Lecture Notes on Data Engineering and Communications Technologies 191

Yap Bee Wah Dhiya Al-Jumeily OBE Michael W. Berry *Editors*

Data Science and Emerging Technologies Proceedings of DaSET 2023



Lecture Notes on Data Engineering and Communications Technologies

Volume 191

Series Editor

Fatos Xhafa, Technical University of Catalonia, Barcelona, Spain

The aim of the book series is to present cutting edge engineering approaches to data technologies and communications. It will publish latest advances on the engineering task of building and deploying distributed, scalable and reliable data infrastructures and communication systems.

The series will have a prominent applied focus on data technologies and communications with aim to promote the bridging from fundamental research on data science and networking to data engineering and communications that lead to industry products, business knowledge and standardisation.

Indexed by SCOPUS, INSPEC, EI Compendex.

All books published in the series are submitted for consideration in Web of Science.

Yap Bee Wah · Dhiya Al-Jumeily OBE · Michael W. Berry Editors

Data Science and Emerging Technologies

Proceedings of DaSET 2023



Editors Yap Bee Wah School of Information Technology UNITAR International University Petaling Jaya, Malaysia

Michael W. Berry University of Tennessee Knoxville, TN, USA Dhiya Al-Jumeily OBE Faculty of Engineering and Technology Liverpool John Moores University Liverpool, UK

ISSN 2367-4512ISSN 2367-4520 (electronic)Lecture Notes on Data Engineering and Communications TechnologiesISBN 978-981-97-0292-3ISBN 978-981-97-0293-0 (eBook)https://doi.org/10.1007/978-981-97-0293-0

© The Editor(s) (if applicable) and The Author(s), under exclusive license to Springer Nature Singapore Pte Ltd. 2024

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

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

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

This Springer imprint is published by the registered company Springer Nature Singapore Pte Ltd. The registered company address is: 152 Beach Road, #21-01/04 Gateway East, Singapore 189721, Singapore

Paper in this product is recyclable.

Conference Organization

Patron

Sahol Hamid Bin Abu Bakar, Vice-Chancellor, UNITAR International University, Malaysia

Honorary Chairs

Witold Pedrycz, University of Alberta, Canada Mohd Uzir Mahidin, Department of Statistics, Malaysia

General and Founding Chairs

Yap Bee Wah, UNITAR International University, Malaysia Dhiya Al-Jumeily OBE, Liverpool John Moores University, UK

International Advisory Committee

Mario Koeppen, Kyushu Institute of Technology, Japan Heri Kuswanto, Institute Teknologi Sepuluh Nopember Mohammed Bennamoun, University of Western Australia, Australia Chidchanok Lursinsap, Chulalongkorn University, Thailand

Proceeding Editors

Yap Bee Wah, UNITAR International University, Malaysia Dhiya Al-Jumeily OBE, Liverpool John Moores University, UK Michael W. Berry, University of Tennessee, USA

Finance and Registration Committee

Chong Kim Loy, UNITAR International University, Malaysia Normaiza Binti Mohamad, UNITAR International University, Malaysia Simranpreet Kaur Hansaram, UNITAR International University, Malaysia

Secretary

Sharifah Nurul Aina Binti Sayed Burhanudin, UNITAR International University, Malaysia

Technical Program Committee

Abdulaziz Al-Nahari, UNITAR International University, Malaysia (Chair) Azlin Ahmad, Universiti Teknologi MARA, Malaysia (Co-chair) Aitizaz Ali, UNITAR International University, Malaysia Hadi Naghavipour, UNITAR International University, Malaysia Jamila Mustafina, Kazan Federal University, Russia Sulaf Assi, Liverpool John Moores University, UK Umar Farooq Khattak, UNITAR International University, Malaysia Tahir Mehmood, UNITAR International University, Malaysia Wasiq Khan, Liverpool John Moores University, UK Abir Jaafar Hussain, Liverpool John Moores University, UK Anazida Binti Zainal, Universiti Teknologi Malaysia, Malaysia Wan Fairos Wan Yaacob, Universiti Teknologi MARA, Malaysia Sverina Azlin Md Nasir, Universiti Teknologi MARA, Malaysia Mohd Fadzil Hassan, Universiti Teknologi PETRONAS, Malaysia Sapto Wahyu Indratno, Institut Teknologi Bandung, Indonesia Farhad Nadi, UNITAR International University, Malaysia Syarifah Bahiyah Rahayu Binti Syed Mansoor, Universiti Pertahanan Nasional Malaysia, Malaysia

Program Book Committee

Noor Lees Ismail, UNITAR International University, Malaysia Iznora Aini Binti Zolkifly, UNITAR International University, Malaysia Rohaizah Abd Latif, UNITAR International University, Malaysia Jan Lunn, Liverpool John Moores University, UK

IT Committee

Danny Ngo Lung Yao, UNITAR International University, Malaysia Noor Azma Binti Ismail, UNITAR International University, Malaysia Paridah Binti Daud, UNITAR International University, Malaysia Intan Izzatul Fariza Binti Rossli, UNITAR International University, Malaysia Ker Boon Chin, UNITAR International University, Malaysia Muhamad Hadri Bin Mohd Hassan, UNITAR International University, Malaysia Mohd Aiman Mohd Sani, UNITAR International University, Malaysia Ahmad Ruzaini Bin Rahim, UNITAR International University, Malaysia

