	41939	45982	46523	125903
Forestry_LIDAR_Landsat	8372	31281	37406	33537	55194	70284	136907	244006
Forestry_LIDAR_Spot	4598	6932	6932	17374	58422	73037	71913	129843
sigmeareal	165	165	165	165	202	236	732	1164
soil_quality	3031	6385	8038	4617	9767	12184	23214	68372
VegetationCondition	81368	165625	192474	136516	294664	348911	368187	663801
vegetationV2	82363	86915	132580	137050	315931	344118	435360	1247710
water-quality	4009	7836	8736	6513	12895	16901	13614	51671

Table 2. Number of leaves in OPTs with different parameter values.

	$\epsilon = 0.1$		$\epsilon = 0.2$		$\epsilon = 0.5$		$\epsilon = 1$	
	$d = 0.5$	$d = 0.9$	$d = 1$	$d = 0.5$	$d = 0.9$	$d = 1$	$d = 0.5$	$d = 0.9$
collembolaV2	9	20	20	34	39	39	39	71
edm1	10	13	13	14	26	30	36	91
Forestry_Kras	28	62	80	44	95	109	87	111
Forestry_LIDAR_IRS	11	20	21	28	53	57	54	110
Forestry_LIDAR_Landsat	7	11	12	12	21	25	60	97
Forestry_LIDAR_Spot	5	5	5	12	52	67	60	107
sigmeareal	1	1	1	1	3	4	12	19
soil_quality	5	11	12	7	18	22	47	87
VegetationCondition	14	31	33	28	49	60	71	104
vegetationV2	10	12	15	18	33	35	51	95
water-quality	9	24	26	15	36	42	41	109

Table 3. Number of option nodes in OPTs with different parameter values.

	$\epsilon = 0.1$		$\epsilon = 0.2$		$\epsilon = 0.5$		$\epsilon = 1$			
	$d = 0.5$	$d = 0.9$	$d = 0.5$	$d = 0.9$	$d = 0.5$	$d = 0.9$	$d = 0.5$	$d = 0.9$		
collembolaV2	0.68	1.17	1.12	1.71	1.56	3.07	3.23	8.39	11.83	6.23
edm1	0.22	0.23	0.23	0.23	0.24	0.34	0.41	0.3	0.44	0.42
Forestry_Kras	449.71	781.2	653.2	1380.88	1065.94	2083.34	2117.13	1920.28	2069.53	2110.88
Forestry_LIDAR_IRS	1.13	1.89	2.64	3.7	7.89	10.02	12.8	7.29	11.58	11.39
Forestry_LIDAR_Landsat	6.04	13.55	14.15	22.11	51.93	92.44	104.78	74.78	121.95	124.78
Forestry_LIDAR_Spot	0.9	1.19	2.59	11.17	11.84	18.36	19.03	13.03	18.37	18.75
sigmeareal	0.58	0.52	0.44	0.35	7.1	5.99	6.09	6.05	5.19	1.12
soil_quality	1.0	1.5	1.25	2.0	12.02	19.25	19.35	10.47	21.25	20.01
VegetationCondition	20.15	40.44	30.65	64.23	80.19	144.58	154.24	121.42	174.21	166.67
vegetationV2	24.45	24.55	39.35	84.73	116.34	324.78	381.27	273.54	373.55	375.06
water-quality	0.57	0.74	0.65	0.98	1.02	5.83	5.73	2.31	11.96	6.07

Table 4. Learning times for OPCTs with different parameter values.

Tomaž Stepišnik is a PhD student at the Jožef Stefan International Postgraduate School and a young researcher at the Department of Knowledge Technologies, Jožef Stefan Institute. He completed his MSc studies in mathematics at the University of Ljubljana in 2016 and was awarded a Young researcher grant by the Slovenian national research agency. His research interests include the study and development of machine learning algorithms and their application. He was on a research visit at the University of Auckland, NZ in 2018 for 3 months. He reviews for various conferences on the topics of Machine Learning and Data Mining (KDD, DS, ECML PKDD, ...).

Aljaž Osojnik is a researcher at the Department of Knowledge Technologies, Jožef Stefan Institute, Ljubljana, Slovenia. He completed his MSc and BSc in mathematics at the Faculty of Mathematics and Physics, University of Ljubljana, Slovenia. He completed his PhD in 2017 at the Jožef Stefan International Postgraduate School, Ljubljana, Slovenia on the topic of structured output prediction on data streams. His research interests are related to online learning, learning from data streams and privacy-preserving data mining.

Sašo Džeroski received his Ph.D. degree in computer science from the University of Ljubljana in 1995. He is Head of the Department of Knowledge Technologies at the Jožef Stefan Institute and full professor at the JS International Postgraduate School (both in Ljubljana, Slovenia). He is also a visiting professor at the Phi-Lab of ESRIN, European Space Agency (Frascati, Italy). His research interests fall in the field of artificial intelligence (AI), focusing on the development of data mining and machine learning methods for a variety of tasks - including the prediction of structured outputs and the automated modeling of dynamical systems - and their applications to practical problems. In 2008, he was elected fellow of the European AI Society for his "Pioneering Work in the field of AI and Outstanding Service for the European AI community". In 2015, he became a foreign member of the Macedonian Academy of Sciences and Arts. In 2016, he was elected a member of Academia Europaea. He is currently serving as Chair of the Slovenian AI Society.

Dragi Kocev is a researcher at the Department of Knowledge Technologies, JSI. He completed his PhD in 2011 at the JSI Postgraduate School in Ljubljana on the topic of learning ensemble models for predicting structured outputs. He was a visiting research fellow at the University of Bari, Italy in 2014/2015. He has participated in several national Slovenian projects, the EU funded projects IQ, PHAGOSYS and HBP. He was co-coordinator of the FP7 FET Open project MAESTRA.

Received: September 28, 2019; Accepted: May 20, 2020.

A Robust Reputation System using Online Reviews^{*}

Hyun-Kyo Oh¹, Jongbin Jung², Sunju Park³, and Sang-Wook Kim⁴

¹ Samsung Electronics
hyunkyo.oh@samsung.com

² Stanford University
jongbin@stanford.edu

³ Yonsei University
boxenju@yonsei.ac.kr

⁴ Hanyang University
wook@agape.hanyang.ac.kr

Abstract. Evaluating sellers in an online marketplace is an important yet non-trivial task. Many online platforms such as eBay and Amazon rely on buyer reviews to estimate the reliability of sellers on their platform. Such reviews are, however, often biased by: (1) intentional attacks from malicious users and (2) conflation between a buyer’s perception of seller performance and item satisfaction. Here, we present a novel approach to mitigating these issues by decoupling measures of seller performance and item quality, while reducing the impact of malignant reviews. An extensive simulation study shows that our proposed method can recover seller reputations with high rank correlation even under assumptions of extreme noise.

Keywords: reputation, reviews, attacks.

1. Introduction

One of the major challenges for online marketplaces such as eBay.com is that of accurately measuring the reliability of sellers on their platform [3, 10, 14, 17, 25]. The most common implementation of this task takes the form of a reputation system in which buyers are tasked with evaluating their interaction with sellers on some common scale (e.g., a 5-star rating [15, 26–29].) These ratings are then aggregated and presented to future buyers as a proxy for a seller’s quality. While such measures have substantial influence on buyer behavior and overall marketplace dynamics, determining how reliable and robust they are to bias is still an open question. The impact of this problem is only getting larger as the presence and significance of online platforms in society increase. With the addition of more complex multi-agent systems and multi-faceted online marketplaces which often span global economies—such as Uber or AirBnB, the question of whether a reputation rating system can be robust to corruption and bias, and what the framework for measuring such robustness could be is becoming more important than ever. [2] Our work aims to address this question by first introducing a novel method for reliably measuring reputation from potentially corrupted ratings, whether maliciously intended or not, and subsequently proposing a general simulation model for quantifying the robustness of various reputation measurement strategies.

* Corresponding author: Sang-Wook Kim (wook@agape.hanyang.ac.kr)

Reputation systems aim to leverage the wisdom of crowds [11, 18], assuming that all participants understand and agree upon the common goal of transparently measuring the quality of a seller. This assumption, however, is flawed, as pointed out by recent studies that identify adversarial behavior in buyer ratings [28]. Various types of cheating behavior have been identified, both in theory and in practice, along with recommendations for how to account for such behavior in aggregating reviews [5, 7, 10, 19, 20, 23, 24].

In addition to the threat of malicious reviews, another—perhaps more subtle—issue for reputation systems is that the common goal of a review may not be immediately obvious to reviewers (e.g., buyers.) By evaluating an interaction, the reviewer is necessarily conflating multiple aspects of a transaction, only one of which is the seller’s performance. For example, a buyer could be extremely satisfied with an item, but find the seller’s competence in communication and execution problematic. In such a case, the buyer’s scoring of the transaction, whether high or low, will be at best a biased measure of seller performance and item quality. We address this issue by first modeling a buyer’s review as a combination of evaluating two factors: seller performance and item quality. Taking advantage of the fact that multiple sellers offer similar items and that one seller often offers multiple items, we propose an iterative method, which we call *RATING SEPARATION*, for teasing out each of the two factors that confound buyer reviews.

In order to create a comprehensive and robust reputation system, we further propose *INTEGRITY WEIGHTING*, a novel approach to mitigate the adverse effects of dishonest reviews. As the name suggests, the idea is to estimate the level of trustworthiness of each review, based on theories of buyers’ cheating behavior. Each review is subsequently re-weighted according to the estimated trustworthiness. We show that, while existing approaches mainly focus on a subset of possible attacks, *INTEGRITY WEIGHTING* is robust against a large pool of attack types and patterns that have been identified in literature.

Finally, we develop a simulation framework for evaluating the efficiency of a reputation system in online marketplaces. Compared to existing marketplace simulations [4, 12, 16], our model allows for the conflation of multiple factors in a buyer’s review. Using this framework, we evaluate the proposed reputation system and find it to be more robust and reliable in measuring seller reputation compared to existing methods.

In summary, contributions of this paper are threefold. First, we propose *RATING SEPARATION*, a method for disentangling, from a single score given by a buyer, the ratings for a seller and the item sold. Second, we present *INTEGRITY WEIGHTING*, a scheme for mitigating the risk of malignant agents on the platform. Third, we present a novel and comprehensive simulation approach for evaluating various policies and systems on an online marketplace platform. Using this framework, we are able to evaluate the practical efficacy of complex reputation systems—such as the ones we propose here. Additionally, this simulation framework allows us to better investigate existing methods, and identify their strengths and potential shortcomings.

2. Related work

Many online platforms have their own reputation systems evaluating the reliability of products or sellers by aggregating buyer reviews. However, these reputation systems are intrinsically vulnerable to malicious users who intentionally give unjustified reviews to

products or sellers. Numerous studies have been conducted to improve the robustness of reputation systems by mitigating the influence of such malicious users.

Online platforms can be largely categorized as either single-agent or multi-agent systems. In a single-agent system, a single provider curates a collection of products—such as movies on IMDb.com—and users are tasked with rating their experience with each product. Since there is a single provider, buyer ratings have a clear mapping to each product, and buyers are often limited to rating each product at most once. Existing studies of such single-agent systems give attention to eliminating anomalous (potentially malicious) ratings based on statistical analysis of the distribution of buyer ratings or ranking/grouping users based on their rating patterns in order to derive the weighted mean of ratings given by the users [6, 8, 21, 22].

In a multi-agent system, numerous sellers can provide multiple goods and services—as on Amazon.com. Unlike single-agent systems, in a multi-agent system, buyers can evaluate both the seller and their product. And as a consequence of repeated interactions, it is possible for buyers to provide more than one rating to the same seller, over a variety of products.

Both single-agent and multi-agent systems share the same goal of diminishing the risk of malicious ratings. However, due to the different graph structure between buyers and sellers being more complex, state-of-the-art strategies that work well for single-agent systems are often insufficient for multi-agent systems. One of the often discussed challenges for multi-agent systems is that of identifying malicious buyers (attackers) who—in coordination with associated sellers—aim to artificially manipulate the reputation of target sellers. This can take the form of either increasing the reputation of partnered sellers, or decreasing the reputation of competing sellers. As online platforms become more commonplace and complex, deceptive rating strategies have also evolved, making it harder to identify attackers. In response, most existing studies focus on not only filtering out statistically insignificant ratings but also finding suspicious buyer-seller relationships by detecting malicious behavioral patterns in their rating systems [1, 5, 6, 9, 13, 15, 20, 24, 26–30]. A less discussed issue for rating systems in multi-agent systems, however, is that of confounded ratings. As discussed in the previous section, benign users can still corrupt a reputation system by conflating their evaluation of a seller and a product. While previous studies deal with the issue of malicious users, we have yet to find studies that address the subtle, yet important, issue of confounded buyer ratings. In addition to addressing the more traditional concerns of malicious ratings, our approach aims to achieve robustness against the threat of such ambiguous ratings as well.

3. Proposed methods

Overall, we propose RATING SEPARATION AND INTEGRITY WEIGHTING (RS&IW), a system for retrieving reputations that are robust to confounded ratings and adversarial behavior. The first part of the system, RATING SEPARATION, decomposes the possibly confounded rating into individual components of seller and item ratings. The second part of the system, INTEGRITY WEIGHTING, extends the work of Oh et al. [23] to quantify the trustworthiness of each rating. In the following sections, we present the details of each method.

3.1. Decoupling ratings

A common goal for online marketplace providers is to identify the reputation, trustworthiness, and quality of active sellers and items that are traded on their platform. Formally, given a set of sellers $S = \{s_1, s_2, \dots\}$ and items $M = \{m_1, m_2, \dots\}$, a platform provider—such as eBay.com—would like to recover some function $\rho_S : S \rightarrow \mathbb{R}$ that ranks sellers and $\rho_M : M \rightarrow \mathbb{R}$ that ranks items. Since seller and item rankings are not readily measurable, platform providers will often estimate rankings by asking buyers to rate their interactions via some common scale (e.g., 5-star ratings). In other words, given a set of buyers $B = \{b_1, b_2, \dots\}$, platform providers observe a set of scores $Y = \{y_{s,m}^{(b)}\}$ when a buyer b rates their interaction with seller s to purchase item m . One important, yet subtle issue in this setting is that the observed scores $y_{s,m}^{(b)}$ do not directly measure values of either the seller ratings ($\rho_S(s)$) or item quality ($\rho_M(m)$), but some function of the two.

To motivate our approach, we first consider $y_{s,m}$, the unobserved *true* score corresponding to the evaluation of seller s with regard to item m . Next, we model this true score as a function of two terms, seller performance ($\rho_S(s)$) and item quality ($\rho_M(m)$):

$$y_{s,m} = f(\rho_S(s), \rho_M(m)). \quad (1)$$

Note that any single observation $y_{s,m}^{(b)}$ is only a noisy approximation of $y_{s,m}$, since different buyers will combine the two components in a different way.⁵ The first part of our proposed method involves an iterative clustering of the observed $y_{s,m}^{(b)}$ to decouple and estimate each component $\rho_S(s)$ and $\rho_M(m)$.

Initial clustering and estimation of ρ_S First, define the set $B_{s,m} \subset B$ as the set of buyers who have rated an interaction of purchasing item m from seller s (i.e., $B_{s,m} = \{b_k \in B \mid \exists y_{s,m}^{(b_k)} \in Y\}$). For each seller-item pair $(s, m) \in S \times M$, we initially estimate $y_{s,m}$ via the sample mean

$$\bar{y}_{s,m} = \frac{1}{|B_{s,m}|} \sum_{b \in B_{s,m}} y_{s,m}^{(b)}. \quad (2)$$

Next, we utilize the fact that multiple sellers can, and often do, offer the same items, to cluster sellers and estimate $\rho_S(s)$. For a specific seller s_i , let $M_{s_i} = \{m_j \in M \mid \exists y_{s_i,m_j}^{(b)} \in Y\}$ be the set of all items for which the seller s_i has been rated. Similarly, for some specific item m_j , let $S_{m_j} = \{s_i \in S \mid \exists y_{s_i,m_j}^{(b)} \in Y\}$ be the set of all sellers who have been rated with item m_j . We initially construct $K = |M|$ clusters of sellers, \mathcal{C}_S , by collecting sellers who happened to have ratings for the same item m . Formally, we write,

$$\mathcal{C}_S(k) = \{s_i \mid s_i \in S_{m_k}\}.$$

⁵ We also note that each buyer ratings $y_{s,m}^{(b)}$ could contain an additional bias component that depends on the specific buyer b . For example, some buyers may intentionally give higher or lower ratings, independent of seller or item quality, to satisfy their idiosyncratic goals. We specifically address this issue in the second part of our method, INTEGRITY WEIGHTING, which is presented in Section 3.2.

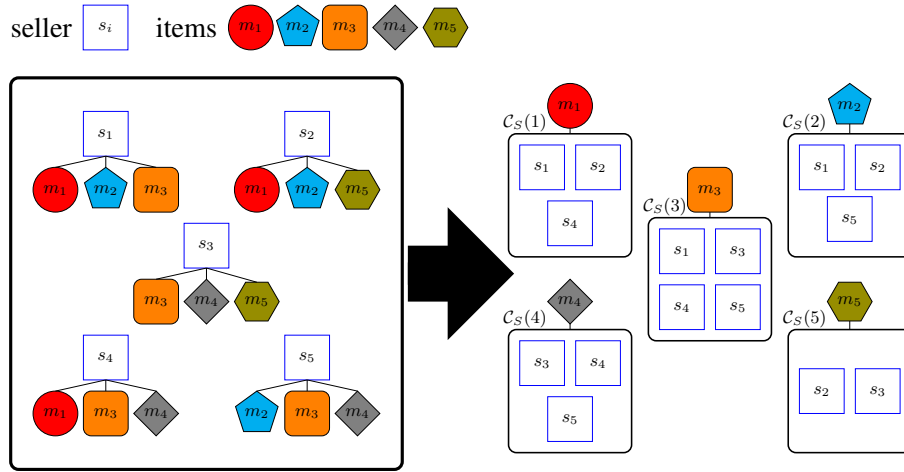


Fig. 1. An example of initial seller clustering

Further, let \mathcal{L}_{s_i} be the set of parameters k such that $\mathcal{C}_S(k)$ contains seller s_i as an element. In other words,

$$\mathcal{L}_{s_i} = \{k \mid s_i \in \mathcal{C}_S(k)\}. \quad (3)$$

An example of this initial clustering is presented in Fig. 1. Fig. 1 represents a platform of five sellers, $\{s_1, s_2, \dots, s_5\}$, and five items, $\{m_1, m_2, \dots, m_5\}$. As a result, sellers are initially organized into five clusters based on the items they offer: $\mathcal{C}_S(1) = \{s_1, s_2, s_4\}$, $\mathcal{C}_S(2) = \{s_1, s_2, s_5\}$, $\mathcal{C}_S(3) = \{s_1, s_3, s_4, s_5\}$, $\mathcal{C}_S(4) = \{s_3, s_4, s_5\}$, and $\mathcal{C}_S(5) = \{s_2, s_3\}$. Correspondingly, while not illustrated in Fig. 1, we can write out the sets of clusters that include each seller as $\mathcal{L}_{s_1} = \{1, 2, 3\}$, $\mathcal{L}_{s_2} = \{1, 2, 5\}$, $\mathcal{L}_{s_3} = \{3, 4, 5\}$, $\mathcal{L}_{s_4} = \{1, 3, 4\}$, and $\mathcal{L}_{s_5} = \{2, 3, 4\}$.

Next, we estimate the ranking of each seller $\rho_S(s)$ by averaging relative ratings within each cluster. Define $e_k : \mathcal{C}_S(k) \rightarrow \mathbb{R}$ as a scoring function for some seller $s_i \in \mathcal{C}_S(k)$ with respect to each item m_k , *relative to all other sellers in $\mathcal{C}_S(k)$* . Specifically, we define,

$$e_k(s_i) = \bar{y}_{s_i, m_k} - \frac{1}{|\mathcal{C}_S(k)| - 1} \sum_{s_j \in \mathcal{C}_S(k) \setminus s_i} \bar{y}_{s_j, m_k}. \quad (4)$$

Then, for each seller, we subsequently estimate ρ_S by computing

$$\hat{\rho}_S(s) = \frac{1}{|\mathcal{L}_s|} \sum_{k \in \mathcal{L}_s} e_k(s) \quad \forall s \in S. \quad (5)$$

In other words, a seller's rating, decoupled from item quality, is estimated by taking the average of all the relative scores achieved across clusters. Note that the range of $\hat{\rho}_S$ will vary depending on the range of the original scale implemented in the platform for recording buyer feedback and ratings. However, for the purpose of quantifying rankings amongst a

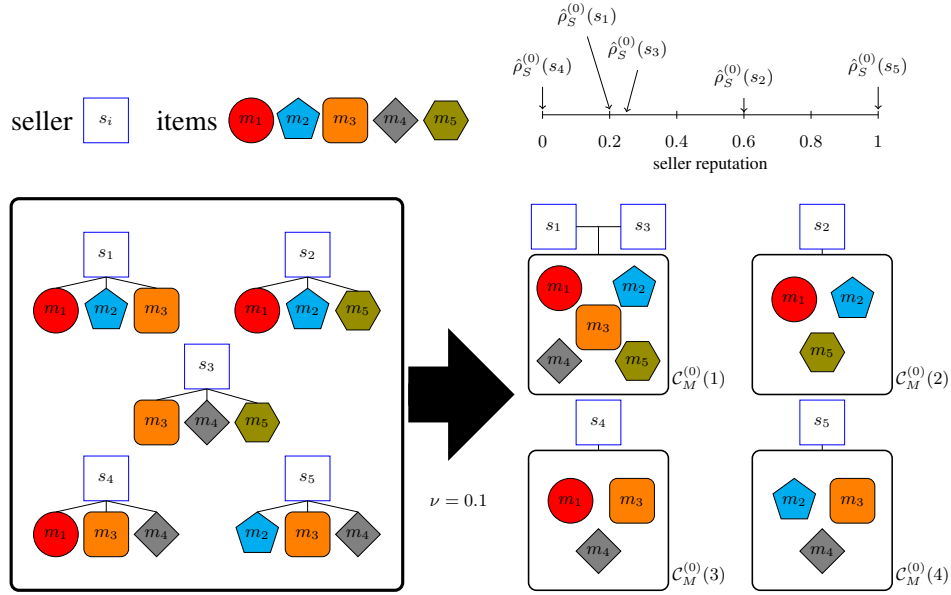


Fig. 2. An example of item clustering based on seller reputation

set of sellers, we can always normalize $\hat{\rho}_s$ to be within some desired range. For the following sections and the experiment described in Section 4, we use min-max normalization to restrict the range of $\hat{\rho}_S$ values to be in $[0, 1]$.

Iterative clustering and estimation of ρ_M Once we have the initial estimates $\hat{\rho}_S$, we can use these values to further cluster items in a similar manner, and subsequently estimate ρ_M . To formalize this iterative approach, we first denote an estimate of ρ_S and ρ_M at the t^{th} iteration as $\hat{\rho}_S^{(t)}$ and $\hat{\rho}_M^{(t)}$, respectively. Hence, our initial estimate from (5) is denoted $\hat{\rho}_S^{(0)}$, and similarly we let $e_k^{(0)}$ be our initial values of e_k computed in (4). At the t^{th} iteration, we create $K \in \mathbb{N}$ clusters of items by first grouping sellers such that seller s_i and seller s_j are in the same group if $|\hat{\rho}_S^{(t)}(s_i) - \hat{\rho}_S^{(t)}(s_j)| < \nu$, where ν is a parameter for determining the granularity and size of clusters and K is determined as a consequence of the distribution of $\hat{\rho}_S^{(t)}$. We define \mathcal{P}_k , the set of sellers in group k , such that $|\hat{\rho}_S^{(t)}(s_i) - \hat{\rho}_S^{(t)}(s_j)| < \nu$ for any $s_i \in \mathcal{P}_k$ and $s_j \in \mathcal{P}_k$. Then, items in the k^{th} cluster are defined as the items that have been ranked for the sellers who are in group \mathcal{P}_k . Formally, we write,

$$\mathcal{C}_M^{(t)}(k) = \{m_j \mid m_j \in M_{s_i} \forall s_i \in \mathcal{P}_k\}.$$

Similar to (3), let $\mathcal{L}_m^{(t)}$ be the set of parameters k such that $\mathcal{C}_M^{(t)}(k)$ contains item m as a member. In other words,

$$\mathcal{L}_m^{(t)} = \{k \mid m_j \in \mathcal{C}_M^{(t)}(k)\}.$$

Continuing our illustrative example from the previous section, Fig. 2 presents a numerical example of such item clustering, where $\nu = 0.1$. Based on the numerical values of $\hat{\rho}_S^{(0)}$, presented on the upper-right scale, sellers s_1 and s_3 are grouped together, while the other sellers form singletons, resulting in $K = 4$ clusters. Without loss of generalization, we can arbitrarily assign numbers $1 \leq k \leq 4$ to each cluster, defining sets $\mathcal{P}_1 = \{s_1, s_3\}$, $\mathcal{P}_2 = \{s_2\}$, $\mathcal{P}_3 = \{s_4\}$, $\mathcal{P}_4 = \{s_5\}$; and clusters $\mathcal{C}_M^{(0)}(1) = \{m_1, m_2, m_3, m_4, m_5\}$, $\mathcal{C}_M^{(0)}(2) = \{m_1, m_2, m_5\}$, $\mathcal{C}_M^{(0)}(3) = \{m_1, m_3, m_4\}$, and $\mathcal{C}_M^{(0)}(4) = \{m_2, m_3, m_4\}$. Correspondingly, the clusters that include each item is stored as $\mathcal{L}_{m_1}^{(0)} = \{1, 2, 3\}$, $\mathcal{L}_{m_2}^{(0)} = \{1, 2, 4\}$, $\mathcal{L}_{m_3}^{(0)} = \{1, 3, 4\}$, $\mathcal{L}_{m_4}^{(0)} = \{1, 3, 4\}$, and $\mathcal{L}_{m_5}^{(0)} = \{1, 2\}$.

Similar to (2), we compute a within-cluster mean $\bar{y}_{m,k}^{(t)}$ for each item m in cluster k by taking the average rating over each buyer and seller within the cluster. In other words,

$$\bar{y}_{m,k}^{(t)} = \frac{\sum_{s \in \mathcal{P}_k} \sum_{b \in B_{s,m}} y_{s,m}^{(b)}}{\sum_{s \in \mathcal{P}_k} |B_{s,m}|}. \quad (6)$$

Let $z_k^{(t)} : \mathcal{C}_M^{(t)}(k) \rightarrow \mathbb{R}$ be a scoring function for some item $m \in \mathcal{C}_M^{(t)}(k)$, relative to all other items in $\mathcal{C}_M^{(t)}(k)$. In particular, we define

$$z_k^{(t)}(m_i) = \bar{y}_{m_i,k}^{(t)} - \frac{1}{|\mathcal{C}_M^{(t)}(k)| - 1} \sum_{m_j \in \mathcal{C}_M^{(t)}(k) \setminus m_i} \bar{y}_{m_j,k}.$$

Then, for each item, we subsequently estimate ρ_M at iteration t by computing

$$\hat{\rho}_M^{(t)}(m) = \frac{1}{|\mathcal{L}_m^{(t)}|} \sum_{k \in \mathcal{L}_m^{(t)}} z_k^{(t)}(m) \quad \forall m \in M.$$

In other words, the quality of an item, decoupled from a seller's performance rating, is estimated by taking the average of all the relative scores achieved by that item across clusters. As in (5) for $\hat{\rho}_S$, the range of $\hat{\rho}_M^{(t)}$ will vary. For the following sections and the experiment described in Section 4, we use min-max normalization at each iteration t to restrict the range of $\hat{\rho}_M^{(t)}$ values to be in $[0, 1]$.

Iterative clustering and estimation of ρ_S Given values of $\hat{\rho}_M^{(t)}$, we can further improve our estimate of ρ_S by taking additional iterations. An iteration of computing $\hat{\rho}_S^{(t)}$ is very similar to the initial estimation procedure we describe for (5), with the primary difference being in how clusters $\mathcal{C}_S^{(t)}(k)$ are defined for $t > 0$.

Specifically, at the t^{th} iteration for $t > 0$, we create $K \in \mathbb{N}$ clusters of sellers by first grouping items such that item m_i and item m_j are in the same group if $|\hat{\rho}_M^{(t-1)}(m_i) - \hat{\rho}_M^{(t-1)}(m_j)| < \nu$. The set of items in group k , \mathcal{Q}_k , is defined such that $|\hat{\rho}_M^{(t-1)}(m_i) - \hat{\rho}_M^{(t-1)}(m_j)| < \nu$ for any $m_i \in \mathcal{Q}_k$ and $m_j \in \mathcal{Q}_k$. The set of sellers in the k^{th} cluster for $t > 0$ are then defined as

$$\mathcal{C}_S^{(t)}(k) = \{s_i \mid s_i \in S_{m_j} \forall m_j \in \mathcal{Q}_k\}, \quad t > 0.$$

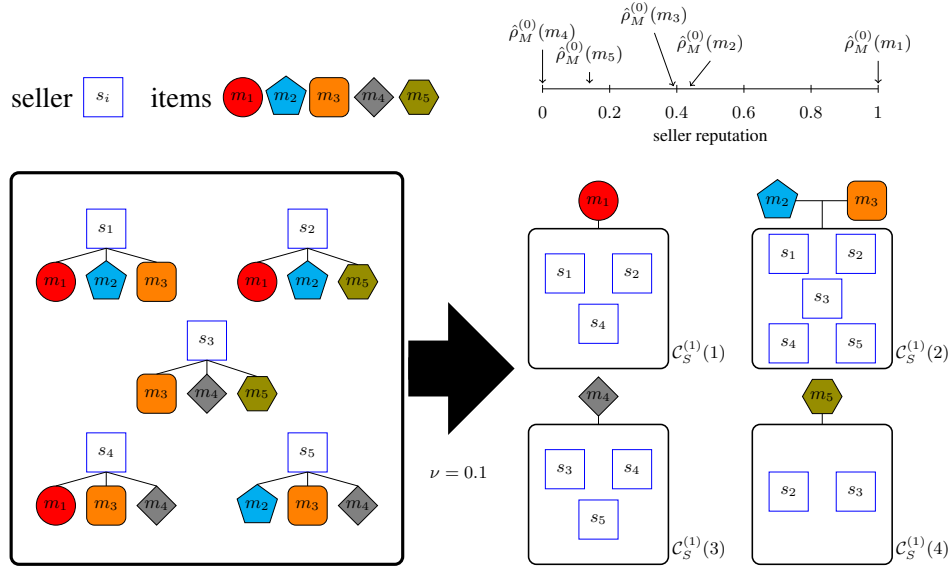


Fig. 3. An example of seller clustering based on item reputation

The set $\mathcal{L}_s^{(t)}$ is trivially defined similar to \mathcal{L}_s in (3).

To continue our example, Fig. 3 illustrates such a clustering of sellers with $\nu = 0.1$ at $t = 1$. Based on the numerical values of $\hat{\rho}_M^{(0)}$, presented on the upper-right scale, items m_2 and m_3 are grouped together, while the other items form singletons, resulting in $K = 4$ clusters. Again, we can assign numbers $1 \leq k \leq 4$ to each cluster, defining sets $\mathcal{Q}_1 = \{m_1\}$, $\mathcal{Q}_2 = \{m_2, m_3\}$, $\mathcal{Q}_3 = \{m_4\}$, $\mathcal{Q}_4 = \{m_5\}$; and clusters $C_S^{(1)}(1) = \{s_1, s_2, s_4\}$, $C_S^{(1)}(2) = \{s_1, s_2, s_3, s_4, s_5\}$, $C_S^{(1)}(3) = \{s_3, s_4, s_5\}$, and $C_S^{(1)}(4) = \{s_2, s_3\}$. Correspondingly, the clusters that include each seller is stored as $\mathcal{L}_{s_1}^{(1)} = \{1, 2\}$, $\mathcal{L}_{s_2}^{(1)} = \{1, 2, 4\}$, $\mathcal{L}_{s_3}^{(1)} = \{2, 3, 4\}$, $\mathcal{L}_{s_4}^{(1)} = \{1, 2, 3\}$, and $\mathcal{L}_{s_5}^{(1)} = \{2, 3\}$.

Within-cluster mean $\bar{y}_{s,k}^{(t)}$ for each seller s in cluster k is computed by taking the average rating over each buyer and item within the cluster. In other words,

$$\bar{y}_{s,k}^{(t)} = \frac{\sum_{m \in \mathcal{Q}_k} \sum_{b \in B_{s,m}} y_{s,m}^{(b)}}{\sum_{m \in \mathcal{Q}_k} |B_{s,m}|}, \quad t > 0. \quad (7)$$

Finally, $e_k^{(t)}$ and $\hat{\rho}_S^{(t)}$ for $t > 0$ are defined similar to the initial case of $t = 0$, following (4) and (5), but replacing the initial estimates $\bar{y}_{s,m}$ with the within-cluster average $\bar{y}_{s,k}^{(t)}$.

Complete algorithm for RATING SEPARATION As a stopping condition of the iterative algorithm, we define a tolerance parameter ε . After completing iteration $t > 0$, and given estimates $\hat{\rho}_S^{(t)}$ and $\hat{\rho}_M^{(t)}$, the algorithm is to advance to the next iteration $t + 1$ until $|\hat{\rho}_S^{(t)} - \hat{\rho}_S^{(t-1)}| < \varepsilon$ and $|\hat{\rho}_M^{(t)} - \hat{\rho}_M^{(t-1)}| < \varepsilon$. The overall procedure presented in this section is formally summarized in Algorithm 1.

Algorithm 1: Rating Separation (RS)

Input: set of buyers B , set of sellers S , set of items M , set of ratings Y , clustering range ν , convergence tolerance ε

Output: estimated quality scores for each seller and item $(\hat{\rho}_S, \hat{\rho}_M)$

$t \leftarrow 0$;

repeat

// Seller rating separation

if $t = 0$ **then**

| cluster sellers by item to build $\mathcal{C}_S^{(0)}$;

else

| cluster sellers by item score $\hat{\rho}_M^{(t-1)}$ and ν to build $\mathcal{C}_S^{(t)}$;

end

foreach $\mathcal{C}_S^{(t)}(k) \in \mathcal{C}_S^{(t)}$ **do**

foreach $s \in \mathcal{C}_S^{(t)}(k)$ **do**

| compute $e_k^{(t)}(s)$;

end

end

foreach $s \in S$ **do**

| compute $\hat{\rho}_S^{(t)}(s)$;

end

// Item rating separation

cluster items by seller score $\hat{\rho}_S^{(t)}$ and ν to build $\mathcal{C}_M^{(t)}$;

foreach $\mathcal{C}_M^{(t)}(k) \in \mathcal{C}_M^{(t)}$ **do**

foreach $m \in \mathcal{C}_M^{(t)}(k)$ **do**

| compute $z_k^{(t)}(m)$;

end

end

foreach $m \in M$ **do**

| compute $\hat{\rho}_M^{(t)}(m)$;

end

$t \leftarrow t + 1$;

until $t > 0$, $|\hat{\rho}_S^{(t)} - e_{t-1}^*| < \varepsilon$, $|\hat{\rho}_M^{(t)} - z_{t-1}^*| < \varepsilon$;

3.2. Mitigating adversarial reviews

A known issue with many buyer rating systems is that malicious actors may negatively affect the accuracy of scores through various cheating behavior. [5, 7, 10, 19, 20, 23, 24] Here, we mitigate such risk by proposing a method to score each review in terms of an estimated measure of trustworthiness—or *integrity*, which is then used to weigh each observed rating. Our proposed measure of trustworthiness considers three components: engagement, diversity, and anomaly. We calculate each component for every buyer, based on observed rating behavior across item categories. Formal definitions of each component are presented below.

Concretely, we define $\mathcal{G}^\ell \subset M$, a subset of items in category ℓ , as a collection of items which satisfy some predetermined criteria⁶. For example, \mathcal{G}^1 might be the collection of all electronics, while \mathcal{G}^2 might be all items classified as furniture. Then, let $B^\ell \subset B$ be the set of all buyers who have rated items in \mathcal{G}^ℓ . Similarly, we use $Y^\ell \subset Y$ to denote the subset of all ratings that were observed for items in category ℓ , while $Y^{\ell[b]} \subset Y^\ell$ further denotes the subset of ratings for items in category ℓ that were given by user b . In other words, we define

$$\begin{aligned} B^\ell &= \{b \in B \mid \exists y_{s,m_j}^{(b)} \in Y, m_j \in \mathcal{G}^\ell\} \\ Y^\ell &= \{y_{s,m_j}^{(b)} \in Y \mid m_j \in \mathcal{G}^\ell\} \\ Y^{\ell[b]} &= \{y_{s,m_j}^{(b)} \in Y \mid m_j \in \mathcal{G}^\ell, b \in B^\ell\}. \end{aligned}$$

Engagement A common measure for quantifying the trustworthiness of users on a typical platform is user engagement. For example, scores provided by a buyer who is highly engaged in the platform—purchasing items and rating interactions on a regular basis—is considered more reliable than that from a one-time visitor. Here, engagement is operationalized as the relative frequency of ratings given by each buyer b within a category ℓ :

$$\alpha_{b,\ell} = \frac{|Y^{\ell[b]}|}{|Y^\ell|} - \frac{|Y^\ell|}{|B^\ell|},$$

where $|Y^\ell|/|B^\ell|$ is the average number of ratings provided by each buyer in category ℓ . The corresponding user engagement weights $\alpha_{b,\ell}$ are further normalized to be within the range $[0, 1]$, using min-max normalization for each category ℓ .

Diversity Another consideration for a buyer's trustworthiness is the concentration of ratings. Conceptually, a buyer is considered more trustworthy if they interact with, and rate, a variety of different sellers, as opposed to repeatedly rating a small number of sellers. Thus, we quantify diversity as the proportion of unique sellers that the buyer has rated over all ratings the buyer has given within that category. Formally, let $S^{\ell[b]} \subset S$ be the subset sellers who have received a rating from user b , for at least one item in \mathcal{G}^ℓ . In other words

$$S^{\ell[b]} = \{s_i \mid \exists y_{s_i,m_j}^{(b)} \in Y, m_j \in \mathcal{G}^\ell\}.$$

Then, the diversity weight $\beta_{b,\ell}$ for a buyer b corresponding to item category ℓ is calculated as

$$\beta_{b,\ell} = \frac{|S^{\ell[b]}|}{|Y^{\ell[b]}|}.$$

Note that this quantification of diversity is relative to the total number of ratings made by the buyer. For example, a buyer who provided only one rating ($|Y^{\ell[b]}| = 1$) is considered to have high diversity ($\beta_{b,\ell} = 1$). As the buyer rates more transactions, $\beta_{b,\ell}$ will decrease whenever the buyer rates a seller that they have already rated previously. Similar to engagement weights, we normalize the corresponding diversity weights $\beta_{b,\ell}$ via min-max normalization within each item category \mathcal{G}^ℓ .

⁶ In this study, we use item categories as defined by the lowest-level grouping of items on eBay (<https://www.ebay.com/v/allcategories>).

Algorithm 2: Integrity weighting (IW)

Input: set of buyers B , set of ratings Y , set of item categories \mathcal{G}
Output: integrity-adjusted ratings, $\hat{y}_{s,m}^{(b)}$
// Compute integrity weights
foreach item group ℓ **do**
 foreach $b \in B^\ell$ **do**
 $w_{b,\ell} \leftarrow \alpha_{b,\ell} \times \beta_{b,\ell} \times \gamma_{b,\ell};$
 foreach $y_{s,m}^{(b)} \in Y^{\ell[b]}$ **do**
 $\hat{y}_{s,m}^{(b)} \leftarrow w_{b,\ell} \times y_{s,m}^{(b)};$
 end
 end
end

Anomaly We are also concerned with how much a buyer's rating of an item deviates or conforms to that of the general consensus of other buyers, which we refer to as anomaly. To quantify anomaly, we first consider the standardized distance of a buyer's rating for each item, from the overall distribution of ratings for that item. For any given item m , let μ_m and σ_m denote the average and standard deviation of ratings that the item received across all buyers and sellers. Then, for each rating $y_{s,m}^{(b)}$ given by buyer b for item m , we compute the normalized distance from the mean as:

$$\delta_{s,m}^{(b)} = \left| \frac{y_{s,m}^{(b)} - \mu_m}{\sigma_m} \right|,$$

where smaller values of $\delta_{s,m}^{(b)}$ indicate that the ratings given by buyer b for item m is similar and consistent with ratings given by other buyers for that same item. Then, $\gamma_{b,\ell}$, the anomaly weight for buyer b in item group ℓ is computed by taking the average of $\delta_{s,m}^{(b)}$ for all items $m \in \mathcal{G}^\ell$:

$$\gamma_{b,\ell} = \frac{1}{|Y^{\ell[b]}|} \sum_{Y^{\ell[b]}} \delta_{s,m}^{(b)}.$$

As with engagement weights and diversity weights, anomaly weights $\gamma_{b,\ell}$ are subsequently normalized via min-max normalization within each item category \mathcal{G}^ℓ .

Integrity weighted ratings Given the normalized weights $\alpha_{b,\ell}$, $\beta_{b,\ell}$, and $\gamma_{b,\ell}$ for engagement, diversity, and anomaly, respectively, we can compute a comprehensive integrity weight for each buyer b within item category \mathcal{G}^ℓ as

$$w_{b,\ell} = \alpha_{b,\ell} \times \beta_{b,\ell} \times \gamma_{b,\ell}.$$

Then, for an observed rating $y_{s,m}^{(b)}$ where $m \in \mathcal{G}^\ell$, we can compute an integrity-adjusted rating—where the observed rating is weighted by the estimated integrity of user b within category \mathcal{G}^ℓ as

$$\hat{y}_{s,m}^{(b)} = w_{b,\ell} \times y_{s,m}^{(b)}.$$

This procedure is formally summarized in Algorithm 2.

3.3. A comprehensive reputation score

Finally, we can compute reputation scores for sellers and items that are robust to confounded and adversarial ratings by combining RATING SEPARATION from 3.1 and INTEGRITY WEIGHTING from 3.2. This is achieved by replacing the raw ratings $y_{s,m}^{(b)}$ with their integrity weighted counter parts, $\hat{y}_{s,m}^{(b)}$, in (2), (6), and (7) of RATING SEPARATION.

4. Experiment

To evaluate the efficacy of the methods proposed, we further present a novel and comprehensive simulation framework. The contributions of our new approach are two fold. First, in contrast to existing literature we directly model item-level transactions. This enables our framework to distinguishing between a buyer's rating of sellers versus satisfaction of a specific item. Second, the proposed framework incorporates a comprehensive model of plausible adversarial behavior, allowing us to test how robust our reputation scoring systems are to numerous realistic attack scenarios.

4.1. A simulation framework for online marketplaces

The simulation framework we propose involves four components: three entities—items, sellers, buyers—and a model for how the different entities interact—transactions. A major advantage of our approach is that by explicitly modeling items, in addition to buyers and sellers, we can further capture the realistic dynamics that take place in online marketplaces.

An online marketplace is characterized by the number of items, buyers, and sellers on the platform. For our experiment, we consider two parameter regimes: small-size and large-size marketplaces. For the small-size marketplace, we set 1,000 items, 500 sellers, and 5,000 buyers. For the large-size marketplace, we set 2,000 items, 1,000 sellers, and 10,000 buyers. For each setting, we simulate 300 days of marketplace activity, where each buyer is limited to one transaction per day. We describe each component of the simulation in detail below.

Items An item is parameterized by its quality and categorization. The quality of an item is represented by a continuous score in the range $[0, 1]$. We allow for multiple hierarchical item categories.

For the purpose of our simulations in this study, we limit the hierarchy of item categories to three levels—top, middle, and bottom, which we find sufficient to represent many typical online marketplace categorizations realistically. In our experiments, we set three top-level categories, each of which consists of five mid-level subcategory. Each mid-level category is further classified in to six bottom-level subcategories. In total, there are 90 different unique item categories. We further assign items uniformly across different categories, so that the number of items available in each category is similar.

Sellers Sellers are parameterized by their capability and the items that they offer. We assume that seller capabilities follow a truncated normal distribution, bounded in $[0, 1]$, with mean 0.5 and standard deviation 0.25. Higher capability scores correspond to faster

delivery and better service, while lower capability scores correspond to late delivery and poor service.

An important characteristic of sellers, which is not often captured in existing simulation frameworks, is the variety of items that they offer. For example, while some sellers may focus on selectively offering only a few items in major categories, others may choose to offer a wide-selection of items across multiple categories. By modeling items as entities, and parameterizing sellers by the items they offer, the simulation framework we propose is capable of representing this diversity.

In our experiment, we assume that sellers offer items in one major category, along with items from up to three minor categories. To operationalize this assumption, for each seller we first sample one major item category, from which they offer between three and six items. Then, we sample between zero and three minor item categories, from which one to six items are subsequently sampled.

Buyers Buyers are parameterized by item categories of interest and the level of interest for each category. Each buyer is randomly assigned to 3 to 6 item categories of interest. The level of interest for each item category is assigned a continuous value in the range $[0, 1]$. We assume that buyers are more likely to purchase items in categories for which they have a higher level of interest. After each transaction, buyers will leave a single score rating as a function of item quality and seller capability.

Transactions Transactions represent the event in which a buyer purchases an item from a seller, and provides a rating. Each buyer is assigned a random purchase cycle, between zero and three days, which represents how often the buyer will participate in a transaction on the marketplace. Buyers are more likely to purchase items from sellers who offer items for which they have higher levels of interest in. This could result in unrealistically high-frequency transactions between the same buyer-seller pairs. To mitigate this issue, we require a *repurchase waiting time* of three, five, or ten days for the same buyer-seller pairs to have a repeat transaction.

4.2. Simulating malicious ratings

To evaluate the robustness of a reputation system in the presence of adversarial buyers, we model the behavior of malicious ratings. Based on existing literature [5, 15, 20, 24, 27, 29], we categorize adversarial buyers by three behavioral patterns and six attack strategies. The three behavioral patterns and six attack strategies are presented in Tables 1 and 2, respectively, along with a short description and relevant literature reference.

Any attacker will adopt one behavioral pattern and an attack strategy, allowing for a total of 18 possible attacker types. Compared to existing work, which only consider a limited subset of these 18 possible pairs, here we investigate the performance of a rating system under all 18 types of attacks. This is achieved by modeling each type of attack behavior and strategy within the simulation framework presented in Section 4.1. In our experiment, we parameterize the intensity of attacks on a platform as the attack rate—the proportion of all ratings that are malicious. We compare results for varying attack rates, from 10% to 90%, in 10% increments.

Table 1. *Three categories of adversarial behavior patterns.*

Pattern	Description	Reference
Basic	Attackers consistently exhibit adversarial behavior—granting high ratings to conspiring sellers or low ratings to rival sellers.	[5, 7, 15, 20, 24, 26–29]
Camouflage	Attackers attempt to camouflage their adversarial intent by strategically mixing justified ratings with malicious ones. Under this behavioral scheme, attackers typically exhibit benign behavior in early interactions, and transition to adversarial activities at later stages.	[29]
Whitewashing	Attackers behave under the basic scheme, while subsequently creating multiple accounts on the platform to mitigate detection and create an illusion that their malicious ratings are socially validated.	[29]

Table 2. *Six categories of attack strategy. Each strategy is assigned a number which we use as a reference in the text.*

# Name	Description	Reference
1 Ballot stuffing (BS)	Attempting to boost the reputation of conspiring sellers by giving maximum ratings	[5, 20, 27]
2 Bad mouthing (BM)	Attempting to hurt the reputation of rival sellers by giving minimum ratings	[5, 20, 27]
3 BS & BM	Employing a mix of both ballot stuffing and bad mouthing	[5, 20, 27]
4 r -high shifting	Attempting to boost the reputation of conspiring sellers by giving ratings that are r points higher than the average ratings	[15, 20, 24]
5 r -low shifting	Attempting to hurt the reputation of rival sellers by giving ratings that are r points lower than average	[15, 20, 24]
6 r -high/low shifting	Employing both r -high shifting and r -low shifting strategies to boost reputation of conspiring sellers while simultaneously reducing the reputation of rival sellers	[15, 20, 24]

Note that the last three strategies— r -high, r -low, and r -high/low shifting—require that the attackers assign some distribution of benign buyer ratings for the target sellers.

5. Results

We use Spearman’s rank correlation coefficient to measure and compare how well a reputation system can recover the true capability rankings of sellers.⁷ Our results are presented in two parts. First, we investigate the efficacy of RATING SEPARATION (RS). To do so, we compare RS with a naive baseline approach of computing a seller’s reputation via simple average of observed ratings. Second, to test whether INTEGRITY WEIGHTING

⁷ Note that here, we focus on seller reputation, but the proposed methods and simulation framework could also be used for estimating item quality via trivial extension of this work.

(IW) and the combined approach of RATING SEPARATION AND INTEGRITY WEIGHTING (RS&IW) is truly robust to adversarial rating activity, we compare performance of each method to existing mitigation techniques.

5.1. RATING SEPARATION (RS) performance

First, we evaluate the performance of RATING SEPARATION in recovering true seller rankings. Because RATING SEPARATION in itself does not mitigate against adversarial ratings, for this section we focus on a simulated platform that assumes no malicious attacks.⁸ We compare performance under two different assumptions of marketplace parameters, as described in Section 4.1. As a baseline, we compute a naive measure of seller reputation by taking the average of all ratings that a seller received.

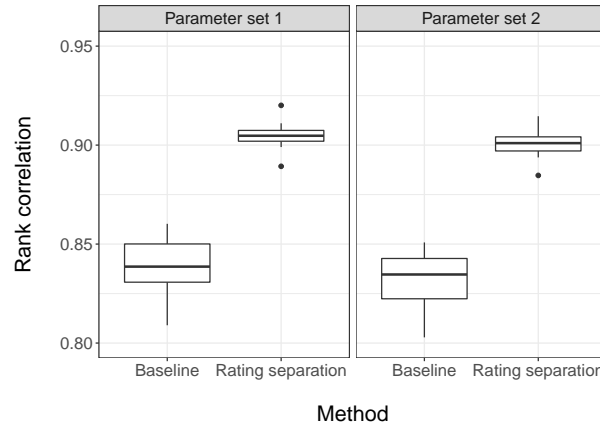


Fig. 4. Simulation results comparing a single iteration of RATING SEPARATION versus the baseline. Each box plot summarizes the results of 10 simulations. The y-axis shows the rank correlation between the estimated seller reputation and true seller capability, for each method. Column panels show the two simulation parameter settings. Overall, RATING SEPARATION consistently recovers the true ranking of seller capability more reliably than the baseline, with less variance across each trial.

In Fig. 4, we compare the rank correlation between estimated seller reputation and true seller capabilities for the baseline and RATING SEPARATION using just a single iteration. The box plot represents the distribution of rank correlation performance achieved for each method, across 10 simulations each.

We find that for every simulation trial, RATING SEPARATION achieved consistently higher rank correlation compared to the baseline. RATING SEPARATION also was more consistent in better recovering true seller rankings, demonstrated by the low variance in performance across simulations, compared to the baseline.

⁸ We investigate robustness of our proposed methods to attacks in Section 5.2.

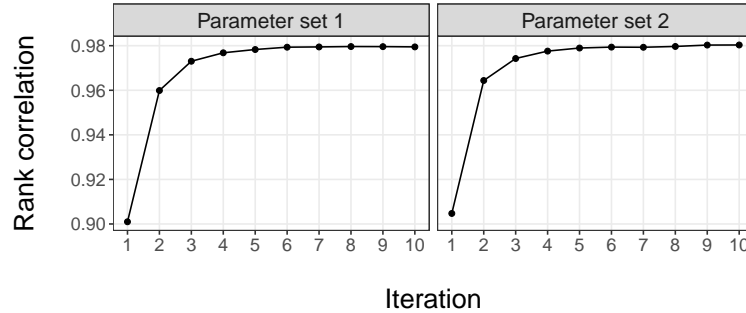


Fig. 5. Rank correlation between the estimated seller reputation computed via RATING SEPARATION and true reputation as a function of the number of iterations. Column panels show the two simulation parameter sets.

While Fig. 4 shows that just a single iteration of RATING SEPARATION can achieve superior performance compared to the baseline, the iterative nature of RATING SEPARATION allows for further improvement. In Fig. 5, we show that the performance of RATING SEPARATION can be substantially improved by just 3 additional iterations, at which point the rank correlation is close to perfect at about 0.98. This represents a tremendous improvement, considering that the baseline approach, at best, recovers seller rankings with about 0.85 correlation.

5.2. Performance with adversarial ratings

Next, we evaluate the efficacy of INTEGRITY WEIGHTING in mitigating the harms of adversarial ratings. The two approaches of using just INTEGRITY WEIGHTING and using both RATING SEPARATION AND INTEGRITY WEIGHTING are evaluated. We compare performance to three existing methods from previous literature: BRS [27], PA [28], and iCLUB [20]. A baseline approach, which does not explicitly adjust for potential adversarial behavior is included as well. As mentioned in Section 4.2, we conduct simulations for 18 possible attack types, unique pairs of three behavior patterns and six attack strategies. For each attack type, we vary the rate of attack between 10% and 90%, in 10% increments. The results are presented in Fig. 6.

From Fig. 6, we first note that the methods we propose (IW and RS&IW) consistently outperform all alternative methods in every setting that we test. Notably, we find that the baseline method, which assigns equal weight to all ratings and does not account for any malicious behavior, performs better than more sophisticated methods under some conditions.

Overall, as the rate of attacks increases, the performance of all methods in all settings decrease, albeit in varying degrees. One interesting finding is that iCLUB typically achieved either *negative* correlation, or the highest performance among the four benchmark approaches. This suggests that while iCLUB can be a high-performing method under specific assumptions of adversarial behavior, it is not generally reliable.

Under either a bad mouthing strategy (Pattern 2) or r -low shifting strategy (Pattern 5), other methods for mitigating adversarial ratings typically do no better than a naive

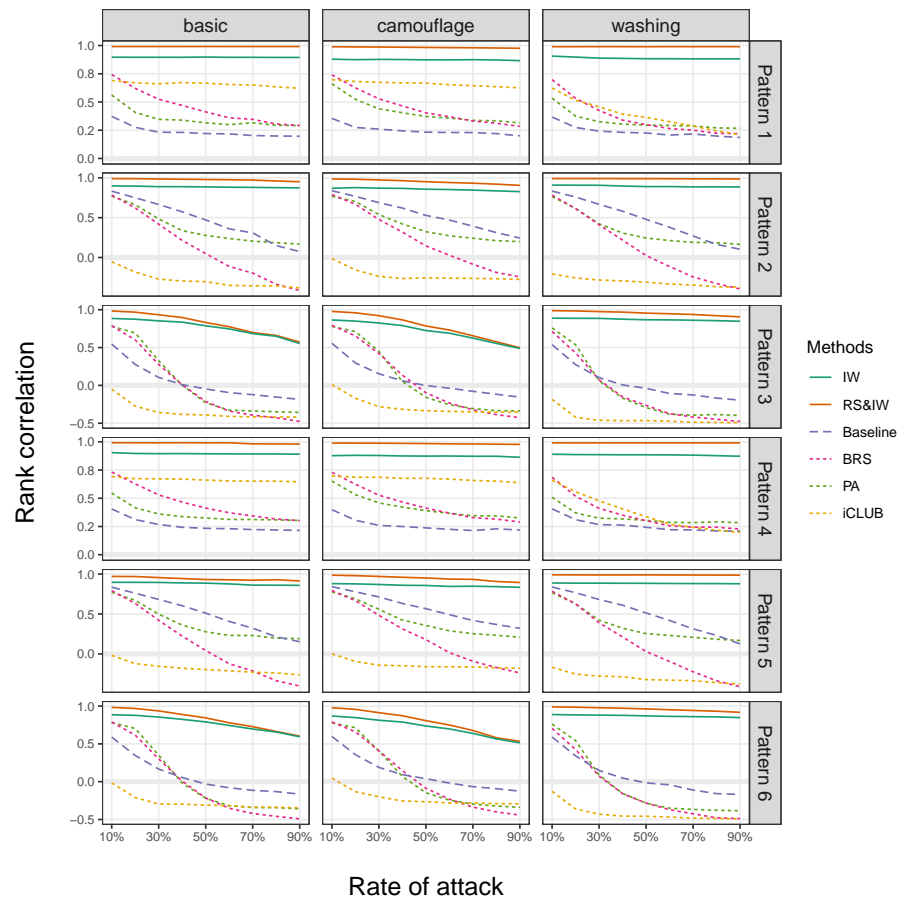


Fig. 6. Comparison of multiple methods for mitigating malicious reviews. The x-axis shows the proportion of attacks that are assumed in each simulation, while the y-axis shows the rank correlation between the estimated seller reputation and true reputation, for each method. Column panels show different attack types and row panels show different attack patterns. Overall, the proposed methods (RS and RS&IW) are able to recover the true reputation more reliably than any existing method across all simulated circumstances.

baseline approach. This indicates that existing methods are tailored to certain types of attack strategies, and do not perform well against a wide range of attacks, in general.

6. Conclusions

In increasingly complex online marketplaces that involve the interaction of multiple agents, evaluating the quality and characteristics of each agent is becoming more important. This paper addresses the issue of the confounded buyer ratings, as well as malicious ratings, in reputation systems and proposes RATING SEPARATION AND INTEGRITY WEIGHTING

(RS&IW), a system for providing agent reputations that are robust to confounded ratings and various types of cheating behavior. Through extensive experiments, we showed that our reputation system can both disentangle scores for sellers from the confounded ratings and are robust to numerous realistic attack scenarios generated by incorporating a comprehensive model of plausible adversarial behaviors.

While, in the interest of clarity and consistency, we have focused our work in this paper on the concrete problem of identifying seller rankings and mitigating malignant buyer behavior, the methods we propose could be extended to a broader family of problems in a more general context of multi-agent platforms. One possible extension would be to apply RATING SEPARATION in matching markets, where participating agents report numerous confounded signals with regard to the quality of other entities. For example, in ride sharing applications, RATING SEPARATION could be applied on ratings to decouple rider satisfaction of driver (e.g., personal, vehicle) and route (e.g., traffic conditions, travel time) characteristics. Or in a three-sided market, such as food delivery services, RATING SEPARATION could be extended to disentangle courier and restaurant ratings from an eater's single score.

Besides seller performance and item quality, item price is the one of the main factors having an influence on conforming a single score. Sometimes, buyers could give a generous score for a seller, even though both this seller's performance and his item quality are not satisfactory, because the price is discovered as the lowest one in an online platform. In a further study, we plan to develop a framework to accurately disentangle this price effect from a user's single score for measuring better purified seller reputation.

Acknowledgments. This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government(MSIT) (No. NRF-2020R1A2B5B03001960), the Next-Generation Information Computing Development Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Science, ICT (NRF-2017M3C4A7069440), and by Institute of Information & Communications Technology Planning & Evaluation(IITP) grant funded by the Korea government(MSIT) (No. 2020-0-01373, Artificial Intelligence Graduate School Program (Hanyang University))

Bibliography

- [1] Bidgoly, A., Ladani, B.: Modeling and Quantitative Verification of Trust Systems Against Malicious Attackers. *The Computer Journal* 59(7), 1005–1027 (2016)
- [2] Biega, A.J., Gummadi, K.P., Weikum, G.: Equity of attention: Amortizing individual fairness in rankings. In: *The 41st international acm sigir conference on research & development in information retrieval*. pp. 405–414 (2018)
- [3] Cabral, L., Hortacsu, A.: The dynamics of seller reputation - theory and evidence from eBay. *The Journal of Industrial Economics* LVIII(1), 54–78 (2010)
- [4] Chandrasekaran, P., Esfandiari, B.: A model for a testbed for evaluating reputation systems. In: *Proceedings of the 10th International Conference on Trust, Security and Privacy in Computing and Communications*. pp. 296–303. IEEE (2011)
- [5] Dellarocas, C.: Immunizing Online Reputation Reporting Systems Against Unfair Ratings and Discriminatory Behavior. In: *Proceedings of the 2nd ACM conference on Electronic Commerce*. pp. 150–157 (2000)
- [6] Fan, Z.P. and Xi, Y., Liu, Y.: Supporting consumer’s purchase decision: a method for ranking products based on online multi-attribute product ratings. *Soft Computing* 22, 5247–5261 (2018)
- [7] Fang, H., Zhang, J., Sensoy, M., Thalmann, N.M.: SARC : Subjectivity Alignment for Reputation Computation (Extended Abstract) Categories and Subject Descriptors. In: *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*. pp. 1365–1366 (2012)
- [8] Gao, J., Zhou, T.: Evaluating user reputation in online rating systems via an iterative group-based ranking method. *Physica A: Statistical Mechanics and its Applications* 473, 546–560 (2017)
- [9] Ghiasi, H., Brojeny, M., Gholamian, M.: A reputation system for e-marketplaces based on pairwise comparison. *Knowledge and Information Systems* 56, 613–636 (2018)
- [10] Houser, D., Wooders, J.: Reputation in auctions: Theory, and evidence from eBay. *Journal of Economics and Management Strategy* 15(2), 353–369 (2006)
- [11] Howe, J.: The Rise of Crowdsourcing. *Wired Magazine* (14) (2006), <http://archive.wired.com/wired/archive/14.06/crowds.html>
- [12] Irissappane, A.A., Jiang, S., Zhang, J.: Towards a comprehensive testbed to evaluate the robustness of reputation systems against unfair rating attacks. In: *User Modeling Adaptation and Personalization Workshops* (2012)
- [13] Jiang, W, X.Y.G.H.W.C.Z.L.: Multi agent system-based dynamic trust calculation model and credit management mechanism of online trading. *Intelligent Automation & Soft Computing* 22(4), 639–649 (2016)
- [14] Jøsang, A., Ismail, R., Boyd, C.: A survey of trust and reputation systems for online service provision. *Decision Support Systems* 43(2), 618–644 (2007)
- [15] Jøsang, A., Ismail, R., Jsang, A., Ismail, R.: The Beta Reputation System. *Proceedings of the 15th Bled Electronic Commerce Conference* 160, 324–337 (2002)
- [16] Kerr, R., Cohen, R.: TREET: The Trust and Reputation Experimentation and Evaluation Testbed. *Electronic Commerce Research* 10(3), 271–290 (2010)

- [17] Kramer, M.A.: Self-selection bias in reputation systems. In: Etalle, S., Marsh, S. (eds.) *IFIP International Federation for Information Processing*, vol. 238, chap. Trust Mang, pp. 255–268. Boston: Springer (2007)
- [18] Leadbeater, C.: *We-think*. Profile books (2009)
- [19] Lim, E.P., Nguyen, V.A., Jindal, N., Liu, B., Lauw, H.W.: Detecting product review spammers using rating behaviors. In: *Proceedings of the 19th ACM international conference on Information and knowledge management - CIKM '10*. pp. 939–948 (2010)
- [20] Liu, S., Zhang, J., Mao, C., Theng, Y., Kot, A.: iCLUB: an integrated clustering based approach to improve the robustness of reputation systems. In: *Proceedings of 10th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*. pp. 1151–1152 (2011)
- [21] Ma, L., Pei, Q., Xiang, Y., Yao, L., Yu, S.: A reliable reputation computation framework for online items in E-commerce. *Journal of Network and Computer Applications* 134, 13–25 (2019)
- [22] Oh, H., Kim, S., Park, S., Zhou, M.: Can You Trust Online Ratings? A Mutual Reinforcement Model for Trustworthy Online Rating Systems. *IEEE Transactions on Systems, Man, and Cybernetics* 45(12), 1564–1576 (2015)
- [23] Oh, H.K., Kim, S.W., Park, S., Zhou, M.: Trustable aggregation of online ratings. In: *Proceedings of the 22nd ACM international conference on Conference on information & knowledge management - CIKM '13*. pp. 1233–1236 (2013)
- [24] Regan, K., Poupart, P., Cohen, R.: Bayesian reputation modeling in e-marketplaces sensitive to subjectivity, deception and change. In: *Proceedings of the National Conference on Artificial Intelligence*. vol. 21, p. 1206 (2006)
- [25] Resnick, P., Zeckhauser, R., Swanson, J., Lockwood, K.: The value of reputation on eBay: A controlled experiment. *Experimental Economics* 9(2), 79–101 (2006)
- [26] Teacy, W.T., Patel, J., Jennings, N.R., Luck, M.: TRAVOS: Trust and reputation in the context of inaccurate information sources. *Autonomous Agents and Multi-Agent Systems* 12(2), 183–198 (2006)
- [27] Whitby, A., Jøsang, A., Indulska, J.: Filtering Out Unfair Ratings in Bayesian Reputation Systems. *Icfain Journal of Management Research* 6(2), 106–117 (2005)
- [28] Zhang, J., Cohen, R.: Evaluating the trustworthiness of advice about seller agents in e-marketplaces: A personalized approach. *Electronic Commerce Research and Applications* 7(3), 330–340 (2008)
- [29] Zhang, L., Jiang, S., Zhang, J., Ng, W.K.: Robustness of Trust Models and Combinations. *Trust Management* VI 374, 36–51 (2012)
- [30] Zhou, X., Murakami, Y., Ishida, T., Liu, X., Huang, G.: ARM: Toward Adaptive and Robust Model for Reputation Aggregation. *IEEE Transactions on Automation Science and Engineering* 17(1), 88–99 (2020)

Hyun-Kyo Oh received his B.S., M.S. and Ph.D. degree in Electronics and Computer Engineering from Hanyang University, Seoul, Korea at 2008, 2010 and 2016. He visited the Department of Computer Science of Carnegie Mellon University as a visiting scholar in 2013. He worked with the Knowledge Computing Group at Microsoft Research Asia as a research intern from 2014 to 2015. In 2016, he joined Samsung Electronics, where he currently is a lead data scientist working on creating a new machine learning and deep

learning platform that deals with the health of V-NAND flash memory and the performance of Solid State Disk (SSD). Now, he is also a visiting researcher at Institute for Software Research at Carnegie Mellon University.

Jongbin Jung received a Ph.D. in Computational Social Science and Decision Analysis from Stanford University. He is primarily interested in using quantitative methods and data analytics to help improve and evaluate human decisions. Jongbin studied operations research (M.S.) and business administration (B.B.A.) at Yonsei University (Seoul, South Korea).

Sunju Park received her B.S. and M.S. in computer engineering from Seoul National University, Seoul, Korea, and the Ph.D. degree in computer science and engineering from the University of Michigan, Ann Arbor, MI, USA. She has served on the faculties of Management Science and Information Systems, Rutgers University, NJ, USA. She is a professor of Operations, Decisions and Information at School of Business, Yonsei University, Seoul, Korea. Her current research interests include analysis of online social networks, multiagent systems for online businesses, and pricing of network resources. Her publications include *Computers and Industrial Engineering*, *Electronic Commerce Research*, *Transportation Research*, *IIE Transactions*, the *European Journal of Operational Research*, the *Journal of Artificial Intelligence Research*, *Interfaces*, *Autonomous Agents and Multiagent Systems*, and other leading journals.

Sang-Wook Kim received the B.S. degree in computer engineering from Seoul National University, in 1989, and the M.S. and Ph.D. degrees in computer science from the Korea Advanced Institute of Science and Technology (KAIST), in 1991 and 1994, respectively. In 2003, he joined Hanyang University, Seoul, Korea, where he currently is a professor at the Department of Computer Science and Engineering and the director of the Brain-Korea21-Plus research program. He is also leading a National Research Lab (NRL) Project funded by the National Research Foundation since 2015. From 2009 to 2010, he visited the Computer Science Department, Carnegie Mellon University, as a visiting professor. From 1999 to 2000, he worked with the IBM T. J. Watson Research Center, USA, as a postdoc. He also visited the Computer Science Department at Stanford University as a visiting researcher in 1991. He is an author of more than 200 papers in refereed international journals and international conference proceedings. His research interests include databases, data mining, multimedia information retrieval, social network analysis, recommendation, and web data analysis. He is a member of the ACM and the IEEE.

Received: November 22, 2019; Accepted: May 27, 2020.

Fuzzy Logic and Image Compression Based Energy Efficient Application Layer Algorithm for Wireless Multimedia Sensor Networks

Arafat SENTURK¹, Resul KARA¹, and Ibrahim OZCELIK²

¹ Department of Computer Engineering, Engineering Faculty, Duzce University
Duzce, Turkey
{arafatsenturk,resulkara}@duzce.edu.tr

² Department of Computer Engineering, Faculty of Computer and Information Sciences, Sakarya University
Sakarya, Turkey
ozcelik@sakarya.edu.tr

Abstract. Wireless Sensor Networks (WSN) are the networks that can realize data processing and computation skills of sensor nodes over the wireless channel and they have several communication devices. Wireless Multimedia Sensor Networks (WMSN) are the networks composed of low-cost sensor nodes that transmit real-time multimedia data like voice, image, and video to each other and to sink. WMSN needs more energy and bandwidth than WSN since they transmit a larger amount of data. The size of the data transmitted by the sensor nodes to each other or the sink becomes an important factor in their energy consumption. Energy consumption is a fundamental issue for WMSN. Other issues that affect the progress of WMSN are limited bandwidth and memory constraints. In these networks, for which the node battery lives are important sources, the limited sources must be effectively used by decreasing the transmitted data amount by removing the redundant data after proper processing of the environmental data. A new algorithm is developed to minimize the energy consumption during image data transmission between sensor nodes on WMSN, and so, make the nodes use their most important source, battery life effectively in this study. This algorithm is named as Energy-aware Application Layer Algorithm based on Image Compression (EALAIC). This algorithm makes use of the top three image compression algorithms for WMSN and decides instantly to which one is the most efficient based on three parameters: the distance between the nodes, total node number, and data transmission frequency. In this way, the sensor node battery lives are used efficiently. The performance analysis of the developed algorithm is also done via Network Simulator – 2 (NS – 2) and it is compared by the existing algorithms in terms of energy rate (consumed energy/total energy) and PSNR (Peak Signal to Noise Ratio).

Keywords: Energy efficiency, Fuzzy logic, Image compression, Wireless multimedia sensor networks.

1. Introduction

Wireless Sensor Networks (WSN) become widespread easily because of some features such as being more appropriate to natural conditions than wired communication, be-

ing more reliable, having a flexible structure, producing solutions with low cost, self-organization, using the energy in a balanced way inside the network, and ease of installation. The purpose of WSN is to give access to data at any time and from everywhere easily. They do this function by collecting, processing, analyzing, and transmitting data [4], [28]. WSN has drawn the attention of the researchers since it can easily solve the practical and theoretical difficulties. Another reason for the increasing popularity is using small devices for transmitting data to far distances after some basic operations are done in physical environments. Besides, WSN has been used frequently in many military and unmilitary applications since it avails to measure some physical phenomena like temperature, pressure, moisture, and location [50].

