



**SEMANTIC TEXT CLUSTERING USING ENHANCED VECTOR SPACE MODEL
USING NEPALI LANGUAGE**

A Dissertation

Submitted to
The Central Department of Computer Science and Information Technology, Institute of
Science and Technology
Tribhuvan University

In Partial Fulfillment of the Requirements for the degree of
Master of Science in Computer Science and Information Technology

By
Chiranjibi Sitaula
January 2012

Under the Supervision of
Prof. Dr. Shashidhar Ram Joshi
(IOE, TU)
Co-Supervisor
Mr. Bikash Balami

Date:.....

Recommendation

We hereby recommend that the dissertation prepared under my supervision by **Mr. Chiranjibi Sitaula** entitled “**Semantic Text Clustering using Enhanced Vector Space Model using Nepali Language**” be accepted as fulfilling in part requirements for the degree of Masters of Science. In my best knowledge this is an original work in computer science.

Dr. Shashidhar Ram Joshi
Professor
Institute of Engineering, Pulchowk,
Nepal
(Supervisor)

Mr. Bikash Balami
Lecturer
Central Department of Computer Science
and IT, T.U, Nepal
(Co-Supervisor)



Tribhuvan University
Institute of Science and Technology
Central Department of Computer Science and Information Technology

We certify that we have read this dissertation work and in our opinion it is satisfactory in the scope and quality as a dissertation in the partial fulfillment for the requirement of Master of Science in Computer Science and Information Technology.

Evaluation Committee

Dr. Tanka Dhamala
Asst.Professor
Head, Central Department of
Computer Science and Information
Technology
Tribhuvan University

(External Examiner)

Date: -----

Dr. Shashi Dhar Ram Joshi
Professor
Institute of Engineering
Pulchowk, Nepal
(Supervisor)

(Internal Examiner)

Date: -----

Acknowledgement

This dissertation has been prepared under **Central Department of Computer Science and Information Technology (Tribhuvan University)**, Kirtipur and the study was done for the “Semantic Text Clustering using Enhanced Vector Space Model using Nepali Language”. I am very grateful to our department for giving me an enthusiastic support. I deeply extend my heartfelt acknowledgement to my respected teachers and dissertation advisor **Prof. Dr. Shashidhar Ram Joshi**, Head of Department of Electronics and Computer Engineering, Institute of Engineering, Pulchowk, for his wholehearted cooperation, encouragement and string guidelines throughout the preparation of this study. With this regard, I wish to extend my sincere appreciation to respected Head of Department of Central Department of Computer Science and Information Technology, **Assoc. Prof. Dr. Tanka Dhamala** for his kind help. I am very grateful and thankful to all the respected teachers of CDCSIT, TU, for granting me broad knowledge and inspirations within the time period of two years.

My special thanks goes to my co-supervisor **Mr. Bikash Balami** who directly and indirectly extended his hands in making this thesis work a success. Similarly, I would like to thank to all of my friends of my class who helped me directly and indirectly to complete my thesis.

As we know that “there is an eraser at the top of the pencil”, so mistakes are done by human beings of which we are also members. Although, I don’t experience any enjoyment, since I did have a very hard time while completing this project, I have done my best. There may be some errors in our project, so any suggestion regarding the mistake of this project will be always welcomed.

Abstract

The vector space model is popular method for the clustering process while doing research in the field of text mining. The main reason of its popularity is its less computational overhead and simplicity. Classical vector space model can not be used for semantic analysis purpose because it simply uses syntactic model for clustering. In order to cluster the documents in sentence level, how individual keyword plays the important roles in the text clustering, is studied in this work. For this, an Enhanced method is proposed which can easily outperform classical vector space model due to the involvement of fuzzy set approach. The classical Vector Space Model is enriched with fuzzy set so as to form the Enhanced Vector Space Model in text clustering. In order to give the semantic text clustering, fuzzy set plays crucial role in addition to classical Vector Space Model.

TABLE OF CONTENTS

Details No.	Page
Acknowledgements	i
Abstract	ii
Table of Contents	iii
List of Figures.	vi
List of Tables.....	vii
Chapter 1	
INTRODUCTION-----	1-3
1.1 Motivation-----	2
1.2 Problem Statement-----	3
1.3 Research Question-----	3
1.4 Organization of thesis-----	3
Chapter 2	
LITERATURE REVIEW-----	4-6
2.1 Text Similarity Study in Literatures-----	4
2.1 Application Area of text similarity-----	5
2.3 Testing Similarity of words-----	6
Chapter 3	
MEASURING TEXT SIMILARITY-----	7-12
3.1 Fuzzy Set-----	8
3.2 Membership function-----	8
3.3 Vector Space Model-----	8
3.4 Similarity Measure Function-----	9
3.5 Weight Functions-----	12
Chapter 4	
THE PROPOSED SYSTEM-----	13-21
4.1 Fuzzy Set-----	13

4.2 Membership function-----	14
4.3 Vector Space Model-----	15
4.4 Enhanced Vector Space Model-----	17
4.5 Cosine Similarity-----	18
4.6 Clustering-----	18
Chapter 5	
RESULTS AND DISCUSSION-----	22-26
5.1 Analysis using Random Index and Confusion Matrix-----	22
5.2 Analysis using Number of clusters produced-----	25
Chapter 6	
CONCLUSION-----	27-28
6.1 Summary of the study-----	27
6.2 Limitations and future works-----	28
REFERENCES-----	29
APPENDICES-----	30-50

LIST OF TABLES

TABLES	Contents	Page
	Table 4.3(a) Document to be clustered-----	15
	Table 4.3(a) Document to be clustered-----	15
	Table 4.3(a) Document to be clustered-----	16
	Table 4.3(a) Document to be clustered-----	16
	Table 5.1(a):Experimental result for single keyword-----	22
	Table 5.1(b):Experimental result for multi keyword document-----	22

LIST OF FIGURES

FIGURES	Contents	Page
Figure 3.1:	Document Vector representation	3
Figure 3.4.1:	w1 and w2 are the text vectors, the black line between w1 and w2 is the L1 Norm.	9
Figure 3.4.2:	L2 Norm between two vectors	10
Figure 3.4.3:	Showing the angle between two text documents w1 and w2	11
Figure 4.1	fuzzy set	14
Figure 4.2	Gaussian membership functions	14
Figure 4.6(a)	Complete algorithms for proposed model	18
Figure 4.6(b)	operational flow of the proposed model	19
Figure 4.6(c)	Vector generated from programming	20
Figure 4.6(d)	Table formed after calculating distance between query and document vector	20
Figure 4.6(e)	raw clusters after performing threshold 0.5	21
Figure 4.6(f)	Final cluster	21
Figure 5.1(a):	Output obtained for single keyword document	23
Figure 5.1(b):	Output obtained for multi word document	24
Figure 5.2(a)	relation between first 24th documents and number of cluster produced	25
Figure 5.2(b)	relation between second 24th documents and number of cluster produced	25
Figure 5.2(a)	relation between third 24th documents and number of cluster produced	26
Figure 5.2(b)	relation between forth 24th documents and number of cluster produced	26