Logistics Committee

Mohd Farizudin Mohammad Fauzi, UNITAR International University, Malaysia Muhammad Shahrir Aizat Muhammad Shuhaili, UNITAR International University, Malaysia Roziyana Binti Bahrudin, UNITAR International University, Malaysia Nurlisnawati Binti Mohd Hijazi, UNITAR International University, Malaysia Syed Munir Barakbah Bin Syed Faozi Barakbah, UNITAR International University, Malaysia Mohd Amar Bin Mohd Mokhtar, UNITAR International University, Malaysia Farah Nazurah Binti Zainal, UNITAR International University, Malaysia

Publicity and Strategy Committee

Hadi Naghavipour, UNITAR International University, Malaysia (Chair) Wan Zakiyatussariroh Wan Husin, Universiti Teknologi MARA, Malaysia Bander Ali Al-Rimy, Universiti Teknologi Malaysia, Malaysia Izzatdin Abdul Aziz, Universiti Teknologi PETRONAS, Malaysia Norshahriah Binti Abdul Wahab, Universiti Pertahanan Nasional Malaysia, Malaysia

Corporate Committee

Stella Chua Ching Yee, UNITAR International University, Malaysia Marc Kevin Natusch, UNITAR International University, Malaysia Mohamad Shah Andrew Ibrahim, UNITAR International University, Malaysia Nikki Poh Li Yi, UNITAR International University, Malaysia

Industry Committee

Badrie Abdullah, UNITAR International University, Malaysia Visnuvarthen A/L. Sakayam, UNITAR International University, Malaysia

International Scientific Committee

Adel Al-Jumaily, Charles Sturt University, Australia (Chair) Naomi Bt Salim, Universiti Teknologi Malaysia, Malaysia Ku Ruhana Ku Mahmud, Universiti Utara Malaysia, Malaysia Siddhivinayak A. Kulkarni, MIT-World Peace University, Pune, India Simon Fong, University of Macau, China Richard Millham, Durban University of Technology, South Africa Layth Sliman, Efrei, Paris, France

Preface

This volume constitutes the proceedings of the Second International Conference on Data Science and Emerging Technologies (DaSET 2023) held from December 4 to 5, 2023, on a virtual platform. DaSET 2023 aims to provide a platform bringing together experts from academia, industries, government, and professional bodies to share recent trends in Artificial Intelligence and Emerging Technologies for Data-Driven Decisions. The theme of the conference is "Towards Green Artificial Intelligence and Sustainable Solutions."

DaSET is committed to creating a forum that brings academic and industry practitioners to share and establish collaborations toward impactful innovative research for community development, business success, and economic prosperity. This conference is an international conference in collaboration with UK Malaysia University Consortium, Universiti Teknologi MARA, Universiti Teknologi Malaysia, Centre for Data Science (CerDaS), Universiti Teknologi PETRONAS, Institut Teknologi Sepuluh Nopember, Indonesia; Chulalongkorn University, Thailand; Charles Sturt University, Australia; Institut Teknologi Bandung, Indonesia; National Defence University of Malaysia and Data Analytics and Collaborative Computing Group, University of Macau, China. We also appreciate the strong support from Microsoft Malaysia, Malaysia Digital Economy Corporation (MDEC), Cybersecurity Malaysia, and Statworks (M) Sdn Bhd.

From a total of 80 submitted papers, 40 were selected after a rigorous review process for oral presentation, and the Best Paper Awards were given for each track. The authors and presenters for these 40 papers represented 10 different countries. We thank all the reviewers and Springer Editors for their time spent reviewing the papers.

We are very honored to have Dato' Sri. Dr. Mohd Uzir Mahidin, Chief Statistician of Malaysia to officiate the opening of DaSET2023. We are privileged to have Prof. Witold Pedcryz, University of Alberta, for his special keynote address. We are proud to have eight distinguished international and local keynote speakers: Prof. Dr. Muhammad Khurram Khan, King Saud University, Saudi Arabia; Prof. Dr. Schahram Dustdar, Tu Wien, Vienna University of Technology, Austria; Prof. Dr. Seifedine Kadry, Noroff Education AS, Norway; Prof. Dr. Naomie Salim, Universiti Teknologi Malaysia, Assoc. Prof. Dr. Chin Kim On, Universiti Malaysia Sabah, Dato' Dr. Amirudin Abdul Wahab, Cybersecurity Malaysia, Ms. Puteri Anis Aneeza binti Zakaria, Statworks Group, and Mr. Raja Segaran, MDEC, Malaysia. All the distinguished speakers shared various data science and emerging technologies perspectives and projects which are beneficial for academics and industry practitioners.

We would like to thank Professor Emeritus Tan Sri Dato' Sri. Ir. Dr. Sahol Hamid Bin Abu Bakar, Vice-Chancellor of UNITAR International University for his great leadership, advice, and support of local and international academic activities to foster collaborations that lead to the exchange of knowledge and skills for research with impactful outcomes for social and economic prosperity.

We also thank the Series Editor, Springer, Lecture Notes on Data Engineering and Communications Technologies, for the opportunity to organize this guest-edited volume. We are grateful to Mr. Aninda Bose (Senior Publishing Editor, Springer India Pvt. Ltd.) and Mr. Radhakrishnan Madhavamani for the excellent collaboration, patience, and help during the preparation of this volume. We are confident that the volume will provide insightful information to researchers, practitioners, and graduate students in the areas of data science, artificial intelligence, and emerging technologies which are important in this digital information era. Last but not least, we thank all the DaSET 2023 committees for working tirelessly to ensure a successful conference.