Wireless Multimedia Sensor Networks (WMSN) are the networks that are composed of low cost sensor nodes that transmit real-time image, video or audio data to each other. Audio and visual data may exist together in a single device or they may exist separately. Besides, WMSN can store real-time multimedia data after retrieving from several sensor nodes [34]. WMSN leads to the most hopeful developments in real-time multimedia monitoring. But having a larger amount of data for transmission when compared with WSN makes the life and efficient energy usage important problems for WMSN [27]. In addition, energy efficiency has been studied extensively in sensor networks [44],[45],[9],[46],[16]. [31] focuses on the energy-aware image transmission in WMSN for flood images. In this paper, the researchers applied image segmentation and the useless parts of the images were removed before transmission. This approach may achieve good performance in flood images in which mostly the trees and sky are at the top of the image and the flood water are at the bottom and overlay a considerable part of the image. Middle part of the image, transition between the land and water, as accepted as the most important part of the image. For that kind of images segmentation can be applied and the regions out of ROI can be removed and then this downsized image can be transmitted. But for example Lena image that we used in our experiments consists of many details on the whole image and the proposed segmentation method does not work for this kind of cases. Energy consumption is the fundamental issue for WMSN. The sensor nodes in the network consume most of their energy during data retrieval and the performed operations. Therefore, proper design is required to maximize the lifetime of wireless sensor networks. Other issues that affect the development of WMSN are limited bandwidth and memory constraints [4],[29].

Delivery of multimedia data requires higher bandwidth than traditional sensor networks and this situation affects the design of WMSN [15]. Designing low complexity encoders and decoders besides high compression ability is an important application layer problem in WMSN for which energy preservation is rather necessary [15]. Some algorithms can be developed to process raw data retrieved from the environment and only the important data can be transmitted instead of unnecessary data. Thus, the size of transmitted data can be reduced and energy consumption can be minimized [33]. Reliability, congestion control, latency minimization, and error check are the requirements provided in the transmission layer in terms of quality of service. Quality of service-based routing and forwarding schemas are the research topics of the network layer. MAC layer protocols should contain priority-based and latency aware timing algorithms [15],[55]. This layer controls the biggest part of radio communication, and so, it plays an important role in energy efficiency and reducing latency [22].

1.1. Image compression algorithms used for WMSN

There are many types of research related to energy efficiency in WMSN and these researches can be classified according to the layer they correspond to. In the application layer, especially the improvements in image compression algorithms constitute an essential part of the researches. The image compression algorithms used for WMSN are given in the following part.

Low Energy Image Compression Algorithm (LEICA) developed by Sun *et al.* [41] claims to minimize energy consumption by shrinking the data amount without any deformation in image quality. LEICA first divides the image into two layers as background and foreground. This layered structure is named as ROI (Region of Interest). The important content falls into the foreground and the relatively important image is in background. The unimportant background image is compressed more than the important foreground image. In this way, both image qualities are preserved and the compression ratio is increased, and so, the energy efficiency is provided.

The first image compression algorithm used in WMSN is Embedded Zero Tree Wavelet (EZW) [39] is modified by Said *et al.* to make encoding and decoding faster and it is named as SPIHT (Set Partitioning In Hierarchical Trees) [35]. SPIHT converts the pixel data of the image in a tree structure and collects the most related pixels in a single set. Then, zero-valued trees are encoded in a single value. Combining this feature of wavelet transform and SPIHT achieves a high compression ratio in images.

Thanks to the algorithm developed in [2], JPEG2000 formatted images are considered to be transmitted in a low-energy fashion without any deformation in image quality. Unnecessary bits inside the image are removed in this study for less-energy data transmission. But it is understood that the sensor nodes in WMSN are not suitable for JPEG2000 formatted image compression in terms of hardware and it consumes more energy than expected. Mulugeta *et al.* have proposed a secure forwarding protocol for sensor networks in [30]. Minimizing the energy consumption in image compression, processing and transmission are aimed at this study.

Discrete Wavelet Transform (DCT) is an image compression technique that is frequently used in WMSN. DCT based image compression techniques provide proper compression efficiency. Since encoding is done after the image is decomposed into small data blocks, it can also be used in low-memory applications [3].

The network was divided into two layers in [43] as camera nodes and common nodes which fulfil the compression duty. Images were divided into blocks and the largest-energy nodes were selected, and then wavelet transform was applied separately. Several nodes do the compression of a subpart of the image at the same time and parts are then combined in order to maintain load balancing in the network. Besides, unlike LEACH they considered the residual energies of the nodes and the average energy consumption for cluster head selection.

Image was transmitted by three methods in [17]. Transform coding, conventional compressed sensing and a new compressed sensing method used together. In the new compressed sensing method they used inverse binary DCT method. It is shown that the proposed approach is more energy efficient than DCT based other approaches.

In [47], the authors developed an object detection based energy efficient transmission method. They first determine the object presence and send image segments instead of sending whole image. The receiving node constructs the whole image based on the frames

it took earlier. Their purpose is to use the approach in real time object tracking. They did not give the details of the experiments like the distance between the nodes, size of the area, number of hops for a transmission, *etc.*

Compression of video data was handled in [20]. Repeated features between the frames and redundant features were eliminated to reduce bit rate of the video stream. Authors evaluate the frames in terms of their inter ayer features and remove redundant features. They do not evaluate each frame seperately. Therefore, this approach can not be used for the compression of images as in our case.

1.2. Energy-aware protocols for layered architecture

In the following part of the study, energy-aware protocols which are proposed for WMSN are explained based on a layered architecture. Energy-aware studies are given starting from the physical layer to application layer.

New technology with less latency and higher data rate are needed for multimedia communication. Even IEEE 802.15.4 provides high data efficiency in radio inactive status, its energy consumption is rather more. In [40], energy-saving IEEE 802.15.4 is handled. The method adjusts activeness according to traffic conditions in radioactive status and provides energy saving with less latency. In [24], a method based on radio interval and frame process adjustment with RED (Random Early Detection) queue management is proposed. The method aims to improve the performance of WMSN by reducing end to end latency and packet losses.

MAC protocols work for satisfying the quality of service requirements in the application layer. MAC layer has a differing important when compared with other layers because of its responsibilities like reducing collusion and number of retransmissions, minimizing energy consumption and interferences, and maximizing reliability and consistency [15]. In [37], a CSMA-CA (Carrier Sense Multiple Access with Collision Avoidance) based MAC protocol is proposed. Traffic classes are given priority according to the jumped node numbers. Also, to reduce collision possibility it allows interfering to report messages. In [6], both TDMA (Time Division Multiple Access) and contention based CSMA used together for scheduling different traffic flows concerning priority. Service difference is ensured by the queue management schema. High-priority traffic has more chances to access the medium and the latency is reduced. Some MAC protocols include additional processes that require sniffing also in a sleep state to make high-priority data has a greater chance [18], [12].

The fundamental difficulties in routing protocols are to optimize network performance and provide energy preservation while satisfying the quality of service requirements [15]. Quality of service metrics in the routing layer is end to end delay, packet loss, the total number of hops, bandwidth, connection quality, and energy. A geographic routing protocol has been proposed in [36]. End to end delay is ensured by evaluating the position relative to the base stations, remaining energy and queue length in the next node. In [7], a routing approach that maximizes network lifetime by balancing energy consumption in the nodes is presented. Service differentiation is used to access the base stations in the acceptable delay interval. In the other study, the cluster heads use ant colony optimization to find the optimal path through base station [8]. In [51], different service requirements are considered for different traffic flows. These are some service constraints like delay, packet loss, remaining energy in the nodes, and required number of hops for every path.

The method works in the way that the base station broadcasts a query indicating its needs and all subadjacent nodes check their sources and added to the path if they have more resources than required. In [21], energy efficiency in WMSN was considered in terms of finding optimum routing paths and managing network load with respect to energy reserves of the nodes. Their routing protocol saves the overall energy by preventing useless data transmission using an energy threshold. Similarly, [25] also deals with finding optimal paths in a WMSN for QoS and balanced energy distribution. They do not consider to reduce packet size for multimedia transmission in WMSN. [1] focuses on the efficient transmission of video over WSN. It proposes a routing protocol and transmits packets based on their priority through the scored paths. There are also many efforts in recent literature for energy aware routing [19], AODV routing protocol was enhanced in [11] considering not only the number of hops in a transmission, but also the residual energies and the congestion issues for an energy efficient routing.

The transport layer plays an important role in the communication over sensor networks with delay constraints. Congestion detection and reduction, end to end delay, and reliability requirements must be satisfied besides energy constraints [15]. Different traffic flows occur in WMSN and each flow may have its own reliability, error and delay constraints. All these constraints must be satisfied for every traffic flow separately. Transport layer protocols in traditional WSN are not suitable for WMSN. No service differentiation, flexible delay intervals, the importance of reliability of the event, not the application are the features that distinguish WMSN from WSN [15].

Pixel-based Wyner-Ziv encoder is proposed as a proper application layer operation [4], [14]. However, this encoder causes delays since it requires many feedbacks from the decoder part. In [56], pixel based Wyner-Ziv the encoder is applied in a sliced structure with automated rate selection. Sending only necessary parity bits and so efficient usage of bandwidth is purposed. Encoder takes the data quality in the end side and sends an appropriate number of parity bits for a successful encoding. In a different study, a two-phase protocol has been proposed [48]. In the first phase, compression gain which is obtained when interrelated cameras encode together is measured by an entropy based discreteness metric. Then, an optimal clustering hierarchy is established for encoding. Therefore, the global compression gain is maximized besides the reliability in the decoder part. Maximization of compression is achieved by ensuring the selection of the clusters with minimum entropy, and reliability is achieved by ensuring all camera nodes are covered at least in two of selected clusters.

Insufficiency of studies in the application layer and the idea to develop more energy-efficient approaches by combining realized improvements with the improvements in other layers have to lead develop an energy-aware method to work in the application layer. Thus, the network lifetime will be prolonged by minimizing the amount of data to be transmitted and selecting the most suitable compression algorithm according to the network condition (distance between nodes, the total number of nodes and data transmission frequency).

The main contributions of the proposed energy aware approach are as follows:

- Instant determination of image compression algorithm among the most efficient compression methods used in WMSN based on the network condition,
- Acceptable PSNR value is checked to satisfy QoS requirements of WMSN,
- Compression ratio is not fixed, i.e. it can be increased if acceptable PSNR is preserved,

- Fuzzy rule based approach is used during the determination of compression algorithm,
- Comprehensive experiments are conducted to investigate the consumed energy and residual energy of the network based on several parameters which are distance between the nodes, data transmission frequency, and number of nodes. Also, PSNR values are calculated for different compression ratios.
- A comprehensive overview of the literature is realized and the obtained results are compared by the state-of-the-art to show the performance of the proposed method.

In the first part of the paper, WSN and WMSN are mentioned briefly. The image compression algorithms used for WMSN are mentioned and energy aware studies related to WMSN are given based on layered architecture. In the second part, the image compression algorithm to be used in the developed algorithm is decided through the referenced research. Next, operation principles of EALAIC are given. In the third part, the performance criterions of WMSN are considered and the parameters of the simulation are mentioned. Afterward, a performance evaluation of EALAIC is done based on energy rate, end to end delay and total energy consumption via the comparison with the other referenced works. The results are given in the fourth part and future works are proposed.

2. Energy - Aware Application Layer Algorithm Based on Image Compression

Image compression methods in WMSN are divided into two groups as lossy and lossless methods. Lossy methods provide important advantages since they have less encoding/decoding time, better compression ratio, and usability in energy-limited applications. Image size in lossless compressed images is considerably greater than the lossy one, and so, the sensor nodes need more power and bandwidth when transmitting images. Therefore, lossless image compression methods are not preferred in WMSN [53]. Three image compression algorithms that give the best solutions in terms of energy efficiency in WMSN stand out in [38]. These algorithms are DCT, SPIHT, and LEICA.

In this part of the paper, the energy-aware method, EALAIC which is based on DCT, SPIHT, and LEICA to be used for WMSN is given. The parameters of EALAIC and how they are obtained are specified below.

A. Obtaining data about the distance between nodes: The sensor nodes in EALAIC have data about the distance between each other via a routing protocol, Modify-Adhoc On Demand Distance Vector Routing (M-AODV) which is obtained by changing “hello” packet in Adhoc On Demand Distance Vector Routing (AODV) routing protocol.

B. Determination of Number of Nodes and Data Transmission Frequency: The number of required nodes and data transmission frequency is determined in the network establishment phase. The reasons why the distance between nodes, number of nodes, and data transmission frequency are given in Section 2.1.

C. Determination of Compression Algorithm via Fuzzy Logic Method: EALAIC first processes the determined parameters for the selection and decides the algorithms to be used. Then, all parameters are considered among the determined algorithms and the most efficient one is decided. The detailed information is given in the following part of the section.

D. Determination of PSNR Value Loss and Its Conservation: PSNR value of the compressed image is calculated and the predetermined value loss interval is tried to be conserved. Therefore, both higher compression is obtained and the quality of the system is preserved. Required operations and sample analyses are in the following.

The data obtained by the above steps are processed by the decision tree method and converted into rule-based data. The selection of the most efficient image compression algorithm is done according to the data obtained by the fuzzy logic method. These operations are explained in detail in Section 2.1.

Reducing the transmission and acquisition energies spent by sensor nodes is aimed. Since the transmission and acquisition energies will be reduced, sensor nodes use their battery lives efficiently and the network lifetime is intended to be prolonged. Also, the PSNR value of the compressed image is kept in the acceptable interval for wireless communication and the image quality is aimed to be preserved.

The most efficient image compression is decided thanks to EALAIC and the energy spent by sensor nodes will be minimized. In the energy calculation operations, both the energy consumed by compression algorithms and the energy needed by EALAIC are considered. The structure of EALAIC is given in Figure 1.

Pseudo code of EALAIC is the following:

```

ALGORITHM EALAIC (cp, D, Ds)
//input : cp, respose to hello packet, the image size
//output: compressed image
//cp_x, cp_y, neighbour node coordinates
//S: data transmission frequency, Ds: number of nodes in network
node = node who sends hello packet
cp_x = take X coordinate inside respose packet
cp_y = take Y coordinate inside respose packet
node_x = take GPS_X coordinate
node_y = take GPS_Y coordinate
distance = sqrt((node_x-cp_x)^2+(node_y-cp_y)^2)
i = take_image
P = psnr(i) //PSNR value of image i
D = 5 //Compression ratio
Sa = algorithmSelection(S,Ds,distance)
//selection of compression algorithm
While do
Sg = COMPRESSimage(D,Sa,i) // %D compressed image according to Sa
Sg_P = psnr(Sg) //PSNR value of compressed image
if (P-Sg_P )< 25 ve (P-Sg_P) > 20
return i //send data to node with cp_x,cp_y coordinate
else if (P-Sg_p) >= 25
D = D + 5
else D = D - 5

```

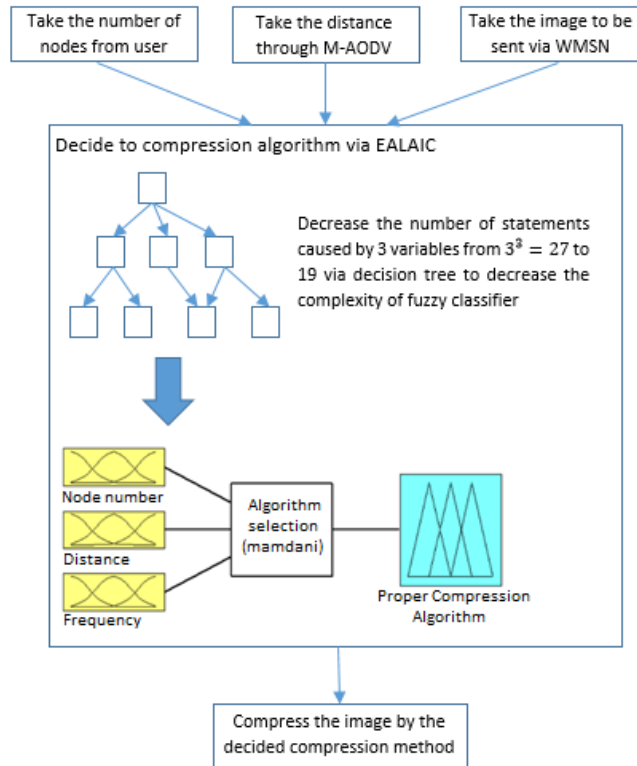


Fig. 1. Structure of EALAIC

2.1. Determination of Compression Algorithm via Fuzzy Logic Method

EALAIC has used a fuzzy logic method to determine the most efficient image compression algorithm. First, the numbers of states are minimized by the decision tree method and the most suitable method is decided by fuzzy logic method.

The parameters of EALAIC are specified based on the system needs. These parameters and the reasons why they are determined are:

- **A.** Distance between the sensor nodes in the network affects their transmission/acquisition energies [32]. This situation is taken into consideration and “distance between nodes” is determined as a parameter. The sensor nodes in the system may be mobile. Thus, the distance between nodes may change at any time. The linguistic variable of membership function is determined as “low” for the cases when the sensor nodes whose coverage area is 200 meters are close to each other, as “high” for the cases the distance is in the border, and “medium” for the cases the distance is neither “low” nor “high”.
- **B.** The Number of sensor nodes over the network affects the total energy consumption. Energy waste is inevitable when there exist a greater number of nodes [26], [13]. Besides, the number of nodes also affects the coverage area of the application. Re-

quired node number can be determined by considering the coverage area of the sensor nodes [13]. “Number of nodes” is determined as a parameter through these inferences. The coverage area of the system is stated by the user in the network establishment phase to obtain the number of nodes to be used for the network. The linguistic variables of fuzzy membership function formed by this calculated value are determined as “low”, “medium”, and “high”.

- C. Sensor nodes transmit the received image to the next sensor node in the coverage area after they perform some operations on the image. Transmission of this image data in a certain time interval is important in terms of energy. Because all nodes used in every transmission spend transmission/acquisition energy [52]. For this reason, “data transmission frequency” is used as a parameter in the system and the linguistic variables of fuzzy membership function are specified as “low”, “medium”, and “high” according to the user request.

Sensor nodes are made to select and use the most efficient algorithm thanks to EALAIC which is developed with the mentioned parameters. Energy consumption during compression (operation complexity), compression amount, quality loss in the output data, and the distance-based evaluation are effective in the selection of the most suitable compression algorithm. For example, in the distance based evaluation, a summation of the operation complexity and the energy consumed to transmit compressed image is calculated for every compression method and the one with less energy consumption is preferred. The method may have high the compression ability, but its operation complexity may also be high. Consuming a high amount of energy may be unnecessary when the distance to transmit data is short. Transmitting data after a simpler method with less operational complexity can be less costly. Therefore, the preferred compression method changes depending on the distance. The method with less complexity is preferred not to cause latencies in case of the evaluation based on transmission frequency and to minimize these latencies. In case when the number of nodes is variable, for instance when the number of nodes is high, number of hops will be high and every hop will bring additional cost. On top of it, if the operation complexity is also high, then the total cost will increase and so more energy will be consumed. This means operation complexity is effective in the evaluation based on the number of nodes.

There are three fuzzy input variables for the Fuzzy Logic Method. These are the ones determined for EALAIC and they are “number of nodes”, “frequency”, and “distance”. The fuzzy output variable is named as “the algorithm to be used”. The fuzzification operation whose block schema is given in Figure 2 is established by MATLAB 2014a Fuzzy Logic Toolbox.

Determination of the algorithm to be used according to the distance between nodes:

“hello” packet in AODV routing protocol has been changed to calculate the distance between sensor nodes which are in the coverage area of each other and this protocol is named as Modified Adhoc On Demand Distance Vector Routing (M-AODV). The location information obtained via the GPS module on sensor nodes is added to the modified “hello” packet. This packet is delivered to the neighbors and they are informed about their locations. Then, the distance between nodes is calculated using this information. If M-AODV is not used, then geographic algorithms must be used. Geography based algorithms use many control packets and they consume a lot of energy [10], [23].

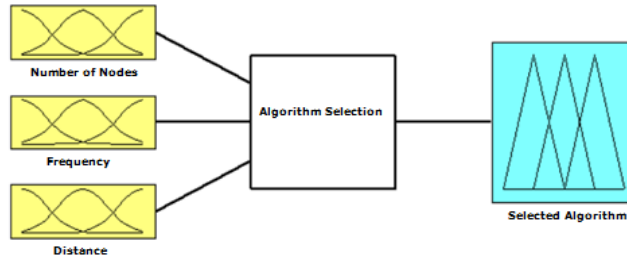


Fig. 2. Block view of fuzzy logic input and output variables

The frame structure of AODV “hello” packet is given in Figure 3 and its modified version is in Figure 4. The distance between nodes is calculated using the equality given in equation 1 thanks to the information provided by GPS module by taking location information of both nodes.

Type	Length	“hello” Message
1 Byte	1 Byte	4 Bytes

Fig. 3. Frame structure of “hello” packet

Type	Length	“hello” Message	Location (x, y)
1 Byte	1 Byte	4 Bytes	8 Bytes

Fig. 4. Frame structure of modified “hello” packet

$$M_{i,j} = \sqrt{(x_j - x_i)^2 + (y_j - y_i)^2} \tag{1}$$

The compression algorithm changes with respect to the distance between nodes. According to the simulation results, “energy rate (remaining energy/total energy)” versus distance between nodes is given in Figure 5 graphically for three compression algorithms.

It is shown by the graphics that DCT is the most efficient algorithm to be used for WMSN in terms of energy when the distance between nodes is low. When the distance between nodes is medium, using LEICA for image compression seems logical and if the distance between nodes is high, SPIHT is the most energy-efficient compression algorithm. However, sometimes the algorithm to be used for compression cannot be decided clearly. That is to say, the energy efficiency of more than one algorithm seems to be high for some distance intervals. When fuzzy logic is used to determine the best one for

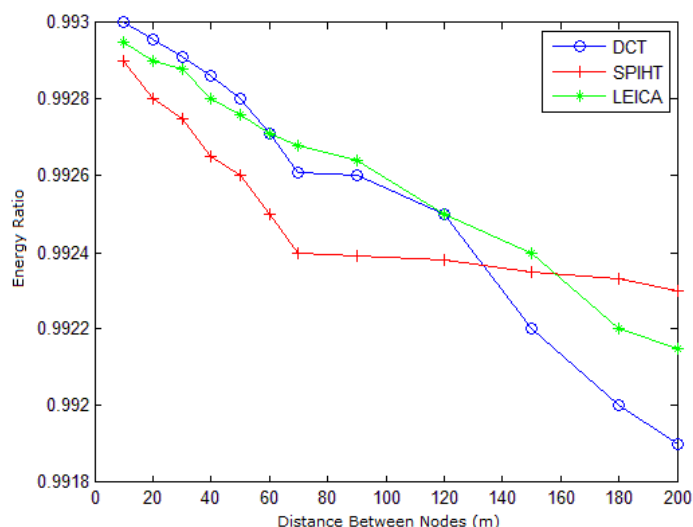


Fig. 5. Energy rate versus distance between nodes

these situations: “Distance” parameter changes between 0 and 200. Three membership functions have been constituted for this input variable and these are named as “low”, “medium”, and “high”. The upper and lower limits of input variable “distance” are given in Table 1 and its graphical representation is in Figure 6.

Table 1. Membership function limits of input variable “Distance”

Linguistic variable	Membership function limit	
Low	0	70
Medium	50	120
High	100	200

As a result of this inference, DCT or LEICA will be used for compression in case “distance” parameter is “low and medium” and LEICA or SPIHT will be used when “distance” parameter is “medium and high”.

Determination of the algorithm to be used according to the number of nodes: The required number of nodes according to the area where they will be set up is known on the average for any system since the area where the WMSN addresses are definite during establishment. The compression algorithm changes according to the number of nodes used for the network while the compressed image data is transmitted. According to the simulation results, “energy ratio (remaining energy/total energy)” versus the number of nodes in the network is given in Figure 7 graphically for three compression algorithms.

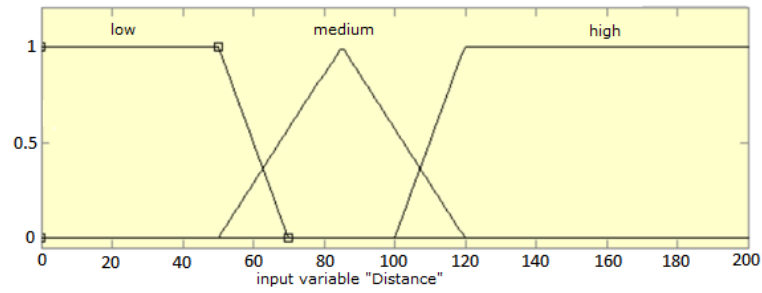


Fig. 6. Graphical representation of membership functions of input variable “Distance”

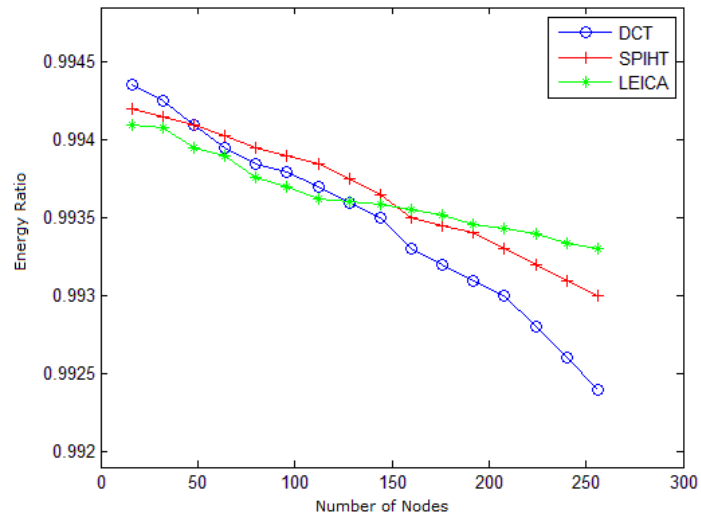


Fig. 7. Energy rate versus distance between nodes

It is shown by the graphics that DCT has the best energy rate when the number of nodes in the network is the lowest. LEICA algorithm seems to have the worst energy rate. However, DCT losses its efficiency as the number of nodes increase and SPIHT seem to be better. LEICA algorithm is at the forefront in the networks with high number of nodes. That is, using DCT is efficient in case the number of nodes in the lowest, using SPIHT is efficient when the number of nodes is medium and using LEICA is logical for the cases in which node number is the highest. As shown in Figure 7, the proper compression algorithm is not decided clearly for node numbers. Fuzzy logic is used for that kind of situation. “Number of nodes” parameter changes between 0 and 256. Three membership functions have been constituted for this input variable and these are named as “low”, “medium”, and “high”. The upper and lower limits of input variable “Number of nodes” are given in Table 2 and its graphical representation is in Figure 8.

Table 2. Membership function limits of input variable “Distance”

Linguistic variable	Membership function limit	
Low	0	64
Medium	50	158
High	144	256

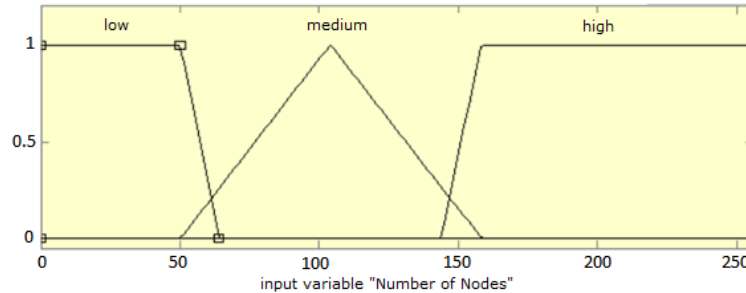


Fig. 8. Graphical representation of membership functions of input variable “Number of Nodes”

As a result of this inference, DCT or SPIHT will be used for compression in case the “number of nodes” parameter is “low and medium” and SPIHT or LEICA will be used when “number of nodes” is “medium and high”.

Determination of the algorithm to be used according to data transmission frequency over the network: Data transmission-acquisition between nodes is realized at certain time intervals. This interval shows the data density in the network and it can be determined optionally during system establishment. The data transmission frequency of the sensors in the network changes between 2 seconds and 19 seconds. According to the simulation

results, data transmission frequency versus “energy rate (remaining energy/total energy)” is given in Figure 9 graphically for three compression algorithms.

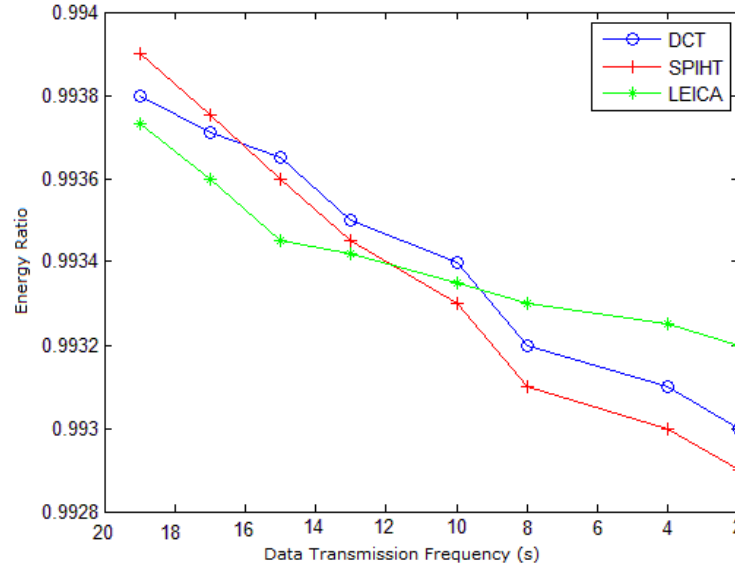


Fig. 9. Energy rate versus distance between nodes

It is shown by the graphics that SPIHT is the most efficient algorithm when the network is in peak time, i.e. data frequency is the highest. DCT algorithm seems to be better when the data frequency becomes lower and LEICA is the most efficient one when the frequency is quite low. However, like the other parameters, in some frequency intervals, there are overlaps and the fuzzy logic method is used to get rid of this problem.

Table 3. Membership function limits of input variable “Frequency”

Linguistic variable	Membership function limit
Low	12 20
Medium	5 14
High	0 7

“Frequency” parameter changes between 0 and 19. Three membership functions have been constituted for this input variable and these are named as “low”, “medium”, and “high”. The upper and lower limits of the input variable “Frequency” are given in Table 3 and its graphical representation is in Figure 10. As a result of this inference, SPIHT or DCT will be used for compression in case the “frequency” parameter is “high and medium” and DCT or LEICA will be used when “frequency” is “medium and low”.

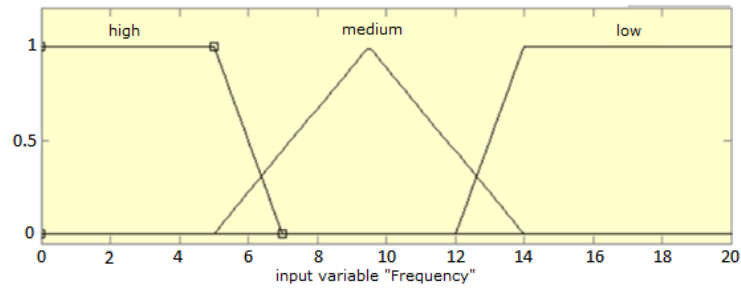


Fig. 10. Graphical representation of membership functions of input variable “Frequency”

Selection of image compression algorithm: “Selected Algorithm” value is obtained using above three input variables. For this purpose, three membership functions are constituted for “Selected Algorithm” output variable. Names of the algorithms (DCT, SPIHT and LEICA) to be used after the method works are given as linguistic variables of membership functions. The upper and lower limits of output variable “Selected Algorithm” are given in Table 4 and its graphical representation is in Figure 11.

Table 4. Membership function limits of input variable “Selected Algorithm”

Linguistic variable	Membership function limit	
Low	0	1
Medium	1	2
High	2	3

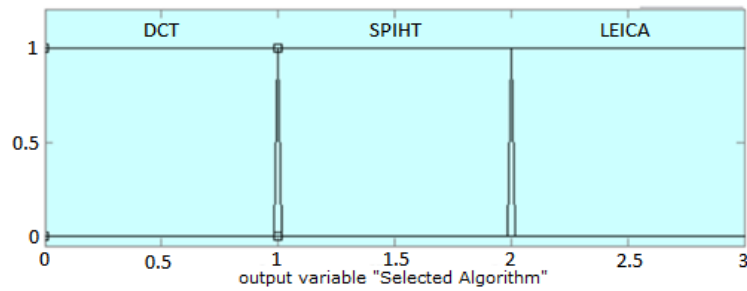


Fig. 11. Graphical representation of membership functions of output variable “Selected Algorithm”

2.2. Determination of PSNR value loss and its Conservation

PSNR value of the compressed image is calculated in the developed system. The acceptable PSNR loss value for wireless communication is about 20-25 dB [42]. If the difference between the original image and the compressed image is in the acceptable loss interval, then the image is transmitted to the next node or the sink. If the difference between the original image and the compressed image is less than or equal to 20 dB, then the compression ratio is increased and the image is transmitted to the next node or the sink after compressing the original image with higher compression ratio. However, if the difference between the original image and the compressed image is more than 20 dB, then the compression ratio is decreased and the image is compressed at lower ratio and sent. Conservation of the quality of the system will be provided by keeping PSNR value of the compressed image in the acceptable interval.

Sample images given below are used for the simulations in Section 3 and the inferences done in this section related to PSNR value loss.

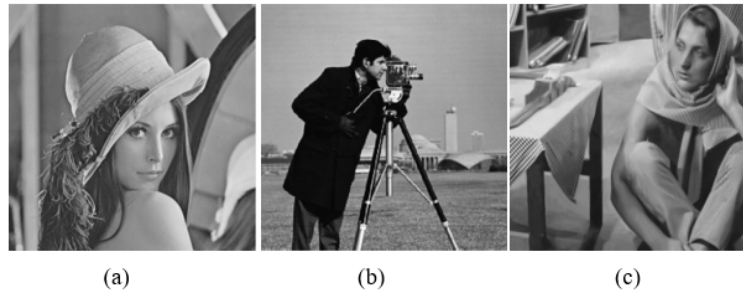


Fig. 12. (a) sample image 1, (b) sample image 2 ve (c) sample image 3

PSNR value loss versus compression ratio graphic of “sample image 1” data for every algorithm is given in Figure 13.

Table 5. Compression ratios with respect to PSNR values of sample images

	CR	PSNR	CR	PSNR	CR	PSNR	CR	PSNR	CR	PSNR
Image 1 (a)	0	98	5	85	10	81	15	77	20	74
Image 2 (b)	0	95	5	78	10	72	15	70	20	67
Image 3 (c)	0	93	5	87	10	74	15	72	20	69

It is shown in Figure 13 that the acceptable PSNR loss value can be reached by applying at most 20% compression ratio for “sample image 1” data. That is, PSNR value of the image decreases from 98 to 74 after compression. So, the PSNR value loss is 24 dB. When the compression ratio is increased to 25, then PSNR value loss will be 28 dB and this is out of the acceptable loss interval. Therefore, 20% compression ratio with 24 dB loss is the most suitable compression ratio.

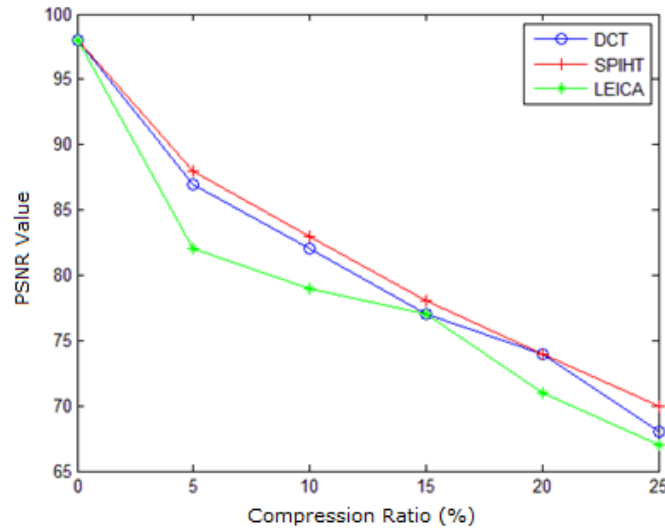


Fig. 13. PSNR value loss versus compression ratio graphic of “sample image 1” for three algorithms

Compression ratios (CR) of three sample image data with respect to their PSNR values are given in Table 5. EALAIC decides the compression ratio to be based on these values.

3. Performance Evaluation of EALAIC

3.1. Performance Criteria for WMSN

WMSN has some constraints like memory, processor, energy consumption, and lifetime. Memory and processor constraints can be overcome by the advances in technology. Enhancing energy consumption and so prolonging network lifetime become possible via the studies carried out [49]. The most important performance criterion for WMSN is energy consumption. More energy and bandwidth is needed to transmit multimedia data [4], [53]. The performance of the method proposed in this study is evaluated based on the following criteria:

- **Energy Consumption:** Energy consumption is evaluated in terms of three variables in the simulations and these are: distance between nodes, number of nodes, and data transmission frequency. Remaining energy after the data is transmitted is proportional to the performance of image compression algorithm used.
- **PSNR Loss Value:** PSNR value of the compressed image is measured in the simulations and image transmission is realized only when the PSNR value of the image decreases 20-25 dB. This interval is the acceptable loss measure for WMSN [42]. Also, the amount of data to be transmitted will be decreased by increasing the compression rate until the acceptable PSNR loss is satisfied. Therefore, the sensor nodes will consume less transmission and acquisition energy.

Calculation of the energy rate is realized under the consideration of the transmission/acquisition energies of the carrier nodes while they carry compressed data towards the sink. That is, a node inside the network takes the image, decides the proper image compression algorithm thanks to EALAIC, and the decided compression algorithm is applied to the image. Then, it sends the image to the next node according to the routing protocol. The neighbor node who takes the compressed data works only as a carrier and forwards the data to the next node. The carrier nodes spend only transmission and acquisition energies meantime. They do not spend computation energy.

Images used for the analysis are given in Section 2.2. Different sample scenarios are made for the performance evaluation and the sensor nodes are made the received images send to the next node under the below conditions. The conditions are:

- Without compression,
- Using only DCT,
- Using only SPIHT,
- Using only LEICA, and
- Using EALAIC.

3.2. Simulation Parameters

Scenarios and simulations are first realized for stationary networks. Then, the simulations are carried out for the networks with mobile nodes whose speeds are 5 m/s and 10 m/s. The nodes are placed by the grid model for stationary network simulation. As for mobile networks, the sensor nodes are first placed according to the grid model and then they are moved randomly. The simulation parameters are given in Table 6 and the images are given in Figure 12.

Table 6. Simulation parameters

Number of nodes	17, 101, 197
Packet size	256 Kbyte
Distance between nodes (meter)	10 - 200
Data transmission frequency (second)	3, 9, 15
Simulation time (second)	120

3.3. Energy Rate

1000-joule energy is assigned to every node as starting energy. The total energy is determined before the simulations according to the number of nodes in the network. The remaining energy is calculated after the simulation time ends. The rate of the remaining energy to the total energy is expressed as “energy rate”. Energy rate is a performance criterion that determines the energy efficiency of the simulation realized. Energy efficiency is determined as good for the cases in which the energy rate is high.

Figure 14 and Figure 15 show the graphics of a stationary network. EALAIC catches the best energy rate by using SPIHT with 10% compression ratio when data transmission

frequency is 15 seconds, by using DCT with 5% compression ratio when data transmission frequency is 9 seconds, and by using LEICA with 20% compression ratio when data transmission frequency is 3 seconds in Figure 14. The corresponding energy rates are given in Table 7.

EALAIC catches the best energy rate by using LEICA with 5% compression ratio when the number of nodes is 17, by using SPIHT with 15% compression ratio when the number of nodes is 101, and by using SPIHT with 10% compression ratio when the number of nodes is 197 in Figure 15. The corresponding energy rates are given in Table 8.

Table 7. Scenario: energy rates of five methods with respect to data transmission frequency in the case that “Distance”-low and “Number of Nodes”-medium

Data Transmission Frequency (sec)	Energy Rate				
	DCT	SPIHT	LEICA	EALAIC	Without Compression
15	0.9918469	0.9919395	0.9919179	0.9921035	0.9887133
9	0.9912076	0.9911705	0.9910011	0.9912076	0.9819545
3	0.9891125	0.989661	0.9897173	0.9901256	0.9727593

Table 8. Scenario: energy rates of five methods with respect to number of nodes in the case that “Distance”-medium and “Data Transmission Frequency”-low

Number of Nodes	Energy Rate				
	DCT	SPIHT	LEICA	EALAIC	Without Compression
17	0.9938447	0.9939762	0.9944439	0.9944439	0.9928777
101	0.9930386	0.9930810	0.9930327	0.9932125	0.9900944
197	0.9925863	0.9926916	0.9923728	0.9927954	0.9889077

Figure 16 and Figure 17 show the graphics of a mobile network (nodes are moving with speeds 5 m/s and 10 m/s). EALAIC catches the best energy rate by using LEICA with 10% compression ratio when data transmission frequency is 15 seconds, by using SPIHT with 5% compression ratio when data transmission frequency is 9 seconds, and by using LEICA with 10% compression ratio when data transmission frequency is 3 seconds in Figure 16. The corresponding energy rates are given in Table 9.

EALAIC catches the best energy rate by using DCT with 10% compression ratio when the distance between nodes is 40 meters, by using DCT with 15% compression ratio when the distance between nodes is 90 meters, and by using LEICA with 5% compression ratio when the distance between nodes is 190 meters in Figure 17. The corresponding energy rates are given in Table 10.

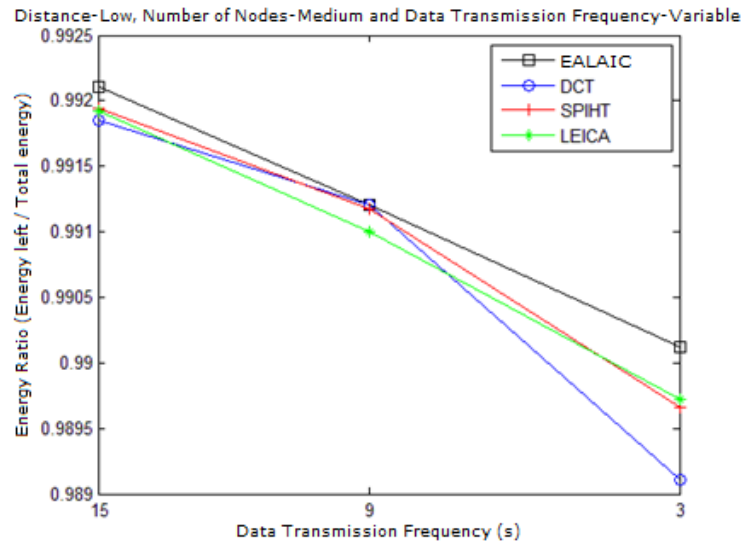


Fig. 14. Scenario: data transmission frequency versus energy rate graphic of four methods in the case that “Distance”-low and “Number of Nodes”-medium

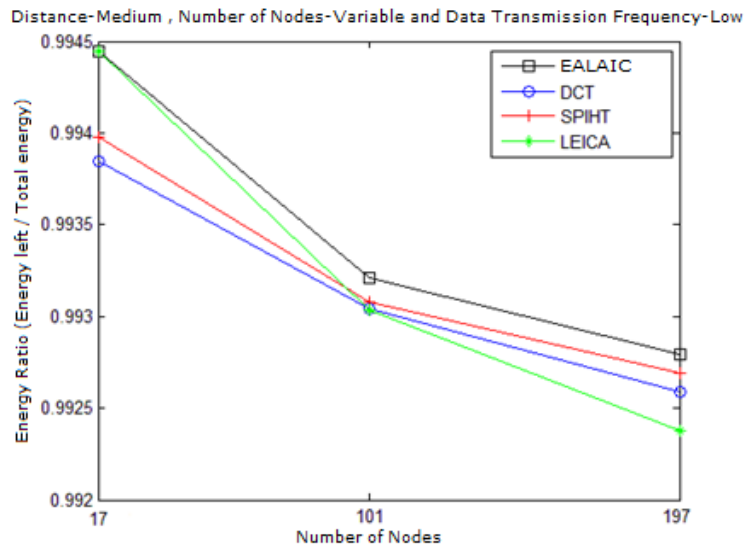


Fig. 15. Scenario: number of nodes versus energy rate graphic of four methods in the case that “Distance”-medium and “Data Transmission Frequency”-low

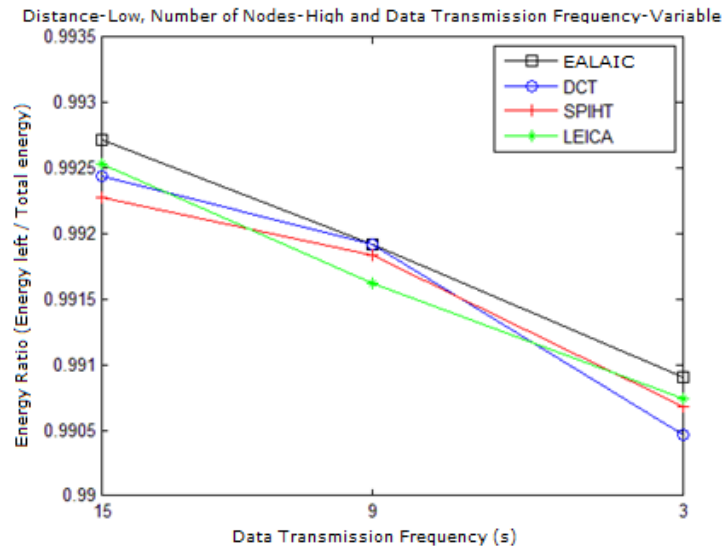


Fig. 16. Scenario: data transmission frequency versus energy rate graphic of four methods in the case that “Distance”-low and “Number of Nodes”-high

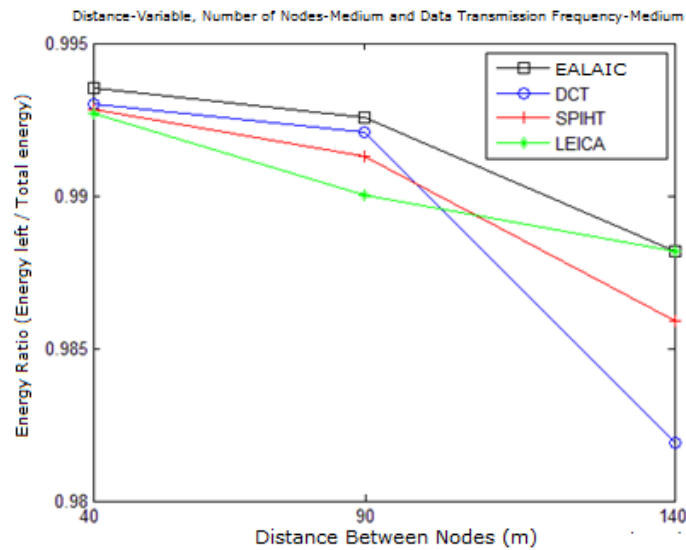


Fig. 17. Scenario: distance versus energy rate graphic of four methods in the case that “Number of Nodes”-medium and “Data Transmission Frequency”-medium

3.4. End to End Delay

Every packet produced in different scenarios is analyzed to obtain end to end delay. In this works, the average of the time to leave the source and reach the target for every packet

Table 9. Scenario: energy rates of five methods with respect to data transmission frequency in the case that “Distance”-low and “Number of Nodes”-high

Data Transmission Frequency (sec)	Energy Rate				
	DCT	SPIHT	LEICA	EALAIK	Without Compression
15	0.9924413	0.9922819	0.9925320	0.9927140	0.9905598
9	0.9919207	0.9918407	0.9916191	0.9919207	0.9885940
3	0.9904621	0.9906785	0.9907414	0.9909013	0.9852178

Table 10. Scenario: energy rates of five methods with respect to distance in the case that “Number of Nodes”-medium and “Data Transmission Frequency”-medium

Distance (m)	Energy Rate				
	DCT	SPIHT	LEICA	EALAIK	Without Compression
40	0.9929828	0.9928449	0.9926847	0.9935465	0.9897579
90	0.9920975	0.9912844	0.9900279	0.9925456	0.9545897
140	0.9819125	0.9859102	0.9881926	0.9881926	0.9325489

which reaches its target correctly is considered. This operation, taking the average time, is repeated separately for each algorithm (DCT, SPIHT, LEICA, and EALAIK) and the packet transmission path is designed in the same manner for every algorithm. That is, in the simulations realized for the same network, the path followed by the packets is the same for every algorithm. The compression times of the image compression algorithms are: DCT in 54.843 ms, SPIHT in 57.968 ms, and LEICA in 54.707 ms.

The time required for a packet to reach its target according to the scenario with 17 nodes is given in Table 11 for each algorithm. EALAIK decided DCT as the most energy-efficient image compression algorithm in this scenario and used it for compression.

Table 11. End to end transmission time of a packet for a scenario with 17 nodes according to the used methods

	Used Method			
	DCT	SPIHT	LEICA	EALAIK
Elapsed Time (ms)	403,800	406,925	403,664	404,488

It is shown in Table 11 that EALAIK has better end to end delay value than SPIHT and worse than DCT and LEICA. The reason is that EALAIK uses the best algorithm it selected at every turn according to the network conditions and the compression times of these algorithms are different. Therefore, having varying values for different scenarios is natural.

3.5. Comparison of the Methods in terms of Energy Consumption

The methods of the referenced studies and the developed method EALAIC are compared in terms of the energy consumption which occurs while the image is compressed and transmitted. The same network model is used in the comparisons (distance between nodes is low, number of nodes is high, and data transmission frequency is variable). The average energy consumption of DCT [5], SPIHT [54], and LEICA [41] methods are calculated. Similarly, the average energy consumption of EALAIC depending on time is calculated using the mentioned network model.

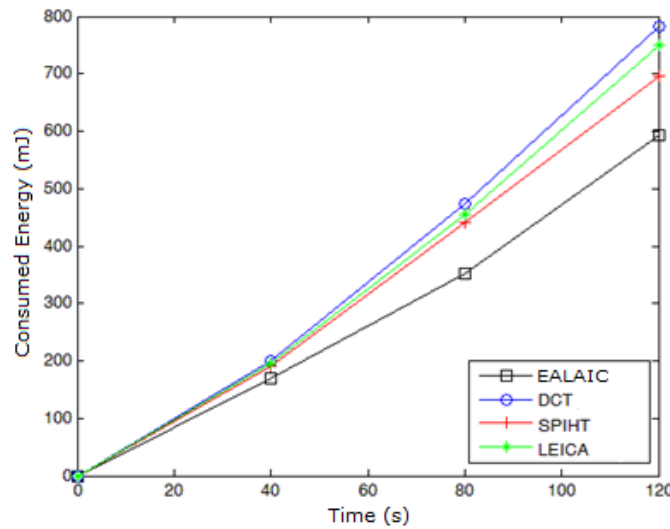


Fig. 18. Comparison of four methods based on the sample scenario depending on time

Figure 18 gives the comparison about time-dependent energy consumption of a node for EALAIC and other methods (DCT [5], SPIHT [54], and LEICA [41]) based on the network model in the sample scenario. The methods have worked for certain times and the average energy consumption of the network is calculated. Accordingly, It is seen that EALAIC consumes 103mJ less energy than the other methods.

This situation can be expressed clearly in terms of the energy consumption rates as: the value 1000 joules can be used approximately 47.6 hours by the least energy consuming algorithm, SPIHT on the average and it can be used by EALAIC for 55.5 hours. A node that uses EALAIC developed according to SPIHT has 7.9 hours more lifetime.

The proposed energy-efficient image compression-based approach was also compared by the state-of-the-art in order to show its performance. The existing papers were compare in Table 12 concerning the consumed energy.

Table 12. Comparison of the state-of-the-art

Reference	Year	Method	Consumed Energy
[43]	2011	Wavelet based image compression	more than 500 mJ
[17]	2015	Compressed sensing based image transmission	350-400 mJ
[31]	2018	Image segmentation	360-380 mJ
Proposed	2020	Adaptive image compression	100-200 mJ

4. Conclusion

A more effective method than the existing popular compression techniques used for WMSN has been proposed in this paper. The solution of the energy efficiency problem is given to be applied in the application layer. The proposed algorithm is different from the others since it decides the compression algorithm to be used through the parameters “distance between nodes”, “number of nodes”, and “data transmission frequency”. Therefore, the most suitable compression algorithm is selected instantly and maximum efficiency is provided. The method decides which algorithm to use according to the instant network conditions and so, it consumes less energy than the traditional methods that uses the same compression method for every condition. EALAIC provides 51% energy save for the stationary networks according to the methods without compression, and 36% energy save for the mobile networks. Also, for stationary networks, it saves 3% energy according to SPIHT, 3.6% save according to LEICA, and 4% save according to DCT. As for mobile networks it saves 2.7% energy according to SPIHT, 3.1% save according to LEICA, and 3.4% save according to DCT. It is seen that EALAIC prolongs the lifetime of every node by 7.9 hours on the average when 1000-joule energy is assigned to every node.

The proposed application layer method related to energy efficiency based on image compression can be combined with the energy-aware methods of other layers and a hybrid network with longer lifetime can be designed in the future works. Additionally, the efficiency of the proposed method can be tested by higher resolution cameras and the method can be generalized considering the resolution as an additional parameter.

References

1. Abd El Kader, M.E.E.D., Youssif, A.A., Ghalwash, A.Z.: Energy Aware and Adaptive Cross-Layer Scheme for Video Transmission Over Wireless Sensor Networks. *IEEE Sensors Journal* 16(21), 7792–7802 (nov 2016)
2. Abraham, A., Narendra, G.K.: Energy Aware Priority-Based Secure Image Transmission in Wireless Sensor Networks. In: 2012 8th International Conference on Wireless Communications, Networking and Mobile Computing. pp. 1–4 (2012), <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=6478513>
3. Ahmed, N., Natarajan, T., Rao, K.: Discrete Cosine Transform. *IEEE Transactions on Computers* C-23(1), 90–93 (jan 1974), <http://ieeexplore.ieee.org/document/1672377/>
4. Akyildiz, I., Melodia, T., Chowdhury, K.: A survey on wireless multimedia sensor networks. *Computer Networks* 51(4), 921–960 (2007), <http://linkinghub.elsevier.com/retrieve/pii/S1389128606002751>

5. Amutha, R., Phamila, Y.: Energy efficient image compression scheme using ZonalDCT for wireless sensor networks. In: IET Chennai 3rd International Conference on Sustainable Energy and Intelligent Systems (SEISCON 2012). pp. 408–414. Institution of Engineering and Technology (2012), <http://digital-library.theiet.org/content/conferences/10.1049/cp.2012.2247>
6. Ben-Othman, J., Diagne, S., Mokdad, L., Yahya, B.: Performance evaluation of a hybrid MAC protocol for wireless sensor networks. In: Proceedings of the 13th ACM international conference on Modeling, analysis, and simulation of wireless and mobile systems - MSWIM '10. p. 327. ACM Press, New York (2010)
7. Ben-Othman, J., Yahya, B.: Energy efficient and QoS based routing protocol for wireless sensor networks. *Journal of Parallel and Distributed Computing* 70(8), 849–857 (2010)
8. Cobo, L., Quintero, A., Pierre, S.: Ant-based routing for wireless multimedia sensor networks using multiple QoS metrics. *Computer Networks* 54(17), 2991–3010 (2010)
9. Din, I.U., Hassan, S., Almogren, A., Ayub, F., Guizani, M.: PUC: Packet Update Caching for energy efficient IoT-based Information-Centric Networking. *Future Generation Computer Systems* (2019)
10. Ehsan, S., Hamdaoui, B.: A Survey on Energy-Efficient Routing Techniques with QoS Assurances for Wireless Multimedia Sensor Networks. *IEEE Communications Surveys & Tutorials* 14(2), 265–278 (2012), <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=5714976>
11. Er-Rouidi, M., Moudni, H., Mouncif, H., Merbouha, A.: A balanced energy consumption in mobile ad hoc network. In: *Procedia Computer Science*. vol. 151, pp. 1182–1187. Elsevier B.V. (2019)
12. Farrag, O., Younis, M., D'Amico, W.: MAC Support for Wireless Multimedia Sensor Networks. In: *GLOBECOM 2009 - 2009 IEEE Global Telecommunications Conference*. pp. 1–6. IEEE (nov 2009)
13. Ghazalian, R., Aghagolzadeh, A., Andargoli, S.M.H.: Energy consumption minimization in wireless visual sensor networks using convex optimization. In: 7th International Symposium on Telecommunications (IST'2014). pp. 312–315. IEEE (sep 2014), <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=7000720>
14. Gürses, E., Akan, Ö.B.: Multimedia communication in wireless sensor networks. *Annales Des Télécommunications* 60(7-8), 872–900
15. Hamid, Z., Hussain, F.B.: QoS in Wireless Multimedia Sensor Networks: A Layered and Cross-Layered Approach. *Wireless Personal Communications* 75(1), 729–757 (mar 2014), <http://link.springer.com/10.1007/s11277-013-1389-0>
16. Haseeb, K., Islam, N., Almogren, A., Ud Din, I., Almajed, H.N., Guizani, N.: Secret Sharing-Based Energy-Aware and Multi-Hop Routing Protocol for IoT Based WSNs. *IEEE Access* 7, 79980–79988 (2019)
17. Hemalatha, R., Radha, S., Sudharsan, S.: Energy-efficient image transmission in wireless multimedia sensor networks using block-based Compressive Sensing. *Computers and Electrical Engineering* 44, 67–79 (may 2015)
18. Hoon Kim, Sung-Gi Min: Priority-based QoS MAC protocol for wireless sensor networks. In: *2009 IEEE International Symposium on Parallel & Distributed Processing*. pp. 1–8. IEEE (may 2009)
19. Hossain, M.S., You, X., Xiao, W., Lu, J., Song, E.: QoS-oriented multimedia transmission using multipath routing. *Future Generation Computer Systems* 99, 226–234 (oct 2019)
20. K, P.K., Govindaraj, E.: Quality enhancement with fault tolerant embedding in video transmission over WMSNs in 802.11e WLAN. *Ad Hoc Networks* 88, 18–31 (may 2019)
21. Kandris, D., Tsagkaropoulos, M., Politis, I., Tzes, A., Kotsopoulos, S.: Energy efficient and perceived QoS aware video routing over Wireless Multimedia Sensor Networks. *Ad Hoc Networks* 9(4), 591–607 (2011)

22. Kumar S., A.A., Ovsthus, K., Kristensen., L.M.: An Industrial Perspective on Wireless Sensor Networks — A Survey of Requirements, Protocols, and Challenges. *IEEE Communications Surveys & Tutorials* 16(3), 1391–1412 (2014)
23. Li, B.Y., Chuang, P.J.: Geographic energy-aware non-interfering multipath routing for multimedia transmission in wireless sensor networks. *Information Sciences* 249, 24–37 (2013)
24. Lin, M.s., Leu, J.s., Yu, W.c., Yu, M.c., Wu, J.l.c.: On Transmission Efficiency of the Multimedia Service over IEEE 802 . 15 . 4 Wireless Sensor Networks. In: *Advanced Communication Technology (ICACT), 2011 13th International Conference on*. pp. 184–189 (2011)
25. Liu, J., Wang, D., Zhao, J.: An energy-efficient distributed adaptive cooperative routing in wireless multimedia sensor networks. In: *2019 IEEE 11th International Conference on Communication Software and Networks, ICCSN 2019*. pp. 163–169. Institute of Electrical and Electronics Engineers Inc. (jun 2019)
26. Liu, X.: An optimal-distance-based transmission strategy for lifetime maximization of wireless sensor networks. *IEEE Sensors Journal* 15(6), 3484–3491 (2015)
27. Ma, T., Hempel, M., Peng, D., Sharif, H.: A Survey of Energy-Efficient Compression and Communication Techniques for Multimedia in Resource Constrained Systems. *IEEE Communications Surveys & Tutorials* 15(3), 963–972 (2013), <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=6226793>
28. Magli, E., Mancin, M., Merello, L.: Low-complexity video compression for wireless sensor networks. In: *2003 International Conference on Multimedia and Expo. ICME '03. Proceedings (Cat. No.03TH8698)*. vol. 3, pp. III–585. IEEE (2003), <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=1221379>
29. Mainwaring, A., Culler, D., Polastre, J., Szewczyk, R., Anderson, J.: Wireless Sensor Networks for Habitat Monitoring. *Proceedings of the 1st {ACM} International Workshop on Wireless Sensor Networks and Applications* pp. 88–97 (2002), <http://doi.acm.org/10.1145/570738.570751>
30. Mulugeta, T., Shu, L., Hauswirth, M., Chen, M., Hara, T., Nishio, S.: Secured two phase geographic forwarding protocol in wireless multimedia sensor networks. In: *GLOBECOM - IEEE Global Telecommunications Conference* (2010)
31. Pal, T., DasBit, S.: A low-overhead adaptive image compression technique for energy-constrained WMSN. *Computers and Electrical Engineering* 70, 594–615 (aug 2018)
32. Palaskar, S.P., Chavhan, N.A.: Design and Implementation of Energy Efficient Node Data Transmission by Using Mobile Sink Node. In: *2014 Fourth International Conference on Communication Systems and Network Technologies*. pp. 79–82. IEEE (apr 2014), <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=6821361>
33. Pudlewski, S., Prasanna, A., Melodia, T.: Compressed-Sensing-Enabled Video Streaming for Wireless Multimedia Sensor Networks. *IEEE Transactions on Mobile Computing* 11(6), 1060–1072 (jun 2012), <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=6188348>
34. Rahimi, M., Baer, R., Iroezi, O.I., Garcia, J.C., Warrior, J., Estrin, D., Srivastava, M.: Cyclops: In Situ Image Sensing and Interpretation in Wireless Sensor Networks. In: *Proceedings of the 3rd international conference on Embedded networked sensor systems - SenSys '05*. pp. 192–204 (2005), <http://doi.acm.org/10.1145/1098918.1098939>
35. Said, A., Pearlman, W.A.: A new, fast, and efficient image codec based on set partitioning in hierarchical trees. *{IEEE} Transactions on Circuits and Systems for Video Technology* 6(3), 243–250 (1996)
36. Savidge, L., Huang Lee, Aghajan, H., Goldsmith, A.: QoS-based geographic routing for event-driven image sensor networks. In: *2nd International Conference on Broadband Networks, 2005*. pp. 68–77. IEEE
37. Saxena, N., Roy, A., Shin, J.: Dynamic duty cycle and adaptive contention window based QoS-MAC protocol for wireless multimedia sensor networks. *Computer Networks* 52(13), 2532–2542 (2008)

38. Senturk, A., Kara, R.: An analysis of image compression techniques in wireless multimedia sensor networks. *Tehnicki Vjesnik* 23(6) (2016)
39. Shapiro, J.: Embedded image coding using zerotrees of wavelet coefficients. *IEEE Transactions on Signal Processing* 41(12), 3445–3462 (1993), <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=258085>
40. Suh, C., Mir, Z.H., Ko, Y.B.: Design and implementation of enhanced IEEE 802.15.4 for supporting multimedia service in Wireless Sensor Networks. *Computer Networks* 52(13), 2568–2581 (2008)
41. Sun, E., Shen, X., Chen, H.: A low energy image compression and transmission in wireless multimedia sensor networks. In: *Procedia Engineering*. vol. 15, pp. 3604–3610 (2011)
42. Thomos, N., Boulgouris, N.V., Strintzis, M.G.: Optimized transmission of JPEG2000 streams over wireless channels. *IEEE Transactions on Image Processing* 15(1), 54–67 (2006)
43. Tian, F., Liu, J., Sun, E., Wang, C.: An energy efficient and load balancing distributed image compression algorithm in WMSNs. In: *Procedia Engineering*. vol. 15, pp. 3421–3427. Elsevier (jan 2011)
44. Toor, A., ul Islam, S., Sohail, N., Akhunzada, A., Boudjadar, J., Khattak, H.A., Din, I.U., Rodrigues, J.J.: Energy and performance aware fog computing: A case of DVFS and green renewable energy. *Future Generation Computer Systems* 101, 1112–1121 (dec 2019)
45. Ul Islam, S., Khattak, H.A., Pierson, J.M., Din, I.U., Almogren, A., Guizani, M., Zuair, M.: Leveraging utilization as performance metric for CDN enabled energy efficient internet of things. *Measurement: Journal of the International Measurement Confederation* 147 (dec 2019)
46. Ullah, U., Khan, A., Zareei, M., Ali, I., Khattak, H.A., Din, I.U.: Energy-Effective Cooperative and Reliable Delivery Routing Protocols for Underwater Wireless Sensor Networks. *Energies* 12(13), 2630 (jul 2019), <https://www.mdpi.com/1996-1073/12/13/2630>
47. Ur Rehman, Y.A., Tariq, M., Sato, T.: A novel energy efficient object detection and image transmission approach for wireless multimedia sensor networks. *IEEE Sensors Journal* 16(15), 5942–5949 (aug 2016)
48. Wang, P., Dai, R., Akyildiz, I.F.: Collaborative Data Compression Using Clustered Source Coding for Wireless Multimedia Sensor Networks. In: *2010 Proceedings IEEE INFOCOM*. pp. 1–9. IEEE (mar 2010), <http://ieeexplore.ieee.org/document/5462034/>
49. Wei, Z., Lijuan, S., Jian, G., Linfeng, L.: Image compression scheme based on PCA for wireless multimedia sensor networks. *The Journal of China Universities of Posts and Telecommunications* 23(1), 22–30 (2016)
50. Xiaojiang Du, X., Hsiao-Hwa Chen, H.h.: Security in wireless sensor networks. *IEEE Wireless Communications* 15(4), 60–66 (aug 2008), <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=4599222>
51. Yan, X., Li, L., An, F.: Multi-constrained routing in wireless multimedia sensor networks. In: *2009 International Conference on Wireless Communications & Signal Processing*. pp. 1–5. IEEE (nov 2009)
52. 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 Transactions on Networking* 23(3), 810–823 (jun 2015), <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=6757018>
53. Zaineldin, H., Elhosseini, M.A., Ali, H.A.: Image compression algorithms in wireless multimedia sensor networks: A survey (2015)
54. ZainEldin, H., Elhosseini, M.A., Ali, H.A.: A modified listless strip based SPIHT for wireless multimedia sensor networks. *Computers & Electrical Engineering* 56, 519–532 (2016)
55. Zeeshan, M., Majid, M., Nizami, I.F., Anwar, S.M., Ud Din, I., Khurram Khan, M.: A newly developed ground truth dataset for visual saliency in videos. *IEEE Access* 6, 20855–20867 (apr 2018)
56. Zhuo, X., Loo, K., Cosmas, J., Yip, P.: Distributed video coding in wireless multimedia sensor network for multimedia broadcasting (2008)

Dr. Arafat Senturk graduated in Computer Engineering from the Faculty of Engineering of Atilim University in 2009 and earned a PhD in Electrical and Computer Engineering from the Enstitute of Science of Duzce University in 2017. He has been worked as a research assistant at Duzce University for six years and working as assistant prof. since 2017 at the same institution. His interests include wireless sensor networks. <https://orcid.org/0000-0002-9005-3565>

Dr. Resul Kara received the B.S. and M.S. degrees in electrical and electronics engineering from Gazi University, Ankara in 1996 and in 1999, respectively and Ph.D. degrees in electronics engineering from Sakarya University, Sakarya, in 2009. He worked as an assistant professor from 2009 to 2012 and as an associate professor from 2012 to 2017 with the Computer Engineering Department in Duzce University, Duzce. Since 2017, he has been a Professor with the same department. His research interests include computer networks, MANETs, wireless communication and network security. <https://orcid.org/0000-0001-8902-6837>

Dr. Ibrahim Ozelik was born in Izmit, Kocaeli, Turkey in 1973. He received the B.S. degree in electric and electronic engineering from Karadeniz Technical University, Trabzon, in 1995 and M.S. and Ph.D. degrees in electronic engineering from Sakarya University, Sakarya, in 1997 and in 2002, respectively. He worked as a research assistant from 1995 to 2002 and as an assistant professor from 2002 to 2012 and as an associate professor from 2012 to 2020 with the Computer Engineering Department in Sakarya University. Since 2020, he has been a Professor with the same department. He is the author of more than 30 articles. His research interests include computer networks, VANETs, industrial communication systems, cyber security and SCADA security. <https://orcid.org/0000-0001-9985-5268>

Received: November 24, 2019; Accepted: May 23, 2020.

Variational Neural Decoder for Abstractive Text Summarization

Huan Zhao, Jie Cao, Mingquan Xu, and Jian Lu

College of Computer Science and Electronic Engineering, Hunan University,
Changsha, Hunan, P.R. China, 410000
{hzhao,jiecao,hb_xmq,jianlu}@hnu.edu.cn

Abstract. In the conventional sequence-to-sequence (seq2seq) model for abstractive summarization, the internal transformation structure of recurrent neural networks (RNNs) is completely determined. Therefore, the learned semantic information is far from enough to represent all semantic details and context dependencies, resulting in a redundant summary and poor consistency. In this paper, we propose a variational neural decoder text summarization model (VND). The model introduces a series of implicit variables by combining variational RNN and variational auto-encoder, which is used to capture complex semantic representation at each step of decoding. It includes a standard RNN layer and a variational RNN layer [5]. These two network layers respectively generate a deterministic hidden state and a random hidden state. We use these two RNN layers to establish the dependence between implicit variables between adjacent time steps. In this way, the model structure can better capture the complex semantics and the strong dependence between the adjacent time steps when outputting the summary, thereby improving the performance of generating the summary. The experimental results show that, on the text summary LCSTS and English Gigaword dataset, our model has a significant improvement over the baseline model.

Keywords: abstractive summarization, sequence-to-sequence, variational auto-encoder, variation neural inferer.

