

UNIVERSITI TEKNOLOGI MARA

**Semantic Analysis of Hadith for Topic Classification
Using Latent Semantic Indexing (LSI)**

Aiman Haziq Bin Ibrahim

**Thesis submitted in fulfilment of the requirements of Bachelor of
Computer Science (Hons.) College of Computing, Informatics and
Mathematics**

February 2024

ACKNOWLEDGEMENT

Alhamdulillah, praises and thank you to Allah because of His Almighty and His utmost blessings, I was able to complete this research within the designated timeline. I would like to extend my sincerest thanks to my supervisor, Roslan Sadjirin, for his guidance and support throughout the project journey. Without him, I would not have been able to complete this project with such success.

I would like to extend my heartfelt gratitude to my lecturer, Madam Ummu Fatimah Binti Mohd Bahrin, for her support and knowledge in the subject of CSP650 that have greatly assisted me in my research. Furthermore, my beloved parents also deserve a special mention for their unwavering support and encouragement during this project. I am forever grateful to them for their love and support.

Finally, I would like to express my appreciation to my dear friends and everyone who has been a part of this journey and has provided me with their unwavering support and encouragement. Without them, this project would not have been possible.

ABSTRACT

The aim of this project is to provide a framework utilizing Latent Semantic Indexing (LSI) to categorize topics in Hadith texts for semantic analysis. Islamic teachings place a high value on the hadith literature, which records the words and deeds of Prophet Muhammad (peace be upon him). To make it simple to access, retrieve, and comprehend pertinent information, Hadith writings must be effectively organized and categorized depending on their topics. The subjectivity, labor-intensive manual categorization, and insufficient capture of semantic links within texts are only a few of the drawbacks of the currently available approaches for Hadith topic classification. To address these challenges, LSI-based framework was proposed that leverages the latent semantic meaning in Hadith texts. LSI captures the underlying semantic relationships between words and enables more accurate topic classification. The research framework consists of six phases, including a preliminary study, requirement analysis, data finding, development, evaluation, and documentation. The data finding involves collecting and preprocessing reliable Hadith datasets. Development focuses on creating an information retrieval system using LSI. The evaluation assesses the system's performance through metrics like cosine similarity, precision, recall, and F1 Score. The experiment assessed the effectiveness of LSI by utilizing ten queries and relevant judgements, precision ranged from 5.4% to 100%, recall from 0% to 65%, yielding an average F1 Score of 19.4%. Finally, documentation encompasses writing a comprehensive report that includes background, methodology, findings, and conclusions.

TABLE OF CONTENTS

CONTENT	PAGE
SUPERVISOR APPROVAL	i
STUDENT DECLARATION	ii
ACKNOWLEDGEMENT	iii
ABSTRACT	iv
TABLE OF CONTENTS	v
LIST OF FIGURES	vii
LIST OF TABLES	ix
CHAPTER ONE	1
1.1 Background of Study	1
1.2 Problem Statement	3
1.3 Objectives	5
1.4 Project Scope	5
1.5 Project Significance	5
1.6 Overview of Research Framework	6
1.7 Conclusion	7
CHAPTER TWO	9
2.1 Overview of Hadith	9
2.1.1 Definition of Hadith	9
2.1.2 The importance of Hadith	10
2.2 Natural Language Processing (NLP)	10
2.2.1 Latent Semantic Indexing Algorithm (LSI)	12
2.3 Similar Application	14
2.4 Similar Algorithm in Various Application	27
2.5 Implication of Literature	44
2.6 Conclusion	44
CHAPTER THREE	46
3.1 Overview of Research Methodology Framework	46

3.2	Preliminary Study	51
3.2.1	Knowledge acquisition	51
3.3	Data Collection	52
3.3.1	Data Description	52
3.3.2	Data Pre-processing	56
3.4	Prototype Design	58
3.4.1	Conceptual Framework	58
3.4.2	Flowchart	66
3.4.3	Interface Design	66
3.4.4	Pseudocode	67
3.5	Development	68
3.5.1	Software Recommendation	69
3.5.2	Hardware Recommendation	69
3.6	Evaluation	70
3.6.1	Cosine Similarity	70
3.6.2	Evaluation Metrics	71
3.6.3	K-dimensional value	72
3.7	Gantt chart	72
3.8	Conclusion	73
CHAPTER FOUR		74
4.1	Data Preprocessing	74
4.1.1	Tokenization	74
4.1.2	Stop Word Removal	75
4.1.3	Lemmatization	77
4.2	Result for Objective 2	78
4.2.1	Input Representation	78
4.2.2	LSI Implementation	80
4.2.3	Prototype Interface	84
4.3	Evaluation	86
4.3.1	Cosine Similarity	87
4.3.2	Precision, Recall and F1 Score	89
4.4.2.1	Experiment on k=5, threshold = 0.5	89
4.4.2.2	Experiment on k=5, threshold = 0.6	96