Petaling Jaya, Malaysia Liverpool, UK Knoxville, USA Yap Bee Wah Dhiya Al-Jumeily OBE Michael W. Berry

About This Book

The book presents selected papers from the Second International Conference on Data Science and Emerging Technologies (DaSET 2023), held online at UNITAR International University, Malaysia, from December 4-5, 2023. This book aims to present current research and applications of data science and emerging technologies. The deployment of data science and emerging technology contributes to the achievement of the Sustainable Development Goals for social inclusion, environmental sustainability, and economic prosperity. Emerging technologies such as artificial intelligence and blockchain are useful for various domains such as marketing, health care, finance, banking, environmental, and agriculture. Innovations in the field of artificial intelligence continue to shape the future of work across nearly every industry. Data Science has a transformative effect on the economy, industry, and society. An important grand challenge in data science is to determine how developments in computational and social-behavioral sciences can be combined to improve wellbeing, emergency response, sustainability, and civic engagement in a well-informed, data-driven society. The topics of this book include, but are not limited to: artificial intelligence, machine and deep learning, statistical learning, and health and industrial applications.

Contents

Artificial Intelligence

A Comparative Study of Lemmatization Approaches for Rojak	2
Language Liu Jun Yoon, Xuan Yi Tan, Khai Yin Lim, Chi Wee Tan, Ling Ern Cheng, and Jenny Tan	3
Multi-aspect Extraction in Indonesian Reviews ThroughMulti-label Classification Using Pre-trained BERT ModelsNur Hayatin, Suraya Alias, Lai Po Hung, and Yuliana Setiowati	17
Artificial Intelligence (AI) Empowered Sign Language Recognition Using Hybrid Neural Networks Ambar Saxena, Nailya Sultanova, Jamila Mustafina, and Noor Lees Ismail	33
The Performance of GPT-3.5 in Summarizing Scientific and News Articles Sabkat Arshad, Muhammad Yaqoob, and Tahir Mehmood	49
Wound Stage Recognition Using YOLOv5Clair Abela and Frankie Inguanez	63
Harvest Palm Tree Based on Detection Through 2D LiDAR Sensor Using Power Equation Luqman Hakim Bin Yusof, Abdulaziz Yahya Yahya Al-Nahari, Danny Ngo Lung Yao, and Normaiza Mohamad	79
Enhancing Security Surveillance Through Business Intelligencewith NVIDIA DeepStreamVishal Pednekar, Nidhi Shettigar, and Sayli Tawhare	91
Fuzzified Hybrid Metaheuristics for QoS-Aware ServiceCompositionHadi Naghavipour, Farhad Nadi, and Ali Aitizaz	105

Machine/Deep Learning

Fraudulent E-Commerce Website Detection Using ConvolutionalNeural Network Based on Image FeaturesNurfazrina Mohd Zamry, Anazida Zainal, Eric Khoo,Mohamad Nizam Kassim, and Zanariah Zainudin	123
A Generic Framework for Ransomware Prediction and Classification with Artificial Neural Networks Saaman Nadeem, Tahir Mehmood, and Muhammad Yaqoob	137
Leveraging Gamification for Engaged Learning in Online Teaching and Learning Experiences Norshahriah Abdul Wahab, A'tifah Hanim Rosli, Syarifah Bahiyah Rahayu Syed Mansoor, Norazliana Akmal Jamaludin, and Siti Hajar Adam	149
Sentiment Analysis Using Large Language Models: A Case Study of GPT-3.5 Farhad Nadi, Hadi Naghavipour, Tahir Mehmood, Alliesya Binti Azman, Jeetha A/P Nagantheran, Kezia Sim Kui Ting, Nor Muhammad Ilman Bin Nor Adnan, Roshene A/P Sivarajan, Suita A/P Veerah, and Romi Fadillah Rahmat	161
Telecom Customer Experience Analysis Using Sentiment Analysisand Natural Language Processing—Comparative StudyAhmed Mohamed Abdou Ahmed, Abdulaziz Al-Nahari,Raghad Al-Shabandar, Chong Kim Loy, and A. H. Mohammed	169
Efficient Time Series of Smoothing and Auto-regressive Forecasting Models for Predicting Police Officer Fatalities in the USA Danush Nagappan, Manoj Jayabalan, Ahmad Alanezi, Farhad Nadi, and Thomas Coombs	181
Multimodal Emotion Recognition Using Attention-Based Model with Language, Audio, and Video Modalities Disha Sharma, Manoj Jayabalan, Nailya Sultanova, Jamila Mustafina, and Danny Ngo Lung Yao	193
Comparative Analysis of Emotion Recognition Using Large Language Models and Conventional Machine Learning Mangu Soujanya Rao, Thomas Coombs, Normaiza Binti Mohamad, Vinay Kumar, and Manoj Jayabalan	211
The Impact of Clustering-Based Sequential Multivariate OutliersDetection in Handling Missing ValuesMety Agustini, Kartika Fithriasari, and Dedy Dwi Prastyo	221