1. Introduction

Text summarization produces a brief summary of the core ideas of the source articles and is different from extractive text summarization ([4], [23],[30], [19]), which selects key sentences or key phrases in the original text to form a summary. Abstractive text summarization builds an internal semantic representation and then uses deep learning techniques to create a summary that is closer to what a human may generate. Most recent models for abstractive text summarization are based on the seq2seq framework with attention ([2],[21],[12],[29]). These seq2seq models consist of an encoder and a decoder; the encoder encodes input text into a semantic representation, and the decoder generates summaries from this representation.

With the development of deep learning, the neural networks-based encoder-decoder models are used in the sequence-to-sequence tasks, such as neural machine translation[28], speech recognition ([32],[33]) and text summarization ([3],[9],[25]). The Seq2seq framework for abstractive text summarization has recently achieved remarkable success and

has become the dominant architecture. In the seq2seq framework, the semantic representation from the encoding end to the decoding end is learned in an implicit way, in which the internal transformation structure of recurrent neural networks (RNNs) is completely determined. [5]. Therefore, the learned semantic representations are poor at capturing all semantic details and dependencies [29],[14]. To address the insufficiency of semantic representations of abstractive text summarization, [20] proposed a generative model to capture the latent summary information, but they considered only a single latent variable for capturing the global semantics of each source text in their generative model, which led to limited representation ability. Furthermore, [14] introduced a seq2seq framework with a deep recurrent generative decoder that considered the recurrent dependencies in their generative model for capturing historical latent variable dependencies. Although this approach can obtain more latent structure information, in practice, long-term sequential recurrent dependencies can result in the loss of previous information and unnecessary noise; hence, this implementation may not be sufficient for capturing strong and complex semantic dependencies between adjacent target words at each time step of the decoding.

To tackle the problem, we present a variational neural decoder (VND) for abstractive text summarization that is more effective at forcing the decoder to make use of latent structure information. We introduce the variational autoencoder (VAE) [11],[24] process to the decoding process and use latent variables to model the complex potential distribution of text semantics at each time step. Drawing inspiration from the current success of the variational RNN (VRNN) [5], we incorporate latent variables into the RNN hidden state. By using latent variables, the VRNN can model the underlying semantics of source or target texts. Our decoding structure consists of variational neural inferers and two RNN layers: a standard RNN layer and a VRNN layer. A variational neural inferer is employed to address the intractable posterior inferer for the latent variables. The standard RNN layer generates a deterministic hidden state, which is employed to model long- and short-term dependencies. The stochastic latent hidden state based on the VRNN layer is used to capture complex and strong potential semantic distributions and is integrated into the summary generation softmax layer to improve the summary generation quality. Specifically, at each time step, we use the stochastic latent hidden state of the VRNN layer in the previous step as the input of the current RNN layer. This implementation integrates the dependencies between the latent variables in adjacent timesteps.

The main contributions of this paper are as follows: (1) we propose a VND model that efficiently captures the complex semantics and strong dependencies between neighboring target words for abstractive text summarization. (2) Experiments on the LCSTS dataset and English Gigaword for the text summarization task show that our proposed model significantly outperforms the baseline models.

2. Related Work

Automatic text summarization is one of the most active research in natural language processing. It produces a concise and smooth summary while preserving key information content and overall meaning [1]. Recently, an increasing number of researchers have employed a neural network framework to natural language processing. Sequence-to-sequence neural networks[29] have been applied to machine translation ([28],[17]), following their success in abstractive summarization ([2],[13]).

Specifically, [25] first proposed a convolution encoder and a recurrent decoder model for the abstractive sentence summarization task, which has achieved significant performance improvement over conventional methods, and provides the benchmark for the Gigaword dataset. [22] replaced the model with a full RNN seq2seq model and achieved outstanding performance. [17] proposed an attention mechanism, which greatly improved the performance of the seq2seq model on abstractive summarization. To address the unknown word problem, [21] proposed a generator-pointer model so that the decoder can generate words in source texts. [6] also solved this problem by integrating a copying mechanism into a seq2seq model. [27] propose an LSTM-CNN based seq2seq model that can construct new sentences by exploring more fine-grained fragments than sentences, namely, semantic phrases. [18] proposed a neural model to improve the semantic relevance between the source contents and the predicted summaries.

Some other work attempts to incorporate Variational auto-encoder for abstractive summarization. The variational auto-encoder is a popular probabilistic generative model ([5],[11]). These models utilize an neural inference model to approximate the intractable posterior, and optimize model parameters jointly with a reparameterized variational lower bound using the standard stochastic gradient technique. Due to its success in various tasks, this method has attracted increasing attention. Although seq2seq-oriented encoder-decoder framework has been developed and has widely used in abstractive text summarization, there are few research works incorporated variational auto-encoder into the text summary system. For example,[20] first proposed a generative model to capture the latent summary information based on the seq2seq framework. [26] presented an unsupervised approach to summarize sentences abstractively using a VAE. Furthermore, [5] extended the VAE into a recurrent framework for modeling complex semantic representations, which is called VRNN. [14] proposed a deep recurrent generative decoder to capture latent structure information.

Inspired by the successful application of variational auto-encoder in related works, we propose variational neural decoder for abstractive summarization. This paper proposes a variational neural decoder model, which introduces a series of continuous latent variables to capture the latent semantics of the content to improve the quality of the summary.

3. Background: Variational Autoencoder

The VAE [11],[24] is a recently introduced latent variable generative model, which combines Variational Inference with Deep Learning. In VAE, a generative network models an observed variable x as a continuous latent variable z , based on which the generate network reconstructs x . Then, the join distribution density function of the generated model is as follows:

$$p_{\theta}(x, z) = p_{\theta}(x|z)p_{\theta}(z), \quad (1)$$

where $p_{\theta}(z)$ is the probability density function of z prior distribution of latent variable, $p_{\theta}(x|z)$ is the conditional probability density function of x when z is known, and θ is the parameter of two density functions. In general, we assume that $p_{\theta}(z)$ and $p_{\theta}(x|z)$ are the standard standard Gaussian distribution that models the generation procedure, which is typically estimated via a deep nonlinear neural network.

Importantly, the VAE models the conditional distribution $p_{\theta}(x|z)$ as a highly flexible function approximator, which makes the inference of the posterior $p_{\theta}(z|x)$ intractable.

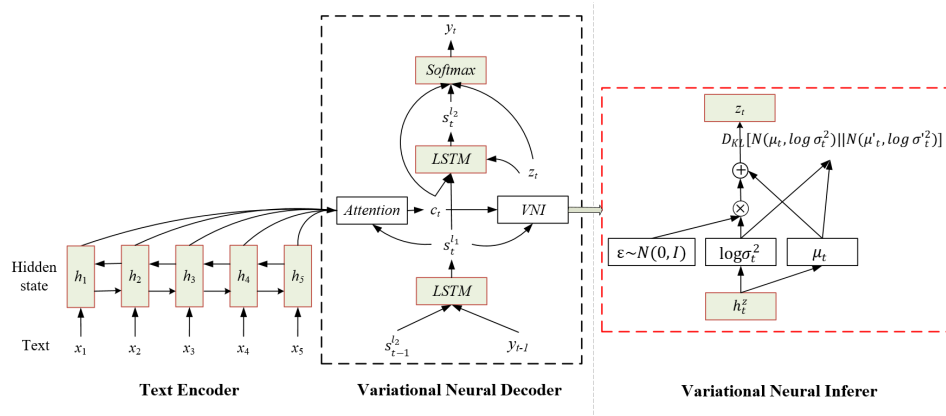


Fig. 1. The overview of the architecture of VAD model.

Thus VAE uses a variational approximation $q_\phi(z|x)$ of the posterior, which introduces the evidence lower bound:

$$\mathcal{L}_{VAE}(\theta, \phi, x) = -D_{KL}(q_\phi(z|x)||p_\theta(z)) + E_{q_\phi(z|x)}[\log p_\theta(x|z)] \leq \log p_\theta(x), \quad (2)$$

where θ and ϕ denote the parameters of the model and $D_{KL}(Q||P)$ is the Kullback-Leibler divergence between two distributions Q and P .

In [11], the approximate posterior $q_\phi(z|x)$ is a diagonal Gaussian $\mathcal{N}(\mu, \text{diag}(\sigma^2))$, whose mean μ and variance σ^2 are the output of a highly nonlinear function of x . The VAE training process maximizes the ELBO, which obtains the optimal parameter selection for the generative model $p_\theta(x|z)$ and inference model $q_\phi(z|x)$. Based on reparametrizing, we compute $z = \mu + \sigma \odot \epsilon$ and rewrite the equation as

$$E_{q_\phi(z|x)}[\log p_\theta(x|z)] = E_{p_\theta(\epsilon)}[\log p_\theta(z = \mu + \sigma \odot \epsilon)], \quad (3)$$

where ϵ is a vector of standard Gaussian variables. Then, the VAE model can be trained through a standard backpropagation technique for stochastic gradient descent.

4. Proposed Model

4.1. Overview

Our model is based on the seq2seq model with attention. The seq2seq model can compress source texts $x = \{x_1, x_2, \dots, x_M\}_N$ into a continuous vector representation with an encoder, and then the decoder generates the summary text $y = \{y_1, y_2, \dots, y_T\}_N$. As shown in Figure 1, VND model mainly contains three neural network based components: the text encoder for encoding text sentences, a variational neural decoder that generate a summary, a variational neural inferer for the posterior and the prior distributions.

4.2. Text Encoder

The text encoder builds meaningful representations of the source sentences. In our model, we use bidirectional long short-term memory (*LSTM*) [8] to encode the source text sequence x from both directions and compute the hidden states for each word, which produces the final hidden state $h = \{h_1, h_2, \dots, h_M\}_N$ from the source text x :

$$\vec{h}_i = LSTM(\vec{h}_{i-1}, x_i, \vec{C}_{i-1}), \quad (4)$$

$$\overleftarrow{h}_i = LSTM(\overleftarrow{h}_{i+1}, x_i, \overleftarrow{C}_{i+1}), \quad (5)$$

$$h_i = [\vec{h}_i; \overleftarrow{h}_i], \quad (6)$$

where \vec{h}_i and \overleftarrow{h}_i are the forward and the backward hidden outputs, respectively, x_i is the input at the i -th time step, N is the number of samples, and C_{i-1} refers to the cell state in the *LSTM* layer.

The transition equations of *LSTM* are defined as follows:

$$I_i = \sigma(W_I x_i + U_I h_{i-1} + b_I) \quad (7)$$

$$F_i = \sigma(W_f x_i + U_f h_{i-1} + b_f) \quad (8)$$

$$O_i = \sigma(W_o x_i + U_o h_{i-1} + b_o) \quad (9)$$

$$C_i = F_i \odot C_i + I_i \odot \tanh(W_c x_i + U_c h_{i-1} + b_c) \quad (10)$$

$$H_i = O_i \odot \tanh(C_i) \quad (11)$$

where \odot stands for element-wise multiplication, σ is the sigmoid function, all $W \in \mathbb{R}^{d \times l}$ and $U \in \mathbb{U}^{d \times d}$ are weight matrices, all $b \in \mathbb{R}^d$ are bias vector.

4.3. Variational Neural Inferer

In order to integrate the stochastic latent variables into our decoder model, we use the variational neural inference to generate the latent variables at each time step of the decoding structure. As described in the background section, the key of variational models is to model the distributions related to latent variables. The variational neural inferer can be divided into two parts: the posterior and the prior distributions. Both posterior and prior distributions are assumed to be multivariate Gaussian distributions with diagonal covariance, but they introduce different parameters. As determined by the ELBO equation, the parameters of the prior are computed by the prior network, which takes the source sentence x and previously generated words $y_{<t}$ as the input. The posterior parameters are also determined from both the source sentence x and previously generated words $y_{<t}$.

Inspired by some ideas in previous works [11],[24]. We use variational neural inference to generate latent variables z in the model, as shown in part variational neural inferer of the Fig. 1, we focus on the use of neural network to simulate the posterior $q_\phi(z_t|x, y_{<t})$ and the prior $p_\theta(z_t|x, y_{<t})$. In this subsection, we use a similar network architecture to that proposed in variational recurrent neural machine translation (VRNMT) [28].

Neural Posterior. Following the standard VAE, we use neural networks for a better approximation in our model. The equation of $q_\phi(z_t|x, y_{<t})$ can be expressed as

$$q_\phi(z_t|x, y_{<t}) = \mathcal{N}(z_t; \mu_t(x, y_{<t}), \sigma_t(x, y_{<t})^2 I), \quad (12)$$

where the μ_t and σ_t denote the variational mean and standard derivation, respectively, which are computed via a neural network based on the observed variables x and $y_{<t}$.

Starting from the VAE, the key to estimating z_t is to calculate the μ_t and σ_t . First, we perform a nonlinear transformation that projects the word embedding y_{t-1} , the deterministic hidden states $s_t^{l_1}$ and the attention content c_t onto our concerned latent semantic space:

$$h_t^z = g(W_z[y_{t-1}; s_t^{l_1}; c_t] + b_z). \quad (13)$$

Then, the above-mentioned Gaussian parameters μ_t and $\log \sigma_t^2$ are calculated through linear regression:

$$\mu_t = W_\mu h_t^z + b_\mu, \quad (14)$$

$$\log \sigma_t^2 = W_\sigma h_t^z + b_\sigma, \quad (15)$$

where W_z , W_μ and W_σ comprise the parameter matrix and b_z , b_μ and b_σ are bias terms. $g(\cdot)$ refers to a nonlinear function. Then, to obtain a representation for the latent variable z_t using reparameterization [24], the latent variables can be expressed as:

$$z_t = \mu_t + \log \sigma_t^2 \odot \epsilon, \epsilon \sim \mathcal{N}(0, I). \quad (16)$$

Intuitively, this reparameterization reduces the gap between $q_\phi(z_t|x, y_{<t})$ and $p_\theta(z_t)$. In other words, it connects these two neural networks. This is important since it enables the stochastic gradient optimization via standard backpropagation.

Neural Prior. The neural model for the prior $p_\theta(z_t|x, y_{<t})$ is the same as that for the posterior $q_\phi(z_t|x, y_{<t})$. Here, we model the prior $p_\theta(z_t|x, y_{<t})$ as

$$p_\theta(z_t|x, y_{<t}) = \mathcal{N}(z_t; \mu'_t(x, y_{<t}), \sigma'_t(x, y_{<t})^2 I). \quad (17)$$

To obtain a representation for latent variable z_t , we first use the same method to employ the latent semantic space:

$$h_t^{z'} = g(W'_z[y_{t-1}; s_t^{l_1}; c_t] + b_z) \quad (18)$$

Then, Gaussian parameters μ'_t and $\log \sigma_t'^2$ in the model are computed by:

$$\mu'_t = W'_\mu h_t^{z'} + b'_\mu, \quad (19)$$

$$\log \sigma_t'^2 = W'_\sigma h_t^{z'} + b'_\sigma, \quad (20)$$

where W'_z , W'_μ and W'_σ comprise the parameter matrix and b'_z , b'_μ and b'_σ are bias terms. Different from the posterior, we directly set z_t as μ'_t , as implemented in [31].

Finally, we integrate the latent variable z_t into decoding our model to enhance the summary generation results, which are described in detail in the following subsection.

Variational Lower Bound. As in the conventional VAE, we learn the generative and inference models jointly by maximizing the variational lower bound with respect to their parameters. We apply ELBO at each t -th time step, and based on factorizations (12) and (17), we have the accumulative ELBO as follows:

$$\mathcal{L}_{VAE} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^T D_{KL}(q_\phi(z_t^{(n)}|x^{(n)}, y_{<t}^{(n)}) || p_\theta(z_t^{(n)}|x^{(n)}, y_{<t}^{(n)})). \quad (21)$$

In the model training, \mathcal{L}_{VAE} is part of the objective function.

4.4. Variational Neural Decoder

The function of the decoder is to generate a series of summary words. As shown in the VND section of Figure 1, at each decoding time step t , the decoder consists of two RNN layers.

The first layer is a standard RNN layer. Given the previously generated word embedding y_{t-1} and the previous stochastic latent hidden state $s_{t-1}^{l_2}$, the standard RNN layer is used to calculate the deterministic hidden state $s_t^{l_1}$ at the t -th time step:

$$s_t^{l_1} = LSTM_1([y_{t-1}; s_{t-1}^{l_2}], C_{t-1}^{l_2}), \quad (22)$$

where the superscript l denotes the decoder LSTM layer.

Next, we apply the dot attention mechanism [17] to obtain the content vector. Then, the deterministic hidden state $s_t^{l_1}$ and the encoder output h_i at each time step t of the process are computed as the attention weight $\alpha_{t,i}$ and the current content vector c_t , respectively:

$$c_t = \sum_{i=1}^m \alpha_{t,i} h_i, \quad (23)$$

$$\alpha_{t,i} = \frac{\exp(e_{t,i})}{\sum_{j=1}^n \exp(e_{t,j})}, \quad (24)$$

$$e_{t,i} = s_t^{l_1 \top} W_e h_i, \quad (25)$$

where W_e is the weight parameter and $\alpha_{t,i}$ is the i -th word of the attention mechanism that assigns weight at time step t .

The second layer is the variational RNN layer. In [5], the author extended the VAE into a recurrent framework for modeling complex semantic representations, which was called VRNN. In VRNN, at each t -th time step, the transition function (LSTM cell) computes the next hidden state based on the previous hidden state and the sampled latent random variables. Inspired by [5], we first combine the deterministic hidden state $s_t^{l_1}$, the current content vector c_t , and the latent variable z_t (implemented in subsection 3.3) to construct a new latent semantics vector \hat{s}_t :

$$\hat{s}_t = W_c c_t + W_s s_t^{l_1} + W_z z_t + b_s, \quad (26)$$

where W_y , W_c , W_s , and W_z are the weight matrices and b_s is the bias term. Then, a VRNN layer is used with the latent semantics vector \hat{s}_t to compute the stochastic latent hidden state $s_t^{l_2}$:

$$s_t^{l_2} = LSTM_2(\hat{s}_t, C_t^{l_1}). \quad (27)$$

Finally, to precisely generate summaries, the softmax layer is introduced to generate the target word y_t based on the latent variable z_t , the current context vector c_t and the stochastic latent hidden state $s_t^{l_2}$. We compute the probability distribution over the target word y_t :

$$p(y_t | z_t, y_{<t}, x) = \text{softmax}(W_v \hat{u}_t + b_v), \quad (28)$$

$$\hat{u}_t = g(W_d [z_t; s_t^{l_2}; c_t] + b_d), \quad (29)$$

where W_v and W_d comprise the parameter matrix of the output layer, b_v and b_d are the bias terms, and $g(\cdot)$ is the nonlinear activation function.

4.5. Objective Function

The objective function of our model consists of two terms. The first objective is the variational lower bound \mathcal{L}_{VAE} in Equation 21. This term is the KL divergence between two Gaussian distributions, which can be computed and differentiated without estimation [11]. The second objective is the maximum likelihood estimation of the generated summaries. Given the latent variable z_t at each time-step and source text x , the models generate a summary \tilde{y} . The learning process is to minimize the negative log-likelihood between the generated summary \tilde{y} and the reference y :

$$\mathcal{L}_{Seq} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^T \{p(y_t^{(n)} | \tilde{y}_{<t}^{(n)}, x^{(n)}, z_t^{(n)})\}. \quad (30)$$

Finally, the objective function, which need to be minimized, is formulated as follows:

$$\mathcal{L} = \mathcal{L}_{VAE} + \mathcal{L}_{Seq}. \quad (31)$$

5. Experiments

5.1. Datasets

LCSTS is a large-scale Chinese short text summarization dataset constructed by [9]. The dataset was collected from a famous Chinese social media website called Sina Weibo and consists of more than 2.4 M text-summary pairs. It is split into three parts, with 2,400,591 pairs in PART I, 10,666 pairs in PART II and 1,106 pairs in PART III. All text-summary pairs in PART II and PART II are manually annotated with relevant scores ranging from 1 to 5. We reserve only pairs with scores of no less than 3, leaving 8,685 pairs in PART II and 725 in PART III. In our experiments, we selected PART I as the training set, PART II as the validation set, and PART III as the test set.

English Gigaword is a sentence summarization dataset based on Annotated Gigawords [22], a dataset consisting of sentence pairs representing the first sentence of collected news articles and their corresponding headlines. We use the dataset preprocessed by [25], which contains 3.8M training pairs, and 189K validation pairs and 2K test pairs were randomly selected. In the data processed directly with [25], we find that there are abnormal data pairs in the corpus, e.g., some unreadable or incomprehensible pairs. Therefore, we reprocessed the data set, therein retaining 3.2M training pairs, a 16K-pair validation set, and a 1,520-pair test set.

5.2. Evaluation Metrics

We employ the recall-oriented understudy for gisting evaluation (ROUGE) score[15] as our evaluation metric with standard options. ROUGE measures the quality of a summary by computing overlapping lexical units, such as unigram, bigram, trigram, and longest common subsequence (LCS). Following previous work [9], our evaluation metrics in the experimental results are the F1 scores of ROUGE: ROUGR-1 (unigram), ROUGE-2 (bigram) and ROUGE-L (LCS).

5.3. Experimental Settings

For the LCSTS dataset, to avoid the effect of Chinese word segmentation error, both our encoder and decoder input texts are Chinese characters. The vocabularies are extracted from training sets, and the source texts and the summaries do not share the same vocabularies. We prune resource and target vocabularies of 8K and 5K words, respectively.

For the English Gigaword dataset, we prune the resource and target vocabularies of 40K and 30K words, respectively. The input word embedding is initialized randomly and learned during the optimization process.

In both experiments, we set the word embedding size and the hidden size to 512, and the number of LSTM layers is 3. The batch size is set to 64, and we do not use dropout on this dataset. We implement the beam search and set the beam size to 10. We use the Adam optimizer [10] to learn the model parameters with the default setting $\alpha = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 1 \times 10^{-8}$. The learning rate is halved every epoch. Our experiments are implemented using PyTorch.

5.4. Baselines

To evaluate the performance of our variational neural decoder (**VND**) model, we compare our results with the results of the baselines and state-of-the-art methods on the LCSTS dataset and Gigaword dataset.

- **ABS** and **ABS+** [25] are the Seq2seq model with attention mechanism and hand-crafted features, which are trained on Gigaword to produce summaries.
- **RNN** and **RNN-context**: [9] are the neural network framework, where RNN+context has the attention mechanism.
- **RNN+distract** [3] employ a new attention mechanism by distracting the historical attention in the decoding steps.
- **DRGD** [14] is a deep recurrent generative decoder model, combining the decoder with a variational auto-encoder.
- **ASC+FSC₁** [20] uses a generative method to model the latent summary variables. The generative model first draws a latent summary sentence from a background language model and subsequently draws the observed sentence conditioned on this latent summary.
- **ARL** [7] is based on the seq2seq framework, which applies the adversarial reinforcement learning strategy to bridge the gap between the generated summary and the human summary.
- **CGU** [16] is a seq2seq model with a convolutional gated unit for global encoding.
- **Seq2seq** is our implementation of the attention mechanism-based seq2seq model, which has a similar setting as our model for comparison.

5.5. Results Analysis

Results on LCSTS Corpus. The summary ROUGE-F1 score comparison of different models on the LCSTS dataset is shown in Table 1. The experimental results show that our VND model has a significantly improved score in each ROUGE compared to the base model. First, we compare the model with the seq2seq baseline. The results show that

Table 1. Comparison with other models on the LCSTS test set. **R-1**, **R-2** and **R-L** denote ROUGE-1, ROUGE-2 and ROUGE-L, respectively.

Models	R-1	R-2	R-L
RNN [9]	21.5	8.9	18.6
RNN-context [9]	29.9	17.4	27.2
RNN+distract [3]	35.2	22.6	32.5
DRGD [14]	37.0	24.2	34.2
ARL [7]	39.4	21.7	29.1
CGU [16]	39.4	26.9	36.5
Seq2seq (our impl.)	33.2	22.4	31.8
VND (our model)	42.5	28.5	37.6

the VND model is significantly improved over the Seq2seq model, with ROUGE-1 9.3% higher, ROUGE-2 6.1% higher, and ROUGE-L 5.8% higher. It proves that our model is effective. Then, we also compare the VND model with the remaining base models.

Next, compared with the DRGD model, our VND model exceeds 5.5% on ROUGE-1, 4.3% on ROUGE-2 and 3.4% on ROUGE-L. In the DRGD model, although the author uses a deep recursive generative model to capture the underlying semantics to improve the quality of the summary, its expressive power is limited and no significant improvement can be achieved. Finally, the ARL model uses adversarial reinforcement learning to optimize parameters on the basis of generating an adversarial network. It is currently the most cutting-edge summary generation model, but the VND model still has obvious advantages. The results show that random latent variables and two RNN layer components plus Seq2seq text summarization model can improve the accuracy of text summarization, which shows the effectiveness of variational neural decoder.

In order to further prove that our proposed model can capture more semantic details and enhance the dependency between the time parts before and after decoding to improve the quality of the summary. We give two summary examples of models randomly drawn from the test set. The original text, the original abstract, the abstract generated by the Seq2seq model and the abstract generated by the model in this chapter are shown in Table 2.

In the first example, the Seq2seq model generates a summary containing the phrases "big bank" and "stock market storm," which is tedious and inaccurate. For comparison, the VND model refines such phrases into more concise and natural phrases. In addition, we can see that the abstract generated by the model is closer to the semantics of the reference abstract than the abstract of the Seq2seq model.

In the second example, the main idea of the original text is that China bans tobacco advertising and imposes fines of up to 30,000 yuan on movies and TV series. However, the summary generated by the seq2seq model only contains information about "TV programs". In addition, it also missed the most important information "forbidden", which makes the semantics very different from the original text, and is not coherent and sufficient. In contrast, the summary of our model is more coherent and streamlined. In addition, in the summary generated by the VND model, the number (30000) is the same as the number in the reference text, which means that the integration of random hidden variables into the decoder model can obtain richer semantics and can make the decoding

Table 2. Examples of the generated summaries on the test set of the LCSTS corpus, compared with that of the Seq2seq model and the reference.

<p>Source(1): 央行今日将召集大型商业银行和股份制银行开会，以应对当前的债市风暴。消息人士表示，央行一方面旨在维稳银行间债券市场，另一方面很可能探讨以丙类户治理为重点的改革内容。此次债市风暴中，国家审计署扮演了至关重要的角色。</p> <p>The central bank will convene large commercial banks and joint-stock banks today to deal with the current debt market turmoil. Sources said that on the one hand, the central bank aims to stabilize the inter-bank bond market, on the other hand, it is likely to explore the reform content focusing on the governance of class C households. The National Audit Office played a vital role in the debt market turmoil.</p> <p>Reference: 媒体称央行今日召集银行开会应对当前债市风暴 The media said that the Central Bank called a meeting of banks today to deal with the current debt market turmoil.</p> <p>Seq2seq: 央行召集大型银行和股份制银行应对债市风暴 The Central Bank calls large banks and joint-stock banks to deal with the storm of the bond market</p> <p>VND: 央行今日召集银行开会应对当前债市风暴 The Central Bank convened a meeting of banks today to deal with the current turmoil in the bond market</p>	<p>Source(2): 国务院法制办公布《公共场所控制吸烟条例(送审稿)》：禁止所有烟草广告、促销和赞助；没有设置室外吸烟点的视为全面禁止吸烟；违反《条例》规定，电影电视剧播放吸烟镜头最高罚3万；禁止通过自动售货机等任何方式向未成年人售烟。</p> <p>The legislative affairs office of the state council announced the "regulations on control of smoking in public places (submitted for review): ban all tobacco advertising, promotion and sponsorship; if there is no outdoor smoking point, it will be considered as a total ban on smoking. Those who violate the regulations can be fined up to 30,000 yuan for smoking in movies and TV series. Selling cigarettes to minors through vending machines or any other means is prohibited.</p> <p>Reference: 我国拟全面禁止烟草广告影视剧播吸烟最高罚3万 China intends to totally ban tobacco in advertisements, movies and TV series. The maximum penalty for smoking is 30,000 yuan.</p> <p>Seq2seq: 我国拟规定电视剧播放吸烟镜头最高罚#万 China intends to stipulate a maximum penalty of # for smoking scenes in TV series.</p> <p>VND: 我国拟规定禁止电影电视剧播放吸烟镜头最高罚3万 China intends to stipulate a maximum penalty of 30,000 yuan for smoking scenes in movies and TV series.</p>
--	--

process adjacent The semantic dependence of the time step is stronger, which makes the generated summary coherent and of high quality.

In addition to the Chinese data set, we also provide several examples of randomly generated summary on the English Gigaword data set. Coherent and more descriptive. We believe that this is mainly because the variational neural decoder can pay

Table 3. ROUGE scores of different models on the English Gigawords dataset.

Models	R-1	R-2	R-L
ABS [25]	29.6	11.3	26.4
ABS+ [25]	29.8	11.9	27.0
ASC+FSC ₁ [20]	34.2	15.9	31.9
DRGD [14]	36.3	17.6	33.6
CGU [16]	36.3	18.0	33.8
Seq2seq (our impl.)	33.5	16.8	31.4
VND (our model)	39.0	19.4	35.3

attention to the potential semantic relationship between the source text and the target sentence. The Seq2seq model cannot obtain more semantics, so it is easy to produce too general phrases and does not necessarily involve the original text.

Table 4. Examples of the generated summaries on the test set of the English Gigawords corpus, compared with that of the Seq2seq model and the reference.

Source(1): jordan's crown prince hassan ibn talal arrived tuesday for his first visit to jerusalem and was to pay his condolences to the widow of assassinated prime minister yitzhak rabin.
Reference: jordan's crown prince makes first visit to jerusalem
Seq2seq: jordan's crown prince arrives in jerusalem
VND: jordan's crown prince arrives for first visit to jerusalem
Source(2): a consortium led by us investment bank goldman sachs thursday increased its takeover offer of associated british ports holdings the biggest port operator in britain after being threatened with possible rival bid.
Reference: goldman sachs increases bid for ab ports
Seq2seq: goldman sachs ups takeover offer of ab british ports
VND: goldman sachs ups increases offer of ab ports
Source(3): the german government and red cross have decided to give ###,### dollars in humanitarian aid to victims of the earthquake which devastated unk in northwest iran the embassy here announced wednesday.
Reference: germany gives ###,### dollars in aid for iran quake victims
Seq2seq: germany to give ###,### dollars in humanitarian aid to iran
VND: germany to give ###,### dollars in aid to iran quake victims
Source(4): global banking giant hsbc said on monday that its pre tax profits had risen in the third quarter despite loan write offs in the united states rising to #.# billion dollars lrb #.# billion euros rrb.
Reference: hsbc says profits rise despite rising us bad debts
Seq2seq: hsbc says profits up despite us loan write offs
VND: hsbc says profits rise despite us loan write offs

Results on Gigaword Corpus. Table 3 presents the test ROUGE F1 score on the Gigaword corpus. Our VND model still outperforms all baseline models and achieves F1

scores of 39.0, 19.4 and 35.3 for ROUGE 1, 2, and L. Comparing with the ABS model, our model performs significantly better by an 9.4 ROUGE-1 F1 score. The ASC+FSC₁ and DRGD models has the highest scores, because both of which incorporates the latent variables into the abstractive summarization model, but in different ways. Compared with these two models, our model can still be better, perhaps the design of our model structure allows for a better enhancement of the summarization performance. In CGU model, the author uses a convolutional gated unit for global encoding to get a high score. In future work, we will try to incorporate this component to our model to improve the quality of the summary.

In addition to ROUGE scores, we also provide several randomly selected examples of generated summaries in Table 4. It can observe that the summaries generated by our model are more coherent and descriptive than those generated by the seq2seq model. We suspect that the main reason for this phenomenon lies in variational neural decoder, the seq2seq model only concerns about the hidden semantics relation behind the source text and target sentence, so it tends to generate phrases which are too general to necessarily refer to the source text.

In the experiment, we train and test our model performance on Chinese and English datasets, and the results on both datasets are better than the baseline model. This proves that our model has good generality in different languages, and also proves that our model can indeed capture complex semantics and strong dependencies through a latent variable that explicitly models the underlying semantics of the source texts to improve the quality of the generated summary.

6. Conclusion

This paper focuses on the problem of redundant and incoherent abstract information in the current field of text summary generation. We propose a new text summarization model based on variational neural decoder. This model combines the advantages of the variational recurrent neural network and the variational decoder to learn complex semantics, and uses a double-layer recurrent network to enhance the strong dependency information between the time before and after the summary output to improve the quality of the summary generation. We evaluate our proposed model on the LCSTS and English Gigaword datasets. The experimental results show that our model outperforms state-of-the-art models and prove that our model captures more strong and complex dependencies to ensure that the generated summary has higher quality.

Acknowledgments. This work was supported by the National Key R&D Program of China No. 2018YFC0831800 and the National Science Foundation of China under Grant 61772188.

References

1. Allahyari, M., Pouriye, S.A., Assefi, M., Safaei, S., Trippe, E.D., Gutierrez, J.B., Kochut, K.: Text summarization techniques: A brief survey. CoRR abs/1707.02268 (2017)
2. Bahdanau, D., Cho, K., Bengio, Y.: Neural machine translation by jointly learning to align and translate. In: ICLR 2015 (2015)

3. Chen, Q., Zhu, X., Ling, Z., Wei, S., Jiang, H.: Distraction-based neural networks for document summarization. CoRR abs/1610.08462 (2016)
4. Cheng, J., Lapata, M.: Neural summarization by extracting sentences and words. In: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers (2016)
5. Chung, J., Kastner, K., Dinh, L., Goel, K., Courville, A.C., Bengio, Y.: A recurrent latent variable model for sequential data. In: Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada. pp. 2980–2988 (2015)
6. Gu, J., Lu, Z., Li, H., Li, V.O.K.: Incorporating copying mechanism in sequence-to-sequence learning. In: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers (2016)
7. Hao, X., Cao, Y., Shang, Y., Liu, Y., Tan, J., Li, G.: Adversarial reinforcement learning for chinese text summarization. In: International Conference on Computational Science (2018)
8. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural Computation* 9(8), 1735–1780 (1997)
9. Hu, B., Chen, Q., Zhu, F.: Lcsts: A large scale chinese short text summarization dataset. *Computer Science* pp. 2667–2671 (2015)
10. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. CoRR abs/1412.6980 (2014)
11. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. CoRR abs/1312.6114 (2013)
12. Kouris, P., Alexandridis, G., Stafylopatis, A.: Abstractive text summarization based on deep learning and semantic content generalization. In: Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers. pp. 5082–5092 (2019)
13. Li, J., Zhang, C., Chen, X., Cao, Y., Liao, P., Zhang, P.: Abstractive text summarization with multi-head attention. pp. 1–8 (07 2019)
14. Li, P., Lam, W., Bing, L., Wang, Z.: Deep recurrent generative decoder for abstractive text summarization. In: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017. pp. 2091–2100 (2017)
15. Lin, C.Y.: Rouge: A package for automatic evaluation of summaries. In: Text Summarization Branches Out (2004)
16. Lin, J., Sun, X., Ma, S., Su, Q.: Global encoding for abstractive summarization. In: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 2: Short Papers. pp. 163–169 (2018)
17. Luong, T., Pham, H., Manning, C.D.: Effective approaches to attention-based neural machine translation. In: Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015. pp. 1412–1421 (2015)
18. Ma, S., Sun, X., Xu, J., Wang, H., Li, W., Su, Q.: Improving semantic relevance for sequence-to-sequence learning of chinese social media text summarization. In: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 2: Short Papers. pp. 635–640 (2017)
19. Mehta, P., Majumder, P.: From Extractive to Abstractive Summarization: A Journey. Springer (2019)
20. Miao, Y., Blunsom, P.: Language as a latent variable: Discrete generative models for sentence compression. In: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016. pp. 319–328 (2016)
21. Nallapati, R., Zhou, B., dos Santos, C.N., Gülçehre, Ç., Xiang, B.: Abstractive text summarization using sequence-to-sequence rnns and beyond. In: Proceedings of the 20th SIGNLL

- Conference on Computational Natural Language Learning, CoNLL 2016, Berlin, Germany, August 11-12, 2016. pp. 280–290 (2016)
22. Napoles, C., Gormley, M.R., Durme, B.V.: Annotated gigaword. In: Proceedings of the Joint Workshop on Automatic Knowledge Base Construction and Web-scale Knowledge Extraction, AKBC-WEKEX@NAACL-HLT 2012, Montréal, Canada, June 7-8, 2012. pp. 95–100 (2012)
 23. Neto, J.L., Freitas, A.A., Kaestner, C.A.A.: Automatic text summarization using a machine learning approach. In: Advances in Artificial Intelligence, 16th Brazilian Symposium on Artificial Intelligence, SBIA 2002, Porto de Galinhas/Recife, Brazil, November 11-14, 2002, Proceedings. pp. 205–215 (2002)
 24. Rezende, D.J., Mohamed, S., Wierstra, D.: Stochastic backpropagation and approximate inference in deep generative models. In: Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014. pp. 1278–1286 (2014)
 25. Rush, A.M., Chopra, S., Weston, J.: A neural attention model for abstractive sentence summarization. In: Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015. pp. 379–389 (2015)
 26. Schumann, R.: Unsupervised abstractive sentence summarization using length controlled variational autoencoder. CoRR abs/1809.05233 (2018)
 27. Song, S., Huang, H., Ruan, T.: Abstractive text summarization using LSTM-CNN based deep learning. *Multimedia Tools Appl.* 78(1), 857–875 (2019)
 28. Su, J., Wu, S., Xiong, D., Lu, Y., Han, X., Zhang, B.: Variational recurrent neural machine translation. In: Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, New Orleans, Louisiana, USA, February 2-7, 2018. pp. 5488–5495 (2018)
 29. Sutskever, I., Vinyals, O., Le, Q.V.: Sequence to sequence learning with neural networks. In: Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada. pp. 3104–3112 (2014)
 30. Verma, R., Lee, D.: Extractive summarization: Limits, compression, generalized model and heuristics. *Computacion Y Sistemas* 21(4) (2017)
 31. Zhang, B., Xiong, D., Su, J., Duan, H., Zhang, M.: Variational neural machine translation. In: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016. pp. 521–530 (2016)
 32. Zhao, H.: A retrieval algorithm of encrypted speech based on syllable-level perceptual hashing. *Computer Science and Information Systems* 14, 24–24 (2017)
 33. Zhao, H., Wang, G., Xu, C., Yu, F.: Voice activity detection method based on multivalued coarse-graining lempel-ziv complexity. *Computer Science and Information Systems.* 8(3), 869–888 (2011)

Huan Zhao is a Professor at the College of Computer Science and Electronic Engineering, Hunan University. She obtained her B. Sc degree and M.S. degree in Computer Application Technology from Hunan University in 1989 and 2004, respectively, and completed her Ph.D. in Computer Science and Technology at the same school in 2010.

Her current research interests include speech information processing, embedded speech recognition and machine learning.

Jie Cao received his bachelor’s degree in information and computing Science from Shenyang Ligong University, Shenyang, China, and his master’s degree in computer technology from Hunan University, Changsha, China, in 2017.

His current research interests include natural language processing, abstractive summarization, dialog systems, machine learning, and deep learning.

Mingquan Xu received his bachelor's degree in network engineering from Wuhan Textile University, Shenyang, China, and his master's degree in computer technology from Hunan University, Changsha, China, in 2017.

His current research interests include natural language processing, named entity recognition, machine learning, and deep learning.

Jian Lu received his bachelor's degree in information and computing Science from Dalian University, Shenyang, China, and his master's degree in computer technology from Hunan University, Changsha, China, in 2017.

His current research interests include natural language processing, named entity recognition, dialog systems, machine learning, and deep learning.

Received: January 31, 2020; Accepted: June 1, 2020.

Adaptive E-Business Continuity Management: Evidence from the Financial Sector

Milica Labus¹, Marijana Despotović-Zrakić¹, Zorica Bogdanović¹, Dušan Barać¹, and Snežana Popović²

¹ Faculty of Organizational Sciences, University of Belgrade, Jove Ilića 154,
11000 Belgrade, Serbia
{milica, maja, zorica, dusan}@elab.rs

² School for Computing, Union University, Knez Mihajlova 6/VI,
11000 Belgrade, Serbia
spopovic@raf.edu.rs

Abstract. This paper focuses on business continuity management in organizations that use modern e-business technologies: the Internet, mobile computing, e-services, and virtual infrastructure. The aim is to make the shift from traditional Business Continuity Management (BCM) towards “e-Business Continuity Management” (e-BCM) suitable for modern technological environments. We have defined a comprehensive framework for the implementation of an adaptive e-BCM adjustable to changes in the business environment. The framework consists of practical steps for defining elements of a business continuity management system: business impact analysis, risk assessment, and a business continuity plan. We have implemented and evaluated the framework within three financial organizations. The key finding is that Business Impact Analysis and the continual improvement of the Business Continuity Management System are the driving factors for the effective establishment of an adaptive e-BCM. The proposed framework is general, and can be applied to any organization that uses modern e-business technologies.

Keywords: E-Business Continuity Management, Business Impact Analysis, Business Continuity Plan, ISO 22301.

1. Introduction

E-business is based on modern information and communication technologies (ICT): the Internet, mobile technology, cloud computing, and next-generation networks, as well as new concepts, such as the Internet of things, Big Data, and ubiquitous computing[1]. In addition to their numerous benefits[2], modern e-business technologies pose new risks to the business, and particularly to business continuity[3]–[5] and information security, given the amount, speed of creation, and availability of information. The more dependent an organization is on modern e-business technologies, the more it becomes vulnerable to the impacts of all types of incidents[6]. The use of modern e-business technologies creates business and information systems that are complex and are considered inherently risky[7]. One of the objectives of business continuity in such

circumstances becomes the ability to adapt to such emerging risks, so as to protect these naturally risky systems[7].

Literature data in the field of BCM offers different frameworks and implementation approaches that address the various aspects of BCM. However, a general framework that could define practical steps for the establishment of business continuity among organizations that use modern e-business technologies is lacking. Furthermore, companies' implementation details and results are often kept confidential, hence organizations are struggling to find experiences and practical guidelines that would help them develop their own adaptive e-BCM strategies.

The primary goal of our research is to explore the possibilities of improving business continuity management through the development and evaluation of a framework that will facilitate the establishment of successful e-BCM. Drawing on the definition of business continuity[8], we consider e-business continuity as the ability of an organization to continue the delivery of electronic services at acceptable, predefined levels of reliability and availability following a disruptive incident. Since e-business technologies are rapidly evolving, the proposed framework should be flexible and adjustable to changes in organizations' business environments. Our framework, defined in [9], has been implemented and used in three organizations active in the financial sector in Serbia (see Table 2). This research paper focuses on the evaluation of the proposed framework. The evaluation was performed using basic statistical analysis and Partial Least Squares Structural Equation Modeling (PLS-SEM).

2. Literature Overview

2.1. Business Continuity Management

Business continuity management (BCM) focuses on possible internal and external risks and their impacts on business processes. As such, Standard ISO 22301:2014[8] defines BCM as a holistic management process that identifies potential threats to a company and its impacts on business operations. It provides a framework for building organizational resilience with the ability to effectively respond by safeguarding the interests of its key stakeholders, reputation, brand, and value-creating activities.

The importance of implementing BCM strategies within the life cycles of an enterprise causes evaluation and changing of BCM permanently, both within academia and practitioners.

The relevant literature from publicly available databases offers different frameworks and implementation approaches that address the various aspects of BCM. After analyzing the relevant sources, we have identified the following types of available BCM research papers:

- 1) General guidelines, applicable to arbitrary industries (e.g., [10], [11], [20], [12]–[19]);
- 2) Frameworks enriched by statistical and mathematical quantitative calculations (e.g., [21]–[24])

- 3) Literature that addresses specific industries, business issues, and requirements (e.g., [25]–[28]);
- 4) The impacts of nature and human destructions on cities (e.g., [29]–[31]).

The first group of BCM research papers recommends general frameworks and implementation approaches that are applicable in arbitrary industries. In the study [10], with reference to [8], the author presents a comprehensive review of a holistic framework for BCM, with a guide for its implementation. Moreover, the author concludes that the proper implementation of a BCM must be aligned with all its key products and processes for an enterprise. The paper [11] presents a systematic framework for the BCM implementation based on the concept of an “always on” business. The model explains that today’s organizations are using continuous computing technologies to offer greater availability and reliability to their customers.

Complexity and changes of business processes, roles, relationships and benefits of new technologies force business systems to reinvent existing BCM strategies and adapt them to new conditions and environments [12]. The author suggests the incremental concept, which includes three types of strategies: process-centric, program-centric, and management resilience management. For each of these strategies, the author proposes a list of specific techniques.

IT infrastructure, as an important and critical factor in overall enterprise business, is the subject of interest of BCM as well. One of the first comprehensive frameworks for IT infrastructure, proposed in 2006 [13], emphasizes the importance of aligning IT infrastructure, data, and services within overall business strategy and all components of the BCM strategy. The severity of IT incidents in modern BCM strategies is pointed out in [14]. The proposed framework for information system continuity management (ISCM) includes external requirements, management support, organizational alertness, and embeddedness to perceive and minimize the business impacts of ISCM.

In the second group of BCM papers, the researchers propose frameworks enriched by statistical and mathematical quantitative calculations. In [21], the authors used a framework and interactive model for making a decision within the resource allocation problem. They applied the concept in a gearbox manufacturer enterprise successfully. The same authors present a framework for estimating the business impact analysis (BIA) [22], which uses multi-parameter decision-making techniques. The evaluation of the framework was performed within an auto part manufacturing enterprise. An integrated BCM framework in [23] uses quantitative metrics through the protection, mitigation, emergency, and recovery phases, and applies them to a critical lightning disruptive event at an oil storage tank farm. In [24], the authors propose the use of a mathematical programming model as a decision support tool for the successful development of disaster recovery plans.

The third identified group of BCM papers presents implementation approaches that address specific industry, business issues, or business requirements. In [25], the authors implemented a wide quantitative study in 75 automobile parts makers in disaster-prone regions, and introduced a new term - Supply Chain Cooperation (SCC) - to the Business Continuity Plan (BCP). They confirm that BIA is the cornerstone of BCP, as it has a powerful effect on other BCP factors and components. The paper [26] describes the concepts of disaster recovery and data replication plan in The Medical Record Company. The authors in [27] introduce a new concept called “Area BCP”, which addresses disaster risk management in industrial agglomerated areas as a whole. The authors propose a framework for coordinated damage mitigation measures and recovery

actions, which involves all the business stakeholders in a certain area. The paper [28] focuses on supply chain business continuity issues which, as authors argue, can be applied to any industry.

The fourth type of BCM papers focuses on the impacts of nature and human destructions on cities. Paper [29] focuses on the risk assessment process (RA) that is, in addition to BIA, one of the fundamental components of BCM. The authors propose an enhanced risk assessment framework by adjusting it to BCM needs. They evaluated the framework in a real service organization in charge of disaster management services in the city of Tehran, Iran. Another study [30] focuses on the Smart City infrastructure for the Olympic Games in Japan in 2020. Authors elaborate aspects of security, convenience, maintenance, etc. of the city for the upcoming Games, and how to maintain the continuity of city life after their end.

The majority of papers present practical implementations of frameworks within companies, classified as large ones. It seems that small and medium-sized enterprises (SMEs) have fallen out of interest although they, for example, represent a share of over 95% in the European Union [32]. In [31], the authors discuss the availability of BCM support for SMEs, including the environmental impacts on BCM in risks areas.

After analyzing the four identified types of relevant BCM papers, we have concluded that the practical steps in establishing business continuity in organizations that use modern e-business technologies are lacking. The following common issues have been identified in existing business continuity frameworks and implementation approaches:

1) their focus is solely on definition and planning of business continuity programs [13] or solely on one aspect of business continuity management [22], [24], [26], [29];

2) there is a lack of specific steps in defining business continuity management [18]–[20], [24], [27], and, more specifically, there is a lack of practical steps for business impact analysis and risk assessment [16], [17], which are essential components of BCMS;

3) the approach is too specific for a particular industry or line of business [27], [28]; and

4) the approach is too complicated to implement in an arbitrary organization [22], [24].

We have attempted to address all of these issues using the proposed framework [9] for the establishment and continuous improvement of e-BCM, tailored specifically to organizations that use modern e-business technologies.

2.2. Adaptive e-Business Continuity Management

As previously stated, the proposed framework for adaptive e-Business continuity management was initially introduced in [9]. It was defined with reference to the ISO 22301 standard [8], which is the de facto primary international standard for business continuity management, following the systematic collection, review, and analysis of business continuity literature and existing frameworks. At later stages of our research, the framework was implemented in three organizations from the financial sector.

The framework combines parameters with responses of an organization. The parameters reflect the basic characteristics of the business environment in which the host organization operates [8], [33], [34], as well as the main characteristics of the underlying technologies: ubiquity, global reach, universal standards, richness,

interactivity, information density, personalization/customization, and the usage of social networks [35]. The framework parameters are the following: P1) BCMS objectives, P2) BCMS scope, P3) Specific context of the organization, P4) Threats, P5) Vulnerabilities, P6) Disruption timescale, P7) Financial impact, P8) Operational impact, P9) Processes' criticality levels, P10) Resource Financial Categories, P11) Risk appetite, and P12) Risk Priority Levels.

In accordance with ISO 22301, the response of an organization draws on business continuity management system, whose key elements are defined by four framework components: Business Impact Analysis (C1), e-BCM Risk Assessment (C2), Business Continuity Plan (C3); and Continual BCMS Improvement (C4). Figure 1 illustrates the components of the framework and their most important interdependencies. The organization's response is precisely defined with a set of specific objectives and procedures.

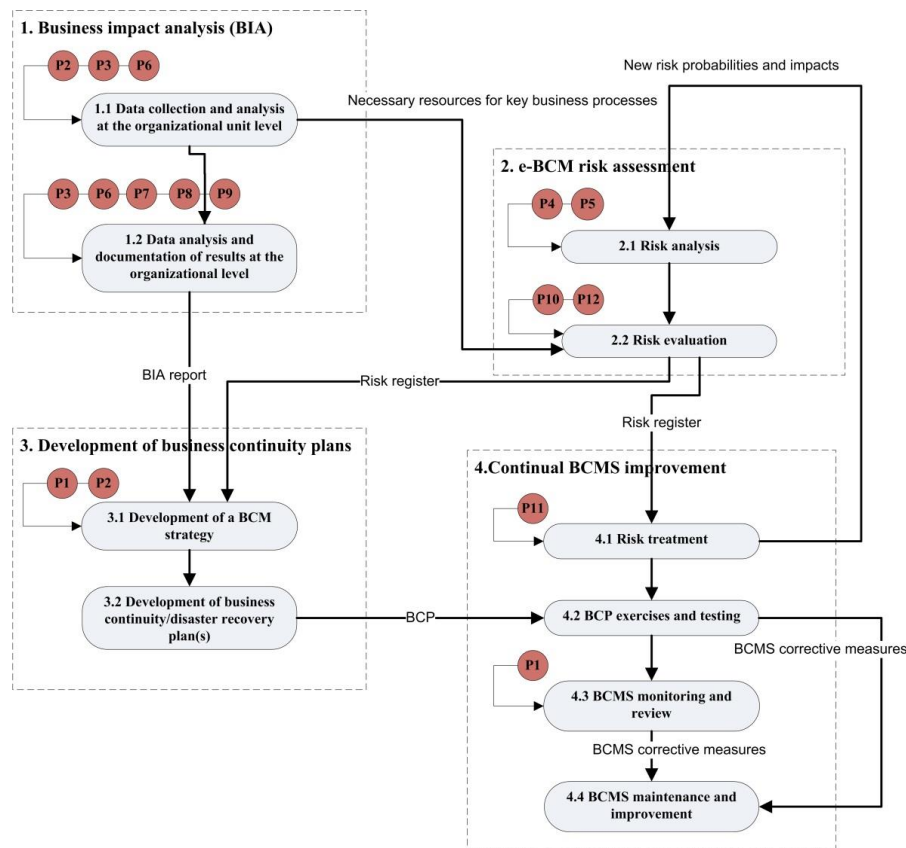


Fig. 1. The framework (parameters in circles; components in dotted-line rectangles; procedures, as part of components, in rounded-corner rectangles) and important interdependencies

Component number one deals with Business Impact Analyses (BIA). Its main objective is to identify the key business processes of the organization, as well as to determine how quickly these processes must recover and begin to provide complete

business functions following an incident. This analysis draws on an assessment of potential financial, operational, regulatory, reputational, and other losses. Using parameter P7, organizations define whether the financial impact on the business will be evaluated, whereas parameter P8 defines all operational impacts that are assessed. Additionally, BIA further identifies the resources necessary for key business processes, and prioritizes recovery activities. As a result, data collection and analysis is performed first at the level of organizational units (procedure 1.1 within this component, see Figure 1). The analysis is performed only on business processes included in the BCMS (P2), and its key stakeholders and legal and regulatory framework of the organization (P3) are taken into consideration. Potential outcomes of disruption are analyzed based on a predefined disruption timescale (e.g., minutes, hours, days), in accordance with the nature of the business (P6).

Results are later analyzed and documented at the organizational level (procedure 1.2, see Figure 1). Criticality levels are identified for all processes based on their maximum acceptable outage (MAO) values (as defined by parameter P9).

The second component deals with e-BCM Risk Assessment, and consists of two procedures: risk analysis (2.1 in Figure 1), and risk evaluation (2.2 in Figure 1). Risk analysis identifies potentially achievable threats (P4) based on the present vulnerabilities of key business processes (P5). Risk evaluation involves assessing risk probabilities and impacts on the resources necessary for the key business processes, and results in calculating the risk values. Important input parameters for risk evaluation are Resource Financial Categories (P10) and Risk Priority Levels (P12). Resource Financial Categories (very high, high, medium, and low) are based on financial values, and are defined specifically for that organization. Risks are further categorized into four priority levels (very high, high, medium, and low) based on risk values, as defined by parameter P12.

The third component deals with the Business Continuity Plan (BCP). Its main goal is to define recovery procedures for key business processes within recovery time objectives, as determined under the BIA. A BCP can comprise one or multiple documents (sub-plans), depending on the geographical, organizational, and other business specifics of an organization. An important aspect of BCP is the recovery of critical information systems and ICT resources, which is often referred to as the Disaster Recovery Plan (DRP). This component consists of two procedures: the Development of a Business Continuity Strategy (3.1 in Figure 1) and the Development of Business Continuity/Disaster Recovery Plans (3.2 in Figure 1). The strategy reflects BCMS objectives (P1), and BCMS scope (P2), and should be based on the results of BIA analyses and risk assessments. It defines prerequisites, and provides all the necessary resources to be made available in case of a crisis situation.

The first three framework components define the most important elements of BCMS. Lastly, the final component deals with continual BCMS improvement. It considers internal and external changes within the organization's environment, and defines corrective measures for the improvement of the BCMS. This component consists of procedures for risk treatment, business continuity plans exercises and testing, monitoring, review, maintenance, and improvement of BCMS. The risk appetite of an organization (P11) is an important parameter for risk treatment (procedure 4.1, Figure 1). Only risks with values above the defined threshold are mitigated. BCP exercises and tests (procedure 4.2, Figure 1) are used to validate BCP content, and to ensure that response and recovery results can be achieved within the defined timeframes [9]. BCMS

monitoring and review (procedure 4.3, Figure 1) includes the evaluation of BCMS performance against defined BCMS objectives (P1). Finally, BCMS maintenance and improvement (procedure 4.4, Figure 1) implements pre-defined corrective measures designed to improve the BCMS. The main objective of this component is to provide an adaptive business continuity management.

For the proposed framework, this research has been focused on the following main research questions: 1) The framework defines clear and effective procedures for the establishment of e-business continuity management; 2) The framework examines the specifics of an organization's business, particularly its use of modern e-business technologies; 3) The framework facilitates the establishment of business continuity in an organization that uses modern e-business technologies; 4) The framework enables adaptive e-BCM in accordance with changes to the organization's environment, and 5) The framework contributes to a positive organizational attitude towards e-BCM. We refer to all these research questions as the core of an Adaptive e-BCM.

In the next chapter, we present the evaluation of the described framework.

3. Research Methods

3.1. Research Context

As stated previously, the proposed framework has been implemented in three organizations active in the financial sector in Serbia (see Table 2). The first implementation was in the Association of Serbian Insurers, a regulatory organization for insurance companies handling the Motor Third-Party Liability (MTPL) insurance in Serbia. The implementation took place in 2013, and the project lasted three months. Since the Association of Serbian Insurers is responsible for the centralized information system for the sale of MTPL policies interviews were conducted with all insurance companies in Serbia during the BIA analysis. The second implementation was a year later, in Milenijum Osiguranje, an insurance company with 300 employees that provides all types of non-life insurance. This project lasted four months. The last implementation was in 2015 and 2016, in the National bank of Serbia (NBS). The project lasted seven months, and included the involvement of over 200 NBS employees at some point. The top management support was critical for the project's success.

All three organizations that implemented the framework agreed to take part in the framework evaluation. However, due to confidentiality reasons, we were explicitly asked not to publish separate evaluation results, but rather to consider all three organizations as a single evaluation sample.

3.2. Evaluation Hypotheses

The main goal of the evaluation is to assess how effective the framework is at establishing adaptive e-BCM. Evaluation hypotheses are based on the theoretical background and assessment of interdependencies between the components.

The first set of evaluation hypotheses tests whether each framework component is instrumental for the establishment of adaptive e-BCM, and how significant its contributions are:

H1a: The Business Impact Analysis component (C1) contributes to the establishment of an adaptive e-BCM.

H1b: The e-BCM risk assessment component (C2) contributes to the establishment of an adaptive e-BCM.

H1c: The component for the development of Business Continuity Plan (C3) contributes to the establishment of an adaptive e-BCM.

H1d: The component encompassing continual BCMS improvement (C4) contributes to the establishment of an adaptive e-BCM.

The second set of evaluation hypotheses pertains to the Business Impact Analysis component (C1). One of the main objectives of BIA analysis is to identify key business processes of an organization and its necessary resources, which is a starting point for risk assessment. The following evaluation hypothesis is, therefore, proposed:

H2a: Implementation of the Business Impact Analysis component (C1) contributes to the implementation of an e-BCM risk assessment component (C2).

Recovery time objectives for key business processes are determined with reference to the results of BIA analysis. These are the Recovery Point Objective (RPO), Recovery Time Objective (RTO), and Maximum Acceptable Outage (MAO). The recovery time objectives define boundaries within which key business process must recover, and constitute the main goals of BCM, as well as a starting point for the development of the Business Continuity Plan document. Hence, the following evaluation hypothesis is derived:

H2b: Implementation of the Business Impact Analysis component (C1) contributes to the implementation of the component for the development of the Business Continuity Plan (C3).

The next set of evaluation hypotheses concerns the e-BCM Risk Assessment component (C2), which is focused on all risks that can lead to the interruption of business processes, and not just security incidents. To account for this fact, the impact parameters (P4 and P5) include predefined lists of potential threats and vulnerabilities for the organization. The potential threats and vulnerabilities are focused on information systems and usage of modern e-business technologies. So defined, risk assessment makes an important contribution to strategic and operational decision-making in the organization [36]. The following evaluation hypothesis is therefore proposed:

H3a: Implementation of the e-BCM risk assessment component (C2) contributes to the implementation of the component for the development of the Business Continuity Plan (C3).

The results of risk assessment inform the treatment of all risks above a certain acceptability level, as identified through the Risk Appetite parameter (P11). The purpose of risk treatment is to take preventive action to address the causes and impacts of potential risk events. Risk treatment is an integral part of the Component for continual BCMS improvement (C4); the following hypothesis is, therefore, proposed:

H3b: Implementation of the e-BCM risk assessment component (C2) contributes to the implementation of the component for continual BCMS improvement (C4).

The component for continual BCMS improvement (C4) is applied at planned intervals and after each important change in the organization's environment. Procedures

that are part of this component include BCP exercises and tests. Short of an actual disruptive incident, BCP exercises and tests are the only means of validating BCP content and ensuring that response and recovery results can be achieved within defined timeframes. Therefore, the final evaluation hypothesis is:

H4: Implementation of the component for continual BCMS improvement (C4) contributes to the component for the development of Business Continuity Plan (C3). All evaluation hypotheses are presented in Figure 2, while Table 1 gives an overview of the literature used.

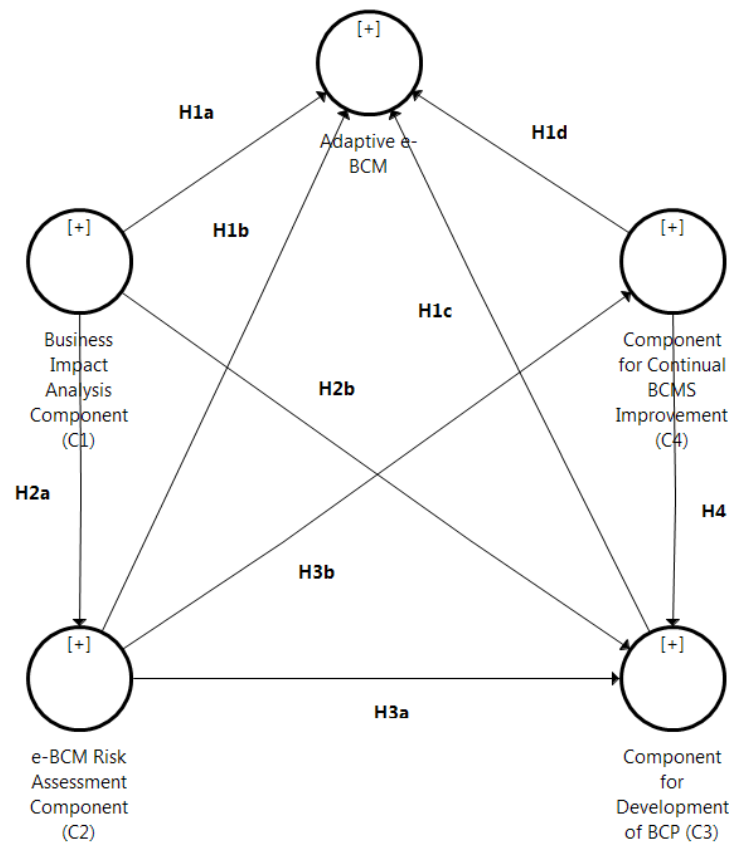


Fig. 2. Evaluation hypotheses (the structural PLS-SEM model)

Table 1. Evaluation hypotheses and theoretical background

Hypothesis	Theoretical background
H1a: C1 -> Adaptive e-BCM	Each framework component is instrumental for the establishment of adaptive e-BCM[9]
H1b: C2 -> Adaptive e-BCM	
H1c: C1 -> Adaptive e-BCM	
H1d: C2 -> Adaptive e-BCM	

H2a: C1 -> C2	Key business processes are identified in BIA analyses[34], [37]; risk assessment is performed on the resources necessary for the key business processes[9], [29], [38]
H2b: C1 -> C3	Recovery time objectives for key business processes are based on the results of BIA analysis[8], [22], [34], [39]
H3a: C2 -> C3	Business Continuity Plan should be based on results of BIA analysis and risk assessment[8], [22], [34]; results of risk assessment are important for strategic and operational decision-making in the organization[36]
H3b: C2 -> C4	Risk treatment is based on the results of risk assessment[8], [29], [38]
H4: C4 -> C3	BCP exercises and tests are essential part of effective business continuity management[8], [34]

3.3. Evaluation Questionnaire and Data Collection

To evaluate the proposed framework and to test evaluation hypotheses, we constructed an evaluation questionnaire (survey) containing a set of evaluation questions (survey items); this questionnaire is provided in Appendix A. All responses were made on a 5-point Likert scale, with the categories being (1) strongly disagree, (2) disagree, (3) neither agree nor disagree, (4) agree, and (5) strongly agree.

Table 2. Organizations that implemented the framework and took part in the evaluation

Organization	No. of employees	Project duration	No. of key business processes	No. of key users who participated in the evaluation
1. Association of Serbian Insurers	20	3 months	8	3
2. Milenijum Osiguranje Insurance Company	300	4 months	31	9
3. National Bank of Serbia	2,000	7 months	115	26
Total no. of evaluations:				38

Survey items used to evaluate the Adaptive e-BCM were based on the key research questions as defined in[9], and discussed in the previous chapter. They were defined in accordance with the international standard[8] and BCMS performance evaluation[40]. In addition, some survey items were based on the following literature:[34], [39], [41] for evaluating the Business Impact Analysis component (C1);[13], [36], [38] for evaluating the e-BCM risk assessment component (C2);[41]–[43] for evaluating the component for

development of Business Continuity Plan (C3); and [44] for evaluating the component for continual BCMS improvement (C4).

Key users from three organizations that implemented the framework took part in the evaluation (see Table 2 for details). All three organizations are from the financial sector in Serbia, and they are all subject to the same standards and rules regarding business continuity. As already stated, due to confidentiality reasons, all three organizations were considered as a single sample. The evaluation was performed between April and June 2017.

3.4. Evaluation Methodology

The evaluation was performed in three steps, as presented in Table 3. Basic statistical analysis, together with statistical significance test (Step 1), was used to determine whether the main research questions hold, and to assess framework effectiveness. A more thorough analysis of the findings of the evaluation (Steps 2 and 3) was conducted using Structural Equation Modeling (SEM), an advanced statistical analysis technique that can identify relationships among measured variables and latent variables (variables that are not directly measured), as well as between latent variables [45]. In our research, we used Partial Least Squares SEM (PLS-SEM), which is a good alternative to Covariance-Based SEM (CB-SEM) for estimating theoretically justified cause-effect relationships in models, especially at small sample sizes [46], [47]. PLS-SEM works efficiently with small sample sizes and generally makes no assumptions about data distributions. We used SmartPLS software (v.3.2.7) for analysis [48].

Table 3. Evaluation Steps

Evaluation step	Evaluation method	Purpose of the evaluation step
1. Evaluation of framework effectiveness	Basic statistical analyses and statistical significance test	Determine whether main research questions hold and whether the framework is effective in implementation
2. Evaluation of the measurement model (outer model)	PLS-SEM standard algorithm	Assess the reliability and validity of the evaluation questionnaire
3. Evaluation of the structural model (inner model)	PLS-SEM blindfolding and bootstrapping	Test evaluation hypotheses

The PLS-SEM analysis consists of the evaluation of corresponding PLS-SEM structural and measurement models.

The structural model (also called the inner model in PLS-SEM) illustrates the research hypotheses and displays the variable relationships that will be examined [45]. We have created the structural model based on our evaluation hypotheses (see Figure 2). In PLS-SEM, variables that are not directly measured are called constructs, and they are represented in models as circles or ovals. The indicators, also called items or manifest variables, are the directly measured variables that contain the raw data, and they are represented in models as rectangles [45].

The measurement model (also referred to as the outer model in PLS-SEM) of the constructs displays the relationships between the constructs and the indicator variables[45]. The measurement model was defined with the survey items explained earlier in this chapter, and is described in more detail in Chapter 4.2 (see Figure 3).

The PLS-SEM models cannot be evaluated with a single goodness-of-fit criterion. Following PLS-SEM guidelines[45], [49]–[53], the study performed a two-stage approach to evaluation, as described in evaluation Steps 2 and 3 (see Table 3).

The following Chapter presents each of the tree evaluation steps in more detail.

4. Research Results

4.1. Evaluation of Framework Effectiveness

In the first evaluation step, a simple statistical analysis was performed based on mean values, standard errors, and standard deviation of all evaluation items. We set the mean value of 4 as the benchmark against which evaluation results may be compared. The basic assumption here was that the mean value of 4 is considered effective in implementing an adaptive e-BCM.

The first part of the questionnaire (items Q.1.1 – Q.1.5) refers to the core of an Adaptive e-BCM. Basic statistical analysis shows that evaluation results support the effectiveness of the core since the average score for each evaluation item is greater than 4 (see “Mean values for Part 1: Adaptive e-BCP” in Appendix A).

The other four parts of the questionnaire address the remaining issues of effectiveness, with questions grouped by framework components. The component for continual BCMS improvement (C4) received the lowest average score, 3.36, as it includes the most advanced activities from the business continuity maturity perspective, and its successful implementation requires some time. Other framework components received average scores greater than or very close to 4: 4.58 for the Business Impact Analysis component (C1), 4.05 for the e-BCM Risk Assessment component (C2), and 3.98 for the component for development of BCP (C3), which suggests that three out of four framework components are effective in implementing an adaptive e-BCM.

To further validate those claims, we have performed the statistical significance test (t-test) on the mean values for each evaluation item (please see Appendix A for details). We have used one-sample t-test with a 95% confidence level in software Stata (v.13). The null hypothesis was: Mean value is equal to 4 ($H_0: \text{Mean}=4$). For each evaluation item we computed test of statistical significance for other alternative hypotheses: Mean value is lower than 4 ($H_a: \text{Mean}<4$), Mean value is not equal to 4 ($H_a: \text{Mean}\neq 4$) and Mean values is greater than 4 ($H_a: \text{Mean}>4$), where the last one is the most relevant for our research. As presented in Appendix A, four out of five items in Part 1: Adaptive e-BCP passed the t-test in the sense that the probability to reject the null hypothesis against the alternative ($H_a: \text{Mean}>4$) is below the significance level of 5%. Only item Q.1.5: “The framework contributes to a positive organizational attitude towards e-BCM” didn’t pass the statistical significance test.

For Business Impact Analysis Component (C1), all evaluation items passed the t-test, which shows that this component is, out of all four components, the most effective in implementation. The second two components partly passed the t-test: two out of six evaluation items for e-BCM Risk Assessment Component (C2), and three out of six items for Component for Development of Business Continuity Plan (C3). These results implicate that those components need further adjustments, especially in the areas of continual improvements, to achieve full effectiveness in practice. For similar reasons, the Component for Continual BCMS Improvement (C4) didn't pass the statistical significance test. Business continuity is an ongoing process[8], and our opinion is that this component needs years of continual practice to achieve the desired level of effectiveness.

4.2. Evaluation of the Measurement Model

After the basic statistical analyses of the findings, we have conducted a more thorough analysis with PLS-SEM. The first part was an evaluation of the PLS-SEM measurement model. Our measurement model (Figure 3) is reflective: each of its measures (items, indicators) represents the effects (or manifestations) of an underlying construct[45]. All constructs have multi-item measures. We evaluated the reliability and validity of the construct's measures by examining PLS-SEM estimates. The results are reported in Figure 3 and Table 4.

Firstly, we assessed internal consistency reliability, i.e., whether the proposed indicators are valid measures of the construct. As argued by [45], the appropriate measure of internal consistency reliability is composite reliability which, compared to alternative measure - Cronbach's alpha, takes into consideration the varying outer loading of the indicator itself and not only the intercorrelations of the observed indicator variables. Table 4 reveals that composite reliability values for all constructs stand at between 0.80 and 0.95, demonstrating high internal consistency reliability[45], [50], [53].

