

**UNIVERSITI TEKNOLOGI MARA**

**A HYBRID OF FUZZY C-MEANS  
CLUSTERING AND LATENT  
DIRICHLET ALLOCATION FOR  
ANALYSING PHILANTHROPIC  
CORPORATE SOCIAL  
RESPONSIBILITY ACTIVITIES**

**NIK SITI MADIHAH BINTI NIK MANGSOR**

Thesis submitted in fulfillment  
of the requirements for the degree of  
**Master of Science**  
**(Computer Science)**

**College of Computing, Informatics, and Media**

**February 2023**

## ABSTRACT

To date, philanthropic corporate social responsibility (PCSR) activities are ad-hoc in nature, where assistance is provided more to basic needs with very little attention to activities that can contribute to eradicating poverty. Based on previous related literatures, it is found that there is no proper categorization and documentation of CSR-related activities. Conventional clustering algorithms are able to extract only a single set of flat topics in which local and global topics are mingled and cannot be separated. Therefore, this research aims to develop hybrid fuzzy document clustering techniques to improve attribute cluster membership values of PCSR activities. This study has extended document clustering technique by integrating the traditional document clustering application with topic modeling approach. This integrated approach is able to produce precise results where it can help infer more coherent themes of PCSR activities. The analysis involved five-year data from the annual reports of 19 CSR-award winning companies in Malaysia where they were converted into a structured format, collated and summarized. Then, text pre-processing for data cleaning was performed followed by identification of the best Latent Dirichlet Allocation (LDA) topic modelling technique that was used to integrate document clustering. Next, documents were then clustered using integration technique of K-Means clustering and LDA as the benchmark method before the proposed integration of Fuzzy C-Means (FCM) clustering and LDA was carried out. Findings from the proposed method in this study revealed thirteen clusters that represent thirteen types of PCSR activities performed by the CSR-award winning companies in Malaysia. This finding was then verified and agreed by experts in CSR on the validity of the obtained PCSR cluster which is able to describe the categories of PCSR activities accordingly. Lastly, the performance of the benchmark method and proposed method was evaluated using model evaluation technique. Based on the evaluation, it was found that K-Mean's performance is better than FCM's performance in terms of computational time. However, FCM is better than K-Means in terms of quality, since FCM produces more quality clusters than K-Means which could produce much better in clustering PCSR activities. As a result, the study shows that the proposed integration method produces much better categories of PCSR activities because FCM clustering is applicable to such incomplete external knowledge and it offers the opportunity to deal with data that belong to more than one cluster at the same time. The findings offer an insight to be considered by companies in strategizing the CSR activities, particularly philanthropic responsibility in ensuring optimum impact to innovatively support the society and contribute towards poverty mitigation.

## **ACKNOWLEDGEMENT**

In the name of Allah, Most Gracious, Most Merciful. All praises be to Allah, the Cherisher and Sustainer of the world, for giving me the strength and opportunity to embark on my MSc and for completing this long and challenging journey successfully.

Firstly, my gratitude and thanks go to my supervisor Dr Syerina Azlin Md Nasir and also my co-supervisor PM Dr Shuzlina Abdul Rahman and Dr Zurina Ismail, for providing me the opportunity to do the study, and giving me all support and guidance, which was essential in completing my MSc. I am extremely thankful to them for providing a pleasant and nice support and guidance.

Finally, this thesis is dedicated to my beloved family especially my parents also my friends for the vision and determination to educate me throughout this journey. This piece of victory is dedicated to all of them. Alhamdulillah.

Thank you.