Contents

Sarcasm Detection in Newspaper Headlines Vishnu Sai Reddy Chilpuri, Saaman Nadeem, Tahir Mehmood, and Muhammad Yaqoob	237
Transformer-Based Named Entity Recognition Model—Tamil Language Karthi Dhayalan, Nailya Sultanova, Jamila Mustafina, and Paridah Daud	251
A Comparative Study of Methods for Topic Modelling in News Articles Swapna D. Rajan, Thomas Coombs, Manoj Jayabalan, and Noor Azma Ismail	269
Application of Deep Learning Algorithms to Terahertz Imagesfor Detection of Concealed ObjectsSoumen Sardar, Sulaf Assi, Iznora Aini Zolkifly, Manoj Jayabalan,Manea Alsaleem, Ammar H. Mohammed, and Dhiya Al-Jumeily OBE	279
Multivariate Comparative Analysis of Statistical and Deep Learning Models for Prediction Hardware Failure Saurabh Gupta, Raghad Alshabandar, Chong Kim Loy, and Ammar H. Mohammed	291
Statistical Learning	
A Case Study via Bayesian Network: Investigating Factors Influencing Student Academic Performance in Online Teaching and Learning During COVID-19 Pandemic Zheng Ning Looi, Poh Choo Song, Huai Tein Lim, and Sing Yan Looi	303
Harnessing the XGBoost Ensemble for Intelligent Prediction and Identification of Factors with a High Impact on Air Quality: A Case Study of Urban Areas in Jakarta Province, Indonesia Wahyu Wibowo, Harun Al Azies, Susi A. Wilujeng, and Shuzlina Abdul-Rahman	319
Modeling Earthquake Catalog in Sumatra by Space–Time Epidemic-Type Aftershock Sequences Model: Combining Davidon–Fletcher–Powell and Stochastic Declustering Algorithms Christopher Andreas, Achmad Choiruddin, and Dedy Dwi Prastyo	335
Small Area Estimation of Mean Years of Schooling Under TimeSeries and Cross-sectional ModelsReny Ari Noviyanti, Setiawan, and Agnes Tuti Rumiati	353
Probabilistic Seismic Hazard Analysis for Sulawesi-Maluku Region of Indonesia Using the Space–Time Epidemic-Type Aftershock Sequence Model Sonia Faradilla, Achmad Choiruddin, and Bambang Widjanarko Otok	369

Contents	5
----------	---

Application of Time Series Regression, Double Seasonal ARIMA, and Long Short-Term Memory for Short-Term Electricity LoadForecastingHafez Afghan and Hidayatul Khusna	385
A Bayesian Network for Classifying and Predicting Ship Collision Iis Dewi Ratih, Ketut Buda Artana, Heri Kuswanto, Emmy Pratiwi, and Muhammad Farhan Nuari	403
Outlier Detection in Simultaneous Equations with Panel DataSuci Ismadyaliana, Setiawan, and Jerry Dwi Trijoyo Purnomo	415
Assessing Departmental Efficiency at Sepuluh Nopember Institute of Technology: A Comparative Study Using Classical and Advanced Data Envelopment Analysis Zakiatul Wildani, M. Naufal Mubarik, Sri Pingit Wulandari, Lucia Aridinanti, and Muhammad Alifian Nuriman	429
Multivariate Adaptive Fuzzy Clustering Means Regression SplinesModel Using Generalized Cross-Validation (GCV) on StuntingCases in Southeast SulawesiMira Meilisa, Bambang Widjanarko Otok,and Jerry Dwi Trijoyo Purnomo	447
Statistical Inferences for Multivariate Generalized Gamma Regression Model Hasbi Yasin, Purhadi, and Achmad Choiruddin	463
Health and Industrial Applications	
W@rk: Attendance Application Framework Using Blockchain Technology Putra Roskhairul Fitri Kaha, Syarifah Bahiyah Rahayu, Afiqah M. Azahari, Mohd Hazali Mohamed Halip, and K. Venkatesan	479
Exploring the Impact of COVID-19 on Individuals' Mental HealthThrough Cluster AnalysisAzlin Ahmad, Siti Nabilah Mohd Abdul Hakim Amir,Ezzatul Akmal Kamaru Zaman, and Abdulaziz Al-Nahari	493
Tracking High Potential Transmission Risk Spots of InfectiousDisease Using Spatial Social Network Analysis and Visualisation(SSNAV) TechniquesIlham Abdul Jalil and Abdul Rauf Abdul Rasam	505
Evaluation of Machine Learning Algorithms for Early Prediction of Liver Disease Sushmitha Geddam, Sulaf Assi, Hadi Naghavipour, Manoj Jayabalan, Abdullah Al-Hamid, and Dhiya Al-Jumeily OBE	521

An Agricultural Information Recommendation Method Based on Matrix Decomposition Knowledge Graph Algorithm Ruipeng Tang, Narendra Kumar Aridas, and Mohamad Sofian Abu Talip	531
Performance of the Auxiliary Information Based Hybrid EWMAChart with Fast Initial ResponsePeh Sang Ng, Huai Tein Lim, Wai Chung Yeong, and Sajal Saha	545
Evaluation of Machine Learning Models for Breast CancerDetection in Microarray Gene Expression ProfilesMohammad Nasir Abdullah and Yap Bee Wah	563
Author Index	577

About the Editors

Professor Yap Bee Wah is the director of the Research and Consultancy Center at UNITAR International University Malaysia. She is the founding and general chair for DaSET2022: International Conference on Data Science and Emerging Technologies and the editor of the proceedings published in Lecture Notes on Data Engineering and Communications Technologies published by Springer. She was the conference chair of the International Conference on Soft Computing in Data Science (2015–2019 and 2021) and an editor of the SCDS conference proceedings published in the Springer CCIS series. She is also one of the editors of the book titled *Supervised and Unsupervised Learning for Data Science* published by Springer Nature Switzerland AG 2020. She actively published papers in ISI and Scopus journals such as *Expert Systems with Applications, Journal of Statistical Computation and Simulation, Communications in Statistics-Computation and Simulation, Journal of Clinical and Translational Endocrinology, and Computers, Materials and Continua.*