Table 4. Summary of Results for Reflective Outer Model

Constructs	Composite Reliability	AVE
Adaptive e-BCM	0.854	0.545
Business Impact Analysis Component (C1)	0.808	0.518
e-BCM Risk Assessment Component (C2)	0.915	0.687
Component for Development of BCP (C3)	0.880	0.648
Component for Continual BCMS Improvement (C4)	0.933	0.737

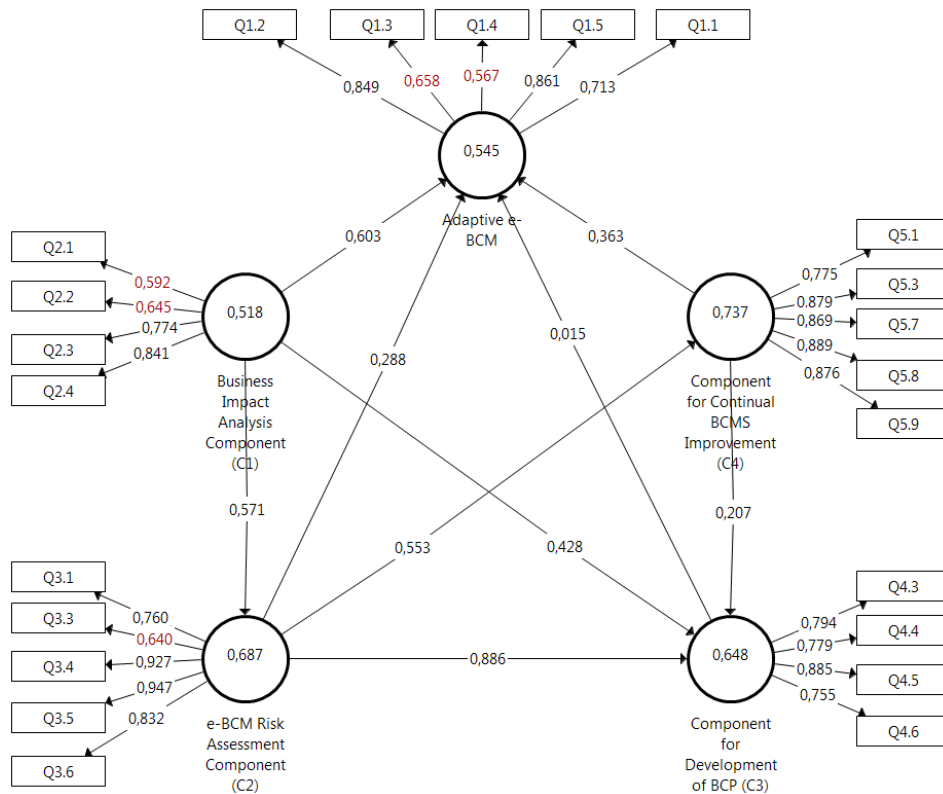


Fig. 3. PLS algorithm results (Inner model: Total Effects, Outer model: Outer Loadings, Constructs: Average Variance Extracted - AVE)

Next, we examined whether convergent validity was established, by measuring the outer loadings of the indicators (also called indicator reliability) and the Average Variance Extracted (AVE) of the constructs[45].

Indicators associated with a construct with high outer loadings have much in common, which is captured by the construct. Out of an initial set of indicators (listed in Appendix A), twelve indicators with outer loadings of between 0.40 and 0.70 were considered for removal from the model, as suggested in[45], by examining whether deleting the indicator would lead to an increase in the composite reliability or AVE above the preferred thresholds. Based on those criteria, six items (Q2.5, Q2.6, Q3.2, Q4.1, Q4.2, and Q5.6) were not included in further analysis. For example, a clear majority of the respondents “strongly agreed” with items “Q2.5 Recovery time objectives are determined.” and “Q4.1 BCMS is adequately defined”.

In addition, three indicators with excessive outer loading values (of around 0.95) were also removed from the model (Q5.2, Q5.4, and Q5.5), as suggested by [45], so that consistency reliability of the corresponding construct could remain below an upper threshold of 0.95.

All the items removed were carefully examined so as not to have major effects on the corresponding construct’s content validity as suggested in[45]. Figure 3 shows that all

remaining indicators have outer loadings above the preferred level of 0.708 or very close to it (sufficiently close values are marked in red).

The results indicated (see Table. 3) that all construct AVE values were greater than the acceptable threshold of 0.5 suggested by [45], which means that the construct, on average, explains more than half of the variance of its indicators.

Lastly, we examined discriminant validity, that is, the extent to which a construct is truly distinct from other constructs by empirical standards. As argued by [45], there are two measures of discriminant validity: cross loadings of indicators, considered a liberal approach[54], and the Fornell-Larcker criterion, a more conservative approach [55]. Discriminant validity was established for all constructs according to the first method, which means that their outer loading on associated constructs are greater than all of their loadings on other constructs for all indicators (that is why this method is called “cross loadings”). Application of the second method, as shown in Table 5, revealed only one minor discrepancy between the e-BCM Risk Assessment Component (C2) and the Component for Development of BCP (C3). Since the correlation between those two variables is only slightly greater than the square root of AVE (approximately 0.037), we considered this acceptable.

Table 5. Discriminant Validity: Fornell-Larcker Criterion

Construct	C1	Adaptive e-BCM	C4	C3	C2
Business Impact Analysis Component (C1)	0.720				
Adaptive e-BCM	0.619	0.738			
Component for Continual BCMS Improvement (C4)	0.359	0.569	0.859		
Component for Development of BCP (C3)	0.437	0.489	0.606	0.805	
e-BCM Risk Assessment Component (C2)	0.571	0.538	0.553	0.841	0.829

4.3. Evaluation of the Structural Model and Hypotheses Testing

As the measurement characteristics of the constructs were proven to be acceptable, we continued with the assessment of the PLS-SEM structural model. The following assessments were conducted:

- Collinearity among sets of constructs;
- Model's predictive accuracy;
- Model's predictive validity; and
- Structural model relationships.

The variance inflation factor (VIF) is a measure of collinearity. In our structural model, all VIF values are below 5, which indicates that there are no collinearity issues among sets of constructs[45].

The coefficient of determination (R^2), as a measure of the model's predictive accuracy, was measured only for endogenous latent variables: Adaptive e-BCP ($R^2 =$

0.525), C2 ($R^2 = 0.326$), C3 ($R^2 = 0.741$) and C4 ($R^2 = 0.305$). The R^2 value for C3 is considered substantial, whereas R^2 values for other endogenous latent variables (Adaptive e-BCP, C2, and C4) are considered moderate [56].

The cross-validated redundancy measure of Q^2 was used to assess the model's predictive validity [57], [58]. Running the blindfolding procedure with an omission distance of seven yielded a Q^2 value of 0.427 for C3, which indicates a large predictive relevance [45]. Other endogenous latent variables have medium predictive relevance: Adaptive e-BCP ($Q^2 = 0.227$), C2 ($Q^2 = 0.187$), and C4 ($Q^2 = 0.175$).

Structural model relationships were assessed using structural model path coefficients, which were obtained after running the PLS-SEM algorithm (see Figure 3). Whether a coefficient is significant ultimately depends on its standard error, which is obtained by means of bootstrapping. The bootstrap standard error allows computing the empirical t value. When the empirical t value is larger than the critical value, we say that the coefficient is significant at a certain error probability (i.e., significance level). The critical value for two tailed tests was 1.96 for a significance level of 5% [45].

The study assessed the structural model through complete bootstrapping with 5,000 samples (significance level 5%, two-tailed test, individual sign changes option). Results are shown in Figure 4 and Table 6.

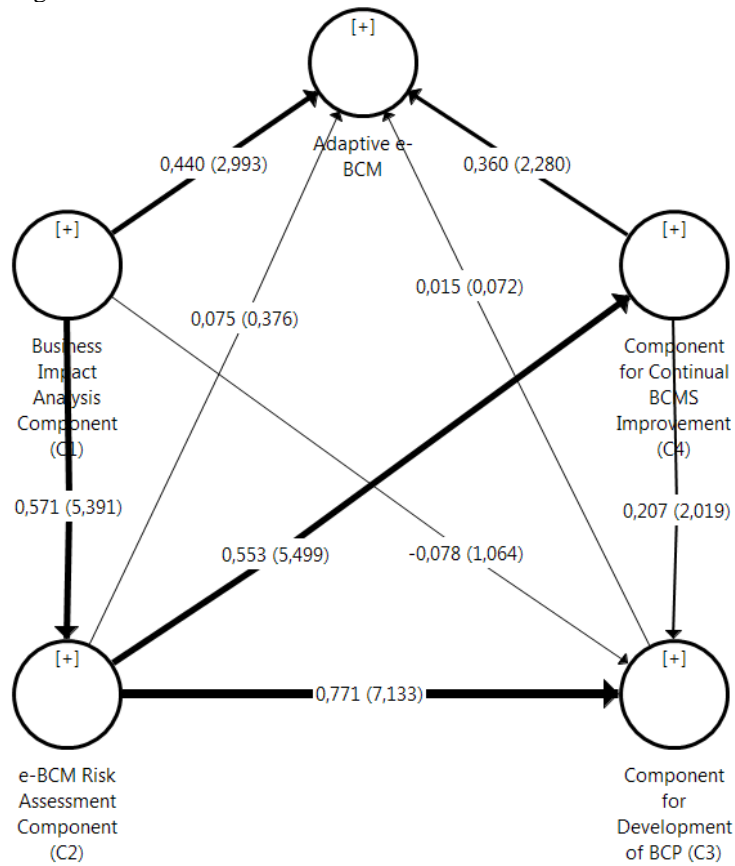


Fig. 4. Bootstrapping Results (Inner model: Path Coefficients and T-Values)

Table 6. Bootstrapping Results

Hypothesis	Original Sample	Sample Mean	Standard Deviation	T Statistics
H1a: C1 -> Adaptive e-BCM	0.440	0.468	0.147	2.993***
H1b: C2 -> Adaptive e-BCM	0.075	0.255	0.200	0.376
H1c: C3 -> Adaptive e-BCM	0.015	0.287	0.214	0.072
H1d: C4 -> Adaptive e-BCM	0.360	0.307	0.158	2.280***
H2a: C1 -> C2	0.571	0.598	0.106	5.391***
H2b: C1 -> C3	-0.078	-0.108	0.073	1.064
H3a: C2 -> C3	0.771	0.770	0.108	7.133***
H3b: C2 -> C4	0.553	0.565	0.100	5.499***
H4: C4 -> C3	0.207	0.207	0.103	2.019***

*** T>1.96 (significance level 5%)

The results revealed a positive and significant effect of the Business Impact Analysis Component (C1) and the Component for Continual BCMS Improvement (C4) on the Adaptive e-BCP, but no significant direct effect of the e-BCM Risk Assessment Component (C2) or the Component for Development of BCP (C3). C1 has the strongest effect on the Adaptive e-BCP ($\beta = 0.440$, $T > 1.96$), followed by C4 ($\beta = 0.360$, $T > 1.96$). Therefore, hypotheses H1a, and H1d hold. The relationship between the e-BCM Risk Assessment Component (C2) and the Adaptive e-BCM is not significant ($\beta = 0.075$, $T < 1.96$), so H1b does not hold. Similarly, the relationship between the Component for Development of BCP (C3) and the Adaptive e-BCM is not significant ($\beta = 0.115$, $T < 1.96$), so H1c does not hold either.

The Business Impact Analysis Component (C1) has a very significant effect ($\beta = 0.571$, $T > 1.96$) on the e-BCM Risk Assessment Component (C2), so hypothesis H2a holds. On the other hand, C1 does not have a significant effect on C3 ($\beta = -0.078$, $T < 1.96$) and, therefore, hypothesis H2b does not hold.

The e-BCM Risk Assessment Component (C2) has a very significant effect on both C3 ($\beta = 0.771$, $T > 1.96$) and C4 ($\beta = 0.553$, $T > 1.96$), so hypotheses H3a and H3b both hold.

Finally, the relationship between the Component for Continual BCMS Improvement (C4) and the Component for Development of BCP (C3) is significant ($\beta = 0.207$, $T > 1.96$), so hypothesis H4 holds.

5. Discussion

Basic statistical analysis reveals that the main research questions hold, and that Business Impact Analysis Component (C1) is the most effective in practice, whereas the other three framework components need some additional adjustments to achieve the desired level of effectiveness.

Evaluation hypotheses tested in the PLS-SEM analysis can be divided into two distinct groups: one that examines direct contributions of framework to the

establishment of adaptive e-BCM (H1a, H1b, H1c, H1d) and the other that examines the framework's internal structure, that is, the interdependencies between the components (H2a, H2b, H3a, H3b, H4).

The adaptive e-BCP, which represents the main research questions, reveals the framework's effectiveness at establishing adaptive e-BCM. The results show that the Business Impact Analysis Component (C1) and the Component for Continual BCMS Improvement (C4) have a positive and significant effect on the implementation of Adaptive e-BCP. For the other two components, the e-BCM Risk Assessment Component (C2), and the Component for Development of BCP (C3), results of PLS-SEM analysis didn't show direct significant effects. This means that, in the proposed framework, BIA analysis and continual BCMS improvements are the factors primarily responsible for the effective establishment of adaptive e-BCM. On the other hand, evaluation results show that component for e-BCM Risk Assessment (C2) has an indirect effect on the Adaptive e-BCM, through the Component for Continual BCMS Improvement (C4), since results of risk assessment are the basis for risk treatment.

The framework shows high internal consistency and logical structure, as all evaluation hypotheses in the second group hold, except for H2b: Implementation of the Business Impact Analysis component (C1) contributes to the implementation of the component for the development of Business Continuity Plan (C3).

The Business Continuity Plan (BCP) is a documented set of procedures that guides an organization to respond, recover, resume, and restore a pre-defined level of operation following a disruptive incident [8]. A BCP ensures that business processes recover rapidly and effectively, that damage to the business is minimized, and that business disruption is managed [17]. It defines the resumption of an organization's key business processes within recovery time objectives, as defined in the BIA [39]. Therefore, the significant direct effect of the Business Impact Analysis Component (C1) on the Component for Development of BCP (C3) is only to be expected. This hypothesis did not hold in our evaluation, which suggests that further improvements should be made to the framework to enhance the connection between these two components. By contrast, the Business Impact Analysis Component (C1) has an indirect impact on the Component for Development of Business Continuity Plan (C3) through the e-BCM Risk Assessment Component (C2), so minor framework adjustments are needed.

Having closely examined all aspects, we concluded that an additional framework parameter should be added, that of Recovery time objectives (P13). Recovery time objectives of key business processes, as defined in BIA, are crucial business continuity objectives that must be addressed by the Business Continuity Plan [34]. This adjustment will define the important relationship between the Business Impact Analysis Component (C1) and the Component for Development of Business Continuity Plan (C3) more precisely.

Regardless of the suggested adjustments, the overall evaluation results are positive, and we can conclude that the framework is effective in establishing adaptive e-business continuity management.

The most relevant practical application of our research is that the proposed framework defines a general context for the establishment and continuous improvement of e-BCM, and is particularly tailored for organizations that use modern e-business technologies. It includes practical steps, grouped into procedures and components, for defining the most important elements of a business continuity management system: business impact analysis, risk assessment, and business continuity plan. In addition, it

guides an organization through continual BCMS improvement which, as we have shown, is one of the main drivers for the effective establishment of adaptive e-BCM.

The small evaluation sample presents the main limitation of our research. However, we consider it as a blueprint for how evaluation should be conducted, once a bigger sample is gathered. We have addressed this limitation by selected PLS-SEM method for evaluation, which works efficiently with small sample sizes. The main threat to the validity of our research is that the proposed framework has only been implemented in organizations from the financial sector, whereas we argue that the framework is general and can be applied in any organization that uses modern e-business technologies. This threat was partly addressed with the implementation of the framework in the National bank of Serbia, within which operates The Institute for Manufacturing Banknotes and Coins as a separate business entity in the field of the graphics industry. Moreover, all three organizations which implemented the framework have the following characteristics common to all organizations that use modern e-business technologies: 1) use of e-services that are critical to business (e.g., online insurance sales, B2B e-services in NBS), 2) cloud IT infrastructure, 3) shared business processes between multiple organizational units, and 4) high standards in information security and data protection.

The framework should be applied every time there is a change in any of the framework parameters, which enables the e-BCM to flexibly respond to changes in the organization's business environment. Corrective BCMS measures are set out in the continual BCMS improvement [8], [34]. Adaptability to changes in the business environment of the organization, together with corrective BCMS measures, allows adaptive e-BCM.

6. Conclusions

Due to the rapid development of e-business technologies in an uncertain world, organizations and their information systems have tremendous risks in the face of potential disaster shocks. In response to that, they need to design adaptive e-Business Continuity Management. Our research is focused on improving business continuity management in organizations that use modern e-business technologies, for which we have introduced the concept of "e-Business Continuity Management" (e-BCM). Our approach was to create a framework for adaptive e-Business continuity management that is sufficiently general and flexible in order to enable any organization to specify their own BCMS: BIA analysis, risk assessment, and business continuity plans. The framework was implemented in three financial organizations of various sizes.

The main contribution of this research paper is the evaluation of the proposed framework based on its application in real business practices. We have evaluated whether the proposed framework is effective in the implementation of adaptive e-Business continuity management in three financial organizations in Serbia. Using the PLS-SEM analysis, we have shown that two out of four framework components, BIA analysis, and continual BCMS improvement, have a statistically positive contribution to the effective establishment of adaptive e-BCM. For the contribution of the remaining two components of the framework, risk assessment and development of Business Continuity Plan, further development and analyses are needed.

Although the evaluation sample is small, we emphasize that one of the institutions that implemented the proposed framework is a central bank, with very complex and challenging specifics. That is why we consider evaluation results relevant, especially in circumstances where there is no similar research available in the open literature. We will continue this research by further developing the framework and implementing it in additional organizations, both financial institutions and those outside the financial industry. Additional implementations will be followed by further detailed evaluation of our approach, with the aim of proving that all framework components have an important contribution to the effective establishment of adaptive e-BCM.

Acknowledgment. Authors are thankful to the Ministry of Education, Science and Technological Development, Republic of Serbia, Grant no.174031.

References

1. Radenković, B., Despotović Zrakić, M., Bogdanović, Z., Barać, D., and Labus, A., *Elektronsko poslovanje* (eng. E-business). Belgrade: Faculty of Organizational Sciences, (2015).
2. Chaffey, D., *E-Business and E-Commerce Management - Strategy, Implementation and Practice*. (2009).
3. Arduini, F. and Morabito, V., Can banks protect themselves for the unthinkable? A Business Continuity Update, *Communications of the ACM*, vol. 53, no. 3, pp. 121–125, (2010).
4. Bergström, J., van Winsen, R., and Henriqson, E., On the rationale of resilience in the domain of safety: A literature review, *Reliability Engineering and System Safety*, vol. 141, pp. 131–141, (2015).
5. Smith, M. and Shields, P.-A., Strategies for IT and communications, in *The Definitive Handbook of Business Continuity Management*, A. Hiles, Ed. John Wiley & Sons, Ltd, (2007), pp. 205–235.
6. Cook, J., A Six-Stage Business Continuity and Disaster Recovery Planning Cycle., *SAM Advanced Management Journal* (07497075), vol. 80, no. 3, pp. 23–68, (2015).
7. Bergström, J., van Winsen, R., and Henriqson, E., On the rationale of resilience in the domain of safety: A literature review, *Reliability Engineering and System Safety*, vol. 141, pp. 131–141, (2015).
8. ISO 22301:2012 Societal security - Business continuity management systems - Requirements. International Organization for Standardization, (2012).
9. Labus, M., Despotović-Zrakić, M., and Bogdanović, Z., Introducing adaptive E-business continuity management, in *Advances in Intelligent Systems and Computing*, vol. 569, pp. 628–637, (2017).
10. Sanchez Dominguez, A. P., *Business Continuity Management: A Holistic Framework for Implementation, Culminating Projects in Information Assurance*, (2016).
11. Bajgoric, N., Business continuity management: a systemic framework for implementation, *Kybernetes*, vol. 43, no. 2, pp. 156–177, (2014).
12. Wong, W. N. Z., Transforming corporate performance: A business continuity management approach, *Organizational Dynamics*, (2018).
13. Gibb, F. and Buchanan, S., A framework for business continuity management, *International Journal of Information Management*, vol. 26, no. 2, pp. 128–141, (2006).
14. Järveläinen, J., IT incidents and business impacts: Validating a framework for continuity management in information systems, *International Journal of Information Management*, vol. 33, no. 3, pp. 583–590, (2013).

15. Herbane, B., Elliott, D., and Swartz, E. M., Business Continuity Management: Time for a strategic role?, *Long Range Planning*, vol. 37, no. 5, pp. 435–457, (2004).
16. Cook, J., A Six-Stage Business Continuity and Disaster Recovery Planning Cycle., *SAM Advanced Management Journal* (07497075), vol. 80, no. 3, pp. 23–68, (2015).
17. Savage, M., Business continuity planning, *Work Study*, vol. 51, no. 5, pp. 254–261, (2002).
18. Herbane, B., Elliott, D., and Swartz, E. M., Business Continuity Management: time for a strategic role?, *Long Range Planning*, vol. 37, no. 5, pp. 435–457, (2004).
19. Tammineedi, R. L., Business Continuity Management: A Standards-Based Approach, *Information Security Journal: A Global Perspective*, vol. 19, no. 1, pp. 36–50, (2010).
20. Lindström, J., Samuelsson, S., and Hägerfors, A., Business continuity planning methodology, *Disaster Prevention and Management: An International Journal*, vol. 19, no. 2, pp. 243–255, (2010).
21. Sahebjamnia, N., Torabi, S. A., and Mansouri, S. A., Integrated business continuity and disaster recovery planning: Towards organizational resilience, *European Journal of Operational Research*, vol. 242, no. 1, pp. 261–273, (2015).
22. Torabi, S. A., Soufi, H. R., and Sahebjamnia, N., A new framework for business impact analysis in business continuity management (with a case study), *Safety Science*, vol. 68, pp. 309–323, (2014).
23. Zeng, Z. and Zio, E., An integrated modeling framework for quantitative business continuity assessment, *Process Safety and Environmental Protection*, vol. 106, pp. 76–88, (2017).
24. (Noel) Bryson, K.-M., Millar, H., Joseph, A., and Mobolurin, A., Using formal MS/OR modeling to support disaster recovery planning, *European Journal of Operational Research*, vol. 141, no. 3, pp. 679–688, (2002).
25. Montshiwa, A., Nagahira, A., and Ishida, S., Modifying Business Continuity Plan (BCP) Towards an Effective Auto-Mobile Business Continuity Management (BCM): A Quantitative Approach, *Journal of Disaster Research*, vol. 11, pp. 691–698, (2016).
26. Lozupone, V., Disaster recovery plan for medical records company, *International Journal of Information Management*, vol. 37, no. 6, pp. 622–626, (2017).
27. Baba, H., Watanabe, T., Nagaishi, M., and Matsumoto, H., Area Business Continuity Management, a New Opportunity for Building Economic Resilience, *Procedia Economics and Finance*, vol. 18, pp. 296–303, (2014).
28. Blos, M. F., Hoeflich, S. L., and Miyagi, P. E., A General Supply Chain Continuity Management Framework, *Procedia Computer Science*, vol. 55, pp. 1160–1164, (2015).
29. Torabi, S. A., Giahi, R., and Sahebjamnia, N., An enhanced risk assessment framework for business continuity management systems, *Safety Science*, vol. 89, pp. 201–218, (2016).
30. Jingye, L. and Takehiro, T., Practical Process for Introducing Smart Business Continuity Management of Smart City in Japan, *Procedia Engineering*, vol. 146, pp. 288–295, (2016).
31. Kato, M. and Charoenrat, T., Business continuity management of small and medium sized enterprises: Evidence from Thailand, *International Journal of Disaster Risk Reduction*, vol. 27, pp. 577–587, (2018).
32. Kaufhold, M.-A. *et al.*, Business Continuity Management in Micro Enterprises: Perception, Strategies, and Use of ICT, *International Journal of Information Systems for Crisis Response and Management*, vol. 10, pp. 1–19, (2018).
33. Zawada, B., The practical application of ISO 22301., *Journal of Business Continuity & Emergency Planning*, vol. 8, no. 1, pp. 83–90, (2014).
34. Zawada, B. and Marbais, G., Implementing ISO 22301, *Avalution Consulting*. (2015).
35. Laudon, K. C. and Guercio Traver, C., E-commerce: business, technology, society. (2012).
36. Aven, T., Risk assessment and risk management: Review of recent advances on their foundation, *European Journal of Operational Research*, vol. 253, no. 1, pp. 1–13, (2016).
37. ISO 22313:2012 Societal security - Business continuity management systems - Guidance. International Organization for Standardization, (2012).
38. Nosworthy, J. D., A Practical Risk Analysis Approach: Managing BCM Risk, *Computers & Security*, vol. 19, no. 7, pp. 596–614, (2000).

39. Sikdar, P., Alternate Approaches to Business Impact Analysis, *Information Security Journal: A Global Perspective*, vol. 20, no. December 2014, pp. 128–134, (2011).
40. Wong, W. N. Z. and Shi, J., Performance evaluation, in *Business Continuity Management System: A Complete Guide to Implementing ISO 22301*, KoganPage, (2015).
41. Watters, J., Disaster Recovery, Crisis Response, and Business Continuity A Management Desk Reference. Apress, (2013).
42. Faertes, D., Reliability of supply chains and business continuity management, *Procedia Computer Science*, vol. 55, no. Itqm, pp. 1400–1409, (2015).
43. Hawkins, S. M., Yen, D. C., and Chou, D. C., Disaster recovery planning: a strategy for data security, *Information Management & Computer Security*, vol. 8, no. 5, pp. 222–230, (2000).
44. ISACA, CISA Review Manual 2014. ISACA, (2013).
45. Hair, J. F. J., Hult, G. T. M., Ringle, C., and Sarstedt, M., A Primer on Partial Least Squares Structural Equation Modeling (PLS-SEM), vol. 46, no. 1–2, (2014).
46. Hair, J. F., Ringle, C. M., and Sarstedt, M., Partial Least Squares: The Better Approach to Structural Equation Modeling?, *Long Range Planning*, vol. 45, no. 5–6, pp. 312–319, (2012).
47. Hair, J. F. and Sarstedt, M., Innovative and established research methods in family business: Description, illustration and application guidelines, *Journal of Family Business Strategy*, vol. 5, no. 1, pp. 1–3, (2014).
48. Ringle, C. M., Wende, S., and Becker, J.-M., *SmartPLS 3*. Boenningstedt: SmartPLS GmbH, (2015).
49. Gudergan, S. P., Ringle, C. M., Wende, S., and Will, A., Confirmatory tetrad analysis in PLS path modeling, *Journal of Business Research*, vol. 61, no. 12, pp. 1238–1249, (2008).
50. Hair, J. F., Sarstedt, M., Ringle, C. M., and Mena, J. A., An assessment of the use of partial least squares structural equation modeling in marketing research, *Journal of the Academy of Marketing Science*, vol. 40, no. 3, pp. 414–433, (2012).
51. Hair, J. F., Ringle, C. M., and Sarstedt, M., Partial Least Squares Structural Equation Modeling: Rigorous Applications, Better Results and Higher Acceptance, *Long Range Planning*, vol. 46, no. 1–2, pp. 1–12, (2013).
52. Reimann, M., Schilke, O., and Thomas, J. S., Customer relationship management and firm performance: the mediating role of business strategy, *Journal of the Academy of Marketing Science*, vol. 38, no. 3, pp. 326–346, (2010).
53. Sarstedt, M., Ringle, C. M., Smith, D., Reams, R., and Hair, J. F., Partial least squares structural equation modeling (PLS-SEM): A useful tool for family business researchers, *Journal of Family Business Strategy*, vol. 5, no. 1, pp. 105–115, (2014).
54. Hair, J. F., Ringle, C. M., and Sarstedt, M., PLS-SEM: Indeed a Silver Bullet, *The Journal of Marketing Theory and Practice*, vol. 19, no. 2, pp. 139–152, (2011).
55. Fornell, C. and Larcker, D. F., Evaluating Structural Equation Models with Unobservable Variables and Measurement Error., *Journal of Marketing Research (JMR)*. Feb1981, vol. 18, no. 1, pp. 39-50. 12p. 1 Diagram, (1981).
56. Hair, J. F., Ringle, C. M., and Sarstedt, M., PLS-SEM: Indeed a Silver Bullet, *The Journal of Marketing Theory and Practice*, vol. 19, no. 2, pp. 139–152, (2011).
57. Geisser, S., A predictive approach to the random effect model, *Biometrika*, vol. 61, pp. 101–107, (1974).
58. Stone, M., Cross-Validatory Choice and Assessment of Statistical Predictions, *Journal of the Royal Statistical Society*, vol. 36, no. 2, pp. 111–147, (1974).

Milica Labus holds a BSc in Mathematics and Computer Science from the University of Belgrade, a MSc. in Computer Science from Stanford University, and an Executive MBA from the Cotrugli Business School. She is a Ph.D. candidate at the Faculty of Organizational Sciences, University of Belgrade, majoring E-Business. Milica has more than seventeen years of varied experience in the IT field and has worked on many projects in the financial and public sectors.

Zorica Bogdanović, PhD, is an associate professor and the head of Department of E-business at Faculty of Organizational Sciences, University of Belgrade. She is the secretary of IEEE computer chapter C16. Her research interests include e-business, internet of things and internet technologies.

Marijana Despotović-Zrakić, PhD, is a professor and the head of Laboratory for Simulation at Faculty of Organizational Sciences, University of Belgrade. She is the vice-chair of IEEE computer chapter C16. Her research interests include e-business project management, e-business risk management, computer simulation and internet technologies.

Dušan Barać is an associate professor at Faculty of Organizational Sciences, University of Belgrade. His research interests include internet technologies, e-commerce and e-business application development.

Snežana Popović is an assistant professor at the School of Computing, Union University in Belgrade. Her research interest is mainly related to information systems modeling, design and development, service-oriented architecture and databases.

Received: February 02, 2019; Accepted: December 01, 2019

7. Appendix A

Table 7. Measurement items included in the questionnaire and indicators for reflective measurement of model constructs, with statistics

ID	Indicator	Mean	Std. Dev.	95% Conf. Interval	t value	The probability to reject the alternative hypothesis			Outer Loadings	
						Ha: Mean < 4	Ha: Mean != 4	Ha: Mean > 4		
Part1: Adaptive e-BCM, Simple average 4.45										
Q1.1*	The framework defines clear and effective procedures for the establishment of e-business continuity management	4.500	0.647	4.287	4.713	t = 4.7621	Pr(T < t) = 1.0000	Pr(T > t) = 0.0000	Pr(T > t) = 0.0000	0.707*
Q1.2*	The framework examines the specifics of an organization's business, in particular its use of modern e-business technologies	4.500	0.647	4.287	4.713	t = 4.7621	Pr(T < t) = 1.0000	Pr(T > t) = 0.0000	Pr(T > t) = 0.0000	0.85*
Q1.3*	The framework facilitates the establishment of business continuity in an organization that uses modern e-business technologies	4.684	0.620	4.481	4.888	t = 6.8058	Pr(T < t) = 1.0000	Pr(T > t) = 0.0000	Pr(T > t) = 0.0000	0.665

Q1.4*	The framework enables adaptive e-BCM in accordance with changes to the organization's environment	4.421	0.758	4.172	4.670	t = 3.4239	Pr(T < t) = 0.9992	Pr(T > t) = 0.0015	Pr(T > t) = 0.0008	0.565
Q1.5	The framework contributes to a positive organizational attitude towards e-BCM	4.132	1.119	3.764	4.499	t = 0.7248	Pr(T < t) = 0.7634	Pr(T > t) = 0.4732	Pr(T > t) = 0.2366	0.861*
<hr/>										
Part 2: Business Impact Analysis Component (C1), Simple average 4.58										
<hr/>										
Q2.1*	BIA is carried out as a part of the organization's key management activities	4.737	0.724	4.499	4.975	t = 6.2780	Pr(T < t) = 1.0000	Pr(T > t) = 0.0000	Pr(T > t) = 0.0000	0.559
Q2.2*	Key business processes and their interdependences are identified	4.711	0.515	4.541	4.880	t = 8.5037	Pr(T < t) = 1.0000	Pr(T > t) = 0.0000	Pr(T > t) = 0.0000	0.600
Q2.3*	Resources necessary for key business processes are identified	4.474	0.762	4.223	4.724	t = 3.8329	Pr(T < t) = 0.9998	Pr(T > t) = 0.0005	Pr(T > t) = 0.0002	0.76*
Q2.4*	Financial and operational impacts are assessed	4.316	0.842	4.039	4.592	t = 2.3129	Pr(T < t) = 0.9868	Pr(T > t) = 0.0264	Pr(T > t) = 0.0132	0.787*
Q2.5*	Recovery time objectives are determined	4.605	0.638	4.395	4.815	t = 5.8446	Pr(T < t) = 1.0000	Pr(T > t) = 0.0000	Pr(T > t) = 0.0000	0.516
Q2.6*	Criticality levels for key business processes and recovery priorities are determined	4.658	0.534	4.482	4.833	t = 7.5940	Pr(T < t) = 1.0000	Pr(T > t) = 0.0000	Pr(T > t) = 0.0000	0.564

 Part 3: e-BCM Risk Assessment Component (C2), Simple average 4.05

Q3.1*	Business continuity risk assessment process is established	4.237	0.786	3.978	4.495	t = 1.8571	Pr(T < t) = 0.9644	Pr(T > t) = 0.0713	Pr(T > t) = 0.0356	0.765*
Q3.2*	Results of risk assessment are adequately documented	4.421	0.758	4.172	4.670	t = 3.4239	Pr(T < t) = 0.9992	Pr(T > t) = 0.0015	Pr(T > t) = 0.0008	0.549
Q3.3	Risk assessment is performed on a regular basis	3.868	1.018	3.534	4.203	t = - 0.7968	Pr(T < t) = 0.2153	Pr(T > t) = 0.4307	Pr(T > t) = 0.7847	0.640
Q3.4	Risk treatment process is established	3.947	1.038	3.606	4.289	t = - 0.3125	Pr(T < t) = 0.3782	Pr(T > t) = 0.7565	Pr(T > t) = 0.6218	0.915*
Q3.5	Risk treatment methods are adequate	3.921	1.194	3.529	4.314	t = - 0.4075	Pr(T < t) = 0.3430	Pr(T > t) = 0.6860	Pr(T > t) = 0.6570	0.927*
Q3.6	Risk treatment plan is regularly monitored	3.921	0.969	3.602	4.240	t = - 0.5021	Pr(T < t) = 0.3093	Pr(T > t) = 0.6186	Pr(T > t) = 0.6907	0.835*

 Part 4: Component for Development of Business Continuity Plan (C3), Simple average 3.98

Q4.1*	BCMS is adequately defined	4.579	0.683	4.354	4.803	t = 5.2248	Pr(T < t) = 1.0000	Pr(T > t) = 0.0000	Pr(T > t) = 0.0000	0.461
Q4.2*	Business continuity strategy is based on BIA and risk assessment results	4.605	0.679	4.382	4.829	t = 5.4917	Pr(T < t) = 1.0000	Pr(T > t) = 0.0000	Pr(T > t) = 0.0000	0.661
Q4.3	Business continuity strategy defines and provides all necessary resources in the event	4.158	1.027	3.820	4.496	t = 0.9474	Pr(T < t) = 0.8252	Pr(T > t) = 0.3496	Pr(T > t) = 0.1748	0.805*

of a crisis situation

Q4.4*	Business continuity plan is adequately defined	4.289	0.867	4.004	4.574	t = 2.0581	Pr(T < t) = 0.9767	Pr(T > t) = 0.0467	Pr(T > t) = 0.0233	0.838*
Q4.5	Business continuity plan training is adequate	3.211	1.255	2.798	3.623	t = - 3.8765	Pr(T < t) = 0.0002	Pr(T > t) = 0.0004	Pr(T > t) = 0.9998	0.773*
Q4.6	Business continuity plan is regularly tested	3.053	1.314	2.621	3.485	t = - 4.4441	Pr(T < t) = 0.0000	Pr(T > t) = 0.0001	Pr(T > t) = 1.0000	0.629

Part 5: Component for Continual BCMS Improvement (C4), Simple average 3.36

Q5.1	BCMS performance evaluation metrics are established	3.447	0.950	3.135	3.760	t = - 3.5858	Pr(T < t) = 0.0005	Pr(T > t) = 0.0010	Pr(T > t) = 0.9995	0.705*
Q5.2	BCMS performance is regularly evaluated	3.447	1.155	3.068	3.827	t = - 2.9484	Pr(T < t) = 0.0028	Pr(T > t) = 0.0055	Pr(T > t) = 0.9972	0.942
Q5.3	Results of performance evaluation are adequately documented	3.316	1.093	2.956	3.675	t = - 3.8585	Pr(T < t) = 0.0002	Pr(T > t) = 0.0004	Pr(T > t) = 0.9998	0.912*
Q5.4	Adequate BCMS management review is established	3.342	1.097	2.981	3.703	t = - 3.6958	Pr(T < t) = 0.0004	Pr(T > t) = 0.0007	Pr(T > t) = 0.9996	0.959
Q5.5	BCMS management review is regularly conducted	3.342	1.122	2.973	3.711	t = - 3.6156	Pr(T < t) = 0.0004	Pr(T > t) = 0.0009	Pr(T > t) = 0.9996	0.965
Q5.6	BCMS management review provides decision-making on implementation of corrective BCMS measures	3.447	1.032	3.108	3.787	t = - 3.3015	Pr(T < t) = 0.0011	Pr(T > t) = 0.0021	Pr(T > t) = 0.9989	0.648

Q5.7	Corrective BCMS measures are adequately documented	3.447	0.978	3.126	3.769	t = - 3.4830	Pr(T < t) = 0.0006	Pr(T > t) = 0.0013	Pr(T > t) = 0.9994	0.806*
Q5.8	BCMS Improvement Plan is comprehensive and up-to-date	3.184	0.955	2.870	3.498	t = - 5.2685	Pr(T < t) = 0.0000	Pr(T > t) = 0.0000	Pr(T > t) = 1.0000	0.904*
Q5.9	Management is regularly informed of the status of corrective BCMS measures	3.263	0.950	2.951	3.575	t = - 4.7830	Pr(T < t) = 0.0000	Pr(T > t) = 0.0000	Pr(T > t) = 1.0000	0.908*

ID*: indicators that passed the t-test, in the sense that probability to reject the null hypothesis H₀ against the alternative (H_a: Mean > 4) is below the significance level (5%)

Outer Loadings*: indicators with desirable levels of outer loadings

Human Activities Recognition with a Single Wrist IMU via a Variational Autoencoder and Android Deep Recurrent Neural Nets

Edwin Valarezo Añazco¹, Patricio Rivera Lopez¹, Hyemin Park¹, Nahyeon Park¹, and Tae-Seong Kim^{1*}

¹ Department of Biomedical Engineering, College of Electronics and Information, Kyung Hee University, Republic of Korea
{edgivala, patoalejor, hmp9669, nhpark, tskim}@khu.ac.kr

Abstract. Human Activity Recognition (HAR) is an active research field because of its versatility towards various application areas such as healthcare and lifecare. In this study, a novel HAR system is proposed based on an autoencoder for denoising and Recurrent Neural Network (RNN) for classification with a single Inertial Measurement Unit (IMU) located on a dominant wrist. A Variational Autoencoder (VAE) is built to denoise IMU signals which improves HAR by Android Deep RNN. Evaluating our VAE and Android Deep RNN HAR system is done in two ways. First, the system is tested on a PC using discrete epochs of activities of daily living. Our results show that VAE improves Signal-to-Noise Ratio (SNR) of the IMU signal from 8.78 to 17.26 dB. In turn, HAR improves from 89.29% to 95.11% in F1-score and from 90.38% to 95.47% in accuracy. Secondly, the system is tested on an Android device (i.e., smartphone) using continuous activity signals. This is done by transferring the PC HAR system to an Android HAR App (i.e., Android Deep RNN). We have achieved 86.13% and 95.09% in accuracy without and with VAE respectively. Our results demonstrate that HAR can be achieved in real-time on a standalone smart device with a single IMU for lifelogging services.

Keywords: Human Activity Recognition; Denoising Autoencoder; Android Deep Recurrent Neural Networks; Mobile Application.

1. Introduction

Human Activity Recognition (HAR) is defined as a context-aware technology that maps activity data collected by IoT sensors into activity labels [1]. In recent years, HAR has gained strong attention because of its various applications [2], such as fall detection [3] [4], rehabilitation [5] [6], activity recognition in sports [7] [8], sports training analysis

* Corresponding author.

Corresponding author address: Room 719, College of Electronics and Information, Kyung Hee University, 1732, Deogyong-daero, Giheung-gu, Yongin-Si, Gyeonggi-do 17104, Republic of Korea.
Corresponding author office Tel.: +82-31-201-3731.
Corresponding author email.: tskim@khu.ac.kr.

[9] [10], energy expenditure estimation [11], and health care [12] [13]. Two issues have recently been raised regarding advanced HAR. First, in HAR studies, activity noises affect the recognition rate [14]. Then, it is necessary to use filters to denoise the motion signals. Second, most HAR studies are implemented on a PC platform. There is a strong need for a practical HAR system on a portable platform, such as smartphones.

Recently Deep Learning (DL) has been adopted to create a better representation of the signal with less noise via Denoising Autoencoder (DAE) [15] [16]. Some studies tested DAE for instance: in Heyman, et al. [17], Bayesian Feature Enhancement (BFE) and DAE were used to clean-up voice signals for speech recognition using the output of BFE as target data. In Xiong, et al. [18], wavelet transform with the scale-adaptive thresholding method was used to filter out noise and DAE to remove the residual noise from electrocardiogram signals. In both cases, additional denoising stages were used as target data for DAE. Then, the performance of DAE depends on the target. In contrast, a totally unsupervised DAE is possible (i.e., using only autoencoder for denoising). For instance, in Mohammed, et al. [14], Variational Autoencoder (VAE) produced desirable denoising performance for motion artifact from a sensor attached to clothes. Our VAE implementation focuses on denoising activity signals from a wrist Inertial Measurement Unit (IMU).

The current popularity of wearable technology creates an invaluable opportunity to build a portable HAR system. However, building an HAR system using DL on a portable device such as smartphones require some considerations. Most of all, large and complex classification algorithms could not be implemented on a portable device because of the computational burden on limited computing power and memory. There are some recent works trying to implement HAR systems on portable devices. For instance, Lane and Georgiev [19] used a low-power Deep Neural Network (DNN) on a CPU paired with a mobile device. In their system, DNN worked with a cloud system (i.e., activity label was inferred on PC and then sent to the smartphone via the Internet). However, the cloud system with a web connection is limited to not be a standalone HAR system. To overcome the limitations of the cloud service computation, some studies implemented their classification algorithms in a standalone mode (i.e., running on a smartphone) [20] [21] [22] [23] [24]. However, conventional classification approaches are used, instead of DL. For instance, Jongprasithporn, et al. [22] applied thresholds values to classify standing, walking, and running activities. Due to the threshold values were fixed for the three activities, complex activities could not be recognized. The other studies used Naïve Bayes (NB), Support Vector Machine (SVM), Random Forest (RF), and K-Nearest Neighbors (KNN) with hand-extracted features to differentiate sitting, standing, lying, walking, running, and cycling. DL should improve the performance over the mentioned conventional classification approaches and avoid the needed of hand-extracted features [25].

In this study, we have developed a HAR system with a single IMU at a dominant wrist based on VAE and Android Deep RNN. Our contributions focus on solving the problems of noise in HAR signals and implementation of a reliable HAR system on portable devices. To solve the problem of noise in the HAR signals, a denoising stage uses VAE based on Convolution Neural Network (CNN) and Recurrent Neural Network (RNN). VAE enhances the IMU signal by approximating the spatial-temporal features to Gaussian distributions. The importance of the denoising stage lays in the fact that IMUs are prone to motion artifacts and sensor noise. To implement a reliable HAR system we used RNN because it outperforms other algorithms in HAR [25] [26]; to

make the HAR system portable avoiding the limitations of cloud systems (e.g., connection delays), the HAR system was implemented under Android (i.e., Android Deep RNN) whereas most DL RNNs are implemented only for a PC [27] [28] [29]. RNN running on Android increases the scope of our system to any smart device using Android, even if the smart device has not internet connection. Our results show that HAR accuracy increases with the denoised IMU signal compared to the raw signal without VAE. Furthermore, the inference time of our Android Deep RNN on Android smartphone is fast enough to implement real-time lifelogging App.

2. Methodology

The proposed system is represented in Fig. 1. Human activities are sensed with a single IMU worn on one dominant wrist. In our HAR system, VAE denoises a tri-axial accelerometer and tri-axial gyroscope signals. Then the denoised IMU signals are classified into activity labels via Android Deep RNN on the PC and Android smartphone platforms. In this study, we recognize Activities of Daily Living (ADL) including standing, walking upstairs, walking downstairs, walking, running, cycling, and Nordic walking.

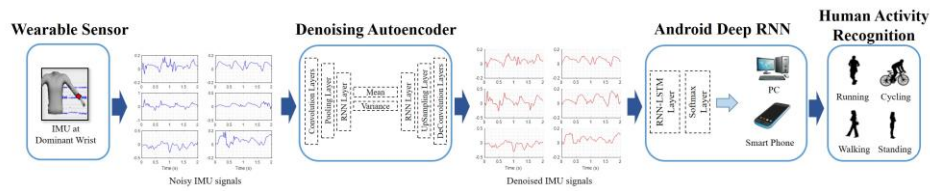


Fig. 1. HAR System workflow.

2.1. ADL Database

In this study, we used a public database (DB), Physical Activity Monitoring for Aging People (PAMAP2) [30] [31]. PAMAP2 includes raw information using three IMUs located on the chest, the wrist, and the ankle, with a sampling frequency of 100 Hz. Each IMU provides the information of tri-axial acceleration (ms^{-2}), tri-axial gyroscope (rad/s), and tri-axial magnetometer (μT). The DB includes data of nine subjects performing 12 ADLs, including ironing, vacuum cleaning, walking, walking up and downstairs, sitting, lying, and standing. High-level dynamic activities have also included, such as Nordic walking, running, cycling, and rope jumping.

For this study, we have selected the following seven ADLs: *Standing* consists of just standing without interaction with another person or standing while the subject is talking. *Walking Upstairs* and *Walking Downstairs* (Ascending and Descending stairs) were carried out in a five-floors building. *Walking* consists of walking outdoors with moderate speed steps. *Nordic Walking* is performed on asphaltic terrain using asphalt pads on the walking poles. *Cycling* is conducted outdoors with a real bike. *Running*

corresponds to jogging outdoors with suitable speed. These ADLs are selected because they should be well reflected on the wrist sensor and most of the subjects perform these activities for all the trials.

2.2. Denoising Autoencoder

Autoencoder (AE) is a feed-forward neural network that reproduces its input as output. It includes an encoder network that generates a featured representation of the input in the latent space. Then the decoder network decodes the feature representation back to the input dimension [32].

VAE is a variant of autoencoder that is able to denoise motion signals. As autoencoder, VAE encodes the input data into a feature-representation vector in the latent space. However, in VAE, the feature representation vector is constrained to Gaussian distributions [33]. Also, VAE uses two loss functions, a combination of Kullback-Leibler divergence (KL) and Mean Square Error (MSE). The first loss function forces the feature-representation vector in the latent space to follow a Gaussian distribution, measuring the relative entropy between the approximate posterior and the prior probability density function [33]. The second loss measures the similarity between the output and the input.

Our VAE model is presented in Fig. 2. It uses a combination of CNN and RNN with Long Short-Term Memory (i.e., LSTM) layers as in [14]. Nevertheless, our VAE implementation has two main differences from the previous approaches [34] [14] [33]. The first difference lies in our VAE architecture. Our VAE has three convolutional layers and three LSTM layers in the encoder. The latent space has two dense layers. The decoder has three LSTM layers and three deconvolutional (ConvTranspose) layers. The second difference is the loss function. It is a weighted combination of KL and MSE similar to Mohammed and Tashev [14]. The weighted loss function is designed by analyzing the effect of KL and MSE in the Signal-to-Noise Ratio (i.e., SNR). Eq. 1 describes the weighted loss function used to train VAE. E_{KL} measures the KL divergence between the prior and posterior probability density functions, described by the feature representation vector z in the latent space. E_{MSE} measures the MSE error between the output and input signals; γ is the weight value, in this work γ value is 5.

$$\mathcal{L} = \gamma * (E_{MSE}[(\hat{y} - y)^2]) + \frac{1}{1000 * \gamma} * (E_{KL}[q(z|x)||p(z)]) \quad (1)$$

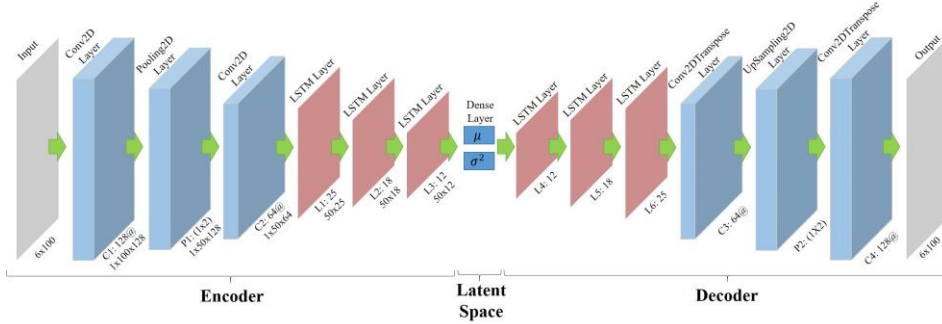


Fig. 2. Variational Autoencoder.

2.3. Android Deep RNN

RNN is characterized by recurrent connections between hidden units. These recurrent connections generate a temporal memory in which the previous state of the network is stored. This lets RNN take decisions based on the previous stage of the network and the current input [26].

Our Android Deep RNN uses the classic LSTM cell proposed in [35]. LSTM produces internal paths regulated by gates, where the gradient can flow for long durations. These regulated paths allow propagation of the error in deep networks. For weight optimization, Stochastic Gradient Descent (SGD) was used. SGD uses a momentum factor to determine how fast the algorithm converges. SGD is described by Eqs. 2 and 3.

$$W_{t+1} = W_t + v_{t+1} \tag{2}$$

$$v_{t+1} = \beta \times v_t - \alpha \times \nabla \mathcal{L}(W_t) \tag{3}$$

where α is learning rate, W_t weight matrix, \mathcal{L} the loss function, and β the momentum factor. The loss function is Negative Logarithmic Likelihood. It is described by Eq. 4, where n corresponds to the number of classes and P_y^i is the probability of the class i .

$$\mathcal{L} = -\frac{1}{n} \sum_{i=1}^n \log(P_y^i) \tag{4}$$

Fig. 3 shows our Android Deep RNN, the main difference with previous works are: first, our Android Deep RNN can run on Android smartphone on standalone mode (i.e., inference on Android device), while most of the previous HAR systems are implemented only for PC [29] [28] [27]. Second, the architecture of our Android Deep RNN, for instance Milenkoski, *et al.* [36] built an HAR system using three LSTM layers with 64 cells each layer. In contrast, our network is lighter than them in terms of number of layers. It has only a single RNN-LSTM layer and a single hidden layer to reduce the memory usage on Android smartphone. The single RNN-LSTM layer has 100 cells

corresponding to 2 seconds data; the hidden layer has 120 nodes; at the end of our network a softmax layer with seven nodes, corresponding to each class.

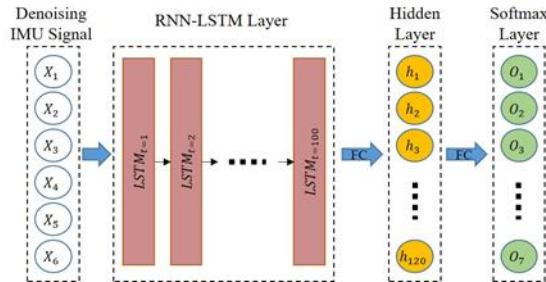


Fig. 3. Android Deep RNN Architecture. The six-channels input matches with the RNN-LSTM layer. The remaining layers are fully connected.

2.4. VAE and Android Deep RNN Implementation

In this study, VAE was implemented on a PC platform using Keras and TensorFlow. The convolutions and the RNN-LSTM layers were implemented to transform the spatial-temporal information of the IMU signals into a family of Gaussian distributions in the latent space with mean μ and variance σ^2 . The mean μ and variance σ^2 are represented as dense layers (see Fig. 2). Mini-batch approach was used to train VAE with a mini-batch size of 100. Adam Optimizer was used with a learning rate of 0.0008. The kernels size and the hidden nodes are same for each encoder-decoder equivalent layer.

The Android Deep RNN was implemented using Deeplearning4J library [37]. The RNN-LSTM layer of our HAR system classifies the data as many-to-one, after an epoch (i.e., 2 seconds data) the model infers a single class label. Truncated Backpropagation Through Time (Truncated-BPTT) was used to train the Android Deep RNN. Truncated-BPTT is a variant of BPTT that restricts the downside to a specific number of hidden nodes (i.e., the gradient is not back-propagated to all hidden nodes). To train our HAR system the mini-batch approach was used with a mini-batch size of 113. Weight initialization was done using a random number generator. Hyperparameters such as learning rate and number of hidden nodes were fixed by analyzing the behavior of these in the loss function. The learning rate was set as 0.09 and the number of training steps was fixed after setting the learning rate and to avoid overfitting. The training of the Android Deep RNN was done using IntelliJ IDE on PC.

An Android app was created using Android Studio IDE with the Application Program Interface (API) 23. API 23 lets the app run on Android 6.0 or higher. The app loads the network architecture shown in Fig. 3 with the trained weights. Then, the inference is made on an Android smartphone without fine-tuning or retraining of the network on smartphone.

The smartphone used in this study is a Samsung Galaxy S7 with the following technical specifications: Android version 7.0, 4GB of RAM, Octa-core processor of

4x2.3 GHz and 4x1.6 GHz. The PC specification includes Processor Intel(R) Core (TM) i5-7500 CPU@ 3.40GHz, 8GB of RAM, and NVIDIA GeForce GTX 1050 Ti.

2.5. Data Preparation

In this study, we utilized activity information from a tri-axial accelerometer and a tri-axial gyroscope on a dominant wrist. We created a discrete epoch dataset, where each epoch is segmented data with a time window of two seconds.

The IMU signals were preprocessed as follows: The activity data was downsampled to 50 Hz, as was suggested in [38] for HAR. The gravity effect was removed by a high-pass Butterworth filter, with a cutoff frequency of 0.25 Hz [39]. The activity data were normalized between -1 to 1. Finally, the sliding window approach was used with an overlap of 50% and a window length of two seconds. The total number of epochs per activity are 610 epochs for standing, 1470 epochs for walking, 700 epochs for running, 1020 epochs for cycling, 1340 epochs for Nordic walk, 640 epochs for ascending stairs, and 500 epochs for descending stairs.

In addition, the performance of our HAR system on a smartphone was evaluated by a real-time simulation with continuous epochs of activities. The activity datasets of discrete and continuous epochs were created from seven subjects in PAMAP2.

2.6. Validations

Validation of our VAE was done by analyzing the SNR improvement. Eq. 5 describes mathematically the SNR. The SNR was calculated by taking the Root-Mean-Square level (i.e., RMS) of the signal of interest involving certain activities and the RMS level of noisy background signals (i.e., no activities performed).

$$SNR = 20 * \log_{10} \frac{RMS(\text{Signal of Interest})}{RMS(\text{Noise})} \quad (5)$$

The following performance metrics are considered to assess our HAR system [40].

$$\text{Recall or Sensitivity} = \frac{TP}{TP + FN} \quad (6)$$

$$\text{Precision or Positive Predictive Values (PPV)} = \frac{TP}{TP + FP} \quad (7)$$

$$\text{Specificity} = \frac{TN}{TN + FP} \quad (8)$$

$$\text{Negative Predictive Values (NPV)} = \frac{TN}{TN + FN} \quad (9)$$

$$F1 - \text{Score} = 2 \times \left(\frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} \right) \quad (10)$$

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (11)$$

To validate the performance of our HAR system two evaluation methods were used. The first evaluation is done using the discrete epochs dataset on a PC. The performance of the HAR system was evaluated through a five-fold. The second validation methodology is a continuous epoch evaluation that runs on an Android smartphone. This test evaluates the feasibility and performance of a real-time HAR system under a standalone mode.

3. Results

3.1. Denoising Performance of VAE

Fig. 4 shows the comparison of some representative epochs from different activities before and after VAE. The effect of VAE is shown as a reduction of the high-frequency components in the denoised activity signals. The overall SNR is improved from 8.78 to 17.26 dB.

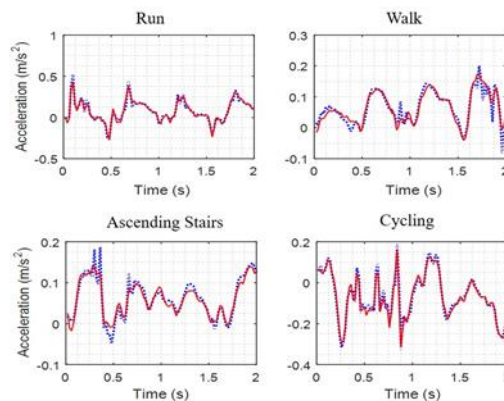


Fig. 4. Raw epochs (dotted in blue) vs. denoised epochs (solid in red) via VAE.

Fig. 5 shows the precision and recall values of our HAR system without and with VAE. With VAE there is an increment in precision for standing from 91.7% to 96.62%, for cycling from 80.07 to 98.85%, for ascending stairs from 86.42% to 88.76%, and for descending stairs from 78.98% to 87.57%. Also, there is an increase in recall for standing from 80.55% to 95.15%, for cycling from 90.94% to 96.62%, for ascending stairs from 80.46% to 89.27%, and for descending stairs from 90.85% to 96.73%. The remain ADLs have a minor improvement in precision and recall.

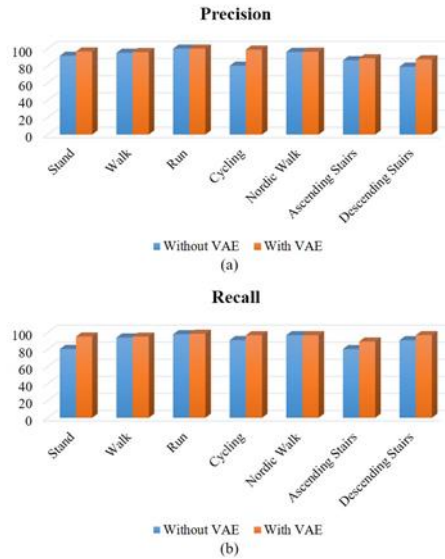


Fig. 5. Comparison of precision and recall with and without VAE. (a) the precision value of our HAR system. (b) the recall value of our HAR system.

3.2. HAR for Discrete Epochs on a PC Platform

The performance of the proposed HAR system without VAE is shown in Table 1 as a confusion matrix and in Table 2 as a summary of the performance metrics for each activity.

Table 1. HAR Confusion Matrix without VAE.

(%)	Stand	Walk	Run	Cycling	Nordic Walk	ASCD ^a Stairs	DESC ^b Stairs
Stand	80.55	0.00	0.00	13.68	1.52	1.22	3.04
Walk	1.26	93.97	0.00	0.50	0.50	1.51	2.26
Run	0.00	1.11	97.78	0.00	0.00	0.56	0.56
Cycling	4.15	0.00	0.00	90.94	0.38	0.75	3.77
Nordic Walk	0.32	0.97	0.00	0.97	96.76	0.65	0.32
ASCD Stairs	2.87	7.47	0.00	3.45	2.30	80.46	3.45
DESC Stairs	1.31	0.65	0.00	2.61	0.00	4.58	90.85

^a Ascending Stairs (ASCD Stairs)

^b Descending Stairs (DESC Stairs)

Table 2. Performance Metrics without VAE.

(%)	Stand	Walk	Run	Cycling	Nordic Walk	ASCD Stairs	DESC Stairs
Precision	91.70	95.17	100	80.07	96.14	86.42	78.98
NPV	95.79	98.30	99.75	98.41	99.33	97.93	99.14
Recall	80.55	93.97	97.78	90.94	96.76	80.46	90.85
Spec^c	98.38	98.65	100	96.11	99.20	98.65	97.76

^c Specificity

The confusion matrix of Table 3 shows the performance of our HAR system with VAE. There is an improvement in classification, especially for standing. The improvement is shown through a reduction in the confusion of standing against the rest of activities and through an increase in the performance metrics.

The values in the main diagonal of the confusion matrix correspond to the recall of each activity. Table 3 shows six of the seven activities (standing, walking, running, cycling, Nordic walking, and descending stairs) having the recall value up to 94.5%. Ascending and descending stairs caused some confusion due to their similarities. For these activities, the recall values are 89.27% and 96.73% respectively. On average, the F1-score increases from 89.29% to 95.11% and accuracy from 90.38% to 95.47%, if VAE is used.

Table 3. HAR Confusion Matrix with VAE.

(%)	Stand	Walk	Run	Cycling	Nordic Walk	ASCD Stairs	DESC Stairs
Stand	95.15	1.21	0.00	0.61	2.12	0.61	0.30
Walk	0.25	94.99	0.00	0.00	0.50	2.01	2.26
Run	0.00	0.55	98.34	0.00	0.00	1.10	0.00
Cycling	1.88	0.00	0.00	96.62	0.38	0.75	0.38
Nordic Walk	0.66	0.99	0.00	0.00	96.70	0.66	0.99
ASCD Stairs	1.13	4.52	0.00	0.56	0.56	89.27	3.95
DESC Stairs	0.65	0.00	0.00	0.00	0.00	2.61	96.73

Table 4 illustrates the performance metrics of our Android Deep RNN with VAE. Comparing to the results without VAE in Table 2, there is an increment in the performance metrics for all activities if VAE is used. On average, the precision increases from 89.78% to 94.88%, NPV from 98.38% to 99.23%, recall from 90.19% to 95.40%, and specificity from 98.39% to 99.24% with the use of VAE.

Table 4. Performance Metrics with VAE.

(%)	Stand	Walk	Run	Cycling	Nordic Walk	ASCD Stairs	DESC Stairs
Precision	96.62	95.95	100	98.85	96.38	88.76	87.57
NPV	98.92	98.59	99.82	99.42	99.34	98.84	99.70
Recall	95.15	94.99	98.34	96.62	96.70	89.27	96.73
Spec	99.26	98.87	100	99.81	99.27	98.77	98.73

3.3. HAR for Continuous Epochs on a Smartphone

Fig. 6 presents the recognized labels for two subjects without VAE, with VAE, and the ground-truth. Without VAE, there is confusion between standing and cycling for almost all epochs in the class of standing.

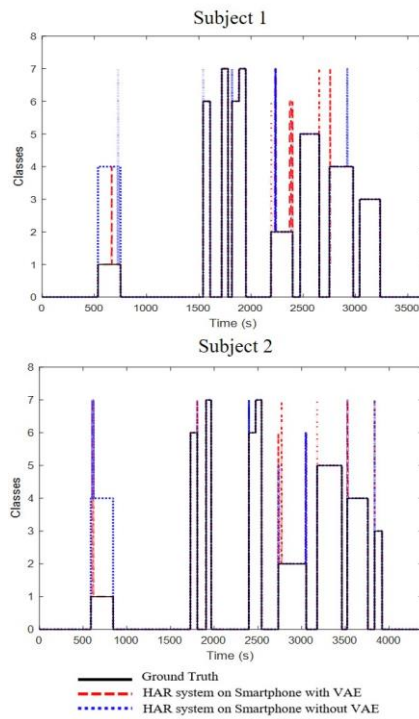


Fig. 6. HAR results via the continuous epochs evaluation on an Android smartphone: Label 0 corresponds to Null, Label 1 Standing, Label 2 Walking, Label 3 Running, Label 4 Cycling, Label 5 Nordic Walking, Label 6 Ascending Stairs, Label 7 Descending Stairs.

Table 5 shows a summary of the accuracy and inference time for all subjects on the Android smartphone. The accuracy per subject is the correctly recognized activities according to Eq. 11. The effect of VAE is reflected with an improvement in accuracy

for each subject. On average, our HAR system increases its accuracy by 8.97% with VAE.

Table 5. HAR results via the continuous epochs evaluation on an Android smartphone.

Subjects	Signal Duration [s]	Accuracy Without VAE [%]	Accuracy With VAE [%]	Inference Time On Smartphone [s]
1	3672	82.64	96.53	709
2	4386	79.65	98.02	847
3	2467	90.33	96.19	477
5	3654	81.54	87.52	707
6	3550	80.54	96.78	686
7	3073	95.30	97.07	594
8	4003	92.89	93.56	774

4. Discussion

This study introduces a HAR system that uses only one IMU on the dominant wrist to recognize ADLs. The proposed system uses VAE to denoise IMU time-series signals. The Android Deep RNN classifies the denoised epochs into activity labels. To make our HAR system practical, our Android Deep RNN was implemented in a standalone mode on a smartphone to perform HAR in real-time using continuous epochs.

Our experiments showed that VAE has a positive effect on our HAR system. VAE reconstructs the IMU signals as less-noisy signals by using the features in the IMU signals because the SNR improves around 9 dB. Overall, the F1-score improves by 6%, and accuracy improves by 5% using VAE.

For statistical significance of performances, a paired T-test, i.e., performance metrics without and with VAE, was performed over Tables 2 and 4. The T-test revealed the statistical significance (p -value ≤ 0.05) for precision ($p=0.046$), recall ($p=0.020$), and NPV ($p=0.40$). Only the specificity shows a p -value > 0.05 with $p=0.068$.

Similar to the discrete epoch evaluation results, in the continuous epoch evaluation results, there is an improvement in performance when VAE is used. This improvement is shown as a decrease in confusion between classes. The average accuracy of our HAR system increases from 86.13% without VAE to 95.09% with VAE. The inference time for a single epoch on Android smartphone is 0.2 seconds. Thus, our Android Deep RNN could be implemented in a real-time mode.

Table 6 summarizes the previous HAR works and their performances including ours. Some of the previous works did not use AE in their HAR systems, such as [41] and [42]. The highest F1-score among the previous works is 93.9%, using Random Forest and hand-extracted features [42]. The highest F1-score using the DL algorithms is 93.7%, using a CNN based classifier [41]. Our approach achieved F1-score of 95.11% using VAE for denoising and Android Deep RNN for HAR. There are only few works that used AE in their HAR systems. For instance, in the work of Mohamed and Tashev [14], AE was used to denoise the IMU signals, then Deep Convolution LSTM (DCLSTM) was used for HAR, achieving an increment of 1% in F1-score. In our case,

F1-score is increased by 6% with VAE. In the work of Almaslukh et al. [43], the HAR classifier was based on a variant of AE called Stacked Autoencoder (SAE). Our VAE differs from SAE in terms of network structures, SAE is a feedforward neural network without convolutional or LSTM layers. Also, SAE was used only for feature extraction and classification similar to the neural networks in [41] and [42], but not for denoising the IMU signals. With SAE, an accuracy of 97.5% was reported instead of F1-score, which is comparable to our accuracy. Note that these summarized works do not represent the fair comparison with the same datasets and models, but the usefulness of AE is clearly demonstrated.

Table 6. Comparison with others HAR studies.

Model	HAR Works	Sensors	AE	Database	F1-Score (%)
KNN	Arif and Kattan [42]	DWS ^d	x	PAMAP	70.80
NN	Arif and Kattan [42]	DWS	x	PAMAP	89.60
DNN	Hammerla et al. [41]	MS ^e	x	PAMAP	90.40
LSTM-F	Hammerla et al. [41]	MS	x	PAMAP	92.90
CNN	Hammerla et al. [41]	MS	x	PAMAP	93.70
Rotation Forest	Arif and Kattan [42]	DWS	x	PAMAP	93.90
DCLSTM	Mohammed and Tashev [14]	MS	✓	Opportunity	90.81
SAE	Almaslukh et al. [43]	SWS ^f	✓	Smartphone-based HAR	97.50*
ADRNN ^g	Proposed HAR System	DWS	✓	PAMAP	95.11

^d Dominant Wrist Sensor (DWS)

^e Multiple Sensor (MS)

^f Single Waist Sensor (SWS)

^g Android Deep RNN (ADRNN)

* Accuracy

5. Conclusion

The reduction of noise in HAR signals improves the SNR as well as the HAR performance. Our VAE successfully denoises motion signals from an IMU at wrist by reducing the noise about 9 dB in SNR. In turn, the HAR performance improves around 6% in all metrics. Furthermore, lifelogging seems feasible on Android smart device, since our Android Deep RNN runs on Android smartphone in standalone mode and the inference time on a smartphone is less than the sensing time. Then, real-time HAR is possible. Finally, our HAR system could be used to develop mobile Apps to monitor in real-time daily or sports activities of daily living.

Acknowledgments. This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MEST) (NRF-2019R1A2C1003713). Edwin Valarezo Añazco gratefully acknowledges to “Escuela Superior Politécnica del Litoral (ESPOL)” and its excellence program “Walter Valdano Raffo”.