# TABLE OF CONTENTS

	<b>Page</b>
<b>CONFIRMATION BY PANEL OF EXAMINERS</b>	<b>ii</b>
<b>AUTHOR'S DECLARATION</b>	<b>iii</b>
<b>ABSTRACT</b>	<b>iv</b>
<b>ACKNOWLEDGEMENT</b>	<b>v</b>
<b>TABLE OF CONTENTS</b>	<b>vi</b>
<b>LIST OF TABLES</b>	<b>ix</b>
<b>LIST OF FIGURES</b>	<b>x</b>
<b>LIST OF SYMBOLS</b>	<b>xi</b>
<b>LIST OF ABBREVIATIONS</b>	<b>xii</b>
<b>CHAPTER ONE: INTRODUCTION</b>	<b>1</b>
1.1 Research Background	1
1.2 Problem Statement	2
1.3 Research Questions	3
1.4 Research Objectives	3
1.5 Significance of Study	3
1.6 Scope and Limitation	4
1.7 Thesis Organization	4
<b>CHAPTER TWO: LITERATURE REVIEW</b>	<b>5</b>
2.1 Introduction	5
2.2 Corporate Social Responsibility (CSR)	5
2.3 Document Clustering	8
2.3.1 Categories of Document Clustering	8
2.3.2 Types of Document Clustering Methods	9
2.3.3 Current Issues in Document Clustering	10
2.3.4 Document Clustering Methods in the CSR Domain	11
2.3.5 Fuzzy-based Approach for Document Clustering	11
2.4 Topic Modeling	17
2.4.1 Overview of Topic Modeling	17
2.4.2 Type of Topic Modeling	17
2.4.3 Topic Modeling Used in Previous Studies	19
2.5 Previous Study on Hybridization of Topic Modeling and Document Clustering	24
2.6 Conclusion	28
<b>CHAPTER THREE: RESEARCH METHODOLOGY</b>	<b>29</b>
3.1 Introduction	29
3.2 Research Framework	29
3.3 Phase 1: Data Collection	30

3.4	Phase 1: Text Pre-processing	33
3.4.1	Convert Text to Lowercase, Remove Numbers, Punctuation and Symbols	33
3.4.2	Remove Stopwords	34
3.4.3	Stemming	35
3.4.4	Lemmatization	35
3.4.5	Remove Words with 3 Characters and Below	36
3.4.6	Tokenization	37
3.4.7	Generating Document Term Matrix (DTM)	37
3.4.8	Building Corpus	38
3.5	Phase 2: Topic Modeling	39
3.5.1	Identify Latent Dirichlet Allocation (LDA) Techniques	39
3.5.1.1	Gibbs Sampling	39
3.5.1.2	Expectation-Maximization	40
3.5.1.3	Variational Bayes Inference	42
3.5.2	Comparing Topic Modeling Techniques	43
3.5.2.1	Topic Coherence Measure	43
3.5.2.2	pyLDAvis	44
3.6	Phase 2: Document Clustering	44
3.6.1	Determine Document Clustering Methods	44
3.6.1.1	K-Means	44
3.6.1.2	Fuzzy C-Means	45
3.6.2	Selecting the Optimal K Clusters using Silhouette Index	47
3.7	Phase 2: Hybridization of Document Clustering and Topic Modeling Technique	47
3.8	Phase 3: Model Evaluation	48
3.8.1	Time Complexity	49
3.8.2	Normalized Mutual Information (NMI)	49
3.9	Phase 3: Obtaining the Local and Global Topics to Describe PCSR Cluster	50
3.9.1	Verify and Validate by Experts	50
3.10	Conclusion	50

## **CHAPTER FOUR: RESULTS AND DISCUSSION ON THE PATTERN RECOGNITION AND DATA MINING METHODS IN CHEMOMETRICS. 51**

4.1	Introduction	51
4.2	Textual Dataset	51
4.3	Text Pre-processing Results	52
4.3.1	Text Cleaning, Stemming and Lemmatize	52
4.3.2	Tokenization	54
4.3.3	Generating Document Term Matrix (DTM) and Building Corpus	55
4.4	Topic Modeling Results	56
4.4.1	LDA Variational Bayes Inferences	56
4.4.2	LDA Gibbs Sampling	58
4.4.3	LDA Expectation Maximization	59
4.4.4	Comparative Evaluation of LDA Techniques	60
4.5	K-Means Clustering Results	62
4.6	Fuzzy C-Means Clustering Results	66
4.7	Model Evaluation	70
4.7.1	Time Complexity	70