Professor Dhiya Al-Jumeily OBE is a professor of Artificial Intelligence and the president of eSystems Engineering Society. His research focus is on developing AI analytics for improving healthcare and the environment fulfilling the United Nations SDGs. His research has been well-recognized and featured in 300+ peer-reviewed articles, 40+ books/book chapters, and attracted over £7.5M. He has successfully supervised 25+ Ph.D. students to completion and has been an external examiner in UK and global universities. He is actively involved as a member of the editorial board/ review committee for numerous international journals. He is the founder/general series chair of the IEEE International Conference on Developments in eSystems Engineering since 2007 and DASET since 2022. He was promoted and appointed by The Queen to the Most Excellent Order of the British Empire, "OBE-Ordinary Officers of the Civil Division" of the said Most Excellent Order for the "Services to Scientific Research".

Professor Michael W. Berry is the co-author and an editor of sixteen books covering topics in scientific computing, information retrieval, text/data mining, and data science. He is the co-editor of the Soft Computing in Data Science volumes from

2015 to 2021 and Data Science and Emerging Technologies 2022 proceedings by Springer. He is also the co-author of popular books published by Society for Industrial and Applied Mathematics: *Understanding Search Engines: Mathematical Modeling and Text Retrieval, Second Edition*, and *Computational Information Retrieval*. He has published over 115 refereed journals and conference publications. He is a member of SIAM, ACM, MAA, ASEE, and the IEEE Computer Society and is on the editorial board of *Foundations of Data Science* (AIMS) and the *SIAM Journal on Matrix Analysis and Applications* (SIAM). He is also a certified program evaluator for the Computing Accreditation Commission (CAC) of the Accreditation Board for Engineering and Technology, Inc. (ABET).

Artificial Intelligence

A Comparative Study of Lemmatization Approaches for Rojak Language



Liu Jun Yoon, Xuan Yi Tan, Khai Yin Lim, Chi Wee Tan, Ling Ern Cheng, and Jenny Tan

Abstract Lemmatization is an important preprocessing step in most natural language processing (NLP) applications where it extracts a valid and linguistically meaningful lemma from an inflectional word. This allows different inflected forms of a word to be grouped into a common root which is the base-form or dictionary-form of a word, known as lemma. Due to the rapid spread of code-mixing languages like the Rojak language that mixes English with Malay, a lemmatizer capable of lemmatizing the language is needed for NLP applications involving this language. Thus, this work proposes a Rojak language lemmatization approach that is able to handle both languages without requiring users to input texts in different language separately. Various methods including rule-based, corpus-based, machine learning, and deep learning-based were experimented and compared using the English Web Treebank (EWT) and Indonesian GSD corpora from the Universal Dependencies (UD) framework. Besides, the effect of POS tags on the performance of lemmatizers was also evaluated based on the accuracy of the train and test sets. From the experiments conducted, the corpus-based approach produced the best results with

L. J. Yoon · X. Y. Tan · K. Y. Lim (🖂) · L. E. Cheng · J. Tan

Department of Computing and Information Technology, Penang Branch, Tunku Abdul Rahman University of Management and Technology, Tanjung Bungah, Malaysia e-mail: limky@tarc.edu.my

L. J. Yoon e-mail: yoonlj-pm20@student.tarc.edu.my

X. Y. Tan e-mail: tanxy-pm20@student.tarc.edu.my

L. E. Cheng e-mail: chengle-pm20@student.tarc.edu.my

J. Tan e-mail: tanj-pm20@student.tarc.edu.my

C. W. Tan

Faculty of Computing and Information Technology, Tunku Abdul Rahman University of Management and Technology, Kuala Lumpur Campus, Kuala Lumpur, Malaysia e-mail: chiwee@tarc.edu.my

© The Author(s), under exclusive license to Springer Nature Singapore Pte Ltd. 2024 Y. Bee Wah et al. (eds.), *Data Science and Emerging Technologies*, Lecture Notes on Data Engineering and Communications Technologies 191, https://doi.org/10.1007/978-981-97-0293-0_1 3

99.90% and 92.27% test set accuracy for Malay and English, respectively, whereas the deep learning-based with POS tag approach produced the worst results of 79.78 and 91.15%.

Keywords Lemmatization \cdot Rojak language \cdot Natural language processing (NLP) \cdot POS tag

1 Introduction

Lemmatizer plays an important role in the preprocessing part in most natural language processing (NLP) applications like information retrieval system (IRS) and machine translation system (MTS). This is because in most NLP applications, extraction of a valid and linguistically meaningful lemma from an inflectional word is required. For example, in MTS, the first phase involves analysis and testing the input text written in the source language for its grammatical correctness. Hence, a lemmatizer is required to perform text normalization to break down inflectional words into root and affixes [1, 2]. By performing lemmatization, different inflected forms of a word are grouped into a common root which is the base-form or dictionary-form of a word, known as lemma, that is a valid root and linguistically meaningful word. For example, if the words "*play, plays,* and *played*" are given to a lemmatizer, "*play*" will be returned by the lemmatizer as a lemma [2]. Since lemmatization is widely used in most preprocessing of NLP applications, a lemmatizer that is able to give a more accurate lemma has to be studied.