References

1. P. Kasnesis, C. Patrikakis and I. Venieris, "Changing the Game of Mobile Data Analysis with Deep Learning," *IT Professional*, no. 99. (June, 2017).
2. A. Avci, S. Bosch, M. Marin-Perianu, R. Marin-Perianu and P. Havinga, "Activity Recognition Using Inertial Sensing for Healthcare, Wellbeing and Sports Applications: A Survey," *Architecture of computing systems (ARCS), 2010 23rd international conference on*, pp. 1-10. (2010)
3. M. Kepski and B. Kwolek, "Embedded system for fall detection using body-worn accelerometer and depth sensor," *Proceedings of the 2015 IEEE 8th International Conference on Intelligent Data Acquisition and Advanced Computing Systems: Technology and Applications, IDAACS 2015*, vol. 2, pp. 755-759. (2015)
4. L. Montanini, A. Del Campo, D. Perla, S. Spinsante and E. Gambi, "A Footwear-Based Methodology for Fall Detection," *IEEE Sensors Journal*, vol. 18, no. 3, pp. 1233-1242. (2018)
5. B.-S. Lin, P.-C. Hsiao, S.-Y. Yang, C.-S. Su and I.-J. Lee, "Data Glove System Embedded With Inertial Measurement Units for Hand Function Evaluation in Stroke Patients," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 25, no. 11, pp. 2204-2213. (2017)
6. D. J. Walker, P. S. Heslop, C. J. Plummer, T. Essex and S. Chandler, "A continuous patient activity monitor: Validation and relation to disability," *Physiological Measurement*, vol. 18, pp. 49-59. (1997)
7. J. Margarito, R. Helaoui, A. M. Bianchi, F. Sartor and A. G. Bonomi, "User-independent recognition of sports activities from a single wrist-worn accelerometer: A template-matching-based approach," *IEEE Transactions on Biomedical Engineering*, vol. 63, pp. 788-796. (2016)
8. A. Anand, M. Sharma, R. Srivastava, L. Kaligounder and D. Prakash, "Wearable Motion Sensor Based Analysis of Swing Sports," *2017 16th IEEE International Conference on Machine Learning and Applications (ICMLA)*, pp. 261-267. (2017)
9. H. Ghasemzadeh and R. Jafari, "Coordination Analysis of Human Movements With Body Sensor Networks: A Signal Processing Model to Evaluate Baseball Swings," *IEEE Sensors Journal*, vol. 11, no. 3, pp. 603-610. (2011)
10. Z. Zhang, D. Xu, Z. Zhou, J. Mai, Z. He and Q. Wang, "IMU-Based Underwater Sensing System for Swimming Stroke Classification and Motion Analysis," *IEEE International Conference on Cyborg and Bionic Systems (CBS)*, pp. 268-272. (2017)
11. D. Ravi, B. Lo and G. Z. Yang, "Real-time food intake classification and energy expenditure estimation on a mobile device," *2015 IEEE 12th International Conference on Wearable and Implantable Body Sensor Networks, BSN 2015*. (2015)
12. N. R. Singh, "Implementation of Safety Alert System for Elderly People Using Multi-Sensors," *International conference of Electronics, Communication and Aerospace Technology (ICECA)*, vol. 1, pp. 282-286. (2017)
13. T. Elakkiya, "Wearable Safety Wristband Device For Elderly Health Monitoring With Fall Detect And Heart Attack Alarm," *Third International Conference On Science Technology Engineering and Management (ICONSTEM)*, pp. 1018-1022. (2017)
14. S. Mohammed and I. Tashev, "Unsupervised Deep Representation Learning to Remove Motion Artifacts in Free-mode Body Sensor Networks," *2017 IEEE 14th International Conference on Wearable and Implantable Body Sensor Networks, BSN 2017*, pp. 183-188. (2017)
15. Y. H. Lai, F. Chen, S. S. Wang, X. Lu, Y. Tsao and C. H. Lee, "A Deep Denoising Autoencoder Approach to Improving the Intelligibility of Vocoded Speech in Cochlear Implant Simulation," *IEEE Transactions on Bio-medical Engineering*, vol. 64, pp. 1568-1578. (2017)

16. D. T. Grozdic and S. T. Jovicic, "Whispered Speech Recognition Using Deep Denoising Autoencoder and Inverse Filtering," *IEEE/ACM Transactions on Audio Speech and Language Processing*, vol. 25, pp. 2313-2322. (2017)
17. J. Heymann, R. Haeb-Umbach, P. Golik and R. Schluter, "Unsupervised Adaptation of a Denoising Autoencoder by Bayesian Feature Enhancement for Reverberant ASR Under Mismatch Conditions," *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, pp. 5053-5057. (2015)
18. P. Xiong, H. Wang, M. Liu, S. Zhou, Z. Hou and X. Liu, "ECG Signal Enhancement Based on Improved Denoising Auto-encoder," *Engineering Applications of Artificial Intelligence*, vol. 52, pp. 194-202. (2016)
19. N. D. Lane and P. Georgiev, "Can Deep Learning Revolutionize Mobile Sensing?," *Proceedings of the 16th International Workshop on Mobile Computing Systems and Applications - HotMobile '15*, pp. 117-122. (2015)
20. I. Roychowdhury, J. Saha and C. Chowdhury, "Detailed Activity Recognition with Smartphones," *2018 Fifth International Conference on Emerging Applications of Information Technology (EAIT)*, pp. 1-4. (2018)
21. A. Vaughn, P. Biocco, Y. Liu and M. Anwar, "Activity Detection and Analysis Using Smartphone Sensors," *2018 IEEE International Conference on Information Reuse and Integration for Data Science*, pp. 102-107. (2018)
22. M. Jongprasithporn, N. Yodpijit and R. Srivilai, "A smartphone-based real-time simple activity recognition - IEEE Xplore Document," *3rd International Conference on Control, Automation and Robotics*, pp. 539-542. (2017)
23. J. J. Guiry, P. Ven, J. Nelson, L. Warmerdam and H. Riper, "Activity recognition with smartphone support," *Medical Engineering and Physics*, vol. 36, pp. 670-675. (2014)
24. E. Bulbul, A. Centin and I. A. Dogru, "Human Activity Recognition Using Smartphones," *2018 2nd International Symposium on Multidisciplinary Studies and Innovative Technologies (ISMSIT)*, pp. 1-6. (2018)
25. J. B. Yang, M. N. Nguyen, P. P. San, X. L. Li and K. Shonali, "Deep Convolutional Neural Networks On Multichannel Time Series For Human Activity Recognition," *IJCAI'15 Proceedings of the 24th International Conference on Artificial Intelligence*, pp. 3995-4001. (2015)
26. E. Valarezo, P. Rivera, J. M. Park, G. Gi, T. Y. Kim and T. Kim, "Human Activity Recognition Using a Single Wrist IMU Sensor via Deep Learning Convolutional and Recurrent Neural Nets," *UNIKOM Journal of ICT, Design, Engineering and Technological Science*, vol. 1, pp. 1-5. (2017)
27. D. Tao, Y. Wen and R. Hong, "Multicolumn Bidirectional Long Short-Term Memory for Mobile Devices-Based Human Activity Recognition," *IEEE Internet of Things Journal*, vol. 3, no. 6, pp. 1124 - 1134. (2016)
28. F. J. Ordonez and D. Roggen, "Deep convolutional and LSTM recurrent neural networks for multimodal wearable activity recognition," *Sensors (Switzerland)*, vol. 16, no. 115. (July, 2016)
29. S. Basterrech and V. K. Ojha, "Temporal Learning Using Echo State Network for Human Activity Recognition," *2016 Third European Network Intelligence Conference (ENIC)*, pp. 217 - 223. (2016)
30. A. Reiss and D. Stricker, "Creating and benchmarking a new dataset for physical activity monitoring," *Proceedings of the 5th International Conference on Pervasive Technologies Related to Assistive Environments - PETRA '12*, vol. 40, pp. 1-8. (June, 2012)
31. A. Reiss and D. Stricker, "Introducing a new benchmarked dataset for activity monitoring," *Proceedings - International Symposium on Wearable Computers, ISWC*, pp. 108-109. (2012)
32. X. Zhang, H. Dou, T. Ju, J. Xu and S. Zhang, "Fusing Heterogeneous Features From Stacked Sparse Autoencoder for Histopathological Image Analysis," *IEEE JOURNAL OF BIOMEDICAL AND HEALTH INFORMATICS*, vol. 20, no. 5, pp. 1377-1383. (2016)

33. H. Li and S. Misra, "Prediction of Subsurface NMR T2 Distributions in a Shale Petroleum System Using Variational Autoencoder-Based Neural Networks," *IEEE Geoscience and Remote Sensing Letters*, vol. 14, pp. 2395-2397. (2017)
34. D. Park, Y. Hoshi and C. C. Kemp, "A Multimodal Anomaly Detector for Robot-Assisted Feeding Using an LSTM-based Variational Autoencoder," *IEEE Robotics and Automation Letters*, vol. 3, pp. 1544-1551. (2017)
35. S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, vol. 9, pp. 1735-1780. (1997)
36. M. Milenkoski, k. Trivodaliev, S. Kalajdziski, M. Jovanov and B. R. Stojkoska, "Real Time Human Activity Recognition on Smartphones using LSTM Networks," *2018 41st International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO)*, pp. 1126-1131. (2018)
37. Deeplearning4j Development Team, "Deeplearning4j: Open-source distributed deep learning for the JVM, Apache Software Foundation 2.0," <http://deeplearning4j.org>. (2017)
38. A. Khan, N. Hammerla, S. Mellor and T. Plotz, "Optimising sampling rates for accelerometer-based human activity recognition," *Pattern Recognition Letters*, vol. 73, pp. 33-40. (2016)
39. V. T. Hees, L. Gorzelniak, E. C. Dean Leon, M. Eder, M. Pias, S. Taherian, U. Ekelund, F. Renstrom, P. W. Franks, A. Horsch and S. Brage, "Separating Movement and Gravity Components in an Acceleration Signal and Implications for the Assessment of Human Daily Physical Activity," *PLoS ONE*, vol. 8, no. 4. (April, 2013)
40. C. Beleites, R. Salzer, and V. Sergo, "Validation of soft classification models using partial class memberships: An extended concept of sensitivity & co. applied to grading of astrocytoma tissues," *Chemometrics and Intelligent Laboratory Systems*, vol. 122, pp. 12 – 22. (2013)
41. N. Hammerla, S. Halloran and T. Plotz, "Deep, Convolutional, and Recurrent Models for Human Activity Recognition using Wearables," *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence (IJCAI-16)*, pp. 1533-1540. (2016)
42. M. Arif and A. Kattan, "Physical Activities Monitoring Using Wearable Acceleration Sensors Attached to the Body," *PLoS ONE*, vol. 10, no. 7. (2015)
43. B. Almaslukh, J. AlMuhtadi and A. Artoli, "An Effective Deep Autoencoder Approach for Online Smartphone-based Human Activity Recognition," *International Journal of Computer Science and Network Security*, vol. 17, no. 4, pp. 160-165. (2017)

Edwin Valarezo Añazco received his B.E. degree in Electronics and Telecommunication Engineering from the Escuela Superior Politécnica del Litoral (ESPOL), Ecuador. He is currently working toward his Ph.D. degree in the Department of Biomedical Engineering at the Bio-Imaging & Brain Engineering Laboratory at the Kyung Hee University, South Korea. His research interest includes Deep Learning, Human Activity Recognition (HAR), real-time applications for HAR, hand gesture recognition, computer vision, reinforcement learning, and autonomous control of anthropomorphic robotic hands.

Patricio Rivera Lopez received his B.E. degree in Electronics, Automation and Control Engineering from the University of the Armed-Forces-ESPE, Ecuador. He is currently working toward his Ph.D. degree in the Department of Biomedical Engineering at Kyung Hee University, Republic of Korea. His research interest includes artificial intelligence, signal processing, and machine learning.

Hyemin Park received her B.S. and M.S. degrees in Biomedical Engineering from Kyung Hee University, South Korea. Her research interests include machine learning such, deep learning, image processing, computer vision, Human Activity Recognition (HAR), and hand gesture recognition.

Nahyeon Park received her B.S. degree in Biomedical Engineering from Kyung Hee University, South Korea. She is currently working toward her M.S. degree in the Department of Biomedical Engineering at Kyung Hee University, Republic of Korea. Her research interests include machine learning such as deep learning, image processing, computer vision, and reinforcement learning with robotic arms.

Tae-Seong Kim received the B.S. degree in Biomedical Engineering from the University of Southern California (USC) in 1991, M.S. degrees in Biomedical and Electrical Engineering from USC in 1993 and 1998 respectively, and Ph.D. in Biomedical Engineering from USC in 1999. After his postdoctoral work in Cognitive Sciences at the University of California at Irvine in 2000, he joined the Alfred E. Mann Institute for Biomedical Engineering and Dept. of Biomedical Engineering at USC as Research Scientist and Research Assistant Professor. In 2004, he moved to Kyung Hee University in Korea where he is currently Professor in the Department of Biomedical Engineering. His research interests have spanned various areas of biomedical imaging, bioelectromagnetism, neural engineering, assistive biomedical lifecare technologies. Dr. Kim has been developing advanced signal and image processing methods, pattern classification, machine learning methods, novel medical imaging modalities, and rehabilitation technologies. Dr. Kim has published more than 300 papers and seven international book chapters. He holds ten international and domestic patents and has received nine best paper awards.

Received: September 20, 2019; Accepted: May 10, 2020

Production of linked government datasets using enhanced LIRE architecture

Nataša Veljković¹, Petar Milić², Leonid Stoimenov¹, and
Kristijan Kuk³

¹ Faculty of Electronic Engineering, Aleksandra Medvedeva 14,
18000 Niš, Serbia

{nataša.veljkovic, leonid.stoimenov}@elfak.ni.ac.rs

² Faculty of Technical Sciences, Knjaza Miloša 7,

38220 Kosovska Mitrovica, Serbia

petar.milic@pr.ac.rs

³ University of Criminal Investigation and Police Studies, Cara Dušana 196,
11000 Belgrade, Serbia

kristijan.kuk@kpu.edu.rs

Abstract. This paper describes the enhanced LIRE (LInked RElations) architecture for creating relations between datasets available on open government portals. The architecture is improved to be applicable on different open government data platforms using minimal configuration at the Data processing layer. We evaluated the applicability of enhanced LIRE, its advantages and disadvantages, which resulted in necessary recommendations for the publication of dataset's metadata to obtain better associations between datasets. Moreover, we introduced a LINDAT indicator that reflects the percentage of linked data in the total possible number of linked data on open government data portals.

Keywords: open government data, linked data, datasets, metadata, dataset relations.

1. Introduction

WEB of data represents related structured data connected through links [1]. Offering open government data (OGD) in the Web of data will ease their discovery and usage. That implies interlinking already published data, which can be a tedious and expensive task for governments, but we argue that benefits prevail. Government moves to the higher degree of openness implicitly by progressing in the domain of open data [2]. Not to forget the final users, who appreciate better accessibility, reusability and easy processing of open data [3].

The real value of OGD is revealed with linking which provides unexpected and unexplored insights into different domains and problem areas [4]. OGD interconnected with Web of Data becomes Linked Open Government Data (LOGD). Processing dataset's metadata to create links makes datasets more comprehensible and leads to uncovering potential faults in metadata definition. Such process contributes to the quality of datasets thus raising the openness and transparency in government.

Government itself is a complex engine running on highly distributed set of institutions. Linking data within and between those institutions enables government to publish data in a modular way, benefiting from a ‘small pieces of loosely joined’ approach to government data [5]. However creating relations between OGD from the government perspective is a costly process, it requires trained and equipped staff for OGD processing, which governments do not have. In that sense, sharing insights on dataset relations between government institutions could help to avoid duplicate work and efforts [6].

User behavior while searching for OGD may help in revealing practical approaches and patterns for linking OGD. OGD users usually browse multiple datasets, while doing so they discover new datasets that may be related to the ones they have already checked. If open data portals provided a way for users to propose relations between datasets, this would create more linked data on the portals without much effort from the government side.

When users build applications with raw OGD they need to inspect government datasets to see if they correspond to their needs. Keeping this in mind, we can ask ourselves: What can be helpful in driving the users toward the dataset they are interested into? We believe the answer is dataset’s metadata. Dataset’s metadata, such as description, format, tags and etc. enable users to inspect if a dataset contains information they need. Metadata helps in generating linkable dataset profiles which is important when relating datasets [7].

This paper analyses how dataset’s metadata can be used for the production and utilization of linked government data. Using the enhanced LIRE (Linked Relations) architecture we demonstrate the process of linking datasets based on the metadata keys and their values. The architecture is flexible and generalized for application on different OGD portals. Moreover, we are contributing to the linked datasets representation on OGD portals by proposing a model for semantic representation of dataset’s relations. Based on the user involvement in the process of consumption of OGD and creating useful applications that depend on OGD, we got an idea to include users in the process of relating datasets. Our approach goes towards the interlinking and integration of OGD. As a proof of concept, we have developed a prototype that enables the users of OGD portals to create linked datasets.

Throughout this article we explained the enhanced LIRE architecture, processes which occur inside the architecture, the specification of data necessary to relate datasets and the semantics of relations. We introduced and evaluated a new measure for reflecting the status of the linking of datasets on OGD portal named LINDAT. Furthermore, we analyzed applicability of the architecture, its advantages and disadvantages, and gave necessary recommendations for publication of dataset’s metadata to obtain better assessment of dataset which we want to link.

The rest of the paper is organized as follows. In Section 2 we analyze current approaches in the area of linking OGD, focusing on those that require metadata utilization or user involvement. Section 3 describes the enhanced LIRE architecture for linking government datasets based on metadata structure. In Section 4 we analyze LINDAT indicator and apply it on CKAN powered OGD platforms to check if open government data portals have linked data. The final section concludes the work done so far and presents ideas for the future work.

2. Related Work

In the following subsections, we present available methods and tools for linking open data in the government domain (LOGD) and differentiation of those approaches to our proposal. We talk about the structure of linked datasets (linksets) on OGD portals, and how our approach contributes to that. To lower the costs associated to datasets exploration and linking, users are often involved in the process of linksets creation, and it was interesting to give an overview of tools which assist in the process and talk about our contribution in this domain.

2.1. Linking government datasets

When publishing OGD, data producers strive to fulfil openness criteria, but they are not so concerned with having the linked criteria fulfilled [3]. LOGD is the practice initially adopted by researchers and third parties who have reused existing open data to create linked open data, by exploring datasets, creating RDFs and publishing them on the Web as new data. Nowadays it happens more often that data producers are interested in publishing linked data as well. An interesting analysis of provision of LOGD data is given by Kalampokis et al. [8]. They show that linked data can be provided in two ways: by publishing on the central open government portal, or on the public agency portal. In the second case, central government portal collects metadata about LOGD datasets published by a public agency in order to make that datasets available on its portal. That is the so-called 'indirect provision of linked data'. This claim is in line with [9] who proposed the architecture for integrating datasets from public agencies via activities and components essential for discovery. Their architecture has two phases, in the first phase, work on downloading and transforming datasets into a common schema language format is conducted, while the second phase addresses semantic heterogeneity with schema matching and statistical analysis of ontology structures. This paper does not take into account the metadata of open datasets, but deals solely with their semantic versions. In relation to our solution, which we will describe in this paper and it is based on the application of metadata, the authors do not extract useful information from "raw" open datasets that can serve to relate them.

On the other hand, some authors claim that successfully linking government datasets requires understanding the data context's sensitive meaning [10]. This is coupled with enabling the semantic interoperability of OGD and provision of metadata that describes them, which is important for creation of value added applications. It is important that there is a standardized level of metadata, because that will allow the harmonization of specific concepts and terminology, interoperability and multilinguality in government open data portals. For that purpose, Assaf, Troncy and Senart [11] have identified a need for a definition of a harmonized model of OGD dataset's metadata which contains sufficient information so that consumers can easily understand and process datasets. That information contains general information, access information, ownership information, provenance information, geospatial information, temporal information, statistical information and quality information with mappings between appropriate data.

Similarly, Kashyap and Sheth's work on the semantic heterogeneity is based on representation of metadata, context and ontologies. They claim that for the proper

linking of OGD, underlying data and capturing of the meaning of domain specific metadata must be considered. This is due to the information overload which arises as a consequence of the heterogeneity of the digital data [12]. In relation to our work, this approach collects different types of metadata independent of type, representation and location, while we utilize metadata from OGD which is not previously curated.

Schmachtenberg et al. [13] presented an overview of relationships between linked datasets in the form of linked open data cloud diagram. Analyzing that, we can notice that two datasets can be linked if there exists at least one RDF link between resources belonging to the datasets. The authors emphasize the use of dataset's metadata because it reveals the origin of datasets and can be used to analyze dataset's quality.

For the integration of OGD datasets into the Web of Data, a group of authors [14] introduce a solution of six stages of dataset integration. Five of these stages named Name, Retrieve, Adjust, Convert, Enhance and Publish are designed to implement an approach for minimization of human effort to incorporate new dataset as linked data, while in the remaining stage the structured and connected descriptions of the initial representations of datasets are added by data modelers. With this approach data structural relations are covered, which supports integration of LOGD from multiple sources. Cverdelj-Fogaraši et al. provides alignment of domain specific metadata ontologies, without modelling relations between OGD datasets [15].

In the work of Scharffe et al. [16] on methods for automated dataset interlinking, different data linking systems are discussed. The authors showed that with an appropriate mapping between dataset's descriptions proper links can be established by denoting the single correspondence between object descriptions. This claim is confirmed by the research of Ellefi et al. [17] who introduced an approach for linking datasets by overlapping dataset's metadata schema definition. For a given dataset, their framework allows identifying the datasets sharing schema with other datasets, which is a useful input for the data linking step. Furthermore, Ngomo and Sherif offer a solution for link discovery between different datasets [18]. Their solution addresses time-efficiency and accuracy challenges which is of central importance when a tool is faced with small amounts of RAM or when is faced with streaming or complex data (e.g., 5D geospatial data). As data alone has little value, to unleash its full potential, it needs to be linked with other referenced data. Datalift tool provided by [19] searches for relationships between local government data and existing public RDF data and enriches that data with links, making them discoverable on the Web of data. Compared to the approach we describe in this paper, the authors integrate ready-made semantic versions of OGD datasets with the rest of the Web of data, while our solution aims to integrate OGD datasets across portals first and then integrate them into the Web of data.

2.2. Representation of linked data on OGD portals

The full potential of OGD has not been realized, and one of the main reasons for that is related to the provision of metadata. Metadata provide documentation, context and necessary background information for interpreting OGD [20], and as such it is often considered as the key enabler for the effective use of linked open data in government domain [21]. To enable the interoperability of OGD and their linking within the Web of Data, they need to be modelled in a semantic way which ensures that data will be

reached by linked data applications. In this respect, short summarization of most used dataset vocabularies to describe linked datasets follows.

The analysis of dataset's metadata structure, consistency and availability conducted by [22], have resulted in development of dcat (data catalog) vocabulary that allows the expression of datasets in the RDF data model. Dcat enables datasets to be queried in RDF, supports the reuse and extension of existing metadata standards such as Dublin Core [23] and compatibility with linked data. It avoids the use of ontologies, because the main purpose of the vocabulary is interoperable data exchange. According to the authors, the goal of this vocabulary is to increase dataset's discoverability enabling applications to easily consume them.

For describing public sector datasets in Europe, a dcat-ap (data catalog - application profile) vocabulary is designed [24]. This vocabulary is intended to represent a selected set of properties from those in dcat and its imported vocabularies, also enables cross-data portal searching and enhances discoverability. It is recommended by the Open Data Support to be the standard for describing linked datasets in Europe.

The VoID vocabulary differentiates linked data publishers, persons or organizations exposing linked data and on the other hand linked data consumers, which may be humans or machines [25] with defining "appropriateness" by following criteria: the content of the datasets, interlinking to other datasets and vocabularies used in datasets. [26]. The designed term 'linkset' deals with interlinking between datasets, especially those that are published on the same portal. With RDF properties: 'void:subset', 'void:target' and 'void:linkPredicate', VoID gives opportunity for dataset interlinking. By modelling dataset's links with VoID we actually connect one or more topics that are originating from a certain source or process and that are accessible on the Web. Fiorelli et al. [27] developed LIME (Linguistic Metadata), an extension of VoID with a vocabulary of metadata about the ontology-lexicon interface. LIME provides the lexicalization of VoID relations in a natural language aligning it with a set of lexical concepts. With LIME, dataset relations intend to be more discoverable, understandable and exploitable.

All of the presented vocabularies describe linked datasets and how they can be accessed through the linked data applications. Nevertheless, the question remains: what about relations between OGD datasets on portals? How to utilize them and reach more related data? In this regard, our solution proposes which elements of vocabularies should be exploited in order to resolve this issues and to make OGD dataset relations available in linked data applications.

2.3. User involvement in linking datasets

As pioneers in creating LOGD, UK and USA government initiatives of linking data have showed that availability of LOGD is costly, and that there is a need for solving this problem or at least to mitigate it. One of the means to reduce the costs is to allow datasets consumers to participate voluntarily in linking OGD. Self-service approach introduced by [28] shifts the burden of interlinking datasets to the data consumers. It encourages them to interlink government datasets without waiting for the government to do so. With the application of Google Refine tool [6] users interlink datasets to the Web of Data space, previously providing Google Refine with SPARQL endpoint or dump file for reconciliation of this datasets against any RDF data available through it. Similar to

this, Li Ding et al. describe linked data ecosystem in which users manage and consume LOGD in connection with online tools, services and societies [3]. LOGD ecosystem is based on converting raw OGD datasets into linked data and their integration with other resources. The work of Li Ding et al. presumes the existence of linked data from OGD data as a base for their further linkage, but does not consider the possibility for the integration of the datasets on the same platform. Our work goes toward this, including users for achieving this goal due to the fact that successfully linking of datasets requires understanding the data's context sensitive meaning [10].

User feedback and application queries can be used to determine whether two datasets can be interlinked [29]. Application queries help filter datasets that are potentially strong candidates for interlinking, whereas user feedback is used as a way to assess the relevance of the candidate datasets. Furthermore, [30] argued that the visual exploration of dataset content in linked open data exploitation can help in identifying links between datasets. This approach is similar to the one we will present in the following sections, but does not take into account the fact that datasets housed on the same OGD platform can also be related mutually. Specifically, it deals with the semantics of linked data without analyzing the attributes of the datasets and getting users attention to the interaction and information use on OGD platform.

To the best of our knowledge, there is no similar work to what we propose, that deals with the production of interlinked government datasets based on its metadata. Some authors [8, 10, 13, 16, 20] discuss linking datasets based on their semantic description without going deeper into the metadata. That does not tackle information hidden in published OGD, which by our opinion can help in linking OGD datasets. While the semantic version of OGD can be linked to the Web of Data space, the semantic interlinking of OGD on same portals remains unexplored. Leme et al. [31] utilize existing dataset relations on open government portals to produce recommendation for dataset linking, but do not take into account metadata describing the dataset such as description, tags, formats, organization and etc., in order to search for possible hidden information that can assist in their connection and consequently their interlinking. Comparing to the [32], the version of architecture that will be presented in following section contains improved model for determining the type of relations between two datasets, applicability on different types of OGD platforms which consequently means that it is interoperable. Moreover, we define and describe a novel measure here, LINDAT indicator, for the purpose of monitoring the linking status of datasets.

3. Enhanced LIRE architecture

Based on the aforementioned approaches for linking OGD, we got the idea to explore dataset's metadata to check whether they can be used to relate to other OGD, and then

to model these relations semantically in order to produce LOGD. For that purpose, we have developed the architecture for managing relations between datasets, and their linking based on user interaction with OGD portal. LIRE architecture enables the interlinking of datasets [32], and can enrich OGD portals with novel functionality. Our architecture assumes the existence of the semantic version of OGD datasets, as the availability of such OGD datasets is embedded in most OGD platforms, to mention few: CKAN¹, DKAN², Opendatasoft³ and Socrata⁴. The prototype exists in a form of a plugin, currently available only for the CKAN platform⁵. The enhanced LIRE architecture is outlined in Fig. 1.

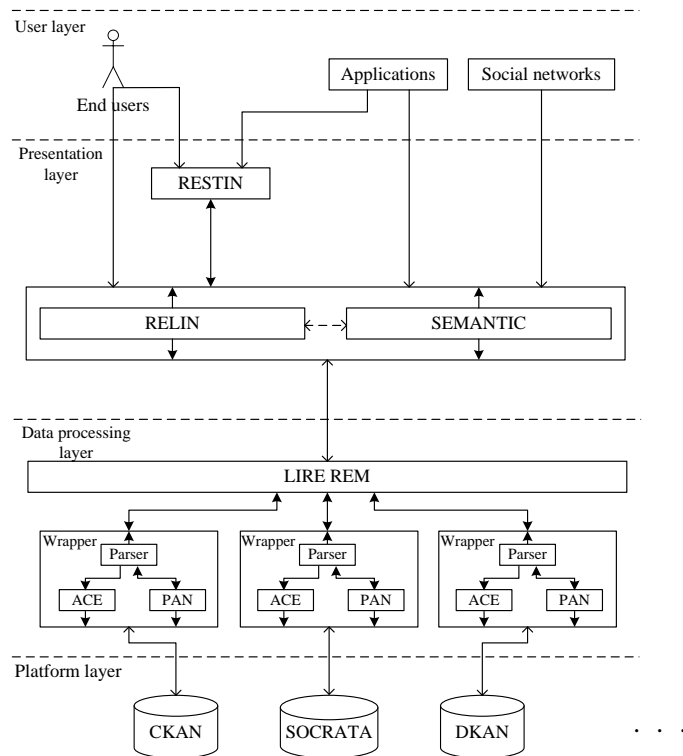


Fig. 1. Enhanced Linked Relations Architecture

¹ <http://docs.ckan.org/en/latest/maintaining/linked-data-and-rdf>
² <https://docs.getdkan.com/en/latest/introduction/catalog-basics>
³ <https://docs.opendatasoft.com/api/explore/v2.html>
⁴ <https://dev.socrata.com/docs/formats/>
⁵ <https://extensions.ckan.org/extension/lire/>

The improvements are reflected in the possibility to apply LIRE on different OGD platforms and in the more precisely defined internal processes for linking datasets. We have grouped some components in order to enable the interoperability of LIRE with different OGD portals.

The enhanced LIRE architecture contains different components (modules) that deal with specific tasks in order to ensure data transformation, data interlinking and creation of dataset's relations:

- WRAPPER collects data from OGD portals. It contains three submodules:
 - *PARSER* deals with requests from REM component. It parses requests and issues relevant commands to PAN or ACE. It works vice versa as well by collecting responses from PAN or ACE and returning the appropriate responses to REM. The specification of information that PARSER sends/receives is given in Section 3.1.
 - *PAN – Portal Analyzer* is a submodule that creates requests to OGD portals to retrieve datasets. Then it checks specific dataset's metadata fields, according to specification received from the PARSER and returns gathered values.
 - *ACE – Action Executor* module executes actions for creating, updating or deleting dataset relations according to PARSER specifications.
- REM (Relations Manager) is responsible for determining the type of relations between datasets and creating and managing dataset's relations. REM examines dataset's metadata received through the PARSER and checks for similarity between datasets for the possible creation of relation. Within this architecture component, a model for relation of OGD is implemented. Given that the proposed model for relation OGD uses the types of relationships that are mutually related, this architecture component does not impose restrictions on the need to use these types of relationships, but leaves the final choice of the type of relationship to the user. The user can choose another type of relationship based on the experience gained in working with OGD on OGD platforms. REM can receive requests and send responses to the components contained in the upper and lower layers of the architecture. REM needs to identify WRAPPER in order to know where to send the applicable requirements and this is done by registering WRAPPER within REM. Moreover, REM creates the specification of request which is sent to the WRAPPER, in order to execute proper actions on OGD platform.
- RELIN (Relations Information) prepares and creates necessary data for the visual preview of datasets. RELIN incorporates specific libraries for visualizing dataset relations and their graphical management. Every dataset is represented with a graphical element that contains information from dataset metadata, for example the description, tag, formats and others. Also, RELIN triggers REM by forwarding him a request for dataset's metadata examination in order to obtain the possible type of relation between datasets.
- SEMANTICS (Semantics of Relations) component performs dataset interlinking by updating the semantic version of OGD datasets i.e. adding links about relations between datasets. Dataset interlinking is modelled by using an RDF graph, described in detail in Section 3.4. The implemented RDF graph model is based on VoID vocabulary, where we exploit linkset feature, because it is naturally intended for modelling of links between datasets.

- *RESTIN (Relations Statistics and Indicators)* offers statistic about linked datasets, such as user created linksets per total linksets on OGD portal, government created linksets per total linksets on OGD portal, linksets per OGD portal datasets, where they all together constitute LINDAT, linked datasets indicator as the indicator of the status of dataset interlinking on a government portal.

3.1. Internal metadata specification

Information provision between REM and WRAPPER occurs by exchanging XML documents. REM issues request to the WRAPPER by sending the XML specification of the information it needs. XML schema specification is given in Fig.2.

```

1 <xs:schema targetNamespace="https://www.w3schools.com" xmlns:xs="http://www.w3.org/2001/XMLSchema">
2   <xs:element name="LireRequest">
3     <xs:complexType>
4       <xs:sequence>
5         <xs:element name="Wrapper">
6           <xs:complexType>
7             <xs:simpleContent>
8               <xs:extension base="xs:string">
9                 <xs:attribute type="xs:string" name="location"/>
10              </xs:extension>
11            </xs:simpleContent>
12          </xs:complexType>
13        </xs:element>
14        <xs:element type="xs:string" name="DatasetSubject"/>
15        <xs:element type="xs:string" name="DatasetObject"/>
16        <xs:element name="TagsMeasure">
17          <xs:complexType>
18            <xs:simpleContent>
19              <xs:extension base="xs:string">
20                <xs:attribute type="xs:string" name="type"/>
21              </xs:extension>
22            </xs:simpleContent>
23          </xs:complexType>
24        </xs:element>
25        <xs:element name="OwnerMeasure">
26          <xs:complexType>
27            <xs:simpleContent>
28              <xs:extension base="xs:string">
29                <xs:attribute type="xs:string" name="type"/>
30              </xs:extension>
31            </xs:simpleContent>
32          </xs:complexType>
33        </xs:element>
34        <xs:element name="GroupMeasure">
35          <xs:complexType>
36            <xs:simpleContent>
37              <xs:extension base="xs:string">
38                <xs:attribute type="xs:string" name="type"/>
39              </xs:extension>
40            </xs:simpleContent>
41          </xs:complexType>
42        </xs:element>
43        <xs:element name="OwnerTagsMeasure">
44          <xs:complexType>
45            <xs:simpleContent>
46              <xs:extension base="xs:string">
47                <xs:attribute type="xs:string" name="type"/>
48              </xs:extension>
49            </xs:simpleContent>
50          </xs:complexType>
51        </xs:element>
52        <xs:element name="GroupTagsMeasure">
53          <xs:complexType>
54            <xs:simpleContent>
55              <xs:extension base="xs:string">
56                <xs:attribute type="xs:string" name="type"/>
57              </xs:extension>
58            </xs:simpleContent>
59          </xs:complexType>
60        </xs:element>
61        <xs:element name="FormatsMeasure">
62          <xs:complexType>
63            <xs:simpleContent>
64              <xs:extension base="xs:string">
65                <xs:attribute type="xs:string" name="type"/>
66              </xs:extension>
67            </xs:simpleContent>
68          </xs:complexType>
69        </xs:element>
70        <xs:element name="CreationDateMeasure">
71          <xs:complexType>
72            <xs:simpleContent>
73              <xs:extension base="xs:string">
74                <xs:attribute type="xs:string" name="type"/>
75              </xs:extension>
76            </xs:simpleContent>
77          </xs:complexType>
78        </xs:element>
79        <xs:element name="DescriptionMeasure">
80          <xs:complexType>
81            <xs:simpleContent>
82              <xs:extension base="xs:string">
83                <xs:attribute type="xs:string" name="type"/>
84              </xs:extension>
85            </xs:simpleContent>
86          </xs:complexType>
87        </xs:element>
88        <xs:element name="TrackingCountTotalViewsMeasure">
89          <xs:complexType>
90            <xs:simpleContent>
91              <xs:extension base="xs:string">
92                <xs:attribute type="xs:string" name="type"/>
93              </xs:extension>
94            </xs:simpleContent>
95          </xs:complexType>
96        </xs:element>
97        <xs:element name="TrackingCountRecentViewsMeasure">
98          <xs:complexType>
99            <xs:simpleContent>
100             <xs:extension base="xs:string">
101               <xs:attribute type="xs:string" name="type"/>
102             </xs:extension>
103           </xs:simpleContent>
104         </xs:complexType>
105       </xs:sequence>
106     </xs:complexType>
107   </xs:element>
108 </xs:schema>

```

Fig. 2. XML schema for REM’s request to WRAPPER

On line 5, REM specifies which WRAPPER will be used, since there can be multiple WRAPPER components for different OGD portals. Further on, lines 14 and 15 REM specify two datasets, the so called subject and object of relation. Lines 16 – 142 specify different schema elements suffixed with *Measure* that will be used for creating measure for linking datasets.

Depicted in Fig. 3., we can see the XML schema for the WRAPPER’s response. Both *DatasetSubject* and *DatasetObject* elements are complex and contain a set of child elements. Lines 7 – 49 contain necessary child elements for DatasetSubject that correspond to requested Measure elements. Same elements must be present for DatasetObject as well (lines from 50 upwards).

```

1 <xs:schema targetNamespace="https://www.w3schools.com" xmlns:xs="http://www.w3.org/2001/XMLSchema">
2 <xs:element name="LireResponse">
3 <xs:complexType>
4 <xs:sequence>
5 <xs:element name="DatasetSubject">
6 <xs:complexType>
7 <xs:sequence>
8 <xs:element name="Tags">
9 <xs:complexType>
10 <xs:sequence>
11 <xs:element type="xs:string" name="TagName" maxOccurs="unbounded" minOccurs="0"/>
12 </xs:sequence>
13 </xs:complexType>
14 </xs:element>
15 <xs:element type="xs:string" name="Owner"/>
16 <xs:element type="xs:string" name="Group"/>
17 <xs:element name="OwnerTags">
18 <xs:complexType>
19 <xs:sequence>
20 <xs:element type="xs:string" name="TagName" maxOccurs="unbounded" minOccurs="0"/>
21 </xs:sequence>
22 </xs:complexType>
23 </xs:element>
24 <xs:element name="GroupTags">
25 <xs:complexType>
26 <xs:sequence>
27 <xs:element type="xs:string" name="TagName"/>
28 </xs:sequence>
29 </xs:complexType>
30 </xs:element>
31 <xs:element name="Formats">
32 <xs:complexType>
33 <xs:sequence>
34 <xs:element type="xs:string" name="FormatName" maxOccurs="unbounded" minOccurs="0"/>
35 </xs:sequence>
36 </xs:complexType>
37 </xs:element>
38 <xs:element type="xs:dateTime" name="CreationDate"/>
39 <xs:element type="xs:string" name="Description"/>
40 <xs:element type="xs:integer" name="TrackingCountTotalViews"/>
41 <xs:element type="xs:integer" name="TrackingCountTotalRecent"/>
42 <xs:element type="xs:integer" name="Fivestar"/>
43 <xs:element type="xs:string" name="Open"/>
44 <xs:element type="xs:string" name="LinkedDataFormat"/>
45 <xs:element type="xs:string" name="MachineProcessable"/>
46 </xs:sequence>
47 <xs:attribute type="xs:string" name="name"/>
48 </xs:complexType>
49 </xs:element>
50 <xs:element name="DatasetObject">
51 ...
52 </xs:element>
53 </xs:sequence>
54 </xs:complexType>
55 </xs:element>
56 </xs:schema>

```

Fig. 3. XML schema for WRAPPER's response to REM

3.2. Workflow: creating linksets

In Fig. 4 we illustrate the basic workflow in the proposed architecture for the case when a user wants to relate two datasets.

When users access an OGD portal via browsers, they make exploration and preview of available datasets and related information depending on their need. During that process, they may find out that some datasets can be related. LIRE architecture is designed to suggest a type of relation to a user based on the examination of dataset's metadata. This examination is taking place when a user tries to relate two datasets. For that purpose, RELIN triggers REM by forwarding a request for dataset's metadata examination. REM then creates the specification of necessary data, where it includes dataset names which will be examined and forwards the information to WRAPPER with a request to gather proper information from the OGD portal. In a few iterations which WRAPPER performs with the OGD portal, a set of information is obtained, which is returned to REM in an appropriate form. When REM receives the requested information, it performs the analysis of the obtained information, as per model

explained in [32], in order to make a suggestion for a possible type of relation for those datasets to a user. The user is given a choice, to apply the recommendation or chose another more appropriate type of relation. As the last iteration in the process, user commits the creation of relation between selected datasets.

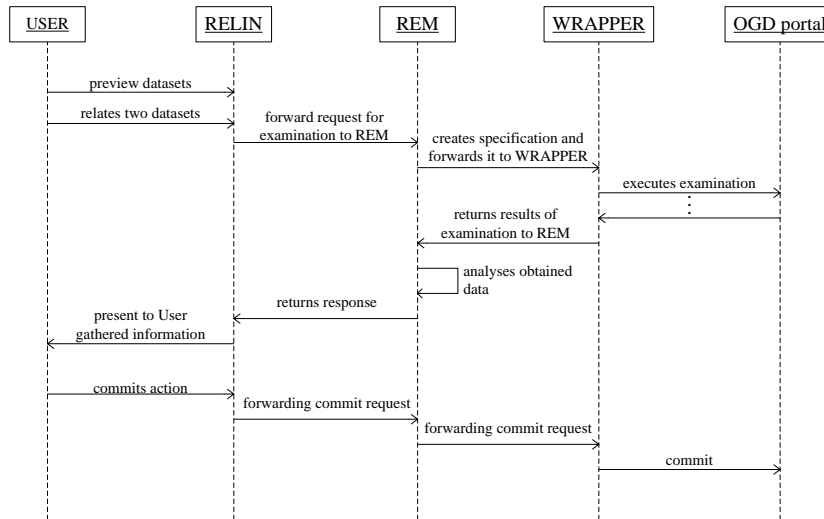


Fig. 4. Workflow for relating two datasets

When users access an OGD portal via browsers, they make exploration and preview of available datasets and related information depending on their need. During that process, they may find out that some datasets can be related. LIRE architecture is designed to suggest a type of relation to a user based on the examination of dataset’s metadata. This examination is taking place when a user tries to relate two datasets. For that purpose, RELIN triggers REM by forwarding a request for dataset’s metadata examination. REM then creates the specification of necessary data, where it includes dataset names which will be examined and forwards the information to WRAPPER with a request to gather proper information from the OGD portal. In a few iterations which WRAPPER performs with the OGD portal, a set of information is obtained, which is returned to REM in an appropriate form. When REM receives the requested information, it performs the analysis of the obtained information, as per model explained in [32], in order to make a suggestion for a possible type of relation for those datasets to a user. The user is given a choice, to apply the recommendation or chose another more appropriate type of relation. As the last iteration in the process, user commits the creation of relation between selected datasets.

The relation itself is added to the dataset’s metadata. Relation metadata keys can be one of the following:

- subject – the first dataset in relation
- object – the second dataset in relation
- type – the type of relation between two datasets
- label – a string that indicates whether the linkset was created by government or user

Relation metadata is nested in dataset's metadata structure under a "relationships" key. This key is initially enabled in CKAN while other OGD platforms require an extension or modification of their existing metadata models to support this (DKAN, Socrata, Opendatasoft) [21].

3.3. The semantics of related dataset

Accessing related OGD datasets by linked data applications goes toward extracting context and meaning of OGD, thus enabling their proper use for the retrieval of information in conjunction with other linked datasets and the Web of Data space. For that purpose, relations between OGD datasets need to be modelled semantically. That is achieved through the RDF description of dataset relations. As already most of OGD portals ensure semantic availability of their datasets [21], it will be only necessary to model dataset relations semantically.

Within the SEMANTICS component of enhanced LIRE architecture, we carry out the semantic modelling of dataset relations by using VoID vocabulary. We exploit the feature "linkset", which is a collection of RDF links, where an RDF triple has a subject and object described in different dataset [26].

In this way we are enabling access to linksets from semantic web applications, but also from OGD portals. By exploring linksets, users can find related data and see more on the topic they searched for.

3.4. Architecture applicability

Although it requires the customization of the WRAPPER component, the main strength of LIRE architecture lies in its application to different OGD portals. OGD portals may be developed in different technologies, and it's up to the portal manager/developer to setup the WRAPPER to communicate with the portal. This is a compromise they need to make for applying the LIRE. This problem can also be solved by encouraging the community of developers to create a custom WRAPPER for the OGD portal, and later on share their effort with others. In this way we will stop accumulating time and costs for this task.

As the application of our architecture requires availability of metadata fields for the examination of datasets and storing dataset relations, OGD portals need to make them available. We can see from Table 1, which gives metadata support on some OGD platforms, that only CKAN supports relations between datasets, while in the SOCRATA some minor modifications are needed. Other platforms don't support relations, because they offer metadata that describe resource data values, but not the metadata about relations. The incorporation of the metadata that describe relations into the dataset structure will enable the usage of our architecture and contribute to the better definition of datasets.

Table 1. Metadata Support in Open Data Platforms

Platform	Dataset Metadata	Resource Metadata	Custom Metadata	Storing Relations
CKAN	YES	YES	YES	YES
DKAN	YES	YES	YES	NO
SOCRATA	YES	YES	YES	NO
OGDI	NO	YES	NO	NO
OGPL	NO	YES	NO	NO
OPENDATASOFT	YES	YES	YES	NO
PLENARIO	NO	YES	NO	NO
JUNAR	NO	YES	NO	NO

The presence of relevant dataset information in the metadata structure of dataset is given in Table 2. We performed an analysis of 8 popular OGD portals powered by four different platforms: CKAN, DKAN, Socrata and Opendatasoft. This analysis aims to validate the model for proposing relations by examining whether the information of importance is presented in the dataset's metadata structure. Defined conditions, named C1-C13 [32], examine the following:

- **C1** - Number of same/similar tags between two datasets
- **C2** - Do they belong to the same organization
- **C3** - Do they belong to the same group
- **C4** - Whether the number of the same/similar tags of the first dataset is greater than (or less than) the number of the same/similar tags in the second dataset organization
- **C5** - Whether the number of the same/similar tags of the first dataset is greater than (or less than) the number of the same/similar tags in the second dataset group
- **C6** - Are they linked via links in extra field
- **C7** - Whether the number of the same/similar resource formats of the first dataset is greater than (or less than) the number of the same/similar resource formats in the second dataset
- **C8** - Whether the first dataset was created after the second
- **C9** - Whether the descriptions of two datasets are similar
- **C10** - Whether the number of total views of the first dataset is less than (or greater than) the number of total views of the second dataset
- **C11** - Whether the number of recent views of the first dataset is less than (or greater than) the number of recent views of the second dataset
- **C12** - Whether the five star index of the first dataset is less than (or greater than) the five star index of the second dataset
- **C13** - Whether they are open

Table 2 shows the percentage of defined data for the whole portal for each condition. It is important to notice that C6 has no values. This indicates that C6 is irrelevant for the determination of the type of relation and we have excluded this condition in the revised model. The remaining of the table, shows the presence of mostly all data. It is interesting to note that portals powered by CKAN have a lower presence of some data relative to other portals.

Table 2. Presence of Relevant Dataset Information in Metadata Structure

Portal	Platform													
		C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	C11	C12	C13
catalog.data.gov	CKAN	84	100	0	100	100	0	88	100	100	100	100	88	53
data.gov.uk	CKAN	28	100	0	100	100	0	9	100	100	0	100	9	9
offenedaten-koeln.de	DKAN	99	0	97	0	99	0	99	99	99	0	0	99	99
www3.unog.ch	DKAN	72	0	0	0	90	0	90	100	95	0	0	90	49
data.edmonton.ca	Socrata	97	100	100	100	0	0	100	99	100	100	100	100	100
data.oregon.gov	Socrata	73	100	98	100	0	0	95	100	79	100	100	95	100
opendurham.nc.gov	Opendatasoft	100	91	0	100	0	0	100	100	75	0	0	100	0
opendata.paris.fr	Opendatasoft	99	93	0	100	0	0	99	99	99	0	0	99	0

To improve dataset relations in OGD portals, we can in the future rely more on data consumers. Whoever browses datasets can be the one creating links between the datasets. We should not wait for the OGD portal's management team to link datasets [6]. This is aligned with civic-sourcing, a particular type of "crowd-sourcing" being adopted as a part of Government 2.0 to harness the wisdom of citizens [33]. As a full automatic approach does not always guarantee good results, the inclusion of consumers improves obtaining linked datasets, which imposes our solution as swift, less demanding and approachable by non-expert users.

4. LINDAT indicator

We introduced an indicator named LINDAT (LINKedDATaset), calculated by RESTIN component of the architecture, which indicates how much linked data is created vs. total possible linked data on the portal. LINDAT is calculated by dividing two measures: Created Linksets (CL) and Total Possible Linksets (TPL), as given in equation (1), and it is expressed in percentage. CL represents a sum of all created linksets on the OGD portal, both by OGD data publishers and data consumers. Before making the data available to stakeholders, government representatives perform internal data linking (GCL – Government Created Linksets). The process probably happens at data import, but can also be done at later stage if required. Data consumers can participate in the linking of open data (UCL – User Created Linksets), because their experience in using open government data is significant and should be taken into account.

$$\text{LINDAT} = \text{CL} / \text{TPL}. \quad (1)$$

Total Possible Linksets is calculated as given in the equation (2), where "n" is the number of datasets that reside on the OGD portal.

$$\text{TPL} = n * (n - 1). \quad (2)$$

LINDAT calculation excludes the possibility of an erroneous data, in a sense that some member of linkset can have missing links. This can happen while adding relations manually, but by using LIRE architecture it is impossible to make this mistake since relations are created automatically. Knowing that linksets are expressed through directed multigraphs, they are completely equivalent to their inverse ones. For example,

a linkset with the triple: A skos:broader B is semantically equivalent to a linkset with the triple: B skos:narrower A [34]. From here it becomes clear that each direction of relation between datasets has a linkset with its own inverse linkset.

4.1. Datasets evaluation using LINDAT

Evaluation of datasets on OGD portals using LINDAT, started with accessing the list of CKAN powered OGD portals around the world⁶. We used custom generated application to search all listed portals and find whether *relationships* metadata key is defined, which would reveal linked relations between datasets. The results of the evaluation are listed in Table 3. Unfortunately only 4 out of 197 portals have used relationships metadata key to link datasets. This extremely low number tells us that not many portals are linking datasets.

Table 3. LINDAT indicator calculation with CKAN powered OGD portals

Portal	Number of datasets	TPL	CL	LINDAT
datahub.io	11462	131365982	340	0.000259%
dados.gov.br	6458	41699306	9	0.000022%
data.gov.ie	8823	77836506	9	0.000012%
dartportal.leeds.ac.uk	25	600	2	0.333333%
data.gov.uk	47139	2222038182	0	0.000000%

Table 3 shows weak interlinking between datasets on every portal, because there is a low number of created linksets. Portal datahub.io has the highest number of defined relationships but very low LINDAT index, because the number of total possible linksets is very high. Portal dartportal.leeds.ac.uk has the best score for LINDAT index, because of the small number of datasets available on the portal. Each of the analyzed OGD portals shows diversity in the number of created relationships with comparing of number of datasets that reside on these portals. The relationships feature is not fully utilized and there should be more attention paid to this issue. There are several approaches [3, 27, 29, 30] that can help to achieve better interlinking between datasets, and lead to higher number of linked datasets on the portal. Using enhanced IRE architecture and applicable prototype application, this number could increase even more since it includes data consumers in datasets interlinking. Their feedback can be used to assess candidate datasets for interlinking and to offer the best suited workflow and best practice to support achieving this aim in collaborative and participatory manner.

⁶ <https://ckan.org/about/instances>

5. Discussion and future remarks

The enhanced LIRE architecture is created for the purpose of managing and creating relations between datasets on OGD portals. It utilizes dataset's metadata to propose relations and to interlink datasets using custom refined model. Model uses dataset's metadata keys in a way it has not been used before, to semantically link datasets and enable their availability in linked data domain. This paper makes a step forward in this direction and utilizes valuable information hidden in metadata with the aim to expose OGD for processing in semantic applications.

Involving data consumers in the process of datasets interlinking helps to increase the number of linksets and contributes to the better accessibility of government data on the portal itself but also in various software applications via publicly available APIs. This is in line with the raised awareness of the user involvement in the process of creation of linked data. In LIRE, users are allowed to identify and reveal possible relations between datasets and to validate them through the proposed model. We like to believe that in this way our architecture contributes to better participation and collaboration between data consumers and data holders (i.e. government). To track the status of dataset's interlinking we have exposed a measure called LINDAT. LINDAT gives the information on how much linksets are added compared to the total number of possible linksets. The total number of possible linksets is directly affected by the total number of datasets on the portal. The value of this indicator is available at any time as it is calculated automatically.

The proposed architecture is applicable on different OGD platforms. Interoperability is ensured by creating custom WRAPPER for each new platform. Connectivity between WRAPPER and the rest of the architecture is achieved via the specification of necessary actions that are issued from components located above the WRAPPER, enabling in that manner the interoperability of the whole architecture. The question that arises is: Should there be a framework for the development of WRAPPER component? This is something that we plan to address in the future, in order to achieve better integrity of the architecture. The future improvements of the architecture could also go in the direction of the examination of the semantic similarity of datasets, based on their names and description. As there are many techniques and tools available for semantic similarity, it will be first necessary to evaluate the applicability of those techniques in short and long texts which exists in dataset metadata description. The research in this area can potentially improve the operability of LIRE application.

While doing the research on dataset's metadata we stumbled upon a question: Should there be a standard for permissible names of metadata's tags, for achieving their unique interpretation? Different OGD platforms differently name the metadata tags [35, 36], giving a task to the developers to develop a technique for their processing based on their notation. The existence of a standard for naming and structure of dataset's metadata will contribute to better processing of this information and speed up the development of end-users applications. A solution to this problem can be found in the development of a mediator, which will process uniquely dataset's metadata from different OGD platforms. Nevertheless, it must be kept in mind that there is no guarantee that metadata will be available on each portal.

References

1. Jacksi, K., Zeebaree, S. R. M., Dimililer, M.: LOD Explor-er: Presenting the Web of Data. *International Journal of Advanced Computer Science and Applications*, Vol. 9, No. 1, pp. 45 – 51. (2018)
2. Veljković, N., Bogdanović-Dinić, S., Stoimenov, L.: Benchmarking open government: An open data perspective. *Government Information Quarterly*, Vol. 31, No. 2, pp. 278-290. (2014)
3. Ding, L., Lebo, T., Erickson, J. S., DiFranzo, D., Williams, G. T., Li, X., Michaelis, J., Graves, A., Zheng, J. G., Shangguan, Z., Flores, J., McGuinness, D. L., Hendler, J. A.: TWC LOGD: A portal for linked open government data ecosystems. *Web Semantics: Science, Services and Agents on the World Wide Web*, Vol. 9, No. 3, pp. 325-333. (2010)
4. Kalampokis, E., Tambouris, E., Tarabanis, K.: Linked Open Government Data Analytics. In *Proceedings of the International Conference on Electronic Government*. Koblenz, Germany, pp. 99-110, 2013.
5. Sheridan, J., Tennon, J.: Linking UK Government Data. In *Proceedings of the Workshop on Linked Data on the Web*. Raleigh, North Carolina, USA. (2010)
6. Maali, F., Cyganiak, R., Peristeras, V.: A Publishing Pipeline for Linked Government. In *Proceedings of the 9th Semantic Web Conference*. Heraklion, Crete, Greece, 778-792. (2012)
7. Assaf, A., Senart, A., Troncy, R.: Roomba: Automatic Validation, Correction and Generation of Dataset Metadata. In *Proceedings of the 24th International Conference on World Wide Web Companion*. Sophia Antipolis, France, 159 – 162. (2015)
8. Kalampokis, E., Tambouris, E., Tarabanis, K.: A classification scheme for open government data: Towards linking decentralised data. *International Journal of Web Engineering and Technology*, Vol. 6, No. 3, pp. 266-285. (2010)
9. Shadbolt, N., O'Hara, K., Berners-Lee, T., Gibbins, N., Glaser, H., Hall, W., Schraefel, M. C.: Linked open government data: Lessons from data.gov.uk. *IEEE Intelligent Systems*, Vol. 27, No. 3, pp. 16-24. (2012)
10. Janssen, M., Estevez, E., Janowski, T.: Interoperability in Big, Open, and Linked Data Organizational Maturity, Capabilities, and Data Portfolios. *Computer*, Vol. 47, No. 10, pp. 44-49. (2014)
11. Assaf, A., Troncy, R., Senart, A.: HDL - Towards a Harmonized Dataset Model. In *Proceedings of the 2nd International Workshop on Dataset Profiling & Federated Search for Linked Data*. Sophia Antipolis, France, 159 - 162. (2015)
12. Kashyap, V., Sheth, A.: Semantic heterogeneity in global information systems: The role of metadata, context and ontologies. In: M. Papazoglou and G. Schlageter (eds.), *Cooperative Information Systems: Tends and Directions 1998*, Academic Press: London, UK. pp. 139-178. (1996)
13. Schmachtenberg, M., Christian, B., Paulheim, H.: Adoption of the linked data best practices in different topical domains. In *Proceedings of the 13th International Semantic Web Conference*. Riva del Garda, Italy, 245-260. (2014)
14. Lebo, T., Erickson, J. S., Ding, L., Graves, A., Williams, G. T., DiFranzo, D., Li, X., Michaelis, J., Zheng, Z., Flores, J., Shangguan, J. G., McGuinness, D. L., Hendler, J. A.: Producing and Using Linked Open Government Data in the TWC LOGD Portal. In *Proceedings of the Linking Government Data Conference*, 51-72. (2011)16. Scharffe, F., Fan, Z., Ferrara, A., Khrouf, H., Nikolov, A.: Methods for automated dataset interlinking. *Datalift Deliverable D4.1*, pp. 34 - 73. (2013)
15. Cverderlj-Fogaraši, I., Sladić, G., Gostojić, S., Segedinac, M., Milosavljević, B.: Semantic integration of enterprise information systems using meta-metadata ontology. *Information Systems and e-Business Management*, Vol 15, No. 2, pp. 257 - 304. (2010)
16. Scharffe, F., Fan, Z., Ferrara, A., Khrouf, H., Nikolov, A.: Methods for automated dataset interlinking. *Datalift Deliverable D4.1*, pp. 34 - 73. (2013)

17. Ellefi, M. B., Bellahsene, Z., Todorov, K., Dietze, S.: Dataset Recommendation for Data Linking: An Intensional Approach. In Proceedings of the 13th European Semantic Web Conference. Heraklion, Crete, Greece, 36-51. (2016)
18. Ngomo, A. C. N., Sherif, M. A.: LIMES - A Framework for Link Discovery on the Semantic Web. *Journal of Web Semantics*, (2018)
19. Kepeklian, G., Bihanic, L., Troncy, R.: Datalift: A platform for integrating big and linked data. In Proceedings of the International Conference on Big Data from Space, Frascati, Italy, 370-373. (2014)
20. Milic, P., Veljkovic, N., Stoimenov, L.: Comparative analysis of metadata models on e-government open data platforms. *IEEE Transactions on Emerging Topics in Computing*, (2018)
21. Zuiderwijk, A., Jeffery, K., Janssen, M.: The potential of metadata for linked open data and its value for users and publishers. *JeDEM-e-Journal of e-Democracy and Open Government*, Vol. 4, No. 2, pp. 222-244. (2012)
22. Maali, F., Cyganiak, R., Peristeras, V.: Enabling Interoperability of Government Data Catalogues. In Proceedings of the 9th International Conference on Electronic Government. Lausanne, Switzerland, 339-350. (2010)
23. DCMI Usage Board. Dublin Core Metadata Initiative: DCMI Metadata Term, <http://dublincore.org/documents/dcmi-terms/>. (2018)
24. EC JoinUP, DCAT Application Profile, https://joinup.ec.europa.eu/asset/dcat_application_profile/home. (2018)
25. Alexander, K., Cyganiak, R., Hausenblas, M., Zhao, J.: Describing linked datasets with void Vocabulary,” <http://www.w3.org/TR/void/>. (2018)
26. VOID: Vocabulary of Interlinked Datasets, Available: <https://www.w3.org/TR/void/>, (2018)
27. Fiorelli, M., Stellato, A., McCrae, J. P., Cimiano, P., Paziienza, M. T.: LIME: The Metadata Module for OntoLex. In: Gandon, F., Sabou, M., Sack, H., d'Amato, C., Cudré-Mauroux, P., Zimmermann, A. (eds.): *The Semantic Web, Latest Advances and New Domains*, Vol. 9088, 321-336. (2015)
28. Cyganiak, R., Maali, F., Peristeras, V.: Self-Service Linked Government Data with dcat and Gridworks. In Proceedings of the 6th International Conference on Semantic Systems. Graz, Austria, 37:1-37:3. (2010)
29. Oliveira, H. R., Tavares, A. T., Lóscio, B. F.: Feedback-based data set recommendation for building linked data applications. In Proceedings of the 8th International Conference on Semantic Systems. Graz, Austria, 49-55. (2012)
30. De Vocht, L., Dimou, A., Breuer, J., Van Compernelle, M., Verborgh, R., Mannens, E., Mechant, P., Van De Walle, R.: A Visual Exploration Workflow as Enabler for the Exploitation of Linked Open Data. In Proceedings of the 3rd International Conference on Intelligent Exploration of Semantic Data. Riva del Garda, Italy, Vol. 1297, 30-41. (2014)
31. Leme, L. A. P. P., Lopes, G. R., Nunes, B. P., Casanova, M. A., Dietze, S.: Identifying candidate datasets for data interlinking. *Lecture Notes in Computer Science*, Vol. 7977, 254-266. (2013)
32. Milić, P., Veljković, N., Stoimenov, L.: Linked Relations Architecture for Production and Consumption of Linksets in Open Government Data. In: Janssen, M., Mäntymäki, M., Hidders, J., Klievink, B., Lamersdorf, W. (eds.). *Open and Big Data Management and Innovation*, Vol. 9373, 221-222, (2015)
33. Nam, T.: The Wisdom of Crowds in Government 2.0: Information Paradigm Evolution toward Wiki-Government. In Proceedings of the 16th Americas Conference on Information Systems. Lima, Peru, 337-348, (2010)
34. SKOS: Simple Knowledge Organization System, Available: <https://www.w3.org/TR/skos-reference/>, (2018)
35. CKAN. CKAN API. <http://docs.ckan.org/en/latest/api/>. (2018)
36. SOCRATA. Socrata API. <http://dev.socrata.com/docs/endpoints.html>. (2018)

Nataša Veljković received the BSc, MSc and PhD degrees in computer science at the University of Niš, Serbia. She is currently working as a Teaching Assistant at Faculty of Electronic Engineering with the Department of Computer Science. Her area of research includes e-government, e-systems, open data, IoT.

Petar Milić received the BSc and MSc degrees in computer science at the University of Priština temporary settled in Kosovska Mitrovica, Serbia and PhD degrees in computer science at the University of Niš, Serbia. He is currently working as a Teaching Assistant at the University of Priština temporary settled in Kosovska Mitrovica, Serbia. His area of research includes e-government, e-systems and web.

Leonid Stoimenov received the BSc, MSc and PhD degrees in computer science at the University of Niš, Serbia. He is a Professor at Faculty of Electronic Engineering at this University. His research interests in computer science include e-systems, GIS, databases, ontologies and semantic interoperability. He is a member of IEEE, IAENG and representative in AGILE association of GIS laboratories in Europe.

Kristijan Kuk received the PhD degrees in computer science at the University of Niš, Serbia. He is a Professor at the University of Criminalistic and Police Studies in Belgrade, Serbia. His research interests in computer science include machine learning, artificial intelligence, pedagogical agent, natural language processing, IT security.

Received: April 20, 2019; Accepted: December 20, 2019

Damaged Buildings Recognition of Post-Earthquake High-Resolution Remote Sensing images based on Feature Space and Decision Tree Optimization

Chao Wang^{1,2,3}, Xing Qiu², *Hui Liu^{4,5}, Dan Li⁴, Kaiguang Zhao³, and Lili Wang⁵

¹ Jiangxi Province Key Laboratory of Water Information Cooperative Sensing and Intelligent Processing,

Nanchang Institute of Technology, Nanchang, China, 330099;

² Key Laboratory of Meteorological Disaster, Ministry of Education (KLME), Nanjing University of Information Science and Technology, Nanjing, China, 210044;

³ College of Food, Agricultural, and Environmental Sciences, The Ohio State University, Wooster, United States, 44691;

⁴ College of Computer and Information Engineering, Hohai University, Nanjing, China, 211100;

⁵ Jiangxi University of Science and Technology, Ganzhou, China, 341000;

Abstract. Earthquake-damaged buildings recognition of the high-resolution remote sensing images has been an indispensable technical means in the post-earthquake emergency response. In view of the difficulties and constraints caused by the lack of pre-earthquake information, this article proposed a novel damaged buildings recognition of high-resolution remote sensing images based on feature space and decision tree optimization. By only using post-earthquake information, the potential building object set is extracted by combining WJSEG segmentation and a group of non-building screening rules. On this basis, an adaptive decision tree number extraction strategy based on the discrimination of classification accuracy by the curve fluctuation is applied. In addition, the spectrum, texture and geometric morphology features are selected according to the feature importance index to form symbolized sets of damaged buildings. Finally, based on the optimized random forest (RF) model, buildings are separated into three categories as undamaged building, partly damaged building and ruin. Experiments on four different datasets show that the overall accuracy all exceed 85% with the proposed method, which is significantly better than the other compared methods in both visual inspection and quantitative analysis.

Keywords: damaged buildings; post-earthquake; high-resolution; feature importance index; decision tree optimization;

1. Introduction

As a fatal disaster, earthquake often occurs with casualty and economic losses. In time and accurate recognition of earthquake-damaged buildings

after the earthquake is of great significance to the rapid assessment of the condition of the disaster, carrying out emergency rescue and post-earthquake reconstruction. Compared with traditional field investigation methods, the recognition of earthquake-damaged buildings based on remote sensing images has the advantages of rapid data acquisition and wide coverage etc., which has become one of the indispensable technical means in the emergency response after the earthquake [1].

With the continuous development of remote sensor technology, the wide application of high-resolution remote sensing images has brought more detailed spatial information that is conducive to a more detailed portrayal of the earthquake-damaged buildings [2]. In this view, earthquake-damaged buildings recognition of high-resolution remote sensing images has a great potential for development and application. The existing recognition methods of earthquake-damaged buildings that only depend on the post-earthquake images has break through its dependence on the pre-earthquake images, thus have higher feasibility in the practical application. However, there is no reference change information which can be extracted from pre-earthquake image, and meanwhile such methods need to face the more prominent phenomenon of “same-object with different spectra” and “same-spectrum with different objects” caused by the increase of the spatial resolution of remote sensing images. That is, the severe challenge of the respective intra-class variance of undamaged building damaged building and other object increases and the inter-class variance decreases. Therefore, constructing more efficient feature space to describe the details of earthquake-damaged buildings accurately, which is the premise and foundation for the recognition of earthquake-damaged buildings. At present, the features which are widely used in the recognition of earthquake-damaged buildings mainly include spectra, texture and geometric morphology [3], [4], [5]. For example, Liu et al. used Morphological Attribute Profiles (MAPs) and Local Binary Pattern (LBP) to extract the geometric and texture features of the image, and then extracted the earthquake-damaged buildings by random forest (RF) classifier [6]; Asli et al. combined with spectral, compactness and smoothness features, proposed a post-earthquake collapsed buildings recognition method [7]. Although combining different kinds of features is beneficial to the multi-dimensional description of earthquake-damaged buildings, the redundant information between different features not only increases computational complexity, but also reduces the recognition accuracy due to the conflict between different features as the evidence of damaged buildings. Therefore, multi-features screening and optimization strategy for building a refined set of features is needed. On this basis, the feature set should be combined with the appropriate classification method to obtain reliable earthquake-damaged buildings recognition results. Currently, RF is a popular integrated classifier, which is applied to the field of earthquake-damaged buildings recognition based on high-resolution remote sensing images. It has the advantages of fewer model parameters and avoiding over-fitting by using the double randomness of samples and features selection [8]. For example, Solomon et al. compared RF with other classifiers, proving the good performance of RF in post-earthquake image classification [9]. Moreover, the rational selection of the number of decision trees is a key factor to improve the performance of RF. It is difficult to obtain reliable classification results when the number of

decision trees is too small, while too large will reduce the efficiency of method implementation, especially when the decision trees exceed a certain quantity, the accuracy of classification fluctuates up and down within a certain range or even has a downward trend. Nevertheless, clear quantitative criteria for decision trees is not given in RF theory [9], and the usual manual assignment way is not only susceptible to subjective factors or fall into local optimized, but also reduces the degree of automation of the classification process.

In view of the above challenges, a damaged buildings recognition of post-earthquake high-resolution remote sensing images based on feature space and decision tree optimization is proposed, and the contributions of this study can be summarized as follows:

(1) In the pre-processing step, a candidate object set extraction strategy based on WJSEG image segmentation and a group of non-building objects screening rules is proposed. It can significantly eliminate the shadows, vegetations, and non-building artificial objects.

(2) An adaptive decision tree number extraction strategy based on the discrimination of classification accuracy by curve fluctuation is proposed. By comparing with the trial and error strategy, it can not only improve the degree of automation, but also achieve ideal classification accuracy.

(3) Based on the optimized classifier, a representative feature set for damaged building description is extracted under the guidance of feature importance index. Through the experiments, this feature set shows outstanding performance in terms of the earthquake-damaged buildings recognition in complex post-earthquake scenes.

This study includes five sections. Section 2 briefly analyzes and summarizes the research progress of earthquake-damaged buildings recognition based on high-resolution remote sensing image. Section 3 elaborates the implementation steps of the proposed method. Section 4 analyses and discusses the experimental results. The conclusion is presented in Section 5.

2. Related Work

In recent years, the recognition of earthquake-damaged buildings based on high-resolution remote sensing images has become a very active research direction, such as Stanford University, the University of Trento, the University of Tokyo, Wuhan University and other research institutes, which are doing related research and achieved many results. According to the different data sources used, the existing methods of earthquake-damaged buildings recognition can be divided into two categories, including multi-temporal (pre-earthquake and post-earthquake) images method and single-temporal (post-earthquake) image method [10].

2.1. Damaged Buildings Recognition Based on Pre-Earthquake and Post-Earthquake Images

This kind of methods recognize the earthquake-damaged buildings by using change detection technology with the pre-earthquake and post-earthquake images. For example, Matsuoka et al. used the difference between the backscatter coefficient and the correlation coefficient of pre-earthquake, and post-earthquake SAR images to obtain the optimized window size [11]. Based on Object Based Image Analysis (OBIA), Faming Huang et al. proposed a method of recognition and loss assessment of earthquake-damaged buildings combined with single temporal image classification and multi-temporal images change detection, which achieved overall accuracy (OA) of 93% in the experiment of IKONOS images [12]. Liu Ying et al. obtained the spectral distance of the object in different temporal images by counting the Histogram Oriented Gradient (HOG), then weighting and fusing the spectral distance and HOG features, finally, the Fuzzy *C*-Means (FCM) method is applied to recognize the earthquake-damaged buildings [4]. Roberta et al. used very high-resolution optical images and existing city maps to identify objects corresponding to buildings, using spectra, textures, and statistical features to classify, and the feasibility of mapping earthquake disaster on a single building scale was evaluated [13].

Because of the introduction of pre-earthquake image, the extracted change information can be used as a key damaged buildings feature, so it is usually possible to obtain high accuracy recognition results. Nevertheless, there are many limitations in the practical application of such methods. First, for many cities, especially developing countries, the lack of reference pre-earthquake image or pre-earthquake image is inchoate to determine whether the changes in buildings are caused by earthquake directly leads to the impracticability of such methods. In addition, the quality of change detection is also significantly affected by the quality of image data, such as radiation and imaging angle differences, etc. Finally, the high accuracy registration between the pre-earthquake and post-earthquake images is also a challenging problem [10]. The above factors seriously restrict the wide application of such methods, while the single-temporal image method has been paid more and more attention by scholars [14].

2.2. Damaged Buildings Recognition Based on Post-Earthquake Images

This kind of methods can only extract features from the post-earthquake image for classification. For example, based on statistical and analytical image gradient, Ye Xin et al. proposed a method based on the local spatial distribution of gradients, and further identified the different types of earthquake-damaged buildings [15]. In recent years, more and more machine-learning based methods have been proposed [16], [17]. However, such methods deeply rely on abundant samples. Christian et al. proposed a series of methods based on Support Vector Machine (SVM) and RF, which include five steps: features extraction, features screening, outlier detection, synthetic

samples generation and supervised classification, and can estimate Seismic Building Structure Types (SBSTs) information [18]. Duarte D et al. proposed a method for recognition post-earthquake damaged buildings with images of different sensors and resolutions based on Convolutional Neural Networks (CNN) and Deep Learning (DP) [19].

With the wide application of various unmanned aerial vehicles and satellite platforms, the timeliness of acquiring high-resolution remote sensing image after the earthquake has significantly enhanced. At the same time, due to this kind of methods do not rely on the pre-earthquake images, it is more in line with the application requirements. It should be noted that the proposed method in this study also belongs to this category.

3. Method

The proposed method mainly includes four steps: potential building set extraction, adaptive selection of the number of decision trees, feature set optimization guided by the importance index, and image classification based on optimized RF model. The specific implementation process is shown in Fig.1:

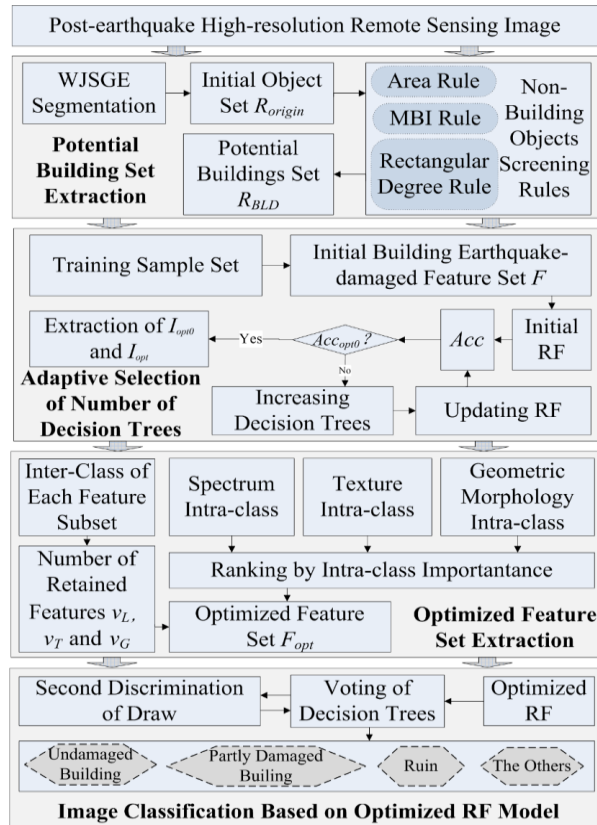


Fig. 1. Flow chart of the proposed method

3.1. Potential Building Set Extraction

● Initial Object Set Extraction

The discrete pixels are firstly divided into geographic objects set with semantic information through image segmentation, thus providing basic analysis units for subsequent earthquake-damaged buildings recognition. For this reason, this study adopts high-resolution remote sensing image segmentation algorithm WJSEG [20], which has the following advantages: WJSEG maintains more complete outlines of geographical objects than the well-known commercial software eCognition, while helping to increase the transparency of the proposed method.

The steps of WJSEG mainly include initial seed region set conduction, secondary extraction of seed region, inter-scale constraint segmentation and regions merging. The specific implementation process of WJSEG can be found in Reference [20]. The initial object set extracted by image segmentation is represented by R_{origin} in this study.

● Non-building Object Screening

Based on R_{origin} , a group of non-building object screening rules is designed in this study, thus avoiding false positives caused by such objects in subsequent processing while reducing the amount of computation. For object R_b in R_{origin} , the specific screening rules are as follows:

(1)Area rule. Count the number of pixels N_{pixels} contained in R_b . Due to the areas of buildings in remote sensing image with different resolution may have large difference, the area rule is set with experience as follows: if $N_{pixels} \leq 80$, then R_b is considered as the small target such as vehicles and noise, and is eliminated.

(2)Rectangular degree rule. The rectangular degree is a parameter that measures the fullness of an object and its smallest external rectangle, which can be represented as $Rd = N_{pixels} / N_{rectangle}$ by calculating the number of pixels N_{pixels} contained in the smallest external rectangle of R_b . The aspect ratio of the smallest external rectangle of R_b is Ar . If R_b meets $Rd < 0.8$ and $Ar > 5$, the object is considered to be a narrow target such as a road, river, etc., and is removed [21].

(3)Morphological Building Index (MBI) rule. MBI takes advantage of the features of pixels belonging to buildings are mostly highlighted in the gray image after top-hat transform and obtains the index value corresponding to a pixel by calculating the multi-scale differential sequence [22]. The higher the MBI, the more likely it is that the pixel belongs to a building, and the calculation formula is as follows:

$$MBI = \frac{\sum_d \sum_s DMP(d, s)}{D \times S} \quad (1)$$

where D and S represent the direction and scale of linear structural elements, respectively; $DMP(d, s)$ is a multi-scale difference morphological sequence. According to the suggestion of document [23], this study sets $D = 8$ and $S = 50$. Based on the MBI values of all pixels in the image, the

separation threshold T_{MBI} is adaptively determined by using the maximum between-class variance method (Otsu) to obtain the proportion $Ratio_{false}$ of non-building pixels in R_b , if $Ratio_{false} > 0.8$, and it is removed.

According to rules (1) to (3), traverse all objects in R_{origin} , and the removed non-building objects will no longer participate in subsequent analysis and discrimination. The remaining objects will form a set of potential buildings R_{BLD} for further of earthquake-damaged buildings recognition.

3.2. Adaptive Selection of Number of Decision Trees

The earthquake-damaged buildings recognition based on machine learning is essentially to transform the problem of target recognition into the problem of image classification through feature extraction of post-earthquake images. The RF classifier used in this study is an integrated classifier based on non-pruning decision tree. Compared with other machine learning classifiers, it has the advantages of higher precision, stronger generalization ability and fewer parameters, which has been widely used in classification research of high-resolution remote sensing images [24]. The number I of decision trees in RF is a key parameter that significantly affects the classification accuracy. Therefore, this study proposes an adaptive selection strategy of number of decision trees, which mainly includes the following five steps:

Step1: In potential building object set R_{BLD} , 20 undamaged buildings, 20 partly damaged buildings, 20 ruins, and 20 other objects are taken by artificial marking to form the training sample set H ;

Step2: Conduct the initial building earthquake-damaged feature set F . In this study, 30 common features in the field of earthquake-damaged buildings recognition including spectrum, texture and geometric morphology are selected to construct the initial building earthquake-damaged feature set F [25], [26], [27]. Among which, the spectral features include R-band mean, G-band mean, B-band mean, R-band standard deviation, G-band standard deviation, B-band standard deviation, R-band contribution rate, G-band contribution rate, B-band contribution rate, Brightness; the geometric morphological features include Area, MajorAxis Perimeter, Eccentric Orientation, MinorAxis, Range, ConvexArea, Diameter, Solidity; and the texture features include grayscale symbiotic matrix Contrast, Homogeneity, Correlation, Entropy, J-value, Roberts operator, Sobel operator, Prewitt operator, Laplacian operator, Canny operator;

Step3: Build the initial RF model with five decision trees and enter a training sample set with all features. The classification accuracy rate is defined as the ratio of the number of correctly classified samples per tree to the number of input samples. The larger the ratio, the closer the classification result is to the real situation. Calculate the classification accuracy of each

tree, and find the mean value of the correct rate of all decision trees

$$Acc_0 = \frac{1}{5} \sum_{i=1}^5 Tree_i, \text{ which } i \text{ is the serial number of the decision tree;}$$

Step4: Build a new RF model after iteratively adding decision tree with five trees as a step size. Using the same steps as Step3, the mean value of the correct rate $Acc_n = (\sum_{i=1}^{5(n+1)} Tree_i) / [5(n+1)]$ is obtained when the iteration number is $n \in [1, 99]$;

Step5: If Acc_n meets $Acc_{n-1} < Acc_n < Acc_{n+1}$, it is considered to be a peak in a sub-interval containing 10 decision trees, and notes $Acc_{opt0} = Acc_n$;

Step6: Continue to calculate the 3 consecutive peak points after Acc_n , which are respectively Acc_{opt1} , Acc_{opt2} , Acc_{opt3} , and if Acc_{opt0} satisfies:

$$Acc_{opt0} \geq \arg \max \{ Acc_{opt1}, Acc_{opt2}, Acc_{opt3} \} \quad (2)$$

the iteration is stopped and the number of decision trees corresponding to Acc_{opt0} is defined as I_{opt0} . Otherwise, proceed to the next step;

Step7: Repeat Step5 and Step6 until the Acc_{opt0} meets formula (2), stop the iteration, and I_{opt0} can be extracted. If Acc_{opt0} cannot be determined within $n \in [1, 99]$, the number of decision trees corresponding to the obtained maximum value of Acc_n is recorded as I_{opt0} ;

Step8: Continue to compare the classification accuracy of I_{opt0} and the four neighboring trees (9 models in total) before and after, and determine the number of optimized decision trees I_{opt} finally extracted according to the maximum classification accuracy.

3.3. Optimized Feature Set Extraction Guided by Feature Importance Index

Based on the optimized selection of the decision tree number of RF, in order to further reduce the redundancy and evidence conflict between the features in the candidate feature set F , this study defines an important degree index of features, and then puts forward the strategy of optimizing the feature set. In a RF model with I_{opt0} decision trees, the importance of all features is firstly calculated. Secondly, the inter-class importance of spectral, texture and geometric morphology features are calculated respectively. On this basis, the intra-class importance of sub-features included in the three types of features is calculated respectively. Finally, the optimized classification feature set of earthquake-damaged buildings is obtained under the guidance of the importance. The specific steps are as follows:

Step1: Calculate the importance of all features. In the training sample set H , the unselected samples after random sampling with return constitute an out-of-bag data set (OOB). The importance of any feature $f_i (f_i \in F, t \in [1, 30])$ to the i th decision tree is calculated by formula (3):

$$W^{(i)}(f_t) = \frac{\sum_{x_j \in \Phi_B} N(l_j = c_j^{(i)})}{|\Phi_B|} - \frac{\sum_{x_j \in \Phi_B} N(l_j = c_{j,f_t}^{(i)})}{|\Phi_B|} \quad (3)$$

where Φ_B represents OOB sample set, x_j and l_j respectively represent any sample in the data out-of-bag and its assigned category label, $c_j^{(i)}$ represents the category label obtained by sample x_j , $c_{j,f_t}^{(i)}$ denotes a category label of sample x_j obtained by replacing the value of the feature f_t with other random values, N is a counting function. After traversing all decision trees, the importance $W(f_t) = \frac{\sum_{i=1}^{I_{opt}} W^{(i)}(f_t)}{I_{opt}}$ of feature f_t to RF classifier can be obtained;

Step2: Calculate the inter-class importance of features. Since the features in F are respectively classified into three categories of spectrum, texture and geometric morphology (hereinafter denoted by subscripts L , T and G respectively), the inter-class importance of the three categories of features, namely spectrum, texture and geometric morphology, can be obtained by summing up the importance of each and the categories to which it belongs respectively, and are respectively recorded as W_L , W_T and W_G . On this basis, the spectral features normalized penalty factor for inter-class redundancy is defined as $\varpi_L = \frac{W_L}{W_L + W_T + W_G}$. By analogy, the normalized penalty factors of inter-class redundancy for texture and geometric morphological features are ϖ_T and ϖ_G respectively;

Step3: Calculate the intra-class importance of features. According to the spectral, texture morphological feature subsets to which the 30 features belong, the intra-class importance is calculated by Step1. On this basis, the features in each feature subset are arranged from high to low according to the intra-class importance;

Step4: For each feature subset, after rounding according to the proportion of redundant normalized penalty factors, only v_L , v_T and v_G features with relatively high importance within the class are respectively retained, thus obtained an optimized feature set F_{opt} with $V = v_L + v_T + v_G$ features.

3.4. Image Classification Based on Optimized RF Model

Based on the extracted F_{opt} and I_{opt} , an optimized RF model is constructed as follows:

$$P(x) = \arg \max_c \sum_{i=1}^{I_{opt}} E(p_i(x) = c) \quad (4)$$

where $P(x)$ represents the classification result; $p_i(x)$ represents the classification result of a single decision tree; c means classification label, $c \in \{\text{Undamaged Building, Partly Damaged Building, Ruin, The Others}\}$. On this basis, voting is carried out according to the classification labels given by each decision tree, and the number of votes is taken as the standard of the final classification label of the sample. If the voting results in a draw, the distance $Dist$ between the sample and the training samples of these categories is discriminated according to formula (5), and the category with smaller $Dist$ is taken as the final classification result of the sample.

$$Dist = \sqrt{\frac{\sum_{v=1}^V (x_{rest(v)} - x_{train(v)})^2}{s_v^2}} \quad (5)$$

In the formula, $x_{rest(v)}$ and $x_{train(v)}$ are the values of the v th feature in the test sample and training sample, respectively; s_v^2 is the variance of the v th feature.

4. Experiment and Analysis

In the experiments, four groups of post-earthquake high-resolution remote sensing images of different sensors are used. Through visual analysis and quantitative accuracy evaluation, the performance of the proposed method is verified by comparison with a variety of advanced methods.

4.1. Experiment Datasets

Dataset 1 and Dataset 2 are GE01 satellite remote sensing images of Yushu in Qinghai Province, China, and the acquisition time is May 6, 2010. The earthquake occurred on April 4, 2010, with the highest magnitude of 7.1. The image includes panchromatic and multispectral (blue, green, red and near infrared) bands with spatial resolutions of 0.41m and 1.65m respectively and a size of 1024×1024 pixels. Pan-sharpened RGB images with a spatial resolution of 0.41m fused by ENVI software are used in the experiment. As shown in Fig. 2 (a) and (b). Dataset 3 and Dataset 4 are QuickBird satellite remote sensing images of Wenchuan in Sichuan Province of China, and the acquisition time is June 3, 2008. The earthquake occurred on May 12, 2008, with the highest magnitude of 8.0. The image includes panchromatic and multispectral (blue, green, red and near infrared) bands with spatial resolutions of 0.6m and 2.4m respectively. Pan-sharpened RGB images with spatial resolution of 0.6m fused by ENVI software are used in the experiments, as shown in Fig. 2 (c) and (d).

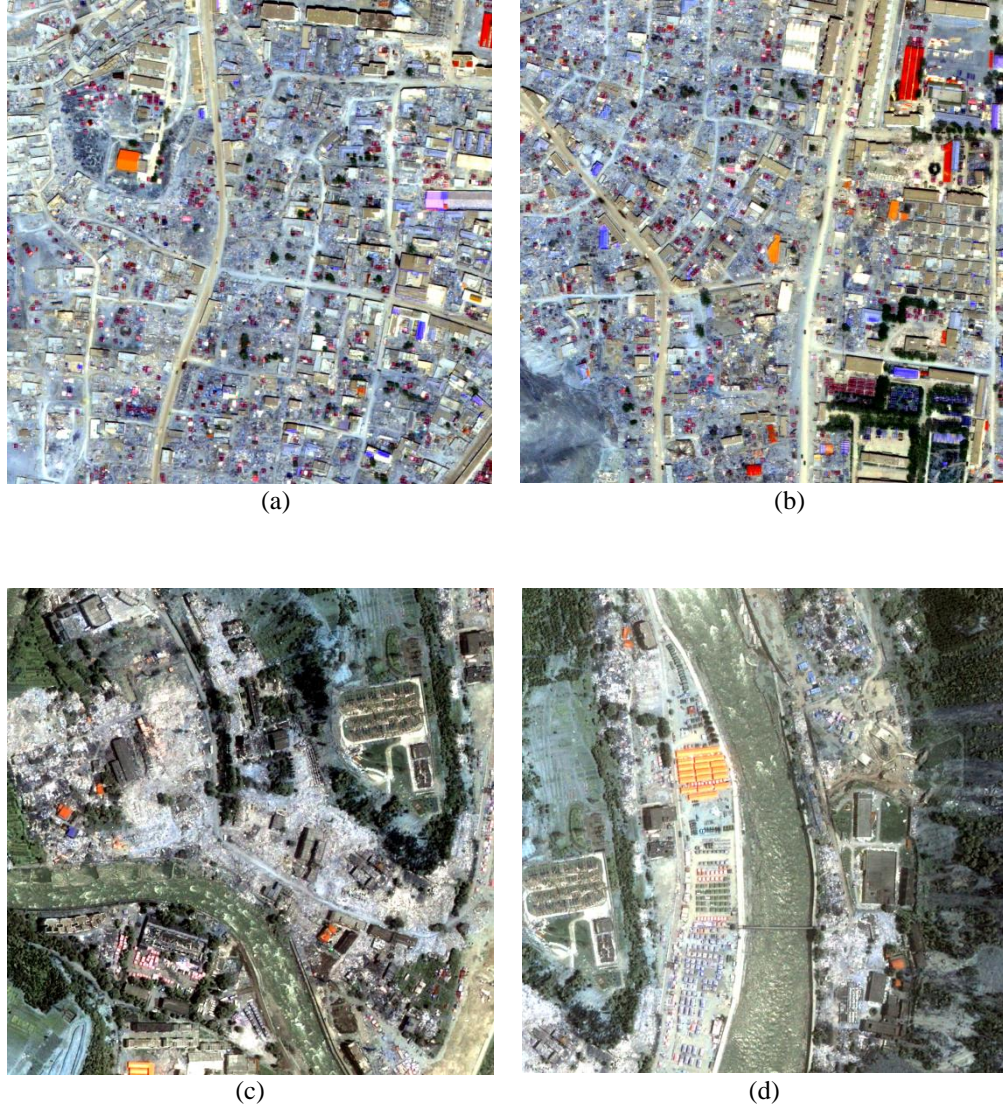


Fig. 2. Experimental datasets: (a) Dataset 1; (b) Dataset 2; (c) Dataset 3; (d) Dataset 4.

The reasons for selecting these four datasets for experiments are as follows: Satellite remote sensing images are one of the main forms of obtaining ground information resources at present. These datasets of different regions and sensors are selected to help analyze the general applicability of the proposed method. Besides, these areas are seriously suffered after the earthquake, including buildings, vegetations, roads, wastelands, rivers, slopes, etc. In addition, there are two types of damaged buildings, including partly damaged buildings and ruins, which is in line with the aim of the identification of damaged buildings with different degrees in this study.

4.2. Extraction Results of Potential Building Sets

● Extraction Results of Initial Object Sets

The image is firstly segmented by WJSEG method, and the initial object set R_{origin} is obtained. In order to facilitate observation, a semitransparent white layer is superimposed on the original image, and then the segmentation result is represented by black pixels and projected into the image, as shown in Fig. 3.



Fig. 3. WJSEG segmentation and initial object set extraction results: (a) Dataset 1; (b) Dataset 2; (c) Dataset 3; (d) Dataset 4.

It can be seen that the segmentation results have completely extracted undamaged buildings and earthquake-damaged buildings, and there is almost no phenomenon of under segmentation. The phenomenon of over-

segmentation mainly exists in river and grassland areas, which are not the region of interesting (ROI) in this study. Therefore, the extracted initial set of objects can provide effective analysis elements for subsequent earthquake-damaged buildings recognition.

● **Screening Results of Non-building Objects**

According to the discrimination rules in section 3.1, the non-building objects in initial object set are screened, and the results are shown in Fig. 4. Among them, the excluded non-building objects are represented by black pixels, and the remaining white objects constitute a potential building set R_{BLD} . At the same time, in order to facilitate analysis, some representative positions and objects in the image have been labeled with different alphabetic symbols, and the same operations are used in the following chapters.

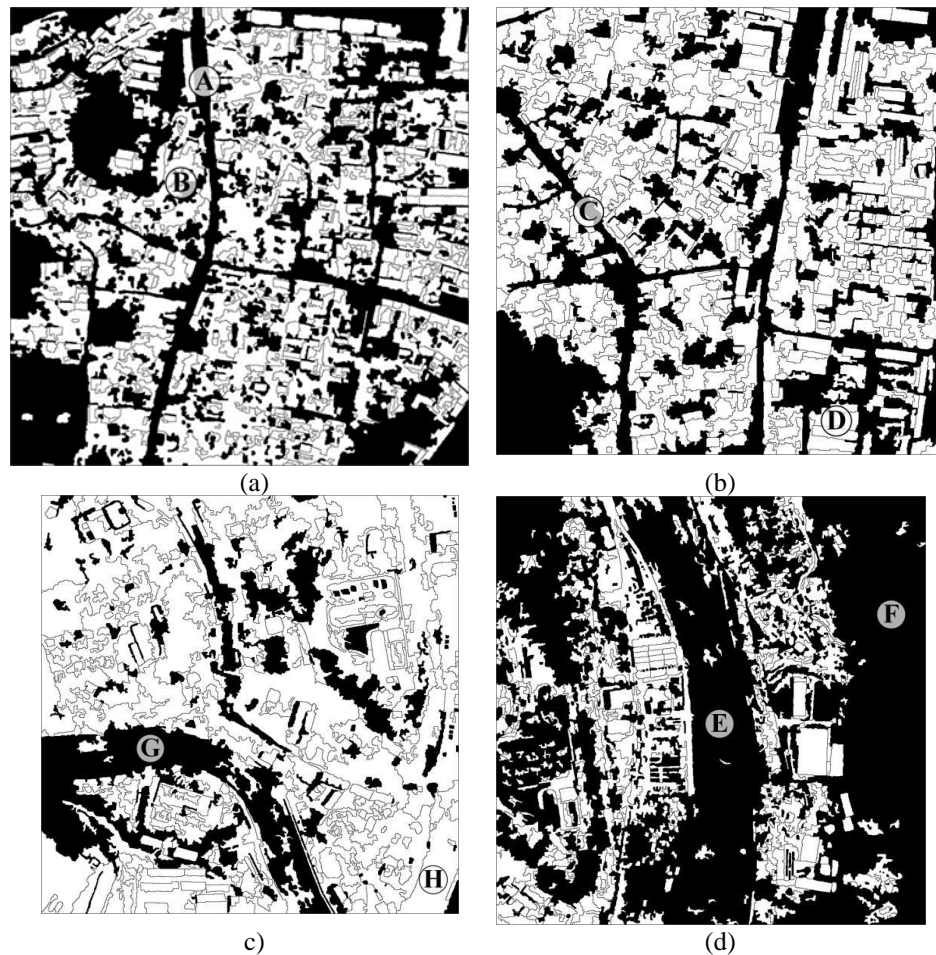


Fig. 4. Results of screening non-building objects: (a) Dataset 1; (b) Dataset 2; (c) Dataset 3; (d) Dataset 4.

As shown in the above figure, roads (e.g., locations A and C) in Dataset 1 and Dataset 2, tents with smaller areas (e.g., location B) in Dataset 1, river courses (e.g., locations G and E) in Dataset3 and Dataset4, and large areas of vegetation region (e.g., location F) in Dataset 4 have been effectively screened. However, there are still some non-building objects, such as wasteland (e.g., location D) in Dataset 2, Grassland and wasteland (e.g., location H) in Dataset 3. Therefore, the proposed non-building objects screening strategy is feasible and effective, but it is also necessary to retain the category of other object in the subsequent classification results.

4.3. Adaptive Selection of the Number of Decision Trees

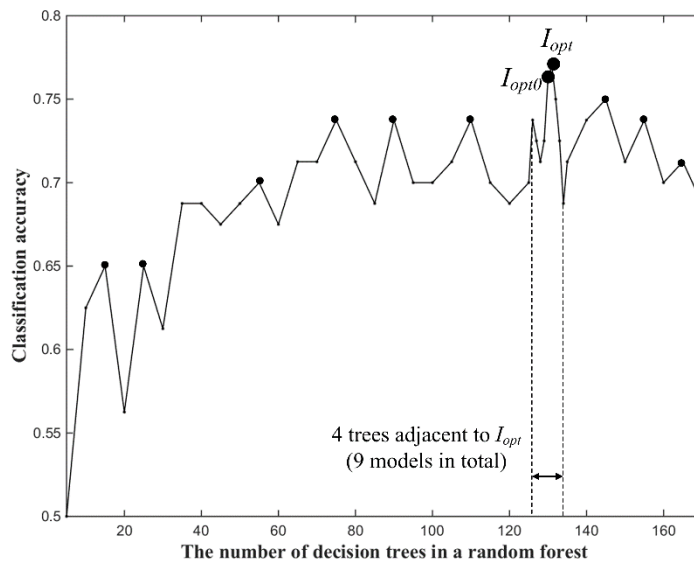


Fig. 5. The adaptively selected number of decision trees in Dataset 1.

The adaptively selected numbers of decision trees corresponding to Dataset 1 to Dataset 4 are 131, 64, 91 and 63 respectively. In the following, Dataset 1 is taken as an example for detailed description of adaptive selection process. As shown in Fig. 5, according to the discrimination rules in section 3.2, the satisfied condition of the selected peak point is $I_{opt0} = 130$. On this basis, an interval is constructed with I_{opt0} as the center, and the classification accuracy of four trees before and after I_{opt0} are calculated respectively. Among them, when the number of decision trees is 131, the classification accuracy rate reaches the maximum, so $I_{opt} = 131$ is taken.

4.4. Optimized Feature Set Extraction Results

On the basis of the adaptive number of decision trees, the optimized feature set extraction results corresponding to four Datasets are shown in Table 1 respectively:

Table 1. Optimized feature set extraction results

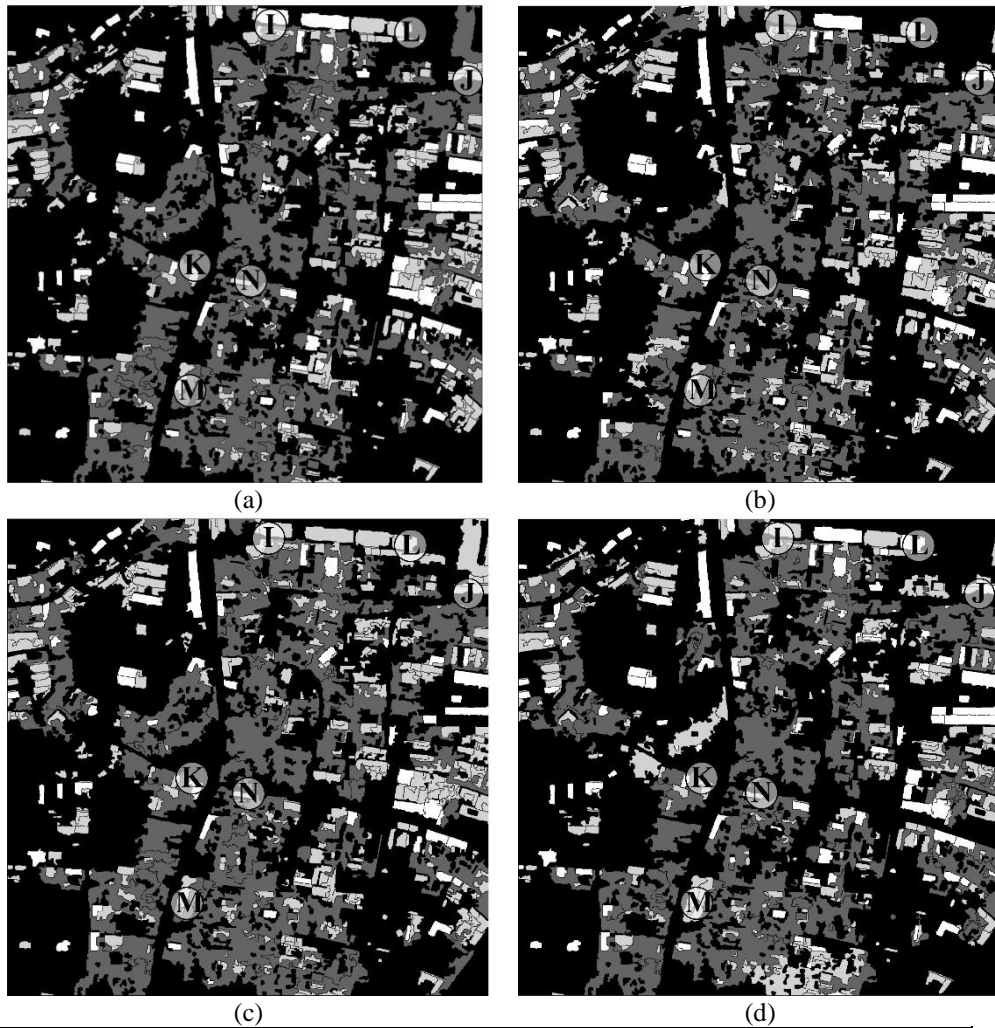
Experimental dataset	Optimized feature set
Dataset 1	G-band standard deviation, B-band standard deviation, R-band contribution rate, MinorAxisLength, Extent, Solidity, GLCM homogeneity, GCLM correlation, GLCM energy, Laplacian operator
Dataset 2	B-band standard deviation, R-band contribution rate, G-band contribution rate, MinorAxisLength, Extent, Perimeter, Solidity, Prewitt, Laplacian operator, Canny operator
Dataset 3	B-band mean, G-band standard deviation, R-band contribution rate, G-band contribution rate, B-band contribution rate, MinorAxisLength, Solidity, Laplacian operator, GLCM correlation, GLCM energy
Dataset 4	B band mean, R band standard deviation, R band contribution rate, G band contribution rate, B band contribution rate, Extent, Solidity, GLCM contrast, GLCM correlation, Laplacian operator

4.5. Comparison Methods and Recognition Results

In order to objectively analyze and verify the performance of this method, traditional RF and another two advanced methods are chosen for comparative experiments. Method 1 [28] is based on the traditional RF. Set the initial building earthquake-damaged feature set F and 500 decision trees as the input of classifier. By comparing with Method 1, it is helpful to analyze the effectiveness of proposed optimization strategy for feature and the number of decision trees selection. Method 2 is an optimized RF based method with the 10-fold cross-validation [29]. More representative training damage samples are chosen for improving the recognition accuracy in this method. The

optimized feature set F_{opt} extracted in this study and 500 decision trees in original literature are used as the input. Method 3 is applied by improved Separability and Thresholds (SEaTH) for feature optimization, and extracting earthquake-damaged buildings based on membership degree. Comparing with the two advanced methods is helpful to objectively evaluate the OA of the proposed method [30], [31], [32]. In addition, since Method 2 is a pixel-level method and it is difficult to directly compare with the results of the object-level classification method in this study, the objects in the initial object set R_{origin} are used instead of pixels as the basic unit of subsequent

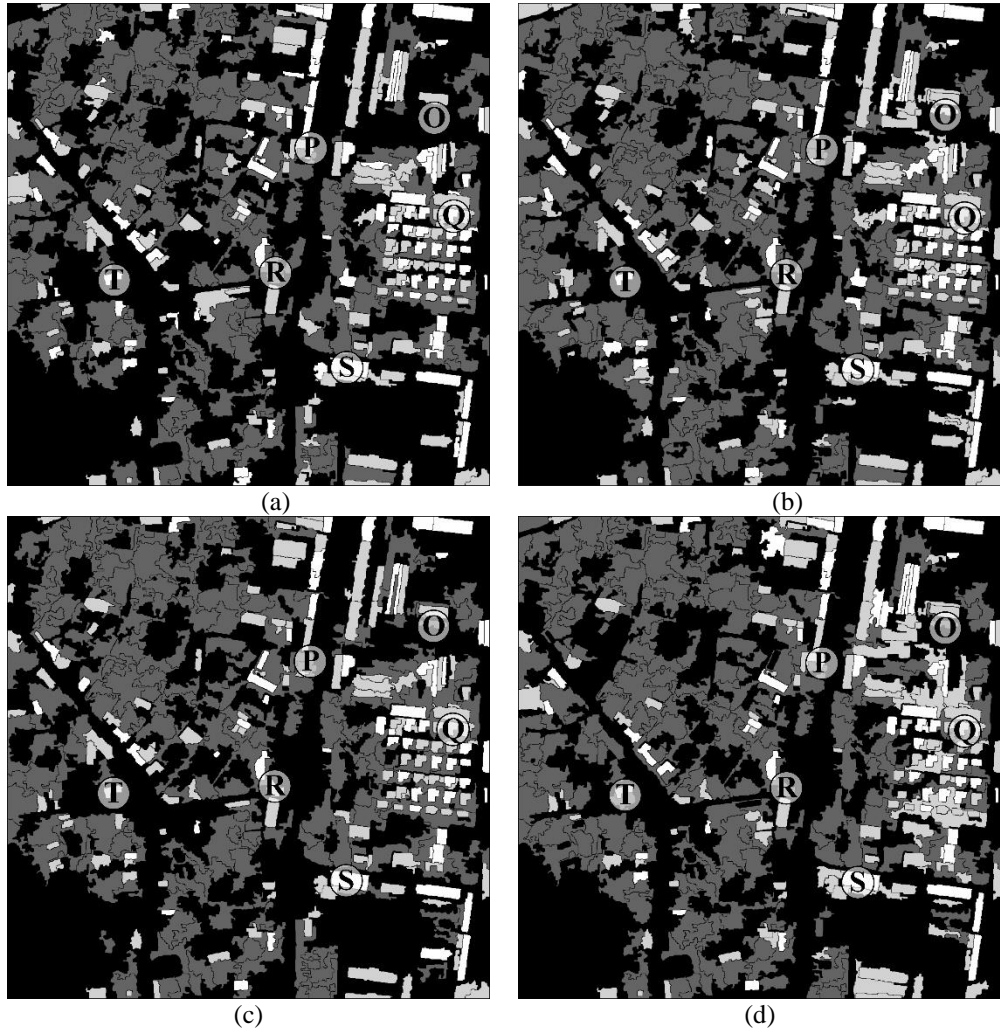
classification. Meanwhile, despite Method 3 is an object-level method, the objects in R_{origin} are also used as the basic unit for insuring the consistency of objects in classification. Besides, in order to avoid the difference in results of earthquake-damaged buildings recognition caused by whether the preliminary screening of non-buildings is taken, all the comparison methods carry out feature extraction and classification on the basis of the potential building set R_{BLD} .



Undamaged Building	Partly Damaged Building	Ruin	The Others

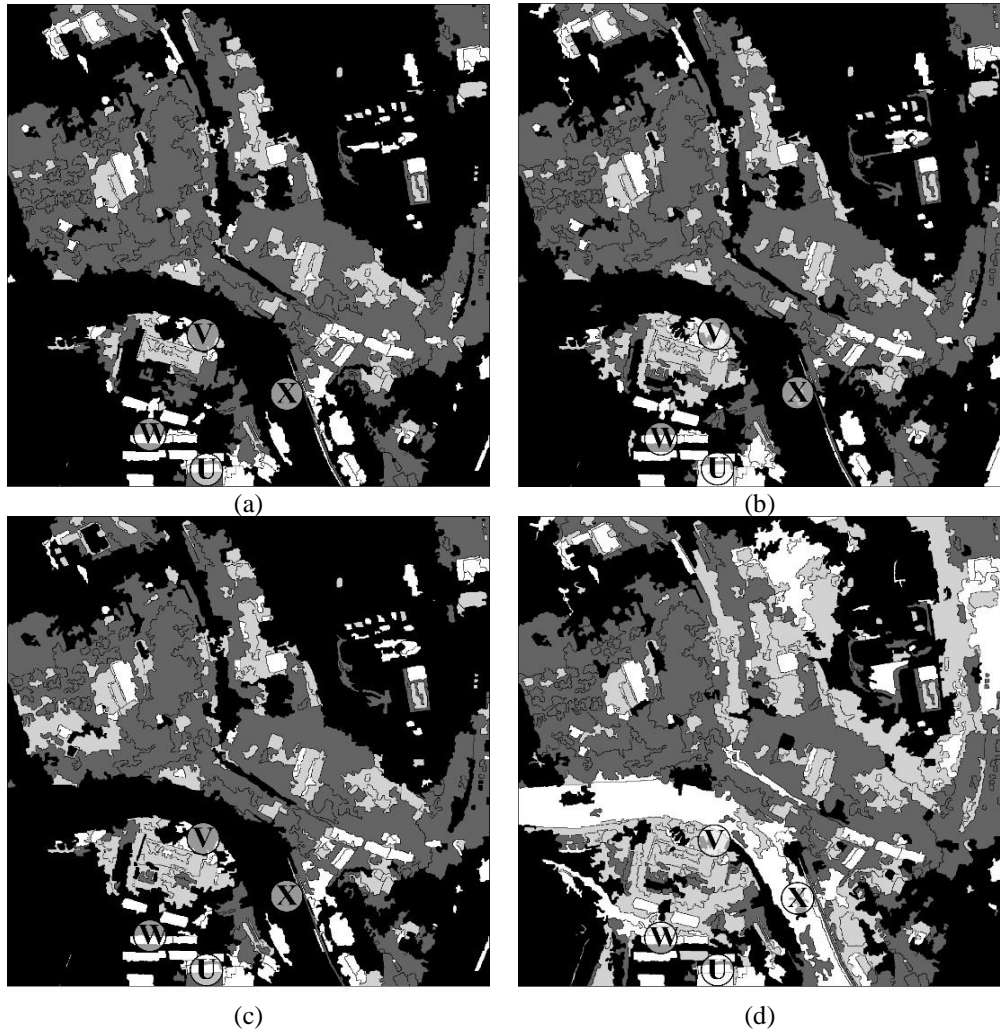
Fig. 6. Recognition Results in Dataset 1: (a) Method of this study; (b) Method 1; (c) Method 2; (d) Method 3.

The results of classification of all comparison methods are including four categories: undamaged building, partly damaged building, ruin and other object, which are respectively represented in different colors. The results of proposed method and three comparison methods are shown in Fig. 6 to Fig. 9.



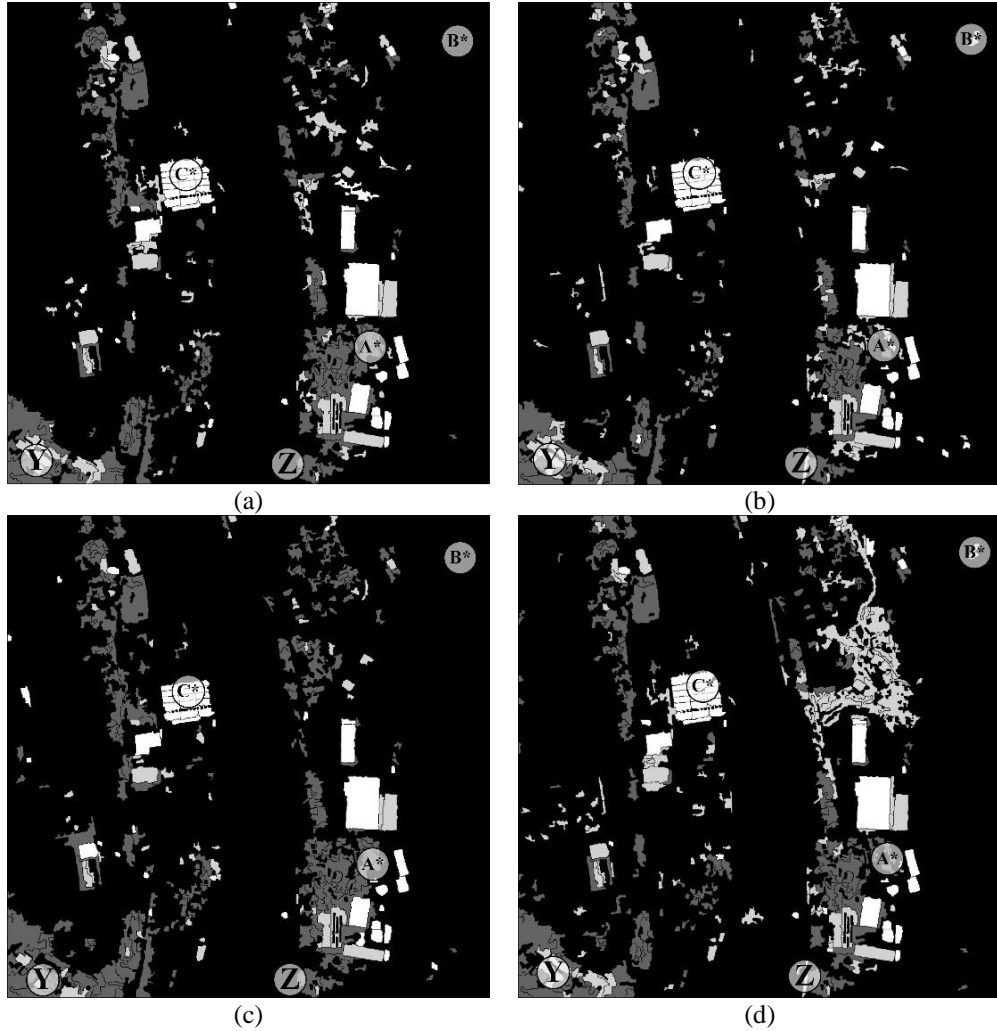
Undamaged Building	Partly Damaged Building	Ruin	The Others

Fig. 7. Recognition Results in Dataset 2: (a) Method of this study; (b) Method 1; (c)



Undamaged Building	Partly Damaged Building	Ruin	The Others

Fig. 8. Recognition Results in Dataset 3: (a) Method of this study; (b) Method 1; (c) Method 2; (d) Method 3

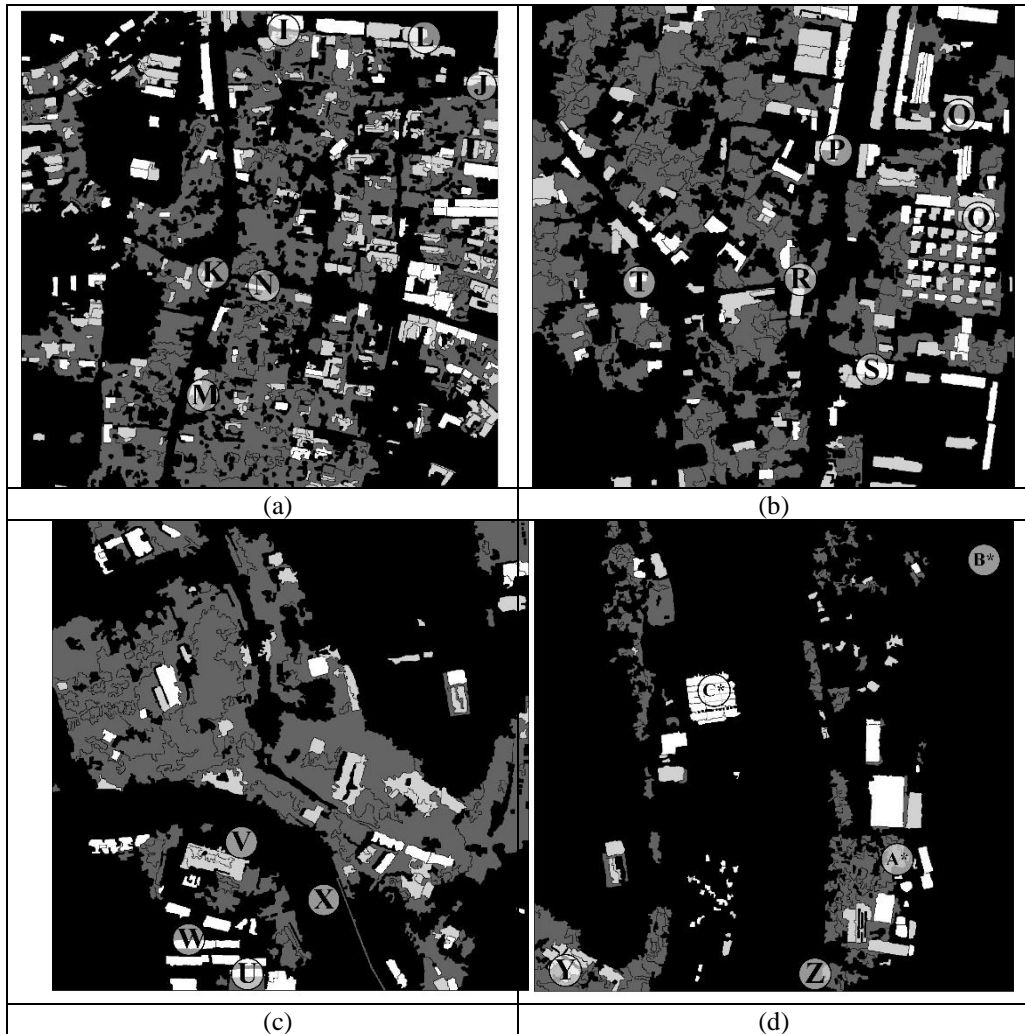


Undamaged Building	Partly Damaged Building	Ruin	The Others

Fig. 9. Recognition Results in Dataset 4: (a) Method of this study; (b) Method 1; (c) Method 2; (d) Method 3.

4.6. Accuracy Evaluation

As the basis of accuracy evaluation, the ground truth maps of four datasets are drawn based on visual inspection and field investigation, as shown in Fig. 10.



Undamaged Building	Partly Damaged Building	Ruin	The Others

Fig. 10. The ground truth maps of four datasets: (a) Dataset 1; (b) Dataset 2; (c) Dataset 3; (d) Dataset 4.

● **Visual Inspection**

By comparing the experimental results with ground truth maps, the recognition effect of the proposed method is significantly better than other three methods, which mainly reflects in: (1) In the four groups of experimental results, the undamaged buildings located at locations I and Q,

the non-buildings located at location O, and the ruin located at location R are all correctly identified only by the proposed method; (2) For the undamaged buildings with regular shape and texture, such as location S, W, Z and C*, the four methods obtain correct discrimination results; However, for partly damaged buildings, for example, only the proposed method and Method 3 obtain the correct results in location K and P, and only the proposed method, Method 2 and Method 3 have made correct judgment for location L and J; (3) For ruins, such as location U and A, only the proposed method and Method 2 have correctly identified them; (4) For non-buildings with similar spectral and shape features to buildings that have not been screened out in potential building sets, such as location V and X, only the proposed method and Method 2 obtain correct results; For some small objects, such as B*, Z, only the proposed method and Method 2 can make correct judgment.

● **Quantitative Analysis**

On the basis of visual inspection, four accuracy indexes including OA, false positive (FP), false negative (FN) and Kappa are used to quantitatively evaluate the accuracy, as shown in Tables 2 to 5.

Table 2. Quantitative accuracy evaluation in Dataset 1

Methods \ Indexes	OA	FP	FN	Kappa
Evaluation criterion	The bigger the better.	The smaller the better.	The smaller the better.	The bigger the better.
Proposed method	91.81%	2.89%	8.19%	0.817
Method 1	90.14%	3.52%	9.86%	0.779
Method 2	90.68%	3.31%	9.32%	0.793
Method 3	88.69%	4.08%	11.31%	0.746

Table 3. Quantitative accuracy evaluation in Dataset 2

Methods \ Indexes	OA	FP	FN	Kappa
Evaluation criterion	The bigger the better.	The smaller the better.	The smaller the better.	The bigger the better.
Proposed method	86.16%	5.08%	13.84%	0.781
Method 1	78.30%	8.45%	21.70%	0.666
Method 2	80.17%	7.61%	19.83%	0.689

Method 3	77.31%	8.91%	22.69%	0.643
-----------------	--------	-------	--------	-------

Table 4. Quantitative accuracy evaluation in Dataset 3

Methods \ Indexes	OA	FP	FN	Kappa
Evaluation criterion	The bigger the better.	The smaller the better.	The smaller the better.	The bigger the better.
Proposed method	85.99%	5.15%	14.01%	0.747
Method 1	80.74%	7.37%	19.26%	0.669
Method 2	81.44%	7.06%	18.56%	0.668
Method 3	74.78%	10.11%	25.22%	0.571

Table 5. Quantitative accuracy evaluation in Dataset 4

Methods \ Indexes	OA	FP	FN	Kappa
Evaluation criterion	The bigger the better.	The smaller the better.	The smaller the better.	The bigger the better.
Proposed method	91.89%	2.86%	8.11%	0.707
Method 1	89.88%	3.62%	10.12%	0.612
Method 2	89.94%	3.60%	10.06%	0.640
Method 3	87.72%	4.46%	12.28%	0.595

From above tables, the OAs of the proposed method can always reach more than 85%, and the four accuracy indexes are superior to the other three comparison methods, which are consistent with the visual inspection results. In the experiments of dataset 1 and dataset 4, the accuracy of the four comparison methods is significantly higher than that of the other two sets of datasets, which is mainly due to the different initial screening results of non-building objects. The specific manifestation is that the proportions of screened non-building objects in the initial set of objects R_{origin} in Dataset 1 and Dataset 4 are 67.84% and 66.6%, respectively, which is significantly higher than that in Dataset 2 (34.66%) and Dataset 3 (38.35%).

Compared with the proposed method, the theoretical reasons underlying the lower accuracy of the compared methods are as follows. Method 1 adopts the unoptimized feature space and classifier, which includes more redundant information and redundancy and evidence conflict. Method 2 adopts the same training sample set and feature space, and the difference exists primarily in the number of decision trees. It has once again proved the necessity of selecting a reasonable number of decision trees. The optimize strategy of feature space in Method 3 is only based on the inter- and intra-class distance between features, while the difference between different datasets should be

considered as introducing the feature importance index proposed in this study.

4.7. Analysis of Relationship between the Number of Decision Trees and OA

In order to further analyze the influence of the number of decision trees on the OA and evaluate the rationality of the number of decision trees adaptively extracted in this study, the variation curves of the OA and the number of decision trees in the four groups of experiments are counted at intervals of 10 in the interval of [50,200], and the maximum, minimum and average of OAs obtained by statistics and the OA obtained in this study are respectively expressed by straight lines of different patterns, as shown in Fig. 11.

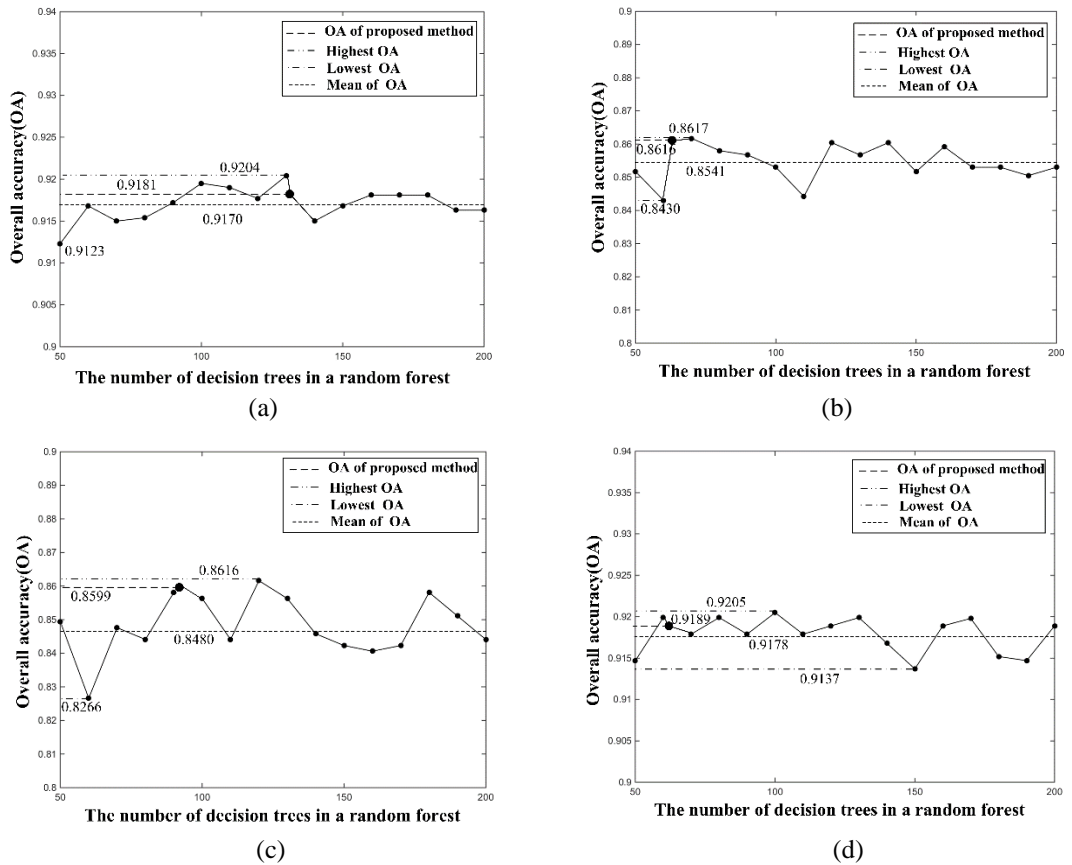


Fig. 11. Relationship between number of decision trees and OA: (a) Dataset 1; (b) Dataset 2; (c) Dataset 3; (d) Dataset 4.

As shown in above figures, the change in the number of decision trees has a significant impact on the OA, so it is necessary to select a reasonable

number of decision trees. On the other hand, although the OA corresponding to the number of decision trees extracted in this study is not the maximum value corresponding to the [50,200] interval, the difference is less than 1% and much higher than the average OA. Therefore, the proposed decision tree number adaptive selection strategy is feasible and effective, which not only improves the degree of automation, but also achieves ideal classification accuracy.

4.8. Analysis of Influence of Feature Combination on OA

In addition to the number of decision trees, this study further analyzes the influence of different types of features and feature combinations on the detection accuracy of earthquake-damaged buildings. With the initial building earthquake-damaged feature set F , the OA obtained by statistics is shown in Table 6 according to different categories and combinations of features.

Table 6. Obtained OA by different features and feature combinations

Features	OA	FP	FN	Kappa
Spectral features	0.8249	0.0661	0.1751	0.6981
Geometrical features	0.6515	0.5856	0.8091	0.1420
Texture features	0.7180	0.5509	0.7863	0.1863
Spectral + Geometrical + Texture features	0.7723	0.4722	0.7285	0.2243
The optimized features in this study	0.8599	0.0515	0.1401	0.7474

As shown in the above table, the feature set extracted in this study is the corresponding to the highest OA. Besides, the OA of spectral features is significantly higher than that of the other two features when spectral, geometrical features and texture features are used alone for classification. The reason is that the inherent geometrical features and texture features of buildings are destroyed after the earthquake, which increases the uncertainty of classification. However, spectral features are not easy to change greatly, so it is more reliable in the specific application field of earthquake-damaged buildings recognition. Moreover, the OA achieved by using three types of features is even lower than that by using spectral features alone. The reason lies in feature redundancy and evidence conflict, and the feature set optimization strategy presented in this study provides an effective solution.

5. Conclusion

Under the premise of the lack of pre-earthquake reference information, this study proposes a method of earthquake-damaged buildings detection based on decision tree and feature optimization. In the experiments of multiple sets of high-resolution remote sensing images from different regions and different sensors, the OA can reach more than 85%, and the FN is less than 6%, which

can provide key and reliable decision support information for post-earthquake emergency response and reconstruction. Its main theoretical contributions are: (1) The proposed feature set optimization strategy guided by the important indicator of feature can provide a feasible solution for the automatic construction of feature space of earthquake-damaged buildings; (2) The proposed feature set screening strategy, combined with the adaptive extraction strategy of the number of decision trees, constructs a novel and efficient optimized RF model for earthquake-damaged buildings recognition.

Author Contributions: Conceptualization, C.W.; methodology, C.W. and X.Q.; software, X.Q.; validation, H.L., X.Q. and H.L.; formal analysis, X.Q. and D.L.; investigation, K.Z. and H.L.; resources, C.W.; writing-original draft preparation, X.Q.; writing-review and editing, C.W.; visualization, C.W. and X.Q.; supervision, C.W., L.W., and K.Z.; project administration, C.W.

Funding: This study is supported by the Open Research Fund of Jiangxi Province Key Laboratory of Water Information Cooperative Sensing and Intelligent Processing (No. 2016WICSIP004), the Six Talent-peak Project in Jiangsu Province (No. 2019-XYDXX-135), the Jiangsu Overseas Visiting Scholar Program for University Prominent Young and Middle-aged Teachers and Presidents (No. 2018-69), the National Natural Science Foundation of China (No. 61601229), the Natural Science Foundation of Jiangsu Province (No. BK20160966).

Conflicts of Interest: All authors have reviewed the manuscript and approved to submit to this journal. The authors declare that there is no conflict of interests regarding of the publication of this article and no “self-citations” included in our manuscript.

References

1. Wu, J., Chen, P., Liu, Y.L.: High Resolution Remote Sensing Information Extraction Method for Earthquake Damaged Buildings. *Geography and Geo-Information Science*, Vol. 3. (2013)
2. Gan, T., Li, J.P., Li, X.Q.: Object-Oriented High-Resolution Remote Sensing Image of Building Seismic Damage Information Extraction. *Engineering of Surveying and Mapping*, Vol. 4, 11-15. (2015)
3. Ji, M., Liu, L., Du, R.: A Comparative Study of Texture and Convolutional Neural Network Features for Detecting Collapsed Buildings after Earthquakes Using Pre-Event and Post-Event Satellite Imagery. *Remote Sensing*, Vol.10, No.11, 1202. (2019)
4. Liu, Y., Li, Q.: Multi-Feature High-Resolution Remote Sensing Images for Seismic Damage Detection of Buildings. *Geomatics & Spatial Information Technology*, Vol. 6, No. 230, 71-74. (2018)

5. Li, Q., Zhang, J.F.: Study on the Identification of Damaged Buildings after Earthquake Based on the Fusion of Different Features. *Journal of Earthquake Research*, Vol. 3, 486-493. (2016)
6. Liu, Y., Cao, G.: Damage Detection Based on Multi-Feature Combination. *Computer Application*, Vol. 9, 2652-2655. (2015)
7. Sabuncu, A.: A Study of Earthquake-Induced Building Detection by Object Oriented Classification Approach. *Egu General Assembly Conference*. (2017)
8. Austin, C., Yang, S., James, C.: Detection of Urban Damage Using Remote Sensing and Machine Learning Algorithms: Revisiting the 2010 Haiti Earthquake. *Remote Sensing*, Vol. 10, 868. (2016)
9. Tesfamariam, S., Liu, Z.: Earthquake Induced Damage Classification for Reinforced Concrete Buildings. *Structural Safety*, Vol. 2, 154-164. (2010)
10. Sui, H.G., Liu, C.X., Huang, L.H., Hua, L.: Application of Remote Sensing Technology in Earthquake-Induced Building Damage Detection. *Geomatics and Information Science of WuHan Universe*, Vol. 7, 1008-1019. (2019)
11. Matsuoka, M., Yamazaki, F.: Use of Satellite SAR Intensity Imagery for Detecting Building Areas Damaged Due to Earthquakes. *Earthquake Spectra*, Vol. 3, 975-994. (2004)
12. Huang, F., Chen, L., Yin, K.: Object-Oriented Change Detection and Damage Assessment Using High-Resolution Remote Sensing Images, Tangjiao Landslide, Three Gorges Reservoir, China. *Environmental Earth Sciences*, Vol. 5, 183. (2018)
13. Anniballe, R., Noto, F., Scalia, T.: Earthquake Damage Mapping: An Overall Assessment of Ground Surveys and VHR Image Change Detection after L'Aquila 2009 earthquake. *Remote Sensing of Environment*, Vol. 210, 166-178. (2018)
14. Du, J.X.: Study and Implementation of Collapse Building Extraction Based on Morphology and Classification. *Shandong University of Science and Technology, QingDao*. (2015)
15. Ye, X., Qin, Q.M., Wang, J.: High Resolution Optical Remote Sensing Images are Used to Detect Earthquake Damage to Buildings. *Journal of Wuhan University (Information Science Edition)*, Vol. 1, 128-134. (2019)
16. Rafiei, M.H., Adeli, H.: A Novel Machine Learning-Based Algorithm to Detect Damage in High-Rise Building Structures. *The Structural Design of Tall and Special Buildings*, Vol.18, No.26, e1400. (2017)
17. Wang, N., Zhao, X., Zhao, P.: Automatic Damage Detection of Historic Masonry Buildings Based on Mobile Deep Learning. *Automation in Construction*, No.103, 53-66. (2019)
18. Estimation of seismic building structural types using multi-sensor remote sensing and machine learning techniques. *ISPRS Journal of Photogrammetry and Remote Sensing*, Vol. 104, 175-188. (2015)
19. Duarte, D., Nex, F., Kerle, N.: Satellite Image Classification of Building Damages Using Airborne and Satellite Image SAMPLES in a Deep Learning Approach. *ISPRS Annals of Photogrammetry, Remote Sensing & Spatial Information Sciences*, Vol. 2. (2018)
20. Wang, C.: A Novel Multi-Scale Segmentation Algorithm for High Resolution Remote Sensing Images Based on Wavelet Transform and Improved JSEG Algorithm. *Optik-International Journal for Light and Electron Optics*. Vol. 19, No. 125, 5588- 5595. (2014)
21. Tao, C., Tan, Y., Cai, H., Bo, D.U., Tian, J.: Object-Oriented Method of Hierarchical Urban Building Extraction from High-Resolution Remote-Sensing Imagery. *Acta Geodaetica Et Cartographica Sinica*, Vol. 39, 39-45. (2010)
22. Huang, X., Zhang, L.: A Multidirectional and Multiscale Morphological Index for Automatic Building Extraction from Multispectral GeoEye-1 Imagery. *Photogrammetric Engineering & Remote Sensing*, Vol. 7, No. 77, 721-732. (2011)

23. Fu, Q.K., Wu, B., Wang, X.Q.: Urban Building Extraction and Height Estimation Based on Morphological Building Index. *Remote Sensing Technology and Application*. (2015)
24. Rodriguezgaliano, V.F., Ghimire,: An Assessment of the Effectiveness of a Random Forest Classifier for Land-Cover Classification. *Isprs Journal of Photogrammetry & Remote Sensing*, Vol. 1, No. 67, 93-104. (2012)
25. Hongyan X., Yan Y.: Detection of Low-Flying Target under the Sea Clutter Background Based on Volterra Filter. *Complexity*, 1-12. (2018)
26. Li, Q., Jiao, Q.S.: Research on Seismic Damage Building Extraction in Beichuan County Based on Ground Lidar Technology Point Cloud. *Science Technology and Engineering*, Vol. 19, 244-249. (2016)
27. Qihao, C., Yuliang, N., Linlin, L.: Buildings Damage Assessment Using Texture Features of Polarization Decomposition Components. *Journal of Remote Sensing*. (2017)
28. Breiman, L.: Random Forest. *Machine Learning*, Vol. 1, 5-32. (2001)
29. Dou, J., Yunus, A.P., Tien, B.D.: Assessment of Advanced Random Forest and Decision Tree Algorithms for Modeling Rainfall-Induced Landslide Susceptibility in the Izu-Oshima Volcanic Island, Japan. *Science of The Total Environment*, Vol. 662, 332- 346. (2019)
30. Hongyan X., Xu Y.: Sound Source Localization Fusion Algorithm and Performance Analysis of a Three-Plane Five-Element Microphone Array. *Applied Sciences*, Vol. 9, No. 12, 2417. (2019)
31. Xu Y., Hongyan X.: Sound Source Omnidirectional Positioning Calibration Method Based on Microphone Observation Angle. *Complexity*, 1-14. (2018)
32. Wang, Z., Liu, C.: Object Level Classification of Low Altitude Remote Sensing Image Information of Unmanned Aerial Vehicle Based on SEaTH Algorithm in Lushan Earthquake. *Journal of Seismological Research*, Vol. 41, No.186, 19-25. (2018)

Chao Wang obtained his Ph.D. degree from Hohai University of China in 2014. He is currently an associate professor in Nanjing University of Information Science and Technology. His research interests include high-resolution remote sensing image processing and computer vision.

Xing Qiu is a graduate student in Nanjing University of Information Science and Technology. Her research interests include high-resolution remote sensing image processing and computer vision.

Hui Liu is an associate professor in Jiangxi University of Science and Technology. His research interests include high-resolution remote sensing image processing and pattern recognition.

Dan Li is a graduate student in Jiangxi University of Science and Technology. Her research focuses on high-resolution remote sensing image processing.

Kaiguang Zhao is an assistant professor of Environmental Modeling and Spatial Analysis in Ohio State University, U.S.A. His research focuses on mapping, monitoring, modeling, and managing terrestrial environments across scales, especially in the context of global environmental changes.

Lili Wang is a graduate student in Jiangxi University of Science and Technology. Her research focuses on high-resolution remote sensing image processing.

Received: August 17, 2019; Accepted: January 02, 2020

Multi-Agent Cooperation Q-Learning Algorithm Based on Constrained Markov Game

Yangyang Ge¹, Fei Zhu^{1,2}, Wei Huang¹, Peiyao Zhao¹, and Quan Liu¹

¹ School of Computer Science and Technology, Soochow University,
Suzhou Jiangsu 215006, China
20184227043@stu.suda.edu.cn, {zhufei, huangwei}@suda.edu.cn,
20195427013@stu.suda.edu.cn, quanliu@suda.edu.cn

² Provincial Key Laboratory for Computer Information Processing Technology,
Soochow University, Suzhou 215006, China

Abstract. Multi-Agent system has broad application in real world, whose security performance, however, is barely considered. Reinforcement learning is one of the most important methods to resolve Multi-Agent problems. At present, certain progress has been made in applying Multi-Agent reinforcement learning to robot system, man-machine match, and automatic, etc. However, in the above area, an agent may fall into unsafe states where the agent may find it difficult to bypass obstacles, to receive information from other agents and so on. Ensuring the safety of Multi-Agent system is of great importance in the above areas where an agent may fall into dangerous states that are irreversible, causing great damage. To solve the safety problem, in this paper we introduce a Multi-Agent Cooperation Q-Learning Algorithm based on Constrained Markov Game. In this method, safety constraints are added to the set of actions, and each agent, when interacting with the environment to search for optimal values, should be restricted by the safety rules, so as to obtain an optimal policy that satisfies the security requirements. Since traditional Multi-Agent reinforcement learning algorithm is no more suitable for the proposed model in this paper, a new solution is introduced for calculating the global optimum state-action function that satisfies the safety constraints. We take advantage of the Lagrange multiplier method to determine the optimal action that can be performed in the current state based on the premise of linearizing constraint functions, under conditions that the state-action function and the constraint function are both differentiable, which not only improves the efficiency and accuracy of the algorithm, but also guarantees to obtain the global optimal solution. The experiments verify the effectiveness of the algorithm.

Keywords: Markov game, Distributed perception, Multi-Agent cooperation, constrained Markov decision process.

1. Introduction

Multi-Agent System (MAS) is a combination of several sub-agents, which decomposes a large complex system into smaller and intercommunicating subsystems that are relatively manageable [1]. MAS is evolved from distributed artificial intelligence, it is applied in various areas such as intelligent robot, traffic control, distributed decision making,

business management, virtual reality and so on [2]. At present, a new mechanism that combines MAS and reinforcement learning is gradually considered to be a research hotspot [3].

Multi-Agent reinforcement learning (MARL) is to apply reinforcement learning algorithm to MAS [4]. Littman in 1900s put forward MARL with Markov Decision Process (MDP) being the contextual framework, which offered a simple mathematical framework for solving most of reinforcement learning problems [5]. MARL possesses certain properties, such as autonomy, distributivity, consistency and so on, and abilities such as learning, reasoning and self-organizing [6]. According to different learning objectives, MARL can be divided into full cooperation task, full competition task and hybrid task [7].

In full-cooperation stochastic game, agents are not making decisions independently, but are cooperating with each other trying to achieve a mutual goal in a parallel way. They share the same reward function and maximize it adopting greedy strategy [8]. Littman M proposed the Team-Q algorithm, which solved the cooperation problem among agents by hypothesizing an optimal union action [9]. Lauer M and Riedmiller M proposed the Distributed-Q algorithm, which solved the cooperation problem among agents without hypothesizing the coordination condition, with an infinite computational cost. It shares the same computation complexity with the single agent Q-learning algorithm [10]. This method, however, is suitable only for deterministic problems with non-negative reward functions. All the algorithms introduced above are limited, they all rely on the precise measurement of the state. Some of them need also the precise measurement of influence to an agent from other agents and may suffer from curse of dimensionality. On the other hand, they ignore the safety problem of agents and other constraint conditions.

In full-competition stochastic game, the agent maximizes its own reward while minimizes reward of its opponents [11]. Minimax-Q algorithm is a full-competition stochastic game that calculates policies and values through minimax principle [12].

In hybrid task, the reward function of the agent is not restricted [13], which is suitable for a selfish agent. In the game many algorithms are only for static tasks based on the concept of equilibrium in game theory [14]. In the process of a hybrid task, an agent needs to know the actions and rewards of other agents. The equilibrium selection problem arises when different agents obtain different policies. Nash Q-learning is a common method for solving this problem [15]. Correlation Equilibrium Q-learning (CE-Q) method solves the equilibrium problem with the concept of correlation equilibrium [16]. Asymmetric Q-learning solves the equilibrium problem using the leader-follower equilibrium. The follower needs not model the Q-table of the leader, but the leader should know how its followers choose their actions [17].