Code-mixing is a widespread phenomenon in Asia that mixes two or more languages. In Malaysia, English has been the focal mixing language practiced with Malay, resulting in the creation of Rojak language [3, 4]. For example, Malaysians tend to write "That is such a canggih gadget, I want one too". Due to the rapid spread of Rojak language, a lemmatizer capable of lemmatizing the Rojak language, namely a combination of Malay and English, is needed [5]. Besides, all existing stemmers and lemmatizers for Malay or Indonesian confuse roots with stems or lemmas. For instance, although the Sastrawi stemmer is a stemmer, it returns roots instead of stems, while Malaya also inherits this problem. This is because they do not consider the parts-of-speech (POS) tags when performing stemming. Conversely, MorphInd is currently the most widely used morphological analyzer for Indonesian regards roots as lemmas. No existing tools provide stems and lemmas [6]. Therefore, an approach capable of lemmatizing inflected Malay and English words, constituting the Rojak language, into its respective lemmas is needed. This study encompasses three primary objectives: (1) to introduce an approach capable of lemmatizing Rojak language, which is able to handle English and Malay simultaneously, without requiring the user to manually identify and separate the language; (2) to compare the performance of rule-based, corpus-based, machine learning, and deep learning-based lemmatization approaches; and (3) to evaluate the impact of POS tags on the performance of machine learning and deep learning-based lemmatizers.

This article is structured as follows. Section 2 analyzes related works. Section 3 describes the experimental set up and proposed methodology. Section 4 details the results and discussions. Section 5 concludes the work.

2 Literature Review

Recently, a lemmatizer that uses a novel dictionary lookup approach has been proposed for the Urdu language [7]. In the approach, the lemma returned depends on the corpus used to train the model. The lemma of a given word was only returned if both the word and its corresponding POS tag were present in the corpus. If the word and POS tag were not found, the lemmatizer would return the word as it was passed in, without any changes. Any differences caused by spelling mistakes or different blank spaces from the user input words would result in a no match in the corpus. The proposed lemmatizer was evaluated with and without POS tagging. The lemmatizer obtained the highest accuracy 76.44% when words without POS tags were lemmatized. On the other hand, the accuracy achieved by words that are passed in together with POS tags obtained a lower accuracy of 66.79%.

A lemmatizer for Icelandic known as Nefnir was proposed in 2019 [1]. Nefnir was developed from rules derived from a morphological database known as Database of Modern Icelandic Inflection (DMII). Suffix substitution rules were derived from the database and used to lemmatize tagged text. New rules were generated to minimize the number of errors in the training set until no further reduction in the error count. In Nefnir, it was assumed all word forms are identical to their lemma. A list of rules was generated for all the errors. The rule that minimizes the number of remaining errors was selected and applied to the training set until the number of errors does not reduce. The criteria for rule generation are that rules are generated only if the rule correctly lemmatizes at least two examples in the training set. The evaluation of Nefnir was performed to determine the accuracy of Nefnir in lemmatizing words with correct POS tags and words that are automatically tagged with POS tagger which was IceTagger. The accuracy achieved by Nefnir with correct POS tags was 99.55%, while for words that are tagged automatically with IceTagger, the accuracy was 96.88%. It was shown that Nefnir accuracy dropped when lemmatization was performed on words that were automatically tagged by POS tagger.

Another rule-based lemmatizer that uses the longest-affix-match approach was proposed for Kannada inflectional words [2]. In the approach, the input word that contains prefix or suffix would be applied with a set of linguistic rules to get the appropriate lemma. Prefixes and suffixes for Kannada inflected nouns and finite verbs were collected manually from Kannada grammar textbooks [8]. A root dictionary was created from Kannada dictionary "Kannada Rathnakosha". The proposed lemmatizer searched for a lemma in the root dictionary, and if the lemma was not found, the lemmatizer would append the obtained lemma to the root dictionary. This further improved the performance of the lemmatizer. The proposed lemmatizer was tested on four datasets with lemmatization performed on official circulars that achieved

an accuracy of 85.72%, newspaper 95.80%, legal documents 97.08%, and All India Radio news 95.39%. The accuracy of the proposed lemmatizer achieved above 85% on different dataset.

In 2020, lemmatization of the Russian language based on machine learning algorithms was proposed [8]. Vectorized word forms obtained from open dictionaries were fitted into various machine learning regression models, which were decision tree, random forest, extra tree, and bagging. Decision tree produced the highest accuracy on the lemmatization of real-world corpora, ABBYY corpus and Open-Corpora corpus, with the accuracy of 75.61% and 70.88%, respectively.

Deep learning sequence-to-sequence approach was proposed in 2021 to perform the automatic Romanian lemmatization [9]. The encoder and the decoder in the sequence-to-sequence model for lemmatization of Romanian words contain a single long short-term memory (LSTM) layer. The encoder and the decoder were enriched with one or two additional LSTM layers to improve the system's accuracy. When the deep learning models were evaluated on Romanian Explicative Dictionary (DEX) dataset, one layer LSTM-based architecture achieved the highest accuracy at both word and character levels with an accuracy of 95.93% and 97.29%, respectively. When POS information was included, the system's accuracy increased by 3.39% at word level and by 2.14% at character level resulting in an accuracy of 99.32% and 99.43%, respectively. The model's accuracy improved when POS information was included.

3 Proposed Methodology

This section describes the proposed framework proposed and the workflow for developing a lemmatizer for Rojak language texts that mixes English and Malay. Figure 1 shows the overall framework of the general process design.

3.1 Data Acquisition

Two corpora, one for the English language and another for the Malay language, were utilized in constructing the lemmatizer designed for the Rojak language. In light of the findings [10] indicating a similarity of over 90% between the Malay and Indonesian lexicons, an Indonesian corpus was used due to its relatively greater availability of resources as compared to Malay language. In this study, the English Web Treebank (EWT) corpus and Indonesian GSD corpus were employed [11–13]. Both the corpora were split into train, development, and test sets in the CoNLL-U file format. The train and development sets were combined to form the train set. Each set consists of sentences made up of words, where various information is provided for each word such as ID, FORM, LEMMA, UPOS, and XPOS. For building the lemmatizers, only the FORM, LEMMA, and UPOS which are universal POS tags

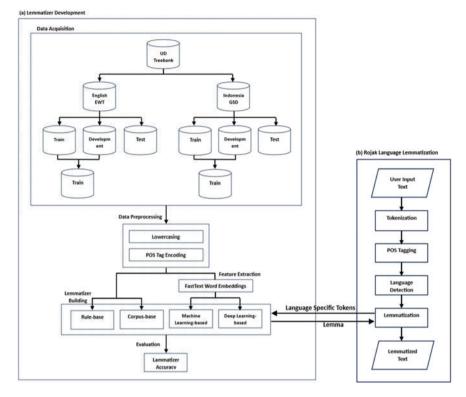


Fig. 1 Overall architecture of the general process design consisting of \mathbf{a} the development of lemmatization and \mathbf{b} the Rojak language lemmatization

are used. UPOS can be used to tag parts-of-speech for all languages, and there are 17 UPOS tags in the UD framework [14]. The number of sentences and words contained in each dataset is as shown in Table 1.

	English			Indonesian		
	Train	Development	Test	Train	Development	Test
Number of sentences	12,544	2001	2001	2001	2001	2001
Number of words	204,612	204,612	204,612	204,612	204,612	204,612

Table 1 Number of sentences and words in each dataset

3.2 Data Preprocessing

Lowercasing was performed in this stage to ensure all words can be mapped or lemmatized despite the different text casings in the corpus.

3.3 Feature Extraction

Feature extraction is required by machine learning-based and deep learning-based lemmatizers. In this process, each word and lemma was transformed into a sequence of numerical values, forming either the feature vector or word embedding [15]. A pre-trained model known as fastText [16] was used to generate word embedding of size 300 for each word in the corpus, which were then fed to the machine learning and deep learning models for training.

3.4 Lemmatizer Building

Lemmatizers for the English and Malay languages were individually constructed using distinct methods. Four different lemmatizers, consisting of rule-based, corpusbased, machine learning-based, and deep learning-based were developed. The training or development process for each method was replicated to create two lemmatizers of the same method for the distinct languages, namely English and Malay.

3.5 Rule-Based

Lowercasing was performed in rule-based lemmatizer on both the train and test sets. The rule-based lemmatizer utilizes information regarding word forms, lemmas, and POS tags to generate rules, as demonstrated in Table 2 for English and Table 3 for the Malay language, enabling it to carry out lemmatization. During the lemmatizer building stage, the initial step involves constructing a corpus that encompasses all word forms and lemmas. This corpus serves as a reference for the lemmatizer to determine whether a generated lemma represents a meaningful or dictionary-form of a word. Subsequently, rules for lemmatizing English and Malay words were formulated based on linguistic knowledge of the languages, observations from the corpora, and insights garnered from prior research. The overall workflow for the rule-based lemmatizer is outlined in Algorithm 1.

POS tag	Rules	Example		
NOUN/PROPN	DUN/PROPN 1. plural form \rightarrow singular form			
VERB	 past tense → present tense continuous tense → present tense perfect tense → present tense singular verb (with s) → plural verb (without s) 	1. played \rightarrow play 2. playing \rightarrow play 3. seen \rightarrow see 4. plays \rightarrow play		
ADJ	1. comparative \rightarrow positive 2. superlative \rightarrow positive	1. greater \rightarrow great 2. greatest \rightarrow great		
ADV	1. better \rightarrow well			
NUM	1. remove comma			
AUX	 third person → first person past tense → present tense 'm, am, 's, is, 're, are, art, was, were, being, been → be 've, has, had, having → have does, did, done, doing → do 'll → will 'd → would 	1. has \rightarrow have 2. had \rightarrow have		
PART	1. n't \rightarrow not			
PRON	1. object \rightarrow subject1. him \rightarrow 2. possessive (with 's) -> not possessive (without 's)2. who's			

 Table 2
 Rules for English rule-based lemmatizer

Algorith	n 1: The execution steps for the rule-based lemmatizer			
1	Input: User input text in Rojak language			
2	Output: Lemmatized words in either English or Malay			
3	Step 1: User inputs sequence of words in Rojak language			
4	Step 2: Convert the received text to lowercase			
5	Step 3: Detect language			
6	Step 4: Check the received POS tag of the received word (perform on both English and Malay text)			
7	Step 5: Perform lemmatization based on the received POS tag of the word where the will be 2 cases:			
8	if the POS tag indicates that no lemmatization should be performed on the word, the word will be returned as the lemma			
9	9 else the word will be lemmatized according to the rules defined for lemmatizing wo with the specific POS tag			
10	Step 6: Check if the output lemma exists in the corpus where there will also be 2 cases:			
11	if the lemma exists in the corpus which means that it is a valid word, return the lemma			

(continued	d)
12	else the received word will be returned as the lemma

The rules for lemmatizing English words were derived from linguistic knowledge and observing the word-lemma pairs in the training set. As there are many irregular verbs and plural forms in English, some irregular words found from the training set are explicitly added to the list of rules for lemmatizing English words to increase coverage and improve accuracy. As for lemmatizing Malay words, the rules are mainly derived from the findings of previous studies supplemented by linguistic knowledge [17].