Traditional MARL ignores the safety problem of the agent which is inescapable in practical use. To solve the problem, in this paper we propose a Constrained Multi-Agent Cooperation Q-learning (CMACQ) algorithm based on constrained Markov game. Compared with traditional methods, this algorithm can ensure the safety of the agent, avoiding dangerous states and their consequences. In this method, before the agent executes an action, it determines all the safe executable actions based on the current state, and chooses the optimal one according to the greedy strategy as well as interactions with other agents. The algorithm ensures the safety of each agent when agents are cooperating to complete the task.

This paper is organized as follows. In section 2, we introduce the related concepts and studies. In section 3, we formalize our model and transform the model into a convex model by linearizing constraint functions, where we exploit the Lagrange multiplier method to obtain the safe action and choose the optimal safe action according to the greedy strategy as well as interactions among agents. In section 4, we compare CMACQ with MACQ in the firefighting-through-multi-agent-cooperation experiment and Deep Sea Fishing experiment. In section 5, a summary of CMACQ and a discussion of future work are presented.

2. Related Work

2.1. Reinforcement Learning

In reinforcement learning tasks, the agent detects the environment and takes actions to obtain the largest long-term cumulative reward which is a valuable encouraging signal [18]. Markov Decision Process framework is used to solve most of the reinforcement learning problems, which is denoted by a tetrad (S, A, P, R) [19] where S is the set of states which contains finite numbers of elements, A is the set of actions, P is the state transition probability and R is the reward function. In MDP, the state transition probability contains actions, which is expressed as follows [20]:

$$P_{ss'}^a = P[s_{t+1} = s' | s_t = s, a_t = a] \quad (1)$$

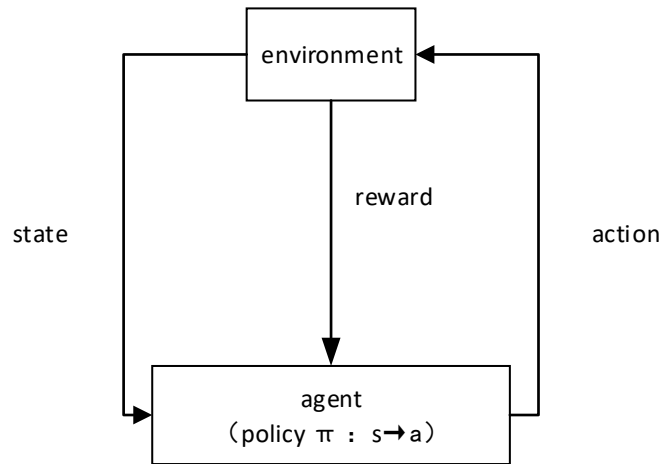


Fig. 1. The illustration of reinforcement learning.

At time t , the agent is in current state s_t , chooses action a according to policy π , receives a feedback from the environment and proceeds to the next state s_{t+1} , and obtains the reward r_t . The policy π , which is divided into deterministic policy and non-deterministic

policy, denotes the mapping from s_t to a_t . In reinforcement learning, the cumulative discounted reward is defined as [19]:

$$R_t = \sum_{i=t}^T \gamma^{i-t} r_i \quad (2)$$

where the discounted factor $\gamma \in (0,1]$, R_t is the reward value from time t to T . The reinforcement learning model is showed as figure 1 [19].

The state-action value function, $Q^\pi(s,a)$, is used to evaluate policies, which denotes the sum of cumulative rewards when the agent chooses action a_t according to the policy π under the current state s_t . The function [19] is shown as:

$$Q^\pi(s, a) = E[R_t | s_t = s, a_t = a, \pi] \quad (3)$$

As the iteration continues, the state-action value will converge to be optimal. Although the optimal policy may not be unique, the optimal state-action value is unique, as shown in equation (4) [19]:

$$Q^*(s, a) = \max_{\pi} E[R_t | s_t = s, a_t = a, \pi] \quad (4)$$

The state-value function $V(s)$ denotes the expectation of all state-action value functions when the agent follows policy π under the state s_t , which is calculated as equation (5) [19]. As the policy iteration continues, the state-value function converges to a unique optimal one. The optimal state-value function is obtained by equation (6) [19]:

$$V^\pi(s) = E[R_t | s_t = s, \pi] \quad (5)$$

$$V^*(s) = \max_{\pi} E[R_t | s_t = s, \pi] \quad (6)$$

Q-Learning algorithm [21] is a typical off-policy algorithm and is one of the most widely used reinforcement learning method. The algorithm defines a Q function, and updates this function using equation (7) [22] and equation (8) [22], that is, TD error, and finally obtains the converged optimal state-action value:

$$\delta_t = r_{t+1} + \gamma \max_a Q_t(s_{t+1}, a) - Q_t(s_t, a_t) \quad (7)$$

$$Q_{t+1}(s_t, a_t) = Q_t(s_t, a_t) + \alpha_t \delta_t \quad (8)$$

where t denotes time step, α_t is learning rate, δ_t is temporal difference error [22], a' is the action taken at next state s_{t+1} . When t tends to infinity, the optimal control policy is obtained [23], [24].

2.2. Multi-Agent Cooperation

With rapid development in areas such as sensor technology, wireless communication technology and computer vision, intelligent robot system was transformed from stand-alone system into multicomputer system. Multi-Agent system is widely applied in which the cooperation between agents played an important role [25]. Therefore, to study the cooperation between two agents is a key procedure when developing Multi-Agent systems [26]. At present, the research is divided into two categories: one is to apply the research of Multi-body behaviors to the cooperation of agents; the other is to concentrating on programming and solving of the problem [27], [28]. Allen et al studied the effect of minor changes on social evolution by assessing the cooperation tendencies

in many different population structures [29]; Mcavoy studied the effect of Evolutionary Game Dynamics on changing the scale of agents by public good game [30]; Engesser et al applied the cognitive programming to Multi-Agent system and solved program tasks in a decentralized way, by expanding the cognitive programming framework and the perspective conversion [31].

Reinforcement Learning is one of algorithms to solve cooperation problems among agents which treats Cooperative Game as the core issue in the cooperation research. Agents cooperate by delivering message with each other to achieve the mutual objective [32]. In game mechanism, mechanisms that promote cooperation are divided into categories of strong-weak reciprocity, network reciprocity and group selection [33], [34]. Elise discussed the internal mechanism of strong reciprocal behaviors: when treachery appeared in an agent group, the strong reciprocal individual would conduct a altruism punishment to the individual who betrayed others, which made betraying individual to be more cooperative [35]; Perolat et al modeled subjects who occupied public resources using Markov Observation Model. The model revealed relations among exclusiveness, sustainability and inequality and demonstrated the solution, which improved efficiency of resource management [36]; Tuyls et al put forward the LDQN algorithm, which imported toleration policy to Deep Q-network that updated passive policy with leniency methods, thus to improve convergence and stability of Multi-Agent cooperation algorithm [37]; Hwang et al combined the multi-agent cooperative Q-learning algorithm with Stochastic recording real-valued unit to solve the problem that the agent is prone to fall into local solution [38].

2.3. Constrained Markov Game

Reinforcement learning, based on MDP, consists of an agent and several states. While MARL is a game of stochastic multi-player cooperation game, including more than one player and state, based on which a Markov game (stochastic) is defined [39]. Markov game is expressed as a tuple $(n, S, A_1, \dots, A_n, T, \gamma, R_1, \dots, R_n)$ [40], where n is the number of players, S is the set of states. The state here is referred as the union state of all players at a certain time instant. $T: S \times A_1 \times \dots \times A_n \times S \rightarrow [0, 1]$ is the transition function, $A_i (i=1, \dots, n)$ is the action set of player i , $\gamma \in [0, 1]$ is the discounted factor, $R_i: S \times A_1 \times \dots \times A_n \times S \rightarrow \mathbb{R}$ is the reward function of player i . In a stochastic game, the transition function probability of next states is determined the current state and the action taken by the player. A reward function $R_i(s, a_1, \dots, a_n, s')$ denotes the reward obtained by the player after the player takes the union action (a_1, \dots, a_n) and transfers from state s to state s' . Similar with MDP, the stochastic game also possesses markov property [41], that is, the next state and reward of the player depends only on the current state and current action of all players. Multi-Agent reinforcement learning model is shown as figure 2 [42].

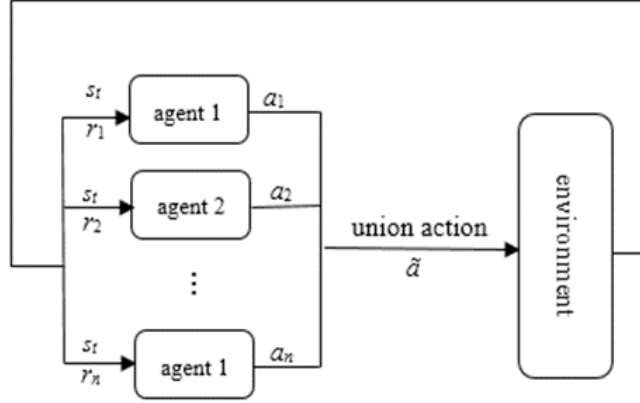


Fig. 2. Model of multi-agent reinforcement learning

Constrained Markov game adds constraint function to the Markov game, which is denoted by a tuple $(n, S, A_1, \dots, A_n, T, \gamma, R_1, \dots, R_n, C)$, where C is the constraint function set. In constrained MDP model, the goal of the agent is changed into maximizing the reward function of the whole plot under condition that the agent satisfies the constraints. The constraint function set is shown as follows:

$$C = \{c_{ij} : S \times A_i \times \dots \times A_n \rightarrow \mathbb{R} \mid j = 1 \dots k\} \tag{9}$$

where k is constant and denotes the number of constraint functions.

2.4. Lagrange Multiplier Method

Adding constraints to the objective function can ensure the safety of the agent in the learning process. Lagrange multiplier method is to find optimal solution in the case of equality constraint. Optimization problem with inequality constraint can be solved using Lagrange multiplier method and Karush Kuhn-Tucker conditions (KKT) [43]. KKT is a sufficient and necessary condition only when the model is convex, otherwise it is only the necessary condition when used to decide whether the solution obtained by Lagrange multiplier method is optimal [45].

$$\begin{aligned} &\max g(y) \\ &s.t. \quad c_j(y) = C_j, \quad j = 1, 2, \dots, k' \\ &\quad \quad c_j(y) \leq C_j, \quad j = k' + 1, \dots, k \end{aligned} \tag{10}$$

The above model is the optimization problem with inequality constraints, where $c_j(y) = C_j$ and $c_j(y) \leq C_j$ are abstract expressions of constraint function, γ is discount factor, $0 < \gamma < 1$, k is constant and denotes the number of constraint funtions. k' is constant and denotes the number of equality constraint functions. The number of inequality constraint functions is denoted by $k - k'$. The optimal solution satisfies $\lambda = 0$ or $c_j(y) - C_j = 0, j = k' + 1, \dots, k$, such that when y satisfies the strict inequality, the

constraint functions are not effective. Only when the constraint functions are equality constraints can the constraint functions be effective. Therefore, the optimization problem with inequality constraints will be turned into the problem with equality constraints and then be solved using Lagrange multiplier method, which simplifies the problem [44]. The model is shown as follows:

$$\max L(y, \lambda) = g(y) - \lambda_j(c_j(y) - C_j) - \lambda'_j(c'_j(y) - C'_j) \quad (11)$$

where the independent variable y is a local optimal solution when satisfying $\nabla_y L(y, \lambda) = 0$ and $\nabla_{\lambda} L(y, \lambda) = 0$ [45], which is obtained using gradient descent method.

When the model is convex, the local optimal solution is the same as the global optimal solution. λ_j is lagrange undetermined multipliers, which denote changes of objective functions in accordance with constraint functions. Since the optimal solution satisfies the constraint $c_j(y) - C_j = 0$, the value of λ_j will not affect the final solution.

3. Constrained Multi-Agent Cooperation Q-learning

Traditional Multi-Agent learning algorithms ignore the safety issues for simplicity while the safety issues have been proved to be crucial for completing the task. To handle the safety problem of Multi-Agent system, we propose Constrained Multi-Agent Cooperation Q-learning algorithms that encode the safety requirement as constraints to ensure a stable Multi-Agent system.

3.1. Model Design

CMACQ model is described by the tuple $(n, S, A_1, \dots, A_n, T, \gamma, R_1, \dots, R_n, C)$ which is shown as figure 3. In the current state, the state-value function of agent i is:

$$\begin{aligned} V_{\pi^i}(s) &= E[R_t | s_t = s] \\ &= E\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | s_t = s\right], i = 1, \dots, n, \forall s \in S \\ \text{s.t. } c_j(s_t, a_t^i) &= C_j, j = 1, \dots, k', i = 1, \dots, n \\ c_j(s_t, a_t^i) &\leq C_j, j = 1, \dots, k, i = 1, \dots, n \end{aligned} \quad (12)$$

Given the current state and the union action of all players, the state-value function of agent i is:

$$\begin{aligned} Q_{\pi^i}(s, a^1, \dots, a^n) &= E[R_t | s_t = s, a_t^1 = a^1, \dots, a_t^n = a^n] \\ &= E\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | s_t = s, a_t^1 = a^1, \dots, a_t^n = a^n\right] \\ & i = 1, \dots, n, \forall s \in S \end{aligned}$$

$$\begin{aligned}
 s.t. \ c_j(s_t, a_t^i) &= C_j, \ j=1, \dots, k', \ i=1, \dots, n \\
 c_j(s_t, a_t^i) &\leq C_j, \ j=1, \dots, k, \ i=1, \dots, n
 \end{aligned}
 \tag{13}$$

In the constrained case, when the optimal policy $\pi^i, i=1, \dots, n$ is obtained, the corresponding optimal state-value function and optimal state-action value function are obtained at the same time, defined as equation (14) and (15).

$$V_{\pi^i}(s) = \max_{\pi^i} V_{\pi^i}(s), \forall s \in S \tag{14}$$

$$Q_{\pi^i}(s, a^1, \dots, a^n) = \max_{\pi^i} Q_{\pi^i}(s, a^1, \dots, a^n), \forall s \in S \tag{15}$$

Index set $\{1, 2, \dots, k'\}$ and $\{k'+1, k'+2, \dots, k\}$ denote the equality constraint index set and inequality constraint index set respectively. State s_t denotes the state of the agent at time t , a_t^i denotes the action chosen by the agent i under the state s_t .

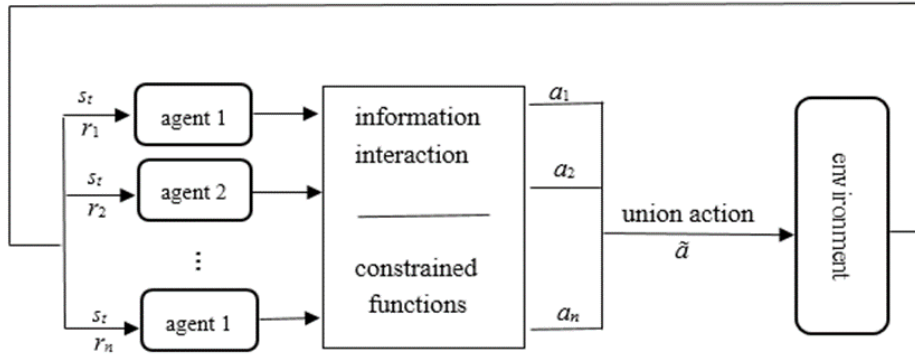


Fig. 3. Model of constrained multi-agent reinforcement learning

In CMACQ, by adding constraint conditions, the state set can be divided into safe state set and unsafe state set and the action set into safe and unsafe ones. By doing this, the safety problem of the agent can be solved at the beginning of the schedule and can also decide whether a state is safe or not using constrained conditions. Feasible region of CMACQ is shown as:

$$\bar{S} = \left\{ s \mid \bar{c}_j(s) = C_j, \ j=1, \dots, k'; \bar{c}_j(s) \leq C_j, \ j=k'+1, \dots, k \right\} \tag{16}$$

$$\bar{A}^i = \left\{ a^i \mid c_j(s, a^i) = C_j, \ j=1, \dots, k'; c_j(s, a^i) \leq C_j, \right. \\ \left. j=k'+1, \dots, k, \ i=1, \dots, n \right\} \tag{17}$$

$\bar{c}_j(s) \leq C_j$ and $c_j(s, a^i) \leq C_j$ are standard forms of inequality constraints concerning state set and action set. If some states satisfies the forms $\bar{c}_j(s) \geq C_j$ and $c_j(s, a^i) \geq C_j$, these forms can be turned into standard forms by multiplying both sides of inequations with -1. For inequality constraints, if there exists $j_0 \in \{k'+1, \dots, k\}$, such that $\bar{c}_{j_0}(s) < \bar{C}_{j_0}$ and $c_{j_0}(s, a^i) < C_{j_0}$, the i_0 th constraint in state s_t is not effective, and can be removed. The effective constraint set is denoted as ξ . Under the effective constraints, CMACQ model can be described as follows:

$$\begin{aligned}
& \arg \max_{a_{t+1}^i} Q^i(s_{t+1}, a_{t+1}^1, \dots, a_{t+1}^n) \leftarrow Q^i(s_t, a_t^1, \dots, a_t^n) \\
& \quad + \alpha \left(R_{t+1} + \gamma Q(s_{t+1}, a_{t+1}^1, \dots, a_{t+1}^n) - Q(s_t, a_t^1, \dots, a_t^n) \right) \\
& \quad s.t. c_j(s_t, a_t^i) = C_j, j \in \xi, i = 1, \dots, n
\end{aligned} \tag{18}$$

The above learning model includes only equality model so as to simplify the model and improve the performance of the algorithm.

3.2. Model Solution

CMACQ algorithm based on constrained Markov game ensures the safety of all agents by adding multi-dimensional constraints, which makes traditional solutions no more suitable. To solve CMACQ model accurately and effectively, in this paper we adopt Lagrange multiplier method to decide all the safe and optimal actions for the agent under the current state. Lagrange multiplier method requires that the objective function and the constraint function are first order continuous differentiable. The objective function satisfies the condition when time t is continuous. The constraint function is not necessarily first order continuous differentiable, which can be solved by linearization of the constrained function. Since the next state of agent is decided by the current state and union action of all agents, it can be concluded that:

$$s_t' = f(s_t, a_t^1, \dots, a_t^n) \tag{19}$$

$$\bar{c}_j(s_t') \doteq c_j(s_t, a_t^i), i = 1, \dots, n, j \in \xi \tag{20}$$

When solving the model, to ensure that the solution is globally optimal, the objective function and constraint function need to be convex. It can be known from the model that the objective function is convex while the constraint function is not. Therefore, the constraint function is linearized. Since linear function is always convex, so is the linearized constraint function. By doing this the globally optimal solution is guaranteed. Linear approximation of constraint function is shown as:

$$c_j(s_t, a_t^i) \approx \bar{c}_j(s_t) + g(s_t; \omega_j)^T a_t^i, i = 1, \dots, n, j \in \xi \tag{21}$$

In the above equation, $g(s_t; \omega_j)$ takes st as input, and output a vector sharing the same dimension with a_t^i , which can be obtained by solving the following equation:

$$\begin{aligned}
& \arg \min_{\omega_j} \sum_{(s_t, a_t^1, \dots, a_t^n, s_t') \in D} \left(\bar{c}_j(s_t') - \left(\bar{c}_j(s_t) + g(s_t; \omega_j)^T a_t^i \right) \right)^2 \\
& \quad D = \left\{ (s_t, a_t^1, \dots, a_t^n, s_t') \right\}, i = 1, \dots, n, j \in \xi
\end{aligned} \tag{22}$$

where set D is composed of tuples $(s_t, a_t^1, \dots, a_t^n, s_t')$ denoting that agent i , under the current state s_t , executes action a_t^i and switches to next state s_t' . Optimal solution of the objective function is contained in set D .

After implementing linear approximation, the learning model is shown as:

$$\begin{aligned}
& \arg \max_{a_{t+1}^i} Q^i(s_{t+1}, a_{t+1}^1, \dots, a_{t+1}^n) \leftarrow Q^i(s_t, a_t^1, \dots, a_t^n) \\
& + \alpha \left(R_{t+1} + \gamma Q(s_{t+1}, a_{t+1}^1, \dots, a_{t+1}^n) - Q(s_t, a_t^1, \dots, a_t^n) \right) \\
& \text{s.t. } c_j(s_t, a_t^j) \approx \bar{c}_j(s_t) + d(s_t; \omega_j)^T a_t^j = C_j, j \in \xi, i = 1, \dots, n
\end{aligned} \tag{23}$$

Utilizing Lagrange multiplier method, the following equations should be solved:

$$\begin{aligned}
& \arg \max_{a^*} \{ Q^i(s_{t+1}, a_{t+1}^1, \dots, a_{t+1}^n) \leftarrow Q^i(s_t, a_t^1, \dots, a_t^n) \\
& + \alpha \left(R_{t+1} + \gamma Q(s_{t+1}, a_{t+1}^1, \dots, a_{t+1}^n) - Q(s_t, a_t^1, \dots, a_t^n) \right) \\
& - \sum_{j \in \xi} \lambda_j \left(\bar{c}_j(s_t) + d(s_t; \omega_j)^T a_t^i - C_j \right) \}, i = 1, \dots, n, j \in \xi
\end{aligned} \tag{24}$$

To avoid falling into locally optimal solution, the linearized constrained function is adopted to guarantee the solution in globally optimal.

3.3. Algorithm Description

In CMACQ, agents cooperate with each other to achieve the goal. During the learning process, under the condition that the constraints are satisfied, the agent chooses an action based on its current state and on observing the action-value function. The model ensures the safety of all agents through constraint function. That a state is safe means that all agents are safe.

Algorithm 1: Constrained Multi-Agent Cooperation Q-Learning

Input: state set S , union action set A , and reward function

Output: safety state sequence and safety union action corresponding to each safety state sequence after training

1: Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$, $\gamma \in (0, 1]$

2: **Initialize:**

a) state-action value function $Q^i(s, a^1, \dots, a^n), \forall s \in S, a^i \in A^i, i = 1, \dots, n$

b) Lagrange multiplier $\lambda_j, j = 1, \dots, k, k \in N^+$

c) parameter ω_j

d) $D = \{(s, a^1, \dots, a^n, s')\}, \forall s, s' \in S, a^i \in A^i, i = 1, \dots, n$

3: **Loop forever** (for each episode):

4: Initialize initial state s

5: **Loop forever** (for each step of episode):

6: Obtain $g(s; \omega_j)$ by solving the formula

$$\omega_j^* \leftarrow \arg \min_{\omega_j} \sum_{(s, a^1, \dots, a^n, s') \in D} \left(\bar{c}_j(s') - \left(\bar{c}_j(s) + g(s; \omega_j)^T a^i \right) \right)^2, \forall a^i \in A^i, i = 1, \dots, n$$

- 7: The constraint function is approximated linearly:

$$c_j \leftarrow \overline{c_j}(s) + g(s, \omega_j^*)^T a^i, \forall a^i \in A^i, i = 1, \dots, n$$
- 8: The Lagrange multiplier method is used to obtain the action

$$a_*^i \leftarrow \arg \max_{a^i} \{Q^i(s, a^1, \dots, a^n) - \sum_{j \in \xi \cup \zeta'} \lambda_j (c_j - C_j)\}, \forall a^i \in A^i, i = 1, \dots, n$$
- 9: Agent i take action $a_*^i, i = 1, \dots, n$, observe $r^i, i = 1, \dots, n$, next state s' ,
- 10:

$$Q^i(s, a_*^1, \dots, a_*^n) \leftarrow Q^i(s, a_*^1, \dots, a_*^n) + \alpha [R_{t+1} + \gamma \max_{a^1, \dots, a^n} Q^i(s', a_*^1, \dots, a_*^n) - Q^i(s, a_*^1, \dots, a_*^n)]$$
- 11: Agent $i, i = 1 \dots, n$ moves to the next state: $s \leftarrow s'$
- 12: **Until** s is terminal
-

Each s of state set S represents a union state $s=(s_1, s_2, \dots, s_n)$, which is different from action set. In $a^i \in A^i$, a^i represents available actions of agent i , and A^i represents all the available actions. According to steps 6-8 in the above algorithm, the constraint function is solved using Lagrange multiplier method so as to decide all the safe and executable actions that each agent can choose under the condition that the long-term cumulative reward is obtained. λ_j is Lagrange multiplier. Step 6 and step 7 describe the linear approximation constraint function to ensure that the constraint functions are differentiable. If the initial constraint function is differentiable, then these two steps are skipped. In step 8, each agent uses Lagrange multiplier to work out the optimal action value under constraints. In step 9, each agent is in safe state after choosing the safe and optimal action. CMACQ ensures the safety of each agent under the condition that the globally optimal solution is obtained.

4. Experiment and Analysis

CMACQ algorithm based on constrained Markov game is suitable for determining a policy for multiple agents under a constrained condition and obtaining the optimal long-term cumulative reward. The constrained algorithm introduced in this paper is mainly to solve safety problem of multiple agents. In the experiment, a dangerous state is defined as a state that causes great damage when chosen by the agent.

To simplify the solution procedure, if problem is discrete and state set and action set are relatively small, the next state can be judged to be safe or not through only the constraint function and can decide all suitable actions for the agent. If the state set and action set are large and are even continuous, the safe state is calculated through Lagrange multiplier method.

To ensure the safety of each agent, CMACQ adopts constraint function to prevent each agent from falling into a dangerous state. The distance between the agent's current state and dangerous state is measured by Manhattan distance [46]. To ensure the safety of the agent, the Manhattan distance between the agent's next state and the dangerous state should be equal or greater than 1. The distance described above is defined as safe distance d , as showed below:

$$d(s_{t+1}, \tilde{S}) > 0 \quad (25)$$

where \tilde{S} is the set of all dangerous states, s_{t+1} is the next state for the agents.

4.1. Firefighting through Multi-Agent Cooperation Experiment

The environment of the firefighting through multi-agent cooperation experiment is a 10×10 grid world with 3 agents starting respectively at (0,9), (9,0), (9,9). In the experiment, agents cooperate with each other to complete firefighting missions, 9 origins of fire are (1,4), (3,3), (4,9), (9,3), (8,9), (7,3), (4,6), (5,2), (9,7). Each agent carries the fire-fighting equipment, and enters into a fire point, quells the fire and receives a reward of 20. Agents starting from (0, 9) and (9, 0) are able to quell four firing places while the agent starting from (9,9) is able to quell 3 firing places. When firefighting materials are used up, agents stop firefighting and go back to starting points. When all 9 firing points are quelled, a whole plot is terminated. Agents cooperate with each other to find out firing points that are not yet quelled. In the process, the agent may encounter 3 dangerous states, located at (4, 3), (8, 2), (5, 6). When the agent enters into a dangerous state, it suffers from permanent damage and receives a reward of -100. To avoid that agents take random walks in the grid world, agents receive a reward of -1 when entering into a new state other than the firing state and dangerous state. By doing this, agents are able to find the shortest path toward firing place in the shortest period of time.

In the experiment, step size α is set uniformly to 0.5, discounted factor γ is set to 1. ϵ -greedy method is adopted to explore actions in the training process so as to avoid the locally optimal solution. Policy parameter ϵ is set to 0.1. The experiment is independently operated for 50 times with 500 plots for each operation.

Figure 4 demonstrates performance of CMACQ. The result shows that agent 3 possesses the highest firefighting speed and cumulative reward for each plot. Agent 2 and agent 1 possess similar cumulative rewards with agent 2 being slightly better which is due to the distribution of firing points and the number of firefighting materials carried by the agent. The overall distribution is closer to the starting point of agent 3 and agent 3 carries only 3 units of firefighting materials, which are 1 unit lesser than those of agent 1 and agent 2, Therefore agent 3 completes the goal in the shortest period of time. Agent 3 need not wander in the grid to search for firing points such that it receives the highest long-term cumulative reward for each plot.

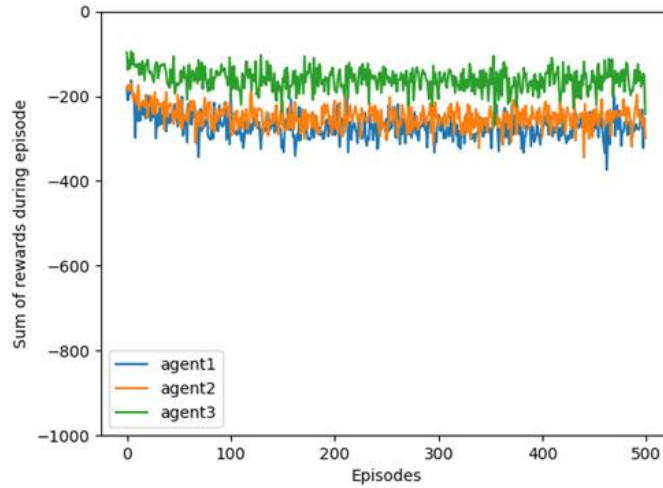


Figure. 4. Performance of CMACQ.

Figure 5 shows long-term cumulative rewards for each plot without safety constraints. The order of rewards for each agent is similar with that of Figure 4, but each reward in Figure 5 is lower than that of Figure 4. The lower rewards are because that agents fall into dangerous states and receive a reward of -100.

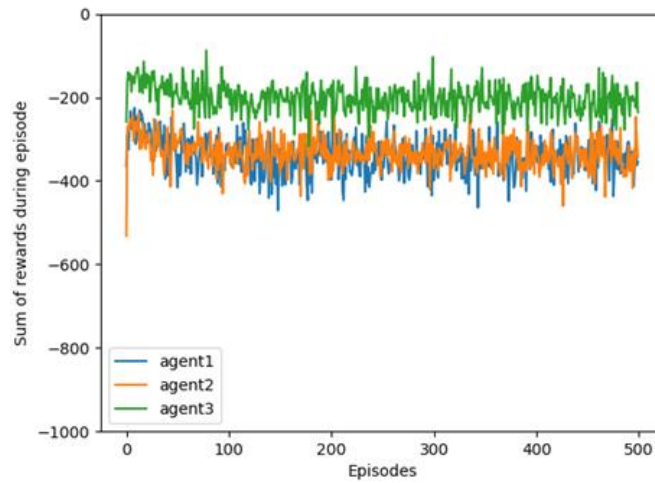


Figure. 5. Performance of MACQ

Figure 6 compares the average long-term cumulative rewards per plot between CMACQ and MACQ [38]. The result shows that CMACQ behaves better than MACQ and that CMACQ enables the agent to avoid dangerous states.

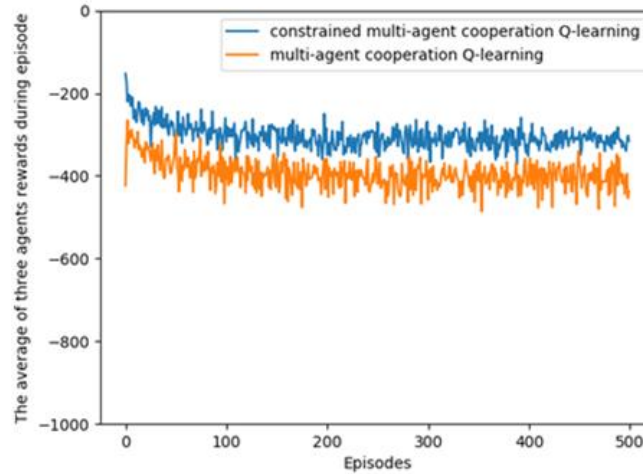


Figure 6. Average long-term cumulative rewards per plot of two algorithms.

4.2. Deep Sea Fishing Experiment

The environment of Deep Sea Fishing experiment is a 12×12 grid world with 2 agents and agents work together to catch fish while steering clear of obstacles in the sea. Agent 1 starts at (0,0) and agent 2 starts at (0,11). In the process, there are 6 dangerous states, located at (3,0), (8,1), (10,4), (11,7), (4,8), (7,11). Dangerous states simulate rocks on the bottom of the sea, sea grass and so on. In the experiment, agents need to avoid hitting rocks and getting entangled in sea grass in order to avoid irreversible damage, continue to fish and return from the sea. States of shoal of fish locate at (7,0), (5,2), (10,5), (4,6), (7,8), (11,9), (5,11). There are one unit of fish in each state of shoal of fish. Resources carried by one agent that can catch 4 units of fish. If one agent catches 4 units of fish, it can no longer fish and leave the sea.

In the experiment, the agent receives a reward of 20 when it catches 1 unit of fish. When the agent enters into a dangerous state, it suffers from irreversible damage and receives a reward of -100. To avoid that the agent takes random walks in the grid world, agents receive a reward of -1 when entering into a new state other than the state of shoal of fish and dangerous state. By doing this, agents are able to find the shortest path toward shoal of fish location in the shortest period of time.

In the experiment, step size α is set uniformly to 0.5, discounted factor γ is set to 1. The ϵ -greedy method is adopted to explore actions in the training process so as to avoid the locally optimal solution. Policy parameter ϵ is set to 0.1. The experiment is independently operated for 50 times with 500 plots for each operation.

Figure 7 shows the experimental results and compares long-term cumulative rewards for each plot solved by CMACQ and MACQ respectively. The experimental results show that CMACQ behaves better than MACQ, because CMACQ enables the agent to avoid dangerous states and ensures the safety of the agent. In the same algorithm, the

difference between the long-term cumulative rewards per plot of agent 1 and those of agent 2 is small, because the danger state and the state of shoal of fish are uniform distribution, while agent 1 and agent 2 enter the experimental environment from the upper left corner and upper right corner respectively.

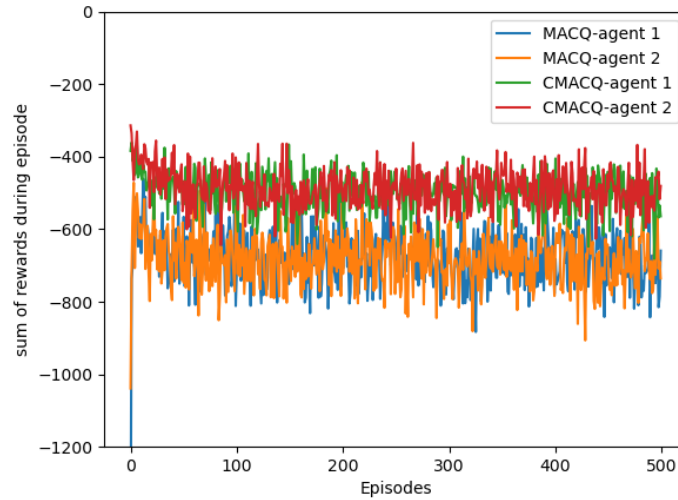


Figure. 7. Comparison of results of Deep Sea Fishing experiment of two algorithms.

The reason why the experimental result of CMACQ algorithm are better than that of MACQ algorithm is that the CMACQ algorithm can avoid the agents from entering the dangerous state and causing irreversible damage, and the safety of agents is guaranteed by the constraint conditions.

5. Conclusion

Multi-Agent system is widely applied in real world in areas such as robot system, distributed decision, traffic control, business management and so on. Reinforcement learning algorithm is a key method for solving Multi-Agent problem. However, traditional algorithms ignore the safety problem of Multi-Agent system. The agent may fall into dangerous states and suffer from great damage. To solve this problem, in this paper we introduce CMACQ algorithm based on constrained Markov game and test this method through the firefighting cooperation experiment. The result shows that CMACQ is able to handle the safety problem.

The CMACQ algorithm presented in this paper guarantees the safety of Multi-Agent system; it can also be applied to the problems where resource or cost is constrained and the area of Multi-Agent cooperation, such as robot system and traffic control, where agents work with each other and are under certain constraints. In future work we will adopt constrained algorithm to solve problems concerning limited resource, minimization of the cost and multiple objectives etc.

Acknowledgment. This work was supported by the National Natural Science Foundation of China (61303108), The Natural Science Foundation of Jiangsu Higher Education Institutions of China (17KJA520004), Suzhou Key Industries Technological Innovation-Pro prospective Applied Research Project (SYG201804), the Program of the Provincial Key Laboratory for Computer Information Processing Technology (Soochow University) (KJS1524), and a Project Funded by the Priority Academic Program Development of Jiangsu Higher Education Institutions.

References

1. Buşoniu, L., Babuška, R., Schutter, B. D.: Multi-agent Reinforcement Learning: An Overview. *Innovations in Multi-Agent Systems and Applications*, Vol. 38, No. 2, 156-172. (2010)
2. Babuška, R., Buşoniu, L., and De Schutter, B.: Reinforcement learning for multi-agent systems. *IEEE International Conference on Emerging Technologies and Factory Automation*. IEEE, 1-7. (2006)
3. Dimeas, Hatziaargyriou.: Multi-agent reinforcement learning for microgrids. *Power & Energy Society General Meeting*. IEEE, 25-29. (2010)
4. Zhang, K., Yang, Z., Liu, H., et al.: Fully Decentralized Multi-Agent Reinforcement Learning with Networked Agents. *Proceedings of the 35th international Conference on Machine Learning*. ICML 2018, Stockholm, Sweden, 1-10. (2018)
5. Littman, M. L.: Markov games as a framework for multi-agent reinforcement learning. *New Brunswick: Machine Learning Proceedings*, 157-163. (1994)
6. Foerster, J., Assael, I., De Freitas, N., et al.: Learning to communicate with deep multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*. Spain: NIPS Press, 2137-2145. (2016)
7. Leibo, J., Zambaldi, V., Lanctot, M., et al.: Multi-agent reinforcement learning in sequential social dilemmas. *Proceedings of the 16th Conference on Autonomous Agents and Multi-Agent Systems*. Singapore: AAMAS Press, 464-473. (2017)
8. Lowe, R., Wu, Y., et al.: Multi-agent actor-critic for mixed cooperative-competitive environment. *Advances in Neural Information Processing Systems*. Los Angeles: NIPS Press, 6379-6390. (2017)
9. Littman, M.: Value-function reinforcement learning in Markov games. *Cognitive Systems Research*, Vol. 2, No. 1, 55-66. (2001)
10. Lauer, M., Riedmiller, M.: An Algorithm for Distributed Reinforcement Learning in Cooperative Multi-Agent Systems. *Seventeenth International Conference on Machine Learning*. Stanford: Morgan Kaufmann Press, 535-542. (2000)
11. Lanctot, M., Zambaldi, V., Gruslys, A., et al.: A unified game-theoretic approach to multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*. Los Angeles: NIPS Press, 4190-4203. (2017)
12. Bing, S., Zhu, H., Wang, J., et al.: Optimize Pricing Policy in Evolutionary Market with Multiple Proactive Competing Cloud Providers. *IEEE International Conference on Tools with Artificial Intelligence*. (2017)
13. Kofinas, P., Dounis, A. I., Vouros, G. A.: Fuzzy Q-Learning for multi-agent decentralized energy management in microgrids. *Applied Energy*, Vol. 219, No. 3, 53-67. (2018)
14. Vidhate, D. A., Kulkarni, P.: Cooperative multi-agent reinforcement learning models (CMRLM) for intelligent traffic control. *2017 1st International Conference on Intelligent Systems and Information Management (ICISIM)*. India: IEEE Press, 325-331. (2017)
15. Hu, Junling, Wellman, et al.: Nash q-learning for general-sum stochastic games. *Journal of Machine Learning Research*, Vol. 4, No. 4, 1039-1069. (2003)

16. Greenwald, A., Hall, K., Serrano, R.: Correlated Q-learning. ICML. Washington: ICML Press, 242-249. (2003)
17. Kononen, V.: Asymmetric multi agent reinforcement learning. International Conference on Intelligent Agent Technology. Canada: IEEE Press, 336-342. (2003)
18. Mnih, V., Kavukcuoglu, K., Silver, D., et al.: Human-level control through deep reinforcement learning. *Nature*, Vol. 518, No. 7540, 529-533. (2015)
19. Sutton, Richard, S. and Barto, Andrew, G.: Reinforcement learning: An introduction. MIT press Cambridge. (2018)
20. Garc, Javier, A., Ndez, F. A.: Comprehensive survey on safe reinforcement learning. *JMLR.org*, 1437-1480. (2015)
21. C. J. C. H. Watkins, and Dayan, P.: Q-learning. *Machine Learning*, Vol. 8, 279-292. (1992)
22. Schaul, T., Quan., Antonoglou, I., et al.: Prioritized Experience Replay. *Computer Science*, 1-21. (2016)
23. Golden, R. M.: Adaptive Learning Algorithm Convergence in Passive and Reactive Environments. *Neural Computation*, 1-28. (2018)
24. Luo, B., Liu, D., Huang, T., et al.: Model-Free Optimal Tracking Control via Critic-Only Q-Learning. *IEEE Transactions on Neural Networks and Learning Systems*, 1-11. (2016)
25. Langergraber, K. E., Watts, D. P., Vigilant, L., et al.: Group augmentation, collective action, and territorial boundary patrols by male chimpanzees. *Proc Natl Acad Sci U S A*, Vol. 114, No. 28, 7337-7342. (2017)
26. Ma, Y., Lu, J., Shi, L.: Diversity of neighborhoods promotes cooperation in evolutionary social dilemmas. *Physica A Statistical Mechanics & Its Applications*, Vol. 468, 212-218. (2017)
27. Rand, D. G., Dreber, A., Ellingsen, T., et al.: Positive interactions promote public cooperation [J]. *Science*, Vol. 325, No. 5945, 1272-5. (2009)
28. Milinski, M., Semmann, D., Krambeck, H. J.: Reputation helps solve the 'tragedy of the commons'. *Nature*, Vol. 415, No. 6870, 424-6. (2002)
29. Allen, B., Lippner, G., Chen, Y. T., et al.: Evolutionary dynamics on any population structure. *Nature*, Vol. 544, No. 7649, 227-230. (2017)
30. Mcavoy, A., Fraiman, N., Hauert, C., et al.: Public goods games in populations with fluctuating size. *Theoretical Population Biology*, Vol. 121, 72-84. (2018)
31. Engesser, T., Bolander, T., Mattm, R., et al.: Cooperative Epistemic Multi-Agent Planning with Implicit Coordination. *Distributed and Multi-Agent Planning*. (2015).
32. Matignon, L., Laurent, G. J., Fort-Piat, N. L.: Hysteretic q-learning: an algorithm for decentralized reinforcement learning in cooperative multi-agent teams. In *Proceedings of IEEE/RSJ International Conference on Intelligent Robots and System (IROS)*. IEEE, San Diego, California, 64-69. (2007)
33. Coultas, J. C., Leeuwen, E. J. C. V.: Conformity: Definitions, Types, and Evolutionary Grounding. *Evolutionary Perspectives on Social Psychology*, 189-202. (2016)
34. Gintis, H.: Strong reciprocity and human sociality. *Journal of Theoretical Biology*, Vol. 206, No. 2, 169-179. (2000)
35. Seip, E. C., Dijk, W. W. V., Rotteveel, M.: Anger motivates costly punishment of unfair behavior. *Motivation & Emotion*, Vol. 38, No. 4, 578-588. (2014)
36. Perolat, J., Leibo, J. Z., Zambaldi, V., et al.: A multi-agent reinforcement learning model of common-pool resource appropriation. *Advances in Neural Information Processing Systems*. Los Angeles: NIPS Press, 3643-3652. (2017)
37. Palmer, G., tuyls, K., Bloembergen, D., et al.: Lenient multi-agent deep reinforcement learning. *Proceedings of the 17th International Conference on Autonomous Agent and Multi-Agent Systems*. Swede: AAMAS press, 443-451. (2018)
38. Hwang K., Lin Y. and Lo C.: Multi-Agent Cooperation by Q-Learning in Continuous Action Domain. *2008 First International Conference on Intelligent Networks and Intelligent Systems*. China, 111-114. (2008)

39. Chalmers, R., Scheidt, D., Neighoff, T., Witwicki, S., and Bamberger, R.: Cooperating unmanned vehicles. AIAA 1st Intelligent System Technical Conference. (2004)
40. Ying-ying, D., Yan, H., and Jing-ping, J.: Self-organizing multi-robot system based on personality evolution. IEEE International Conference on Systems, Man, and Cybernetics. (2002)
41. Sidney, G.: Analysis and Design of Swarm Based Robots Using Game Theory. Ph. D. Thesis, Ottawa, ON: Carleton University. (2009)
42. Nowé A., Vrancx P., De Hauwere YM.: Game Theory and Multi-agent Reinforcement Learning. In: Wiering M., van Otterlo M. (eds) Reinforcement Learning. Adaptation, Learning, and Optimization, Springer, Berlin, Heidelberg, vol. 12, 441-470. (2012)
43. Martyna, J.: Power Allocation in Cognitive Radio with Distributed Antenna System. Lecture Notes in Computer Science, 745-754. (2017)
44. Kuan, T. W., Wang, J. F., and Wang, J. C., et al.: VLSI Design of an SVM Learning Core on Sequential Minimal Optimization Algorithm. IEEE Transactions on Very Large Scale Integration (VLSI) Systems, Vol. 20, No. 4, 673-683. (2012)
45. Farina, F., Garulli, A., and Giannitrapani, A., et al.: Asynchronous Distributed Method of Multipliers for Constrained Nonconvex Optimization. 2018 European Control Conference. Limassol, Cyprus, Vol. 103, 243-253. (2018)
46. Blackburn, S. R., Homberger, C., Winkler, P.: The minimum Manhattan distance and minimum jump of permutations. Journal of Combinatorial Theory, Series A, Vol. 161, 364-386. (2019)

Yangyang Ge is a postgraduate of School of Computer Science and Technology, Soochow University. Her main research interests include safe reinforcement learning. E-mail: 20184227043@stu.suda.edu.cn.

Fei Zhu (corresponding author) got Ph.D. and associate professor in School of Computer Science and Technology, Soochow University. He is a member of China Computer Federation. His main research interests include reinforcement learning, text mining, and pattern recognition. E-mail: zhufei@suda.edu.cn. He is the corresponding author of this paper.

Wei Huang is associate professor in School of Computer Science and Technology, Soochow University. Her main research interests include reinforcement learning and pattern recognition. E-mail: huangwei@suda.edu.cn.

Peiyao Zhao is a postgraduate student in the Soochow University. His main research interest is reinforcement learning. He programmed the algorithms and implemented the experiments. E-mail: 20195427013@stu.suda.edu.cn.

Quan Liu got Ph.D. and professor in School of Computer Science and Technology, Soochow University. His main research interests include intelligence information processing, automated reasoning and machine learning. E-mail: quanliu@suda.edu.cn.

Received: December 20, 2019; Accepted: May 02, 2020

A K-means algorithm based on characteristics of density applied to network intrusion detection*

Jing Xu¹, Dezhi Han¹, Kuan-Ching Li^{2*}, and
Hai Jiang³

¹ Shanghai Maritime University, Shanghai, 201306, China

² Providence University, Taichung 43301, Taiwan

³ Arkansas State University, Jonesboro, Arkansas 72467, USA

Abstract. K-means algorithms are a group of popular unsupervised algorithms widely used for cluster analysis. However, the results of traditional K-means clustering algorithms are greatly affected by the initial clustering center, with unstable accuracy and low speed, which makes the algorithm hard to meet the requirements for Big Data. In this paper, a modernized version of the K-means algorithm based on density to select the initial seed of clustering is proposed. Firstly, Kd-tree is used to divide the hyper-rectangle space, so those points close to each other are grouped into the same sub-tree during data pre-processing, and the generalized information is stored in the tree structure. Besides, an improved Kd-tree nearest neighbor search is used in the K-means algorithm to prune the search space and optimize the operation for speedup. The clustering results show that the clusters are stable and accurate when the numbers of clusters and iterations are constant. Experimental results in the network intrusion detection case show that the improved version of the K-means algorithms performs better in terms of detection rate and false rate.

Keywords: Network security; K-means; Kd-tree; Network intrusion detection.

1. Introduction

In recent years, with the rapid development of new technologies such as cloud platforms, Big Data, and artificial intelligence, Internet has become an indispensable driving force for economic development and social progress [1, 2]. The Internet has

*Corresponding author

changed people's living manners, and now it is gradually covering intelligent transportation, smart cities, and others, which is exceptionally convenient for people's living [3]. However, the exploding number of big data environments of networked information data has put tremendous pressure on data processing systems, and various network security threats have become more severe in anomalies such as viruses, hackers, network attacks, and Internet frauds plague various application users [4]. The way to detect abnormal traffic in the process of using the Internet is significant in preventing further harm [5]. As the name implies, anomaly detection is the detection of abnormal behaviors [6, 7]. Through collecting and analyzing some critical points of the computer networks [8, 9], we find out whether there is any violation of the security policy and sign of being attacked in the network.

In order to prevent network anomaly flow from harming Internet users and causing the loss of personal, company, enterprise, and government information, it is necessary to find a reasonable and effective network anomaly detection method [10]. In this paper, the K-means clustering algorithm is investigated, in which network behaviors are classified, and abnormal traffic is identified through clustering.

Clustering is a type of unsupervised learning [11]. That is, the points with the same property are grouped closer and closer in space. There are several types of clustering algorithms, and can be classified as: (1) Partitioning method: the main idea of the partitioning method is to give the number of initial seeds and the clustering center points, and then solve the problem iteratively according to the distance computation formula to retrieve a final partition, as in K-means, K-medoids [12,13], (2) Hierarchical method: the datasets are solved hierarchically, and the results are obtained by merging and iterating them to meet specific conditions as in HAC [14], CURE [15], (3) Density method: according to the comparison of the density within the range of regional data points against a certain threshold, the classification result can be achieved as in a mean shift clustering algorithm, DPC [16], DBSCAN [17], (4) The grid method: taking grid as the basic clustering unit, the dataset space is divided into a grid for clustering as in STING [18], WAVE-CLUSTER, and (5) Model method: the data matches the constructed model.

In the clustering algorithm of the partitioning method, K-means is a classical algorithm that is widely used in various fields: after analyzing the internal principles and characteristics of the algorithm, Heil, Jannis et al. [19] applied K-means algorithm to soil diffuse reflection spectrum classification; Gulnashin, Fatima et al. [20] studied and compared K-means algorithm in various aspects to classify documents; Means Kimberly et al. [21] applied K-means algorithm to the diagnosis of atrial fibrillation in the emergency department in the medical field. The working principle of a fast and straightforward K-means algorithm is to compute the average value of each cluster with a simple algorithm, and then divide the classes iteratively until all the clusters are covered. At present, K-means algorithm has become one of the crucial methods of machine learning and data mining technology. With the expansion of clustering applications and the rapid growth of different kinds of data, traditional clustering algorithms are difficult to meet people's needs. Even more, investigations have focused on improving the traditional clustering algorithm to meet new development needs.

Aimed at the selection of the initial clustering center, this paper improves the K-means algorithm and applies the improved algorithm to detect abnormal network traffic. The main contributions of this paper can be summarized as follows:

- A K-means algorithm based on density to select the initial clustering center is proposed. Firstly, the Kd-tree structure is used to store the given dataset. Secondly, the initial clustering centers are selected using the data range information stored by the Kd-tree node structure so that the initial clustering centers are distributed in a region with a large data density, the relative distance between the initial clustering centers is reasonable, and the anti-noise ability is stronger. Third, the kd-tree nearest neighbor search algorithm is improved. The idea of this improved algorithm is used in the K-means algorithm. This advantage is used to prune the search space and optimize the operation for speedup. When the amount of data in the given data set is small, a clustering method using correlation nearest neighbor search is proposed for the situation of excessive pruning. Finally, the cluster partition is iteratively computed to retrieve the results, which show the stability of clustering results and accuracy under the condition where the numbers of clusters and iterations are constant.

- The algorithm proposed in this paper is superior to the traditional K-means algorithm and literature [27, 28] in detection rate and false rate in the network abnormal traffic detection scenario with a high data dimension.

The remainder of this paper is organized as follows. In the related work of Section 2, the Kd-tree data structure, and the principle of using Kd-tree to accelerate the K-means algorithm are introduced. The advantages and disadvantages of the K-means algorithm based on Kd-tree improvement are analyzed. In Section 3, the process of the traditional K-means algorithm is reviewed. An improved K-means algorithm based on density selection of initial seed is proposed. Section 4 shows the experimental details and gives the results and analyses. Finally, the conclusions and future work are provided in Section 5.

2. Related work

K-means algorithm, as a clustering method, has been used to detect anomalies in network environments. However, the K-means algorithm still has some defects. Many researchers start from data pre-processing before using clustering to improve the selection of the initial central point set for better accuracy of the K-means algorithm.

Canopy is an unsupervised pre-clustering algorithm proposed by McCallum [22] in 2000. Canopy algorithm only compares the distance between data in the same region and the central point each time. This algorithm reduces the number of comparisons, which reduces the execution time of clustering and improves the computational efficiency of the algorithm. Zhang et al. used a Canopy algorithm with density parameters as the pre-processing process of K-means [23]. The outputs are used as the cluster number and initial central point of the K-means algorithm, which improved the computational efficiency and accuracy of the algorithm. R-tree [24] and Kd-tree [25] are two improved methods based on indexes. The space points contained in each node of R-tree are included by a minimum enclosing rectangle (MBR). The core idea is to aggregate the points close to each other in the tree structure from bottom-up and cover them with a larger MBR in the upper layer until all space points are covered in the root node. Based on this idea, R-tree is widely used in spatial data indexing and partition optimization. Wei Wu et al. proposed a reverse k-nearest neighbor (RkNN) query using the R-tree structure [26]. The algorithm takes the query object as one of its K-nearest

neighbors and uses the new INCH algorithm to solve the evaluation RkNN query and its variant two-color RkNN query on two-dimensional location data.

Kd-tree is simply a multi-dimensional binary tree. The clustering iteration of the K-means algorithm can be improved according to the idea of the nearest neighbor search with the data initialized by Kd-tree. The pruning process is used to eliminate the distance computations and improve the efficiency of the K-means algorithm. Redmond et al. constructed Kd-tree and used its structural characteristics to obtain the density estimation of leaf buckets and used the max-min method to find the initial cluster seeds [27]. Kumar et al. improved the algorithm [28], from the viewpoint that density weight and distance weight are used to replace density estimation and distance estimation when selecting the initial seeds. Next, we discuss Kd-tree in detail and use the Kd-tree nearest neighbor search to accelerate the K-means algorithm.

2.1. Kd-tree

Kd-tree is a binary tree and a data structure for storing a finite set of points from an m -dimensional space. The points with similar partition-dimension can be stored in the same Kd-space. Given a dataset S , a Kd-tree can be constructed: the binary tree is built in a top-down manner, and each sub-tree is contained in a hyper-rectangle, which is containing the projection of maximum and minimum coordinate values of points in each dimension. Computing the variance in each dimension of the sub-tree to obtain the partition-dimension. The median value of partition-dimension is selected to balance both the left and right sub-trees. Then the hyperplane of the axial pair is used to divide the dataset into two subsets for two hyper-rectangles, which are used as the left and right sub-trees of the parent node. Iteration through the partition process continues until the leaf node or a specified number of leaf bucket nodes are divided. The schematic diagram of Kd-tree in two-dimensional space is shown in Fig. 1.

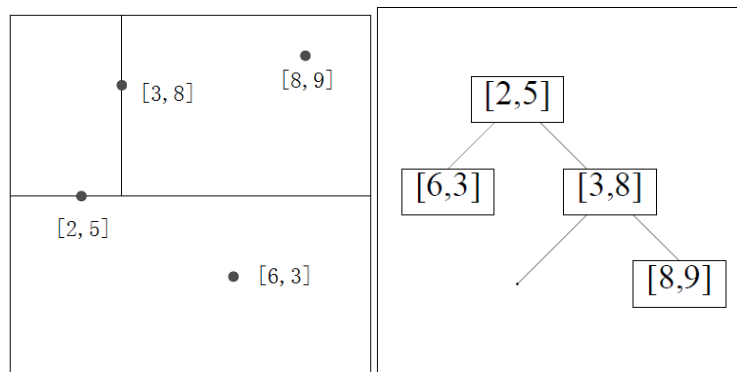


Fig. 1. Construction of Kd-tree with two-dimension data

2.2. Kd-tree nearest neighbor search

K-means algo The reason why Kd-tree can be used to reduce the number of nearest neighbor queries in K-means is that their nodes can represent a large number of points [29]. The node can store the super rectangle information of the dataset, including the Max and Min vectors. The process of finding the nearest neighbor of a node Q: Starting from the root node of dataset S, Kd-tree is accessed from top to the leaf node. The current node is denoted as “nearest neighbor point” and the minimum distance set as *dis*. A traceback program is used to find any points closer to Q in the branch of the access path. The evaluation is based on the fact that the hypersphere with Q as the center and *dis* as the radius intersects with the hyper-rectangle of the branch. The current “nearest neighbor” is updated if there is a data point that is closer to Q until you go back to the root node.

2.3. K-means accelerated by Kd-tree

This K-means algorithm with Kd-tree acceleration stores the data set in Kd-tree during data pre-processing, as displayed in Fig. 2. This algorithm starts using randomly generated cluster centers. Nodes in Kd-tree are accessed from the top down during each iteration. According to the spatial region information, the data of a specific point, and its left and right sub-trees are used to determine whether they belong to a specific central point. If it is not possible to do so, traverse the tree into the left and right sub-trees of this point until all data points are included in clustering.

The factors why Kd-tree can accelerate K-means [30] are because Kd-tree’s nearest neighbor algorithm can prune Kd-tree during the traversal: firstly, search the minimum distance between each cluster center in the candidate cluster center set and the space area represented by the node object, which is obtained by Kd-tree’s nearest neighbor query algorithm. Then, the maximum distance of the space area represented by the node object is denoted as Maxdis, and the maximum distance adopts the maximum value of the space range represented by the node. Finally, the cluster center of the smallest one whose nearest neighbor distance is greater than the maximum distance is cut off. A rectangle can represent the h in the two-dimensional space. The point in Fig. 2 represents the distribution of the subtree dataset of a particular data point, and the rectangle is the data range h.

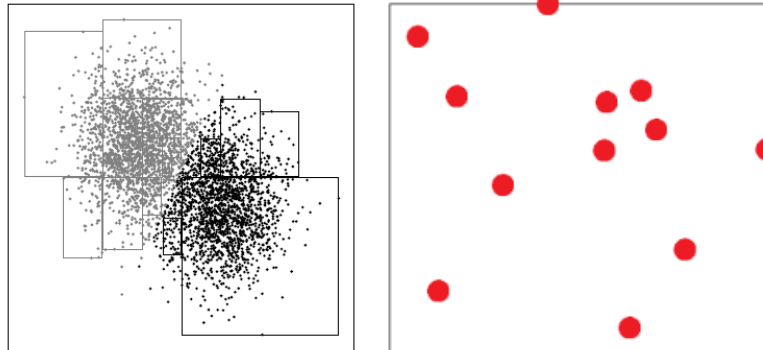


Fig. 2. Two-dimensional dataset rectangle partition

In this paper, the selection of the initial cluster centers is improved. Based on the pre-processing dataset with Kd-tree data structure, a method of correlation weighted distance computation and correlation weighted nearest neighbor search is proposed to improve the selection of clustering centers and iterative clustering computation, respectively. Based on the data structure of Kd-tree, the average density weight of the nodes in the leaf bucket is computed. A point with a maximum density as the first initial center point; then, the remaining $k-1$ initial central points are determined according to the product of the density weight of the buckets and the correlation weighted distance of the selected initial center. Finally, the improved method of correlation distance weighting based on Kd-tree acceleration is used for the iterative computation.

3. K-means algorithm

3.1. Traditional K-means algorithm

As mentioned above, the principle of the K-means algorithm is to obtain clustering by taking the initial central point as the prototype [31]. This algorithm belongs to the partitioning clustering algorithm. The letter K means the number of k used in the randomly selected initial cluster center. The Euclidean distance computation in the Cartesian coordinate system is used to classify the data into the closest cluster. Each data point belongs to a cluster. Each iteration recomputes the center point of each cluster, and the next iteration produces new clustering results. Repeat this step until the center point no longer changes or converges.

The algorithm process is as follows:

Let the input dataset be $S = \{x_1, x_2, \dots, x_n\}$, among them $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})$ where $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})$ is a data object with p dimensional characteristics.

- 1) k initial clustering centers are randomly selected from the sample S .

2) The Euclidean distance $d(x_i, c_i)$ between data points x_i and the clustering center c_i is computed. Through partitioning, data points were classified to the nearest clustering center, and k partitioning classes $L = \{l_1, l_2, \dots, l_k\}$ were obtained.

3) The center c_i of each cluster is computed and the result is used as the new cluster center.

4) Repeat steps 2 and 3 until the maximum number of iterations is reached, or the cluster difference is less than a threshold.

The output are k clusters.

K-means is a classical algorithm to solve clustering problems with the features of simplicity and fast speed. Scalability and efficiency of the algorithm for processing large datasets are maintained. When a cluster is close to the Gaussian distribution, the clustering effect is better.

However, in the K-means algorithm, an initial partition shall be determined according to the initial clustering center. Different initial centers determine the level of clustering accuracy, so that the algorithm is sensitive to the initial values. If the initially selected clustering center has a substantial deviation from the actual clustering center, the clustering effect will be deplorable. The random selection of the center of the initial cluster is not stable. If the cluster contains abnormal points, such as noises and sensitive points of isolated data, the mean deviation will be severe. The Euclidean distances from all sample points to each cluster center need to be computed in each iteration. As the number of iterations increases in the clustering center, the execution speed of the algorithm cannot meet the requirements of Big Data.

3.2. Improved K-means algorithm for selecting the initial center based on density

At this point, a K-means algorithm based on Kd-tree is proposed to select the initial center according to density. To overcome K-means algorithm's unstable clustering results and slow execution, this paper presents a method of selecting the initial clustering centers. Meanwhile, the distance weighting is used in Kd-tree nearest neighbor search and the initial centers at the weighted distance, which avoids the problem of using Euclidean distance to classify samples with similar distances but low correlation into the same category. However, samples with low correlation will cause problems due to the deviation in the same class.

Symbol definitions are shown in Table 1. The algorithm details are as follows:

Table 1. Notations

Symbol	Description
k	Number of cluster center
V	Volume of leaf node
q	Number of leaf buckets

l	Number of leaf nodes in one leaf bucket
c	Cluster center
m	Mean of leaf node
d	Dimensionality of dataset
r	Distance weighting coefficient
ε	Density estimate of leaf node
ρ	Density weight of leaf node
g	Distance estimate of leaf node
β	Correlation distance weight of leaf node

Pearson correlation coefficient [32] is used to weight the Kd-tree nearest neighbor search process and the distance computation used at the initial center point.

$$P = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right). \quad (1)$$

Formula definition: The Pearson correlation coefficient P of two continuous variables (\mathbf{x}, \mathbf{y}) is equal to the product of their covariance $\text{COV}(x, y)$ divided by their respective standard deviations. The value of the Pearson correlation coefficient is between -1.0 and 1.0. Variables close to 0 are considered not correlated, while those close to 1 or -1 are considered strongly correlated.

The distance weighting coefficient is computed as follows:

$$r = 1.0 - P + 0.001. \quad (2)$$

The value 0.001 is set so that the value of the correlation weighting coefficient is not 0.

Data have been standardized and normalized during data preprocessing.

Initial centers selection method uses correlation weighting as follows:

The Kd-tree structure is iteratively established. When the number of data sets is to a certain number, it is no longer divided. This structure is named leaf bucket. As the data is split into a specified number of child nodes, it is no longer divided and stored as a leaf bucket. As with other nodes, the information on the hyper-rectangle space of the dataset is stored on the parent node. The number of leaf buckets is q , and the number of leaf nodes in the leaf bucket is l . First, the average value of the data points in the leaf bucket is computed. Assume that the average value of data point objects in the leaf bucket stored on node x_i is m_i , that is:

$$m_i = \frac{\sum_{i=1}^l x_i}{l}. \quad (3)$$

The Max and min vectors stored in the Kd-tree are used to compute the volume of each leaf bucket:

$$V_i = \prod_{i=1}^d (x_{i\max} - x_{i\min}). \quad (4)$$

The density estimate $\mathcal{E}_i = N_i / V_i$ is computed; where N_i represents the number of data objects in the leaf bucket; the density weight is computed as:

$$\rho_i = \frac{\mathcal{E}_{sum}}{\mathcal{E}_{sum} - \mathcal{E}_i}. \quad (5)$$

where \mathcal{E}_{sum} is the sum of density estimate with $\mathcal{E}_{sum} = \sum_{i=1}^q \mathcal{E}_i$. The average value of the leaf bucket with the most considerable density weight is selected as the first center point C_1 . Then, the correlation weighted distance between the remaining leaf buckets and the center points C_i are computed. The higher the correlation weighted distance, the more significant the difference between the two points and the less the correlation will be.

$$g_i = \arg \min_{k=1 \dots q} [d(m_i - c_k)r]. \quad (6)$$

The Eq. (6) represents the minimum value of the correlation weighted distance between each initial central point and the average value of the leaf bucket. By computing $g_{sum} = \sum_{i=1}^q g_i$, the distance weights are as follows:

$$\beta_i = \frac{g_{sum}}{g_{sum} - g_i}. \quad (7)$$

$$c_k = \{m_i \mid z = \arg \max_{i=1 \dots q} (\rho_i \beta_i)\}. \quad (8)$$

The candidate central point is determined while the density of leaf bucket is large and it is far away from the initial central point. The easiest way to do so is to select it as the next initial central point C_k . However, there is a case where if a leaf bucket itself has a low density but is far from the existing initial central point, the outlier is easily misselected as the initial center, and such a point is called an abnormal point. When running the algorithm, a scheme to remove noise points before the next selection of the center algorithm is executed.

$$Noise = \beta_i / \rho_i. \quad (9)$$

Prior to the execution of each algorithm, the leaf bucket with the maximum noise value of 10% is discarded to obtain a new candidate central point set. The algorithm is as follows.

Table 2. Description of selecting the initial center algorithm

Algorithm description
Input: n-dimensional data set $S = \{x_1, x_2, \dots, x_n\}$; Cluster data quantity k
Output: k clustering centers
begin
Step1: Establish a Kd-tree.

```

Step2:Compute the density weight  $\rho_i$  of each leaf bucket
using Eq. (5).
Step3:The average value of leaf bucket with the largest
 $\rho_i$  value is selected from the candidate center data set
as the first initial center  $c_1$ .
Step4:for j=2 to q
    Compute the distance weight  $\beta_i$  of each leaf
    bucket in the candidate center dataset, and
    use Eq. (7) to compute the z value. Select the
    ith central point according to Eq. (8).
end for
Step5:Eq. (9) is used to compute the noise value Noise of
each leaf bucket in the candidate center dataset.
Discard the top 10% of noise points to generate a new
candidate center dataset. Repeat step 4, until i=k.
Output:  $(C_1, C_2, \dots, C_k)$ 
end

```

3.3. Relevance weighted nearest neighbor search

Kd-tree can accelerate the K-means algorithm, which is introduced in Section 2. This section introduces the relevance weighted nearest neighbor search. Based on the nearest neighbor search idea, data points are allocated to the nearest central point, which avoids the redundant computation of distance in the iterative process of the traditional K-means algorithm; this process saves a lot of time while improving the efficiency of the algorithm. In the improved nearest neighbor search, a correlation weighting is also used in the distance computation, and the dataset is classified into the central point with the strong correlation in the steps requiring a distance computation to make the clustering effect more accurate.

The dataset is traversed from the root node to the leaf node. Let this point be the node to calculate the weighted distance between the leaf node and the candidate central point. The node is classified into the cluster with the smallest distance, and then the backtracking operation is performed. The correlation weighted distance of each current node space area information h and the candidate central point is calculated. If the spatial region information can be used to judge that the data in the sub-tree of this node belong to the central point c , the computation of the distance from the data in the sub-tree of this node to the central point c can be reduced. If it can be judged that none of the data in h belongs to a particular central point c , the computation of the distance from the node data in the h to the central point c can be reduced. When it is determined that all the data in this node sub-tree of this node belong to c , the statistical information stored in the h can be directly classified into the central point c .

Determine the node center in the h , including the following steps:

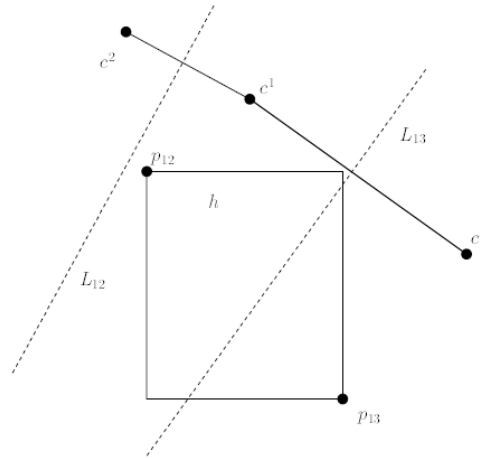


Fig. 3. Domination with respect to a hyper-rectangle

Step 1: determine if there is a minimum distance.

The point with the smallest distance by computing the distance value of $d(c_i, h)$ from the point of dataset C of the candidate centers to the h of the current node is found. If there is one minimum distance, this c_i may be the central point of the Kd-tree or its sub-Kd-tree represented by the hyper-rectangle h . If there is more than one minimum distance, the central point of the Kd-tree or its sub-Kd-tree represented by the hyper-rectangle h does not exist and h is divided into left and right subtrees for search. Through the above analysis and computation, some inferior central points can be excluded in the process of subtree traversals.

Step 2: determine whether c_i is the central point.

In step one, if the result of the operation is that there is a minimum distance, then this c_i may be the central point of h , and further judgment is needed. Take Fig. 3 as an example. To determine whether c_1 is superior to c_3 : orientation $v = (c_3 - c_1)$. If the p is picked as a member of h , the principle is to maximize the inner product of two vectors. As shown in Fig. 3, L_{12} is the decision line between centers c_1 and c_2 . L_{13} is the decision line between c_1 and c_3 . Vector v 's positive and negative dimensions are $(+, -)$. The selection principle of any point p is as large as possible on the X-axis and as small as possible on Y-axis. In this case, p_{12} is the extreme point in h in the direction $c_1 - c_2$, and p_{13} is the extreme point in h in the direction $c_3 - c_1$. The inner product of the two selected vectors is the maximum. Since p_{13} is picked, whose Euclidean distance in a two-dimensional space is closer to c_3 , so c_1 is not better than c_3 . If c_i is superior to other points, it can be judged that c_i is the central point of Kd-

tree or its sub-tree represented by h . Otherwise, c_i is not the central point of Kd-tree or its children represented by h ; Although c_i is not the central point of Kd-tree or its sub-tree represented by h , c_i is better than c_2 , which is excluded from the set of candidate points from the central point of the sub-tree of h .