3.6 Corpus-Based

As for the corpus-based lemmatizer, the training dataset and development database from UD Treebanks for the two different languages were used as training data. Test set was used for evaluating the corpus-based lemmatizer to determine the accuracy of the lemmatizer on unseen dataset. The idea behind this approach is to build a corpus containing all the possible lemmas for a word according to the different POS tags, for each language. After the corpus has been built, it can then be used to look for a matching word form and POS tag to return its corresponding lemma, given an input word and POS tag.

Therefore, to build the corpus, lowercasing was first performed on the UD Treebanks data. The words and lemma from the training set were then used to build a dictionary corpus for the lemmatizer. Words in the training set that are not in the dictionary were added to the corpus together with the POS tag. At the same time, the POS tag of the word was checked if it exists in the corpus. In cases where the POS tag was not found in the corpus, it was appended to the corpus along with its lemma. Hence, to use the lemmatizer, it takes a word and its POS tag as input and then checks the corpus to determine whether the given word and POS tag exist in the corpus. If they exist, the corresponding lemma is returned; otherwise, the received word is returned.

3.7 Machine Learning-Based

In machine learning-based lemmatizer, the train and development sets were combined to form a train set for training the lemmatizer. Features extraction was then performed on the words and lemma through fastText word embedding. In this study, decision tree regression model (DTR) was selected as it managed to yield the highest accuracy during the lemmatization process on the real-world corpora for the Russian language [8]. Thus, DTR was expected to be able to produce outstanding results in Rojak

POS tag	Rules	Example
NOUN	1. anti+ 2. peN+ 3. +an 4. antar+ 5. ke+ 6. +wan 7. per+ 8. +wati 9. ke + tidak+	 l. olahragawan → olahraga pengiriman → kirim perbaikan → baik
ADJ	1. non+ 2. ter+ 3. +an 4. ke+ 5. +nya 6. se+	 tercantik → cantik kecantikan → cantik secantik → cantik
VERB	1. meN+ 2. per+ 3. +kan 4. di+ 5. + i 6. ber+	 mengambil → ambil berduduk → duduk diletak → letak
NUM	1. nya+ 2. ber+	1. bertujuh \rightarrow tujuh 2. tujuhnya \rightarrow tujuh
PRON	1. nya \rightarrow dia 2. ku \rightarrow aku 3. mu \rightarrow kamu 4. kau \rightarrow kamu	

Table 3	Rules for Malay	
rule-based lemmatizer		

language lemmatizer. The DTR takes the word embedding of the inflected word as the input and learns to predict its corresponding output, which is the word embedding of the lemma of the input word. Depending on the type of approach, it can also take the encoded POS tag of the inflected word as input, allowing the DTR to learn from both the word and POS tag.

The DTR is a model based on the decision tree (DT) algorithm that learns using a tree structure that contains a root node, decision nodes, edges, and leaf nodes. It starts from the root node by selecting an attribute and splitting values as its starting point and continuously repeats this process to generate decision nodes, until it reaches its leaf nodes. The leaf nodes contain the values of the final prediction generated by the DT. In this case, the DTR uses the DT algorithm to perform lemmatization that is posed as a regression problem, whereby the predicted lemmas are represented by their feature vectors, which are sequences of continuous values. Using the generated tree, a path leading to the leaf nodes can be found to obtain the predicted values, which can also be expressed in rule form. Hence, the DTR is able to perform predictions and generate results for performing lemmatization [18].

Layers	Parameters	Parameters	
	With POS tag	Without POS tag	
Masking layer	$mask_value = 0$	mask_value = 0	
LSTM layer	units = 317 return sequences = true	units = 300 return sequences = true	
Dropout layer	rate = 0.2	rate = 0.2	
Time distributed layer	layer = Dense units = 300 activation = 'linear'	layer = Dense units = 300 activation = 'linear'	

Table 4 Parameters of the LSTM

3.8 Deep Learning-Based

Similar processes to the machine learning-based lemmatizer were performed. LSTM was employed for the lemmatization task. Similar to the machine learning-based approach, the LSTM also takes word embedding of the inflected word and optionally, encoded POS tags, and lemma as inputs and output respectively for training. With this, it is then able to predict the embedding of the lemma of a given word based on its embedding and optionally, its encoded POS tag.

LSTM is a variant of recurrent neural network (RNN) that retains its chain-like structure. However, the difference between the two is that RNN has difficulties dealing with long term dependencies, which are sequential data that require previous data or context to be retained for a longer period due to the vanishing gradient problem. LSTM overcomes this problem with the use of cell state, which differentiates its recurrent unit architecture from RNNs. It uses three types of gates, namely the forget gate, input gate, and output gate, in the cell state to control addition and deletion of information from the recurrent unit, while allowing information to flow through it using the door mechanism [19]. The parameters of the layers in the LSTM with POS and without POS tags are as shown in Table 4. About 100 epochs with a batch size of 4 were used to train the model. The chosen loss function was cosine_ similarity, and the optimization technique utilized was rmsprop.

3.9 Evaluation

Both the English and Malay lemmatizers built on the different methods were evaluated using the accuracy (Eq. 1). To perform the evaluation, the test sets of the two corpora were used.

$$Accuracy = \frac{\text{Number of Words Correctly Lemmatized}}{\text{Total Number of Words in the Test Set}}$$
(1)