Table 3. Description of the clustering algorithm

Algorithm description
Input: Kd-tree of dataset S ; Initial cluster center C improvedKmeans(S, C): begin Start from the S 's root node Kd-tree.root. if point u is the leaf node The correlation weighted distance of the point was computed and the point was classified into the clustering center with the minimum distance. Return. else for each C compute $d(c_i, \text{node.h}) * r$ if there are multiple minimum d Recursively invoke the improved kmeans algorithm. Call improvedKmeans($\text{node.lchild}, C$). Call improvedKmeans($\text{node.rchild}, C$). Return. else Set the minimum value to ct . Compute whether there is a center superior to ct in C . if(ct better then c_i) Add c_i to C_{out} . if there is only one advantage This point and its child nodes were classified into ct . $C_{new} = C - C_{out}$. Recursively call the improved K-means and pass in the filtered clustering center. Call improvedKmeans($\text{node.lchild}, C_{new}$). Call improvedKmeans($\text{node.rchild}, C_{new}$). Return. Output: k cluster centers end

For small datasets, a clustering method using correlation nearest neighbor search is designed, and the algorithm steps are presented in Table 4. In this method, the Kd-tree structure is performed to find the nearest neighbor.

Table 4. Description of the clustering algorithm

Algorithm description
Input: N-dimensional dataset $S = \{x_1, x_2, \dots, x_n\}$; Cluster data quantity k begin for i=1 to len(S) for j=1 to k The relevance weighted nearest neighbor is used to search for the nearest point x in S to center[j]. Mark it as category j . end for Remove x from S . end for Output: k cluster centers end

3.4. Algorithm analysis

Using Kd-tree leaf bucket structure to compute the density poses some problems. Since the density weight and distance weight dimensions are different, the data density calculations of different datasets are different, and the distance varies greatly. Then, the actual effect of the algorithm could be biased. Therefore, it is necessary to do some processing on these two values to eliminate such bias. Firstly, the influence of density weight is that leaf buckets with low density are underestimated. When selecting the initial center, more leaf buckets with high density are preferred. Thus, most leaf buckets with high density are selected as the initial center. The results ignored some smaller clusters in the actual situation. A method is proposed to solve this problem by using the rank value of density estimation instead of the actual density estimation [27].

Suppose there are 1024 leaf buckets. Compute the density weight of the leaf bucket and arrange it in descending order. The node with the highest density weight is assigned a density level of 1024, followed by 1023, 1022, until the lowest level is 1. The method reduces the influence of the leaf bucket with high-density weight on the center selection to a certain extent. But the density rank method does not reflect the actual density of the node, ignoring the variation of the density weight between the leaf buckets. It is assumed that after the sorting, two adjacent leaf buckets with a significant difference in density weights and two other smaller leaf buckets with different density weights are sorted; thus, the span between ranks is 1 and the degree of discrimination turns to be smaller. When the gap between the leaf bucket density values of the data is small, the gap is enlarged. A hierarchical classification strategy with a large density estimation span is proposed [28], reflecting the gap between the estimated values of the leaf bucket density. It can distinguish the density estimation value more effectively, though this method is random. Besides, there is a drawback that the buckets with the value of low-density are neglected. A mapping method that utilizes a function to slow down density growth is performed to increase the weight of the leaf bucket with a lower density value,

and maintain the leaf bucket with a higher density value. The final density level is divided into points that are not integers in the number of leaf buckets.

3.5. Performance analysis

The time consumption of the proposed algorithm mainly includes the establishment of Kd-tree and clustering, which also contains the initial central point and the partition. The time complexity of constructing a modified Kd-tree is $O(dn \log n)$ [28]. The time complexity of density weight computation in the initial cluster center is $O(dq)$, the same as in [27, 28]. The average value of the leaf bucket is the distance weight and its time complexity is $O(dql)$. The time complexity of computing the distance weight is $O(qk)$, so the time of the 2 to k initial clustering center algorithm is $\frac{1}{2}qk^2 - \frac{5}{2}qk + 3q$. The partition time is the time spent on browsing the tree multiplied by the pruning rate. If the pruning rate is α , the partition time is $O((1-\alpha) \log n)$, ($0 \leq \alpha < 1$).

The clustering time complexity of traditional K-means is $O(kdn)$, and the clustering complexity is $T(n) = O((1-\alpha) \log n)$, same as in [27, 28]. Among them, the computation formula of the pruning rate is $\alpha = \text{Amount of pruned data} / \text{Total data}$. The initial center selected in the improved algorithm is of higher accuracy, so increasing the pruning rate reduces the computation of clustering times in the iteration and thus reduces the time complexity of clustering.

4. Experiment and analysis

4.1. Experimental environment

All the experiments are executed as single-threaded on ThinkStation P920 workstation, configured with a 12-core Intel(R) Xeon(R) Silver 4116CPU 2.10 GHz and 32G RAM. All application programs are compiled and executed using Python3.5 in Pycharm environment running on Ubuntu16.04 OS.

4.2. Experimental data set

The datasets selected for the experiment are Iris and Wine, two classic data sets of UCI Repository, and Nsl-kdd datasets. Iris dataset is a common standard data set for cluster

analysis. It has 150 data and 4 features, and the last one is the data label. The features include the length of the sepals and petals of the iris, and the classification of this plant can be determined through the evaluation of the features. The wine dataset is often used for cluster analysis. There are 178 pieces of data and 13 different numeric features. All the experimental data in this investigation have been pre-processed for standardization.

The Kdd cup 99 [33] is the dataset of the Kdd cup network intrusion theme contest in 1999. Nsl-kdd [34] is an enhancement to the KDD Cup 99 dataset that does not contain redundant and duplicate records in its own set of tests, making it suitable for experiments and tests. Nsl-kdd is the same as the KDD Cup 99 dataset, containing 23 attack behaviors and 42 attribute features.

4.3. Experimental results and analysis

Experiment 1

In order to verify the advantages of the initial center selection of the algorithm in this investigation, traditional K-means clustering and Kd-tree accelerated clustering were used respectively after the completion of the center selection.

Principal component analysis (PCA) is a method to simplify datasets and reduce the data dimension [35]. The main idea of PCA is to reallocate the original d-dimensional features to the k-dimension, which is a brand new orthogonal feature, also known as the principal component [36]. PCA sequentially searches for a set of orthogonal coordinate axes in the original linear space, taking maximum variance as the standard, and projecting the transformed data onto a new set of coordinate axes. The choice of new axes is closely related to the data itself.

The results of the initial cluster centers are shown in Fig. 4. The Iris dataset was projected onto the two-dimensional plane using PCA dimensionality reduction. After being standardized, the data of the Wine dataset is projected onto a two-dimensional plane using the PCA method.

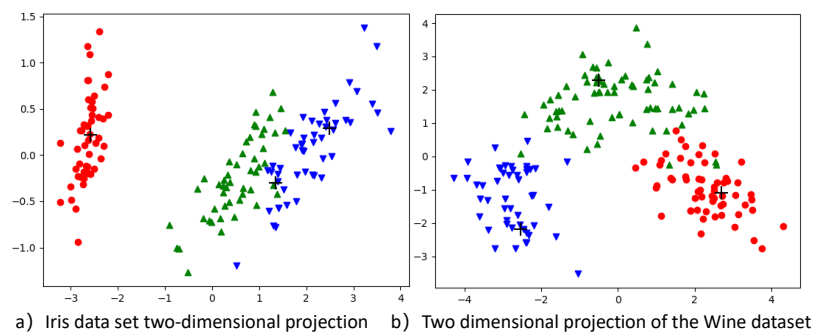


Fig. 4. Data centralized cardiac projection

Three shapes represent different categories of data. The initial central points of clustering are distributed in three different classes.

The Iris and Wine datasets contain a small amount of data. As a result, the traditional Kd-tree accelerated K-means clustering algorithm has no tangible benefit in reducing the execution time, and the high pruning rate will affect the clustering accuracy.

For these two datasets, the algorithm of Table 4 is used. In order to verify the improvement of the initial data center selection method and the correlation weighted nearest neighbor search clustering method, four types of initial cluster center selection methods and three combinations of clustering methods were used. In this experiment, the numbers of nodes 6, 7, and 8 are selected in the leaf buckets. Experimental clustering results in the Iris dataset are presented in Table 5.

Table 5. Iris data set experimental results

The initial center	Clustering method	k	1	2	3	means
K-means	K-means	3	0.593	0.740	0.613	0.649
		4	0.653	0.906	0.827	0.795
		5	0.740	0.760	0.820	0.773
Literature [27]	K-means	3	0.866	0.866	0.866	0.866
		4	0.866	0.866	0.866	0.866
		5	0.900	0.900	0.900	0.900
	Literature [27]	3	0.706	0.706	0.706	0.706
		4	0.740	0.740	0.740	0.740
		5	0.767	0.767	0.767	0.767
Literature [28]	K-means	3	0.887	0.887	0.887	0.887
		4	0.893	0.893	0.893	0.893
		5	0.913	0.913	0.913	0.913
	Literature [28]	3	0.740	0.740	0.740	0.740
		4	0.787	0.787	0.787	0.787
		5	0.806	0.806	0.806	0.806
Proposed	K-means	3	0.900	0.900	0.900	0.900
		4	0.900	0.900	0.900	0.900
		5	0.913	0.913	0.913	0.913
	Proposed	3	0.953	0.953	0.953	0.953
		4	0.940	0.940	0.940	0.940
		5	0.967	0.967	0.967	0.967

The experimental results of the Iris dataset show that the clustering results of the K-means algorithm are unstable. The accuracy of each experiment is uncertain, but the overall accuracy is low. The clustering results of the algorithm proposed in [19, 20] using the traditional K-means algorithm are better than those using the traditional K-means algorithm to select the initial clustering center randomly. However, when using the Kd-tree accelerated K-means algorithm to achieve clustering, there exists a problem that the clustering results are inaccurate due to the high pruning rate.

The method proposed in this search for selecting the initial center has higher accuracy than the traditional K-means algorithm for randomly selecting the clustering center after using the traditional K-means clustering method; though, the clustering effect is better than the algorithm in [27, 28]. The clustering results are also more stable, as the accuracy of the new method of clustering of relevance concerning the nearest neighbor is considerably improved compared to the traditional K-means algorithm. Experiments show that the improved K-means algorithm in this investigation can achieve a better clustering effect. For the same dataset, the same initial clustering center is obtained by setting the same number of leaf buckets l and the same number of

clusters k . Under the condition of constant clustering number and iteration number, the clustering results are stable and accurate.

The Wine dataset contains 13 attributes, and the data with different dimensions among the attributes have considerable differences, which is not favorable to the clustering algorithm. The normalization is used to preprocess the data and transform it into dimensionless scalars to improve the accuracy of the algorithm. Like Iris dataset, four clustering center initialization methods and the combination of three clustering methods are used to verify the accuracy and reliability of the algorithm in this investigation. The numbers of leaf buckets selected in the experiment are 4, 6 and 9. Experimental clustering results in the Wine dataset are presented in Table 6.

Table 6. Experimental results of Wine data set

The initial center	Clustering method	k	1	2	3	means
K-means	K-means	3	0.579	0.557	0.562	0.553
		4	0.512	0.601	0.579	0.564
		5	0.713	0.663	0.612	0.663
Literature [27]	K-means	3	0.657	0.657	0.657	0.657
		4	0.674	0.674	0.674	0.674
		5	0.719	0.719	0.719	0.719
	Literature [27]	3	0.545	0.545	0.545	0.545
		4	0.624	0.624	0.624	0.624
		5	0.679	0.679	0.679	0.679
Literature [28]	K-means	3	0.740	0.740	0.740	0.740
		4	0.767	0.767	0.767	0.767
		5	0.787	0.787	0.787	0.787
	Literature [28]	3	0.567	0.567	0.567	0.567
		4	0.668	0.668	0.668	0.668
		5	0.680	0.680	0.680	0.68
Proposed	K-means	3	0.740	0.740	0.740	0.74
		4	0.781	0.781	0.781	0.781
		5	0.803	0.803	0.803	0.803
	Proposed	3	0.781	0.781	0.781	0.781
		4	0.826	0.826	0.826	0.826
		5	0.871	0.871	0.871	0.871

It can be seen from the experimental results of the Wine dataset in Table 5 that the clustering results of the K-means algorithm are still unstable. The experimental accuracy of the algorithm in this investigation on the Wine dataset is not as high as that of Iris, but the clustering results are relatively more stable than the results obtained by the K-means algorithm. The accuracy of the proposed algorithm in selecting the initial clustering centers is verified.

Experiment 2

In experiment 2, the improved algorithm for network intrusion detection is used to verify the stability and accuracy of the improved algorithm in datasets with large data volumes and high dimensions. The Nsl-kdd dataset contains 125,973 pieces of data with rich characteristic attributes, which provides rich information for research and application, but also increases the workload. The problem can be simplified by analyzing the correlation between variables to reduce the dimension of data. Large

characteristic quantity can reduce the analysis index and minimize the information loss contained in the original index. Taking into consideration that the transformation of closely related variables into fewer new variables, the experimental efficiency is improved.

PCA method is used for dimensionality reduction. As depicted in Fig. 5, different dimension K values in the data and the results obtained by selecting different dimensions are shown that, as K increases, the curve flattens out, and choosing $K=20$ explains 93% of the variation of the initial variable. Therefore, 20 main components were chosen to be conserved and the characteristics of the Nsl-kdd dataset were reduced to the characteristics of the 20-dimensional vectors.

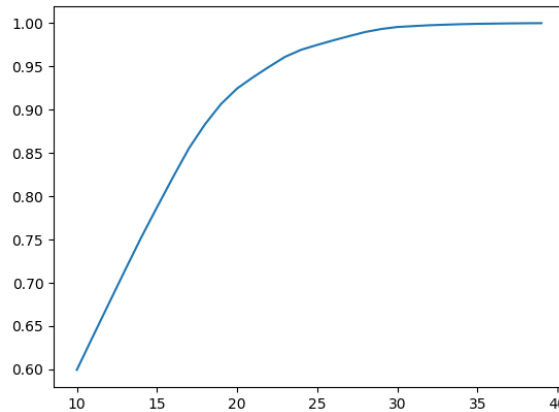


Fig. 5. Data dimension reduction

Experiments were carried out under different clustering numbers k and the number of data l in the leaf bucket. The parameter of the leaf bucket did not exist in the experiment of the traditional K-means algorithm. K-means with the same K value corresponds to several randomized experiments. The results were analyzed, and the detection rate (DR) and false alarm rate (FR) are compared. The detection rate represents the proportion of abnormal records detected correctly in all abnormal records, and the false alarm rate represents the proportion of normal records misdetected as abnormal records in all normal records. In the test of computing running time and computation, several groups of experiments are defined with different K values, in which the iteration times of each algorithm are 10. The experimental results are presented in Table 7.

Table 7. Intrusion detection experiment accuracy results

k	l	d	DR				FR			
			K-means	[27]	[28]	Proposed	K-means	[27]	[28]	Proposed
40	24	20	0.4781	0.8129	0.8545	0.8688	0.1161	0.0619	0.0640	0.0594
50	24	20	0.5173	0.8254	0.8178	0.8713	0.1323	0.0685	0.0788	0.0625
60	24	20	0.5268	0.8052	0.8313	0.8775	0.1355	0.0642	0.0628	0.0633
40	25	20	0.5091	0.8471	0.8664	0.8981	0.1432	0.0721	0.0699	0.0705
50	25	20	0.5314	0.8336	0.8683	0.8913	0.1574	0.0874	0.0767	0.0697

60	25	20	0.5581	0.8594	0.8922	0.9021	0.1711	0.0925	0.0794	0.0720
40	11	20	0.4237	0.7923	0.8339	0.8885	0.1324	0.0674	0.0626	0.0631
50	11	20	0.6140	0.8316	0.8241	0.8091	0.2156	0.0748	0.0683	0.0522
60	11	20	0.5947	0.8191	0.8611	0.8758	0.1978	0.0884	0.0899	0.0823

Table 8. Results of computational efficiency of intrusion detection experiment

k	d	Runtime(s)				Compute distance			
		K-means	[27]	[28]	Proposed	K-means	[27]	[28]	Proposed
40	20	360	189	116	113.	50389200	27646772	23874631	23302891
45	20	407	212	146	139	56687850	29297394	26781201	24547793
50	20	439	228	152	150	62986500	36538875	30100369	28947137
55	20	472	239	163	164	69285150	39909562	33674610	31897414
60	20	496	251	173	177	75583800	40242634	34816779	33471059

Experimental results show that the DR of the algorithm in this investigation is improved, and the FR is low compared to the traditional k-means algorithm and references [19, 20] under different clustering numbers. The algorithm model has a better intrusion detection effect of abnormal traffic.

As presented in Table 8, the center of the cluster converges and the pruning rate of the Kd-tree search space tends to be stable. Experimental results of the Nsl-kdd dataset show that: with the increase of the number of iterations, the algorithm proposed in this study is superior to the traditional K-means algorithm in terms of efficiency and the distance computation in each iteration process is lower than the traditional k-means algorithm. The accuracy of selecting the initial center is improved, and the distance computation in the iteration process can also be reduced. The reliability of the algorithm in this paper is proven when applied in the field of network intrusion detection.

Experiment 3

This experiment intends to verify the advantages of the algorithm in the case of a large volume of high dimensional data. NumPy library in Python is used to generate 20 uniformly distributed random cluster central sets C, where the dimension is 20. Make_blobs from the Sklearn library is then performed to generate random data for the cluster model. 20,000 ~ 300,000 samples are generated each time, with 20 features for each sample and a total of 20 clusters, as presented in Table 9. The test made a comparison between clustering time and tree building time.

Table 9 Random data sets

Data sets	Leaf buckets	Clustering number	Data set size /103	Dimension
S1	18	20	20	20
S2	13	20	30	20

S3	18	20	40	20
S4	12	20	55	20
S5	17	20	80	20
S6	10	20	100	20
S7	16	20	150	20
S8	10	20	200	20
S9	14	20	250	20
S10	17	20	300	20

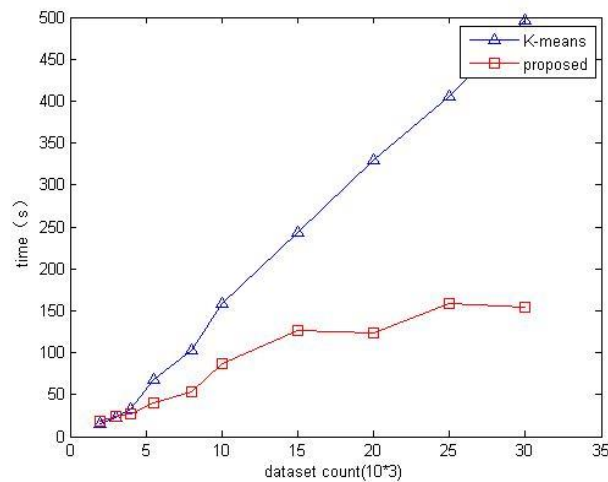


Fig. 6. Experimental results

As displayed in Fig. 6, the K-means algorithm shows that the use of Kd-tree acceleration for datasets with a small amount of data is not significant. There may happen that the efficiency is even lower than the traditional K-means algorithm. The main reason is that there is extra computation when building Kd-tree, and the access overhead also brings time consumption. As the amount of data increases, the effectiveness of the algorithm in this paper is more prominent.

5. Conclusions and Future Work

K-means algorithm is one of the most classic data mining algorithms. Based on the K-means algorithm of unsupervised learning method, aiming at solving the problems of unstable cluster effect caused by the random selection of initial clustering central points and by a low operation efficiency when the data volume is large, a k-means algorithm of selecting clustering center based on density is proposed.

In the improved algorithm, Kd-tree is used as a data pre-processing method to store data, and leaf buckets are performed to organize leaf nodes when building Kd-tree. The initial cluster center was selected by computing the density estimation and density weight of leaf buckets, as well as the correlation weighted distance estimation and distance weight between leaf buckets. In the process of clustering, the relevance

weighted nearest neighbor search is used to prune the search space, to improve the accuracy of clustering and the overall rate of the algorithm operation, and to guarantee the reliability of it. Meanwhile, high efficiency and accuracy are more suitable for detecting abnormal network traffic with large amounts of data under complex and changeable conditions.

From the experiment 1, the results show that the proposed algorithm has higher accuracy. The network intrusion detection results in experiment 2 show that the proposed algorithm has higher accuracy and a lower false alarm rate, where the initial clustering center is more accurate and has a stronger anti-noise ability. Results of experiment 3 show that the proposed algorithm is faster and more user-friendly to datasets containing a large amount of data, and can adapt to the efficiency requirements of computing in the era of Big Data.

As future research, the density-based K-means algorithm can be improved through parallelization, so thus, the speed of execution can be further improved in large datasets. An example of such a parallel environment is Multi-GPU parallel systems [28, 29], where thousands of different types of cores are available in a single computing environment, and the exploitation of the parallelism of the K-means algorithm is crucial for high performance.

References

1. Jiang, H., Chen, Y., Qiao, Z., Li, K. C., Ro, W. W., & Gaudiot, J. L. (2014). Accelerating MapReduce framework on multi-GPU systems. *Cluster Computing*, 17(2), 293-301.
2. Jiang, H., Chen, Y., Qiao, Z., Weng, T. H., & Li, K. C. (2015). Scaling up MapReduce-based Big Data Processing on Multi-GPU systems. *Cluster Computing*, 18(1), 369-383.
3. Cui, M., Han, D., & Wang, J. (2019). An Efficient and Safe Road Condition Monitoring Authentication Scheme Based on Fog Computing. *IEEE Internet of Things Journal*, 6(5), 9076-9084. <https://doi.org/10.1109/JIOT.2019.2927497>
4. Lin, C. H., Xiao-Dong, L. I., Jin, J., Xue-Biao, Y., & Jun, W. U. (2013). DNS traffic anomaly detection based on W-Kmeans algorithm. *Computer Engineering & Design*.
5. Zhang, W., Han, D., Li, K.-C., & Massetto, F. I. (2020). Wireless sensor network intrusion detection system based on MK-ELM. *Soft Computing*(2).
6. Liang, W., Fan, Y., Li, K.-C., Zhang, D., & Gaudiot, J.-L. (2020). Secure Data Storage and Recovery in Industrial Blockchain Network Environments. *IEEE Transactions on Industrial Informatics*, PP(99), 1-1.
7. Liang, W., Li, K.-C., Long, J., Kui, X., & Zomaya, A. Y. (2019). An Industrial Network Intrusion Detection Algorithm based on Multi-Feature Data Clustering Optimization Model. *IEEE Transactions on Industrial Informatics*, PP(99), 1-1.
8. Han, D., Pan, N., & Li, K. (2020). A Traceable and Revocable Ciphertext-policy Attribute-based Encryption Scheme Based on Privacy Protection. *IEEE Transactions on Dependable and Secure Computing*, 1-1. <https://doi.org/10.1109/TDSC.2020.2977646>
9. Han, D., Yu, Y., Li, K.-C., & Mello, R. F. d. (2020). Enhancing the Sensor Node Localization Algorithm Based on Improved DV-Hop and DE Algorithms in Wireless Sensor Networks. *Sensors*, 20(2), 343.
10. Li, L., Chen, X., Jiang, H., Li, Z., & Li, K. (2016, 30 May-1 June 2016). P-CP-ABE: Parallelizing Ciphertext-Policy Attribute-Based Encryption for clouds. The 2016 17th IEEE/ACIS International Conference on Software Engineering, Artificial Intelligence, Networking and Parallel/Distributed Computing (SNPD).

11. Wei, M., Su, J., Jin, J., & Wang, L. (2014). Research on Intrusion Detection System Based on BP Neural Network ((pp. 657-663). https://doi.org/10.1007/978-3-642-40618-8_85
12. Huang, Z. Extensions to the k-Means Algorithm for Clustering Large Data Sets with Categorical Values. *Data Mining & Knowledge Discovery*, 2(3), 283-304.
13. Jiang, H., Chen, Y., Qiao, Z., Weng, T.-H., & Li, K.-C. Scaling up MapReduce-based Big Data Processing on Multi-GPU systems. *Cluster Computing*, 18(1), 369-383.
14. Aliahmadipour, L., & Eslami, E. (2016). GHFHC: Generalized Hesitant Fuzzy Hierarchical Clustering Algorithm. *International Journal of Intelligent Systems*, 31(9), n/a-n/a.
15. Yanfeng, Peng, Di, & Jian. Research of Performance of Distributed Platforms Based on Clustering Algorithm.
16. Rodriguez, A., & Laio, A. Clustering by fast search and find of density peaks. *Science*, 344(6191), 1492-1496.
17. Viswanath, P., & Pinkesh, R. (2006, 20-24 Aug. 2006). 1-DBSCAN: A Fast Hybrid Density Based Clustering Method. 18th International Conference on Pattern Recognition (ICPR'06).
18. Wang, W. (1997). STING: A Statistical Information Grid Approach to Spatial Data Mining. *Proc. of the 23rd Very Large Database Conf.*, 1997.
19. Heil, J., Haring, V., Marschner, B., & Stumpe, B. (2019). Advantages of fuzzy k-means over k-means clustering in the classification of diffuse reflectance soil spectra: A case study with West African soils [Article]. *Geoderma*, 337, 11-21. <https://doi.org/10.1016/j.geoderma.2018.09.004>
20. Gulnashin, F., Sharma, I., & Sharma, H. A New Deterministic Method of Initializing Spherical K-means for Document Clustering.
21. Kimberly, M., Amanda, G., & Tammy, N. Intravenous Continuous Infusion vs. Oral Immediate-release Diltiazem for Acute Heart Rate Control. *Western Journal of Emergency Medicine*, 19(2), 417-422.
22. Mccallum, A., Nigam, K., & Ungar, L. H. (2000). Efficient clustering of high-dimensional data sets with application to reference matching. *Proceedings of the Sixth International Conference on Knowledge Discovery and Data Mining*.
23. Zhang, G., Zhang, C., & Zhang, H. Improved K-means Algorithm Based on Density Canopy. *Knowledge-Based Systems*, S0950705118300479.
24. Guttman, A. (1984). A Dynamic Index Structure for Spatial Searching. *Acm Sigmod Record*, 14(2), 47-57.
25. Bentley, J. (1975). Multidimensional Binary Search Trees Used for Associative Searching. *Communications of the ACM*, Septembar, 1975. vol. 18: pp. 509-517: ill. includes bibliography., 18.
26. Wu, W., Yang, F., Chan, C. Y., & Tan, K. L. (2008). Continuous Reverse k-Nearest-Neighbor Monitoring. (Ed.),^(Eds.). *The 9th International Conference on Mobile Data Management, 2008 (MDM '08)*.
27. Redmond, S. J., & Heneghan, C. A method for initializing the K-means clustering algorithm using kd-trees. *Pattern Recognition Letters*, 28(8), 965-973.
28. Kumar, K. M., & Reddy, A. R. M. (2017). An Efficient k -Means Clustering Filtering Algorithm Using Density Based Initial Cluster Centers. *Information Sciences*, 418.
29. Bris, R. L., & Paul, F. An automatic method to create flow lines for determination of glacier length: A pilot study with Alaskan glaciers. *Computers & Geosciences*, 52, 234-245.
30. Pelleg, D., & Moore, A. (1999). Accelerating Exact k-means Algorithms with Geometric Reasoning. <https://doi.org/10.1145/312129.312248>
31. Hartigan, J. A., & Wong, M. A. (1979). A K-means Clustering Algorithm: Algorithm AS 136 ((Vol. 28, pp. 100-108). <https://doi.org/10.2307/2346830>
32. Wang, X., & Bai, Y. (2016). A Modified MinMax -Means Algorithm Based on PSO. 2016(7), 4606384.
33. Pfahringer, B. (2000). Winning the KDD99 Classification Cup: Bagged Boosting. 1(2), 65-66.

34. Deshmukh, D. H., Ghorpade, T., & Padiya, P. (2015). Improving classification using preprocessing and machine learning algorithms on NSL-KDD dataset. (Ed.),^(Eds.).
35. Roweis, S. (1999). Em algorithms for pca and spca. *Advances in Neural Information Processing Systems*, 10.
36. Tian, Q., Han, D., Li, K.-C., Liu, X., Duan, L., & Castiglione, A. (2020). An intrusion detection approach based on improved deep belief network. *Applied Intelligence*. <https://doi.org/10.1007/s10489-020-01694-4>

Jing Xu is pursuing a master's degree in software engineering at the School of Information Engineering, Shanghai Maritime University, China. Her research interests include Big Data storage and intelligent decision-making, key technologies for DoS attack defense, and wireless sensor network security.

Dezhi Han received the PhD degree from the Huazhong University of Science and Technology. He is currently a professor of computer science and engineering with Shanghai Maritime University, China, and his research interests include network security, cloud computing, mobile networking, wireless communication, and cloud security.

Kuan-Ching Li is a Distinguished Professor in the Dept of Computer Science and Information Engineering (CSIE) at Providence University, Taiwan, where he also serves as the Director of the High-Performance Computing and Networking Center. He is a recipient of awards and funding support from several agencies and industrial companies, as he also received distinguished chair professorships from universities in several countries. He published more than 300 scientific papers and articles and is co-author or co-editor of more than 25 books published by leading publishers. His research interests include parallel and distributed computing, Big Data, and emerging technologies. Professor Li is a Fellow of the IET and a senior member of the IEEE.

Hai Jiang is a Professor in the Department of Computer Science at Arkansas State University, USA. His current research interests include Parallel & Distributed Systems, Cloud Computing, Big Data, and Cryptography. He has published more than 150 research papers in major international journals and conference proceedings. He has served as a U.S. National Science Foundation proposal review panelist and a U.S. DoE (Department of Energy) Smart Grid Investment Grant (SGIG) reviewer multiple times. He is a professional member of ACM and IEEE computer society, also a representative of U.S. NSF XSEDE (Extreme Science and Engineering Discovery Environment) Campus Champion for Arkansas State University.

Received: April 06, 2020; Accepted: May 18, 2020

Cognitive Computation on Consumer's Decision Making of Internet Financial Products Based on Neural Activity Data

Hongzhi Hu¹, Yunbing Tang², Yanqiang Xie³, Yonghui Dai^{4*}, and Weihui Dai^{1*}

¹ School of Management, Fudan University,
Shanghai 200433, China
{hongzhihu, whdai}@fudan.edu.cn

² School of Journalism, Fudan University
Shanghai 200433, China
tangyunbing@fudan.edu.cn

³ Shaoyang Branch, Ping An Insurance Company of China LTD,
Shaoyang 422000, China
775202789@qq.com

⁴ Management School, Shanghai University of International Business and Economics,
Shanghai 201620, China
daiyonghui@suibe.edu.cn

Abstract. Internet finance has become a popular business in today's society. However, different from the physical objects or services sold online, Internet financial products are actually contracts defined by financial terms which make customers bear the possibility of capital losses and liquidity restrictions, but they can obtain profits in the future with some uncertainties. This paper takes consumer's cognition in the decision making of Internet financial products as research circumstances, studies the above issue by conducting an EEG-fNIRS experiment, and proposes an effective cognitive computation method based on neural activity data through BP-GA algorithm. On this basis, a new recommendation approach of Internet financial products is explored according to consumer's typical shared mental model. The computing and testing results indicate that researches of this paper provide promising new ideas and novel methods for the cognitive computation of artificial intelligence and the recommendation of Internet financial products.

Keywords: Internet Finance, consumer preferences, cognitive computation, mental model, neural activity data, artificial intelligence.

1. Introduction

With the rapid development of Internet and advanced mobile technology, Internet finance (ITFIN), a new financial business model, has developed vigorously and is generally accepted by the vast number of consumers in today's society [1]. During the past ten years, Internet finance has provided various new businesses and products such

as third-party payment, crowd funding, P2P (Peer-to-Peer lending), Internet financial portal products and so on [2, 3]. Internet finance can serve many useful purposes, bring more convenient service experience, meet the diverse needs of the consumers, and also improve the operation efficiency of traditional business. Internet technology has revolutionized financial markets, which has a significant impact on reducing the transaction cost, extending financial services, improving time efficiency and bringing a proliferation of innovative Internet financial products and services. However, the wide range of new financial products and services can also make consumers confused when choosing the suitable products. Different from the physical products sold on the Internet, Internet financial products usually involving complex financial terms, which may contain various risks. The selection of Internet financial products is a decision-making process which involves consumer's financial needs, cognitive characteristics, consumption experience, risk preference, etc [4].

Previous studies have shown that consumers' cognition and preferences, as the most important psychological factors in their mental models, play a leading role in the decision-making of their Internet financial products [4, 5]. Therefore, how to analyze the above factors and conduct a successful recommendation to customers has become the focus of researchers. In this respect, various methods and technologies commonly used in traditional e-commerce are also applied to Internet financial products. These methods and technologies are usually based on data mining or questionnaire survey of consumers' online behavior [6, 7]. However, the influencing factors and cognitive characteristics of consumers' decision-making on Internet financial products are quite different from those of traditional products in terms of profitability and risk risks [1, 3, 4, 5]. Without considering the intrinsic neural mechanism underlying the consumer's cognition and behaviors, it is difficult for traditional methods to achieve accurate and reliable results [8, 9].

In this paper, the cognition of consumers in the decision-making process of Internet financial products is studied through a neural experiment, aiming to explore an effective cognitive computing method based on neural activity data, so as to improve the recommendation performance according to different mental models of consumers. It is organized as follows: Section 1 introduces the background and motivation of our research; Section 2 reviews the related research work; Section 3 conducts a neural experiment and proposes the customer's cognitive computing method based on neural activity data; Section 4 analyzes the consumers' mental model and explores a new recommendation approach; and Section 5 undergoes summary and discussion of this paper.

2. Related Work

2.1. Consumer's Consumptive Behaviors of Internet Financial Products

Consumers' behaviors and their influence on financial market belong to the research field of behavioral finance theory. According to the traditional theory, consumers have limited rationality, limited ability of information processing in decision-making, but

with less consideration of the psychological influences on behaviors. New financial behavior theories, such as prospect theory and emotional psychology theory, have brought human's psychology and behaviors into the research framework of finance [10]. Those theories study individual behavior and psychological motivation from a micro-perspective to explain and predict the development of financial market, effectively making up for the shortcomings of traditional financial theory.

Financial consumption is also a psychological process related to consumers' investment beliefs, preferences and decisions. Unlike the physical objects or services that can be experienced immediately, financial products are contracts defined by financial terms, which have the possibility of making customers bear capital losses and liquidity restrictions, but may make some uncertain profits in the future. Different from general products or services, customers of financial products should consider factors related to financial profitability, liquidity, and risk and so on in their decision making [4]. On the other hand, as irrational people, customers' inherent psychological and cognitive biases may lead them to make inappropriate judgments on market trends and product risks [11, 12]. In fact, they are easily misled by the exaggerated marketing strategies of the trading platform in Internet financial transactions. If consumers are unable to understand the existing information provided due to cognitive bias, they may suffer economic losses after making improper decisions. Therefore, consumers' cognition of Internet financial products decision-making is quite distinctive and situational, which needs further research in order to recommend appropriate products according to different preferences.

2.2. Consumer's Cognition of Decision Making in Cyber Space

Through literature review, we found that most of the existing researches on cognitions of consumers' online decision making are based on data mining or questionnaire survey. They obtain consumers' survey data and online behavior data, and establish models to analyze the cognitive characteristics and influencing factors of consumers. For example, Cui et al. studied the consumer satisfaction and trust through the online comment extraction and semantic analysis of emotional words [7]. Wang and LV used data mining and K-means clustering algorithm to analyze the preferences of different consumers for commodities [13]. In addition, text feature extraction technology, principal component analysis, EM clustering and other technologies have also been applied.

Jiang et al. analyzed the stakeholders and related topics according to the differences of information writing style and content characteristics of social media [14]. Social media variables are extracted from two aspects of information activity intensity and sentiment tendency, and a regression prediction model is established to analyze the impact of social media activities on stock behaviors. Experiments on Yahoo Financial Forum show the effectiveness of the proposed method. Besides, some scholars began to build behavior prediction model from user emotional information. These methods provide an effective way to study consumers' cognition based on self-report or behavioral data. However, they may be suffered by the influences of some subjective elements, such as inconsistencies and inaccuracies in description, behavior and intention mismatch and so on [15]. It seems that the customers' cognitions can't be fully and accurately explained if only through historical behavior data analysis or questionnaire

survey, so we need to further explore its internal cognitive mechanism. Therefore, cognitive computing is considered as one of the key technologies to deal with big data challenges and better navigate big data analysis. [16].

The preliminary conceptual framework of cyberspace psychology was proposed by John Suler in his hypertext book of “The psychology of cyberspace” in January 1996 [17]. It is intended to understand the behavior of individuals and groups in cyberspace. Cyberspace psychology is a discipline that studies the psychological phenomena and related laws of people’s online behaviors. In the modern information society, physical space and cyberspace are increasingly closely tied. Users’ understanding, experience, values, and ideology can be mirrored in their behaviors in cyberspace. For the sake of having a better understanding of the relationship between different spaces, Dai studied the evolution mechanism of social public emergencies, and first presented a comprehensive framework of CyberPsychosocial and Physical (CPP) computation based on social neuroscience mechanism in 2014 [18, 19]. CyberPsychological computing (CPC) is developed under the framework of CPP, and it aims to quantitatively study the relationship between the user’s internet behavior and its mental states, so as to establish a computing model. CPC has already been applied in the cognitive interpretation concerning user behaviors in cyberspace. For example, Zhou et al. computed the attention, interest, emotion, and satisfaction of learners according to the learner’s online actions of keyboard and mouse based on the prior knowledge of neural mechanism [17].

2.3. Consumer Mental Model of Product Design and Recommendation

The concept of mental model was firstly proposed by Craik in 1943[20]. It is basically used to describe the inner psychological activities and internal cognitive processes conceived and projected by human thought, which affect people’s comprehension, interpretation and view of their own behavior. As an effective design method, mental model has been well applied in product and service design. For example, it can be used as a design method to improve understanding between product suppliers and consumers [21]. Mental models can also suggest how to react and adapt to new situations. In many cases, knowledge can be better coordinated, communicated and shared among heterogeneous multidisciplinary teams based on the shared mental model [22, 23]. Through different process incentive mechanism, Turumugon *et al.* showed that mental model can make product meet the expectations of different users by adopting “localization” and content analysis method [24].

Mental model theory can be used as a paradigm to identify, measure and model end-users’ privacy attitudes, concerns, intentions and behaviors. It can also provide a valuable structure to investigate connection between privacy related online behaviors and provide valuable insights into privacy decisions [25]. In addition, mental model can identify gaps and misunderstandings in order to develop effective communication. For instance, it can be used as a tool to reveal the basic principles of collision insurance, in which the budget income of consumers is predicted across consumption categories. The results show that consumers’ purchase of insurance is a normal commodity and follows a cognitive model that emphasizes budget constraints [26]. As for the mental models of customers in the decision-making of Internet financial products, there are still many issues worthy of further study. It is difficult to accurately find the cognitions of

consumers if only through data mining of their online behavior data. With the help of advanced neural experiment technology, we can study the cognition of consumers based on neural activity data, extract the typical mental models which may exist in the consumers' decision making of Internet financial products, and then obtain valuable prior knowledge for recommendation.

3. Experiment and Cognitive Computation

3.1. Neural Mechanism of Consumer's Cognition

Based on the existing research results of neuroscience, we proposed the neural mechanism of consumers' cognition and decision-making of Internet financial products, as shown in Figure 1.

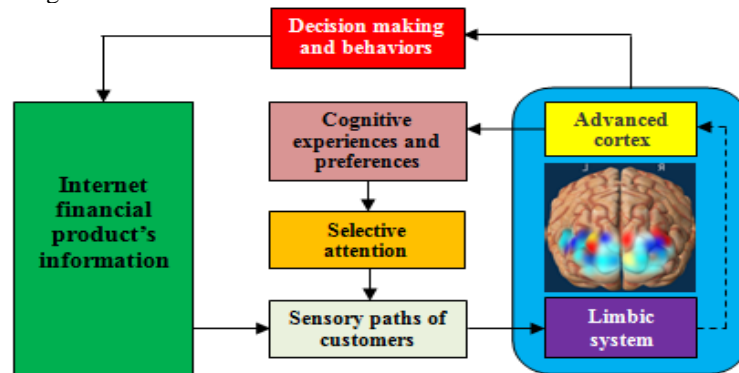


Fig. 1. Neural mechanism of cognition and decision making

When consumers perceive the information of Internet financial products through the sensory path, a series of neural activities from the limbic system to the advanced cortex of consumers' brain may be triggered. First of all, the expressive features and interesting points of financial products will generate preliminary cognition through the limbic system. At last, cognitive experience of external information forms rational cognition by advanced cortex. In this process, consumers' cognitive experience and preference for previous similar products can regulate sensory paths of attention resources and time allocation through the neural activities of advanced cortex, thus affecting their selective attention. The customers' decision making and behaviors are the result of their judgment influenced by the above neural activities in the brain.

The above neural mechanism provides the basis of psychological pattern for consumers to make decisions on the choice of Internet financial products. Among them, the expressive features and interesting points of products can attract consumers' initial attention, but their final decisions making are mostly based on rational cognition. According to the existing studies [1, 2, 4], consumers' cognition of Internet financial products can be summarized into following four aspects: convenience, profitability,

liquidity, and risk. The above cognition will be affected by a variety of consumer personal factors and environmental factors, such as age, gender, education level, experience, product reputation, herd behavior, and so on. Our cognitive computing does not focus on the relationship between consumers' cognition and the above factors, but aims to explore the methods to compute consumers' cognition of convenience, profitability, liquidity, and risk from their neural activity data. There have been a lot of researches related to human cognition, involving the dynamic and complex neural activities of brain function network [28], which is hard to be well described with a clear mathematical model. Therefore, we employ the BP-GA algorithm to realize the cognitive computing as shown in Figure 2.

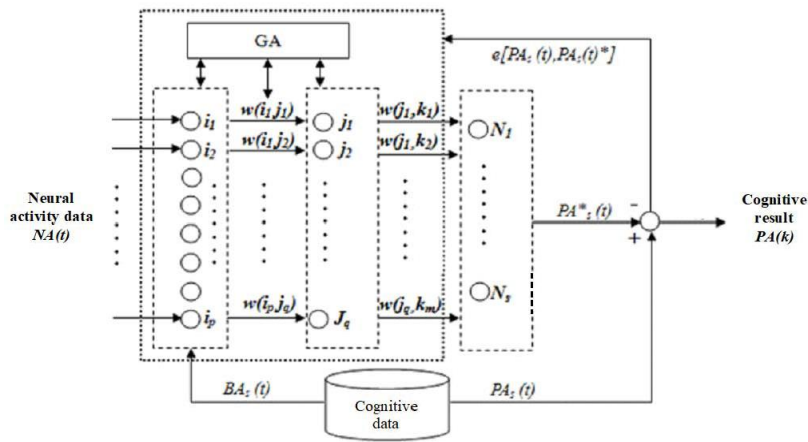


Fig.2. Cognitive computation with BP-GA algorithm

The input variable $NA(t)$ is the multi-dimensional time series data of neural activity, which represents the consumers' cognitive responses in their brains. It consists of a set of feature parameters $f_1(t), f_2(t), \dots, f_M(t)$ extracted from the data collected by neural experimental equipment.

$$NA(t) = [f_1(t), f_2(t), \dots, f_M(t)]. \tag{1}$$

The output variable $PA(k)$ is the result of consumer cognition composed of four parameters PA_1, PA_2, PA_3 and PA_4 , which represent convenience, profitability, liquidity, and risk respectively. In human cognitive research, there are a variety of neural activity data available for application. Among them, EEG (Electroencephalograph) and fNIRS (functional Near Infrared Spectroscopy) data can be collected by portable devices, which are widely used in the actual environment. EEG has superior dynamic performance, and fNIRS can better reflect the spatial characteristics of neural activities by detecting the concentration changes of oxygenated (HbO) and deoxygenated hemoglobin (HbR). In order to improve the accuracy, we use the multi-modal feature parameters of EEG-fNIRS composite data as the input variable $NA(t)$ for cognitive computing.

3.2. Experiment and Cognitive Computation

Alipay launched a financial product “Yu’E Bao” in 2013, which can be operated and transacted by mobile phone. It has been an innovative Internet financial product with huge number of consumers. Similar to Yu’E Bao, we design eight Internet financial products for testing, as shown in Table 1, which reflect different degrees of convenience, profitability, liquidity, and risk. Convenience represents how easy it is for consumers to fully understand product information and complete a transaction operation. Its degrees are set to highest, high, medium, low, or lowest for each product respectively.

Table 1. Internet financial products for test

Product No.	Convenience (Degrees from highest to lowest)	Profitability (Expected annual rate of return)	Liquidity (No-redeemable days)	Risk (Maximum Possible Loss)
1	highest	1% fixed	0	0%
2	high	2% fixed	30	0%
3	medium	3% fixed	180	0%
4	low	5%	7	10%
5	low	20%	30	20%
6	medium	50%	60	30%
7	High	100%	180	80%
8	lowest	200%	360	100%

In our experiment, we set up a mobile phone to display the test signals of the above products. The neural activity data of consumers’ cognition is recorded by EEG-fNIRS device and converted into the input variable of BP-GA computation procedure. Figure 3 shows the experimental data collection and processing.

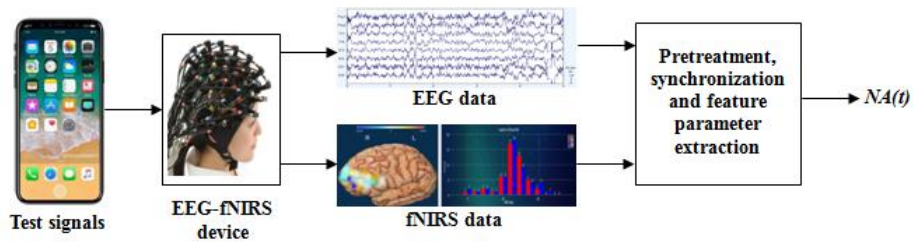


Fig. 3. The experimental data collection and processing

In the cognitive and decision-making study of EEG fNIRS, every specific test task should be carried out according to the strict experimental paradigm. Generally, only one influencing factor can be considered in a task. The purpose of the task is to find out the characteristics and mechanisms of neural activities significantly related to the above influencing factors. However, the goal of our study is to explore an effective method of

consumer cognitive computing from their neural activity data. Therefore, we designed the experimental paradigm as follow:

$$\text{Product information(60s)+Black screen(1s)+[Question(10s)+Answer(10s)]_{1-4} \quad (2)$$

Product information comes from the eight design products in Table 1, but it appears in random order. The test steps of [question (10s) + answer (10s)] are carried out four times. In each cycle, consumers are required to assess the convenience, profitability, liquidity, and risk of products with a score from 1 to 10 respectively. Since consumers' cognition before decision-making is mainly reflected in the thinking process of the question, we only extract the neural activity data during the period of [Question(10s)] to generate the input variables calculated by BP-GA, and take the data during the black screen (1s) as its deviation. Take the consumer evaluation score during the period of set [Answer(10s)] as the output variable calculated by BP-GA. According to the experience of preliminary experiment, the input variables are composed of the following multimodal feature parameters: wavelet decomposition coefficient of EEG signal, mean value, variance, skewness coefficient, kurtosis coefficient and slope of fNIRS signal.

We conduct this experiment with total 38 consumers of Yu'E Bao, including 22 males and 16 females. They are between the ages of 18 and 54 with the educational level from junior college to Ph.D. The consumer's evaluation score of convenience, profitability, liquidity, and risk of each tested product is between 1 and 10, of which 10 represents the highest level and 1 represents the lowest level. At the end of this experiment, the consumers give the degree of purchase intention for each tested product according to their preferences based on the above scoring method. After that, we randomly select the data of 26 consumers for training, and use the remaining data of 12 consumers for computation test. Figure 4 shows the change of its relative error with the number of iterations.

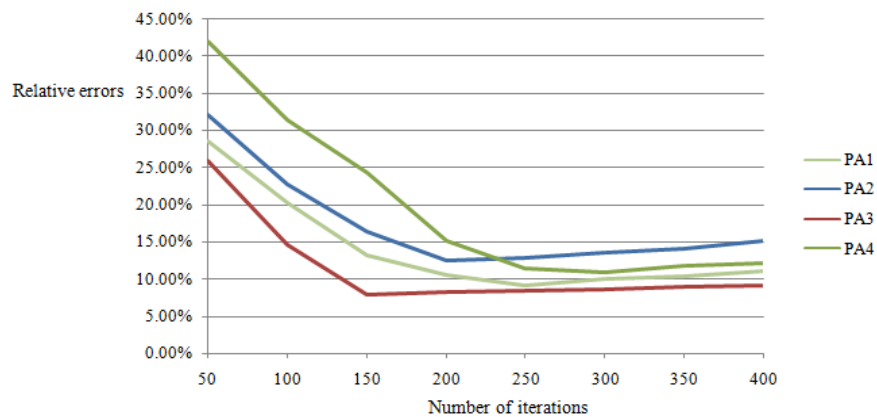


Fig. 4. The change of relative error with the number of iterations

As we can see from Figure 4, the optimal number of iterations for parameters PA_1 , PA_2 , PA_3 and PA_4 are different. Furthermore, Table 2 shows their minimum relative errors at the optimum number of iterations.

Table 2. The minimum relative errors with optimal number of iterations

Parameters	Minimum errors	Optimal number of iterations
PA_1	8.8%	237
PA_2	12.4%	216
PA_3	7.7%	162
PA_4	10.8%	307

Table 2 indicates that the cognitive computing method based on customer's EEG-fNIRS neural activity data can achieve a high accuracy from 87.6% to 92.3%. In the traditional method, the above parameters are usually evaluated by consumer's subjective self-report. Because of the fuzziness and inaccuracy of consumers' feelings and descriptions, it is difficult to accurately reflect the implicit cognitive characteristics of consumers. Our method also provides an objective way to calibrate the cognitive state of consumers.

4. Mental Model and Recommendation Approach

4.1. Mental Model

Consumers' cognitions and their decision making of Internet financial products are usually affected by a series of individual factors such as age, gender, education level, experience, income, personal preference and psychological characteristics. However, consumer psychology research shows that consumers with similar cognitive and decision preferences may have a common mental model. In order to extract the potential sharing mental model of Internet financial product consumers, we explained the structural equation model of variable relationship, as shown in Figure 5.

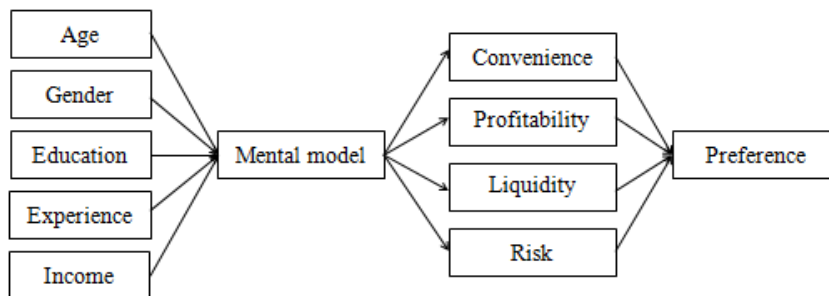


Fig. 5. Structural equation model of variable relationship

The mental model of consumers is determined by age, gender, education, experience, income and other personal factors. It affects consumers' cognition and behavior. In this case, we can see different evaluation of convenience, profitability, liquidity, and risk of Internet financial products, as well as different decision preferences. Among them, the

comprehensive influence of consumer mental model on preference can be summed up as the cognitive characteristics of convenience, profitability, liquidity, and risk. Therefore, we conducted cluster analysis on the above four parameters of the experimental participants. The consumers' evaluation scores of (Convenience), (Profitability), (Liquidity), (Risk) and the Preference scores can be divided into the most appropriate five clusters by the ant colony optimization clustering algorithm [29, 30]. Table 3 shows their clustering and analysis of variance.

Table 3. Clustering and variance analysis

Parameters	Cluster1	Cluster2	Cluster3	Cluster4	Cluster5	F-value	Sig.
PA_1	8.234	8.012	4.132	6.010	3.212	32.774	.000
PA_2	4.455	9.231	8.312	8.281	9.278	58.918	.000
PA_3	8.391	6.034	6.055	7.003	2.329	27.625	.000
PA_4	1.052	8.297	8.291	8.135	9.116	127.479	.000

Table 3 indicates that there are five typical cognitive patterns exist in customers on the decision-making of Internet financial products as follows:

- **Pattern 1:** high convenience, medium profit, high liquidity, and low risk, corresponding to Cluster 1;
- **Pattern 2:** high convenience, high profit, medium liquidity, and high risk, corresponding to Cluster 2;
- **Pattern 3:** low convenience, high profit, medium liquidity, and high risk, corresponding to Cluster 3;
- **Pattern 4:** medium convenience, high profit, medium liquidity, and high risk, corresponding to Cluster 4;
- **Pattern 5:** low convenience, high profit, low liquidity, and high risk, corresponding to Cluster 5.

Actually, the above patterns also reflect the typical shared mental models of customers while they make the decisions of Internet financial products. Although, the similar cognitive patterns may possibly be found by the traditional methods such as data mining of consumer's online behaviors or questionnaire survey, however those patterns in our study are extracted from the customers' neural activity data, which not only verify the existing traditional research findings based on neural mechanism interpretations, but furthermore provide precise descriptions of customers' cognitive characteristics as well as the important prior knowledge for customer's analysis and product's recommendation.

4.2. Recommendation Approach

In the field of online recommendation, a variety of intelligent recommendation methods have been proposed. Especially the recommendation method based on collaborative filtering algorithm has been widely used in financial product recommendation because of its high efficiency and good performance [31, 32]. The basic idea of collaborative filtering is to assume that customers who have similar interests or often choose similar products may have similar preferences [33]. It exists some problems such as cold start, scalability, sparsity and so on, and is in continuous improvement [34, 35]. However, we

believe that the lack of consideration of the internal cognitive mechanism of consumers is an important defect. On this basis, a novel method of Internet financial product recommendation based on consumer mental model is proposed as follows:

- **Step 1:** Classify the categories of recommended products. There are many kinds of financial products, such as banking products, insurance products, stocks, futures, etc., but consumers usually tend to focus on one or several products of interest. Therefore, product category is the first factor to be considered. In fact, the test products designed in our experiment are shown in Table 1, covering all the common categories of Internet financial products, and the recommendation should be put forward based on the analysis of the same category.
- **Step 2:** Calculate the parameters of the recommended product. This step is performed by the following formula:

$$PA_{ij} = c_{ij} * PA_{ib}, i=1,2,3,4; j=1,2,3,\dots,M. \quad (3)$$

In the above formula, PA_{ij} is the degrees of convenience (PA_{1j}), profitability (PA_{2j}), liquidity (PA_{3j}), and risk (PA_{4j}) of the recommended product j , where $j=1, 2, 3,\dots,M$; c_{ij} is the ratio coefficients of product j obtained by comparing with the most similar product of the same category in Table 1; PA_{ib} is the consumers' average score for the above comparison products in our test. Finally, the values of PA_{ij} for all M products considered as recommended should be standardized in the range of 1 to 10, which reflects the characteristics of the products.

- **Step 3:** Identify the customer's shared mental model. For new consumers, this calculation will be based on the historical data of the products purchased by that consumer. Refer to Step 2, first compute the degree PA_{ih} ($i=1, 2, 3, 4; h=1, 2, 3,\dots, N$) of N products purchased by consumers in history. Then, the following Euclidean distance function is used to identify the consumers' shared mental model through pattern recognition:

$$d_m = \|PA_{ih} - PA_{im}\|, i=1,2,3,4; h=1,2,3,\dots,M; m=1,2,3,4,5. \quad (4)$$

Where, d_m is the Euclidean distances between PA_{ih} and the five typical shared mental models, as shown in Table 3 respectively; and PA_{im} is the similar center values of cluster m ($m=1, 2, 3, 4, 5$) in that table. Because for different product categories, the value of cluster center may be different, so PA_{im} should be calculated according to the same product category to which PA_{ih} belongs. Finally, the shared mental model of consumers is determined to be the lowest to the highest possible level d_m .

- **Step 4:** Realize personalized recommendation. The values PA_{ij} obtained in step 2 represent the characteristics of those products to be recommended. The obtained values of d_m in Step 3 reflect the degree of closeness between the new consumers and the tested consumers in the experiment, which are described by Euclid distance of five clustering centers of different typical sharing mental models. Therefore, we

design the personalized recommendation approach for new consumers based on the above relationship, as shown in Figure 6.

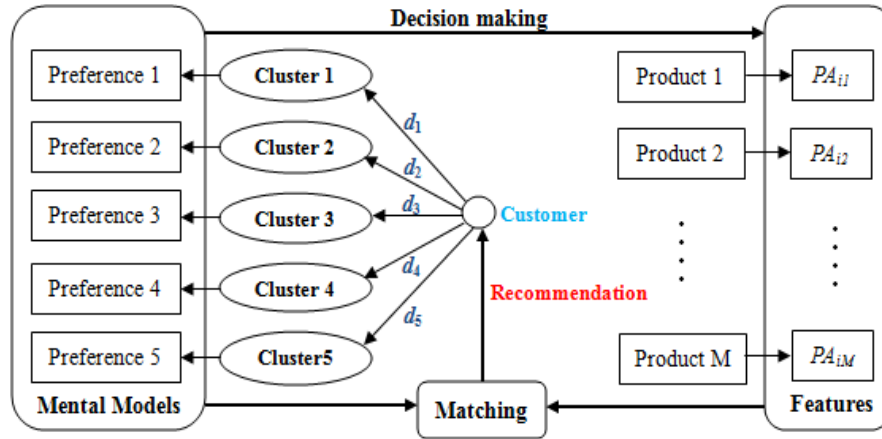


Fig. 6. Recommendation approach for new customer

New consumers’ decision-making on Internet financial products should be affected by the preferences of closest clusters and have a shared mental model. Therefore, personalized recommendation is to find the products that best match the preferences of the above clusters. In order to achieve this purpose, we take the minimum cluster d_m as the closest cluster that the consumer is most likely to belong to, and introduce the following criteria function:

$$RK_j = \sum_{i=1}^4 PA_{ic} * PA_{ij}, j = 1, 2, 3, \dots, M. \tag{5}$$

Where, RK_j is the value of criterion function, PA_{ic} is the center values of the closest cluster, and PA_{ij} is the characteristic values of the product j ($j=1, 2, 3 \dots, M$) considered to be recommended. RK_j is actually the weighted sum of PA_{ij} with coefficients PA_{ic} . The products are recommended according to the descending order of RK_j .

In order to test the performance of the above methods, we make recommendations to 11 new consumers from the actual financial products in the Chinese market. The physical products to be recommended include 12 banking products, 8 insurance products, 30 stocks and 22 futures, which are representative of the product features shown in Table 1. All consumers report that recommendations of our approach are more in line with their preferences than traditional approaches. Table 4 shows the comparison of consumers’ averaged satisfactions between collaborative filtering approach (CFA) and our mental model approach (MMA) respectively.

Table 4. Comparison of the two approaches to the average satisfaction of consumers

Approach	Bank products	Insurance products	Stocks	Futures
CFA	83.1%	78.7%	71.1%	76.8%
MMA	92.6%	88.4%	79.2%	84.1%

Table 4 indicates that our proposed approach (MMA) can achieve about more than 8.7% of the average consumer satisfaction compared with the collaborative filtering method (CFA). However, factors such as the number of recommended products and the time effect of product profitability may affect the performance of recommended methods and consumers' satisfaction. For example, stocks and futures have more products than banking and insurance products, and their returns can be known in a very short time, which may lead to greater fluctuations in consumer satisfaction. Therefore, further rigorous testing and comprehensive analyses are needed to verify the performance of this method.

5. Summary and Discussion

Cognitive computing is an important development direction in the field of artificial intelligence, which involves the acquisition, calculation and expression of human complex cognitive characteristics. Under the background of consumer's decision-making of Internet financial products, this paper conducts the EEG-fNIRS experiment to study the above cognitive characteristics, and proposes a cognitive computing method based on neural activity data through BP-GA algorithm. The results show that it can achieve 87.6% to 92.3% accuracy in the computation of consumer's cognition of product convenience, profitability, liquidity and risk.

In addition, this paper also studies the characteristics of consumer mental model, and explores a new recommendation method based on its typical sharing mental model. The test results show that compared with the traditional method, this method is more in line with the preferences of consumers, and the average satisfaction is about 8.7% higher than collaborative filtering method.

This paper provides new ideas and methods for cognitive computing and Internet financial product recommendation. However, there are many problems that need further study. Future research can be carried out from the following aspects: first, accurate neural learning and prior knowledge utilization of human complex dynamic cognitive characteristics, as well as more effective computational methods; second, accurate description and classification of consumer mental model, and the improvement of the proposed recommendation approach.

Acknowledgements. This research has been supported in part by National Natural Science Foundation of China (No. 71971066), Projects of Social Science Research, Ministry of Education of China (No. 19JZD010, No. 18YJA630019), and International Cooperation and Exchange program, Ministry of Science and Technology of China (No. 4-9/2018). Hongzhi Hu, Yunbing Tang, and Yanqiang Xie are the joint first authors who made the equal contributions. Yonghui Dai and Weihui Dai are the joint corresponding authors of this paper.

References

1. Junping, Q., Qiang, Y., Lilin, G.: On the influencing factors of Internet financial products' usage. *Journal of Intelligence*, Vol.34, No. 1, 179-184. (2015)
2. Ratten, V.: Entrepreneurship, E-finance and mobile banking. *International Journal of Electronic Finance*, Vol.6, No. 1, 1-12.(2012)
3. Tsai, C. H., Zhu, D. S., Jang, Y. M.: A study on the consumer adoption behaviors of Internet Bank. *Journal of Networks*, Vol.8, No. 12,263-268.(2013)
4. Yi, H.: Study on the Influence Factors of Internet Financial Product Selection of Investor. Beijing University of Posts and Telecommunications, Beijing, China. (2015)
5. Zhao, Q., Xue, J., Li, Y.: Research on the Choice of Internet Financial Products Based on Nested Logit Model. *Science and Management*, Vol.37, No.5,38-47. (2017)
6. Jiang, L., Cheng, Y., Li, H, Yan, Y., Wang, X. Q.: A Trust-based Collaborative Filtering Algorithm for E-commerce Recommendation System. *Journal of Ambient Intelligence & Humanized Computing*, 1-12. (2018)
7. Cui, X. L., Hong, Y., Zhao, N.: Consumer Repurchases Intention Mining Analysis Based on Online Review. *Journal of Shandong University*, Vol.50, No. 3, 28-31. (2015)
8. Martino, D. B., Kumaran, D., Seymour, B., Dolan, R. J.: Frames, Biases, and Rational Decision-Making in the Human Brain. *Science (New York, N.Y.)*, Vol.313, No.5787, 684-687. (2006)
9. Wang, T., Wang, C. C., Han, W.H., Ren, J. Dai W. H.: CyberPsychology Computation of Online Shopping Behaviors Based on EEG Signal Analysis, *Asian Journal of Information and Communications*, Vol.9, No. 2, 14-28. (2017)
10. Tversky, K. A.: Prospect theory: an analysis of decision under risk, *Ecomometrica*, No.47, 263-291. (1979)
11. Zhang, Y. F.: Discussion on Cognitive deviation and Protection path of Financial consumers, *Fujian Forum*, No.1, 38-41.(2015)
12. Liu, S. H.: On the Online Financial Services in the Digital Age. Springer London. (2013)
13. Wang, H., Lv, X.: Application of Data Mining in Online Shopping Commodity Feature Analysis. *Journal of Langfang Normal College (Science)*, 35-37. (2015)
14. Jiang, C. Q., Liang, K., Ding, Y., Liu, S. X., Liu, Y.: Stock behavior prediction based on social media. *China management science*, Vol. 23, No. 1, 17-24. (2015)
15. Acquisti, A., Grossklags, J.: Privacy and rationality in individual decision-making. *IEEE Security and Privacy*, Vol.3 ,No. 1, 26–33.(2005)
16. Guo, T.: Cognitive computing better controls big data, *China Computer News (Data Center)*, No. 18, (2013).
17. Zhou, X., Dai, G. H., Huang, S., Sun, X. M., Hu, F., Hu, H. Z., Ivanović, M.: Cyberpsychological computation on social community of ubiquitous learning. *Computational Intelligence and Neuroscience*, Vol.2015, Article ID 812650 , 1-7.(2015)
18. Dai, W.H.: Information Classification, Dissemination and Cognitive Model in Unconventional Emergency, Project Report of National Natural Science Foundation of China. (2015) .
19. Xu, D., Dai, W. H., Zhou, X. Y., Fu, Y., Dai, Y. H., Xu, D.R.: Smart Learning for Cyber Psychosocial and Physical Computation, Proceedings of the 2018 IEEE International Conference on Progress in Informatics and Computing, Suzhou, China,49-56.(2018)
20. Craik, K. J. W.: *The Nature of Explanation*. Cambridge, UK: Cambridge University Press.(1943)
21. Pillan, M., Pavlović, M., He, S.: Mental model diagrams as a design tool for improving cross-cultural dialogue between the service providers and consumers: Case of the Chinese restaurant business in Milan. Proceedings of the International Conference on Cross-cultural Design. Springer, Cham, 78-96. (2018)

22. Zeiler, W.: Morphology to illustrate the mental model of a design team's process. Proceedings of the 19th International Conference on Engineering and product Design Education (E&PDE 2017), Oslo, Norway, 1-6. (2017).
23. Schmidtke, J. M., Cummings, A.: The effects of virtualness on teamwork behavioral components: The role of shared mental models. *Human Resource Management Review*, Vol.27, No.4, 660-677.(2017)
24. Turumugon, P., Baharum, A., Hanapi, R., Kamarudin, N.: Users' mental model pattern for user interface design of mobile shopping Apps, *Journal of Computational and Theoretical Nanoscience*. Vol. 24, No. 2, 1158-1162. (2018)
25. Coopamootoo, P.L.K., Groß, T.: Mental models: an approach to identify privacy concern and behavior. In: Symposium on Usable Privacy and Security (SOUPS), 9–11 July, Menlo Park, CA. (2014)
26. Austin, L. C., Fischhoff, B.: Consumers' collision insurance decisions: a mental models approach to theory evaluation. *Journal of Risk Research*, Vol.13, No.7, 895-911. (2010)
27. Xu, D., Dai, Y. H., Dai, W. H.: Smart Recommendation of personalized cloud music service based on mental model. Proceedings of ACM Turing Celebration Conference-China 2019, Chengdu, China, 247-252. (2019)
28. Jollans, L., Whelan, R., Venables, L., et al. Computational EEG modelling of decision making under ambiguity reveals spatio-temporal dynamics of outcome evaluation. *Behavioural Brain Research*, Vol.321, 28-35. (2017)
29. Dai, W. H., Liu, S. J., Liang, S. Y.: An improved ant colony optimization cluster algorithm based on swarm intelligence. *Journal of Software*, Vol.4, No.4, 299-306. (2009)
30. Dai, W. H., Han, D. M., Dai, Y. H., Xu, D. R.: Emotion recognition and affective computing on vocal social media. *Information & Management*, Vol.52, No.7, 777-788. (2015)
31. Yang, B.: Research on Collaborative Filtering Algorithm for Financial Products Recommendations. Harbin, China: Harbin Institute of Technology. (2019)
32. Chu, X.Q.: Financial Product Recommendation Research Based on Artificial Immune System Combines Collaborative Filtering for Bank Term Deposit. Nanchang, China: Nanchang University. (2015).
33. Linden, G., Smith, B., York, J.: Amazon.com recommendations: Item-to-item collaborative filtering. *IEEE Internet Computing*, Vol.7, No.1, 76-80. (2003)
34. Fu, M. S., Qu, H., Yi, Z., Lu, L., Liu, Y. S.: A novel deep learning-based collaborative filtering model for recommendation system. *IEEE Transactions on Cybernetics*, Vol.49, No.3, 1084-1096. (2019)
35. Ramezani, M., Moradi, P., Akhlaghian, F.: A pattern mining approach to enhance the accuracy of collaborative filtering in sparse data domains. *Physica A: Statistical Mechanics and its Applications*, Vol.408, 72-84. (2014)

Hongzhi Hu is currently a post-doctoral researcher at School of Management, Fudan University, China. She received her Ph.D. in Management Science and Engineering from Fudan University, China in 2015. Her research interests include Internet financial data analysis, consumer behavior and psychology, and social network dynamics. Contact her at hongzhihu@fudan.edu.cn.

Yunbing Tang is currently an associate professor at School of Journalism, Fudan University, China. She received her Ph.D. in journalism and communication from Fudan University, China in 2008. Her research interests include marketing communication, social media analysis, and visual arts. Contact her at tangyunbing@fudan.edu.cn.

Yanqiang Xie is currently a senior staff at Shaoyang Branch, Ping An Insurance Company of China LTD. She has engaged in educational training and financial enterprise management for eight years since she graduated from university in 2012. Contact her at 775202789@qq.com.

Yonghui Dai is currently a lecturer at the Management School, Shanghai University of International Business and Economics, China. He received his Ph.D. in Management Science and Engineering from Shanghai University of Finance and Economics, China in 2016. His current research interests include management information systems, affective computing and artificial intelligence. Yonghui Dai and Weihui Dai are the joint corresponding authors of this paper. Contact him at daiyonghui@suibe.edu.cn.

Weihui Dai is currently a professor at Department of Information Management and Information Systems, School of Management, Fudan University, China. He received his Ph.D. in Biomedical Engineering from Zhejiang University, China in 1996. He serves as a standing director member of the Society of Management Science and Engineering of China, and an executive member of Shanghai Chapter, China Computer Federation (CCF). His recent research interests include neuroartificial intelligence, neuro-management, and man-machine hybrid intelligence. Yonghui Dai and Weihui Dai are the joint corresponding authors of this paper. Contact him at whdai@fudan.edu.cn.

Received: February 29, 2020; Accepted: May 28, 2020

CIP – Каталогизacija y publikaciji
Народна библиотека Србије, Београд

004

COMPUTER Science and Information
Systems : the International journal /
Editor-in-Chief Mirjana Ivanović. – Vol. 17,
No 2 (2020) - . – Novi Sad (Trg D. Obradovića 3):
ComSIS Consortium, 2020 - (Belgrade
: Sibra star). –30 cm

Polugodišnje. – Tekst na engleskom jeziku

ISSN 1820-0214 (Print) 2406-1018 (Online) = Computer
Science and Information Systems
COBISS.SR-ID 112261644

Cover design: V. Štavljanin
Printed by: Sibra star, Belgrade