

DECLARATION

I hereby declare that the thesis entitled “**SMS based FAQ Retrieval using Hybrid Similarity Measure**” which is being submitted to the **Delhi Technological University**, in partial fulfillment of the requirements for the award of degree of **Master of Technology in Computer Science Engineering** is an authentic work carried out by me. The material contained in this thesis has not been submitted to any university or institution for the award of any degree.

Sonal Meena

Department of Computer Engineering

Delhi Technological University,

Delhi.

CERTIFICATE



DELHI TECHNOLOGICAL UNIVERSITY

(Govt. of National Capital Territory of Delhi)

BAWANA ROAD, DELHI-110042

Date: _____

This is to certify that the thesis entitled “**SMS based FAQ Retrieval using Similarity Measure**” submitted by **Sonal Meena (Roll Number: 2K11/CSE/15)**, in partial fulfillment of the requirements for the award of degree of Master of Technology in Computer Science Engineering, is an authentic work carried out by her under my guidance. The content embodied in this thesis has not been submitted by her earlier to any institution or organization for any degree or diploma to the best of my knowledge and belief.

Project Guide

Mr. Manoj Kumar

Associate Professor

Department of Computer Engineering

Delhi Technological University, Delhi-110042

ACKNOWLEDGEMENT

I take this opportunity to express my deepest gratitude and appreciation to all those who have helped me directly or indirectly towards the successful completion of this thesis.

Foremost, I would like to express my sincere gratitude to my guide **Mr. Manoj Kumar, Associate Professor, Department of Computer Engineering, Delhi Technological University, Delhi** whose benevolent guidance, constant support, encouragement and valuable suggestions throughout the course of my work helped me successfully complete this thesis. Without his continuous support and interest, this thesis would not have been the same as presented here.

Besides my guide, I would like to thank the entire teaching and non-teaching staff in the Department of Computer Science, DTU for all their help during my course of work.

SONAL MEENA

ABSTRACT

Mobile technology gave contribution to the progress of media of communication for example: chats, emails and short message services (SMS). The Popularity, utility and simplicity of SMSes is encouraging people to access information via SMSes, accessing information via internet creates hassle, it's not necessary that internet connection is always available. So user can clarify their query, make complaint and get updates of result etc by sending SMS. Accessing information in such a manner makes information access very economic and easy for everybody from rural to metro city people.

“FAQ retrieval” means there is corpora of frequently asked questions, and user sends a query in SMS language to retrieve some information. Such systems finds best match from FAQ corpora for given user defined query written in SMS language. The main problem in SMS language is the noise associated with it. Spelling mistakes, transliteration, phonetic spellings, abbreviations and short forms create difficulties in string matching.

In proposed work, a novel approach has been presented by developing Hybrid similarities which evaluates similarity scores with the questions in the corpus for SMS query. In this way, we can further improve the accuracy of the SMS based FAQ system significantly by refining the results of the system using different hybrid similarity scores.

Table of Contents

COVER PAGE

DECLARATION.....	i
CERTIFICATE.....	ii
ACKNOWLEDGEMENT.....	iii
ABSTRACT.....	iv
Table of Contents.....	v
List of Figures.....	viii
List of Tables.....	xi
List of Equations.....	xi
CHAPTER 1	1
INTRODUCTION.....	1
1.1. MOTIVATION OF WORK.....	3
1.2. RESEARCH OBJECTIVE	4
1.3. RELATED WORK	4
1.4. SCOPE OF WORK.....	5
1.5. ORGANIZATION OF THESIS	6
CHAPTER 2	8
LITERATURE SURVEY.....	8
2.1. BASIC CONCEPTS OF SMS BASED QA SYSTEM.....	8
2.1.1. SMS BASED QUESTION ANSWERING SYSTEM.....	8
2.1.2. TECHNIQUES FOR SMS BASED QA SYSTEM.....	9
2.1.2.1. HUMAN INTERVENTION BASED SYSTEM.....	9
2.1.2.2. NATURAL LANGUAGE PROCESSING BASED SYSTEM.....	10
2.1.2.3. INFORMATION RETRIVAL BASED SYSTEM.....	10
2.1.2.4. FREQUENTLY ASKED AUESTION BASED SYSTEM.....	11
2.2. SMS BASED FAQ RETRIEVAL SYSTEM.....	11
2.2.1. BASELINE.....	11

2.2.2. TEXT NOISE.....	12
2.2.3. SEARCH ALGORITHM.....	13
2.3. PROBLEM FORMULATION.....	14
2.3.1. XML DATABSED.....	15
2.3.2. TOKENIZATION.....	15
2.3.3. STOP WORDS.....	16
2.3.4. STEMMING.....	16
2.3.4.1. PORTER'S ALGORITHM.....	18
CHAPTER 3	22
SIMILARITY MEASURES.....	22
3.1. SOUNDEX SIMILARITY MEASURE	23
3.2. JACCARD'S SIMILARITY MEASURE	26
3.3. COSINE SIMILARITY MEASURE.....	28
CHAPTER 4	30
PROPOSED HYBRID SIMILARITY MEASURES.....	30
4.1. PROPOSED METHODS	31
4.1.1. HYBRID SOUNDEX SIMILARITY MEASURE.....	31
4.1.2. HYBRID S-JACCARD SIMILARITY MEASURE.....	33
4.1.3. HYBRID S- COSINE SIMILARITY MEASURE.....	34
4.1.4. HYBRID JACCARD SIMILARITY MEASURE.....	35
CHAPTER 5	36
IMPLEMENTATION AND EXPERIMENTAL RESULTS	36
5.1. ENVIRONMENTAL SETUP	36
5.1.1. HARDWARE CONFIGURATION.....	36
5.1.2. SOFTWARE CONFIGURATION.....	37
5.2. DATASETS	37
5.3. ANALYSIS AND RESULTS.....	38
5.4. SUMMARY	47
CHAPTER 6	48
CONCLUSION AND FUTURE WORK	48
6.1. CONCLUSION	48

6.2. FUTURE WORK	49
6.2.1. N-GRAMS TECHNIQUE.....	49
6.2.2. INVERSE BIGRAM FREQUENCY.....	49
6.2.3. CACHING RESULTS.....	49
6.2.4. EXTENSION OF WORK FROM MONOLINGUAL TO MULTILINGUAL.....	50
REFERENCES	51
APPENDIX A- CODING	
APPENDIX B- SOUNDEX ALGORITHM	

List of Figures

Figure 2.1. Portesr' Algorithm.....	21
Figure 4.1. Soundex Matching.....	32
Figure 5.1. FAQ Format.....	37
Figure 5.2. SMS Format.....	38
Figure 5.3. Graph for T1 with Maximum Value.....	40
Figure 5.4. Graph for T1 with S-Jaccard	40
Figure 5.5 Graph for T1 with S-Cosine	41
Figure 5.6 Graph for T1 with Hybrid Jaccard.....	41
Figure 5.7 Graph for T2 with Maximum Value.....	42
Figure 5.8 Graph for T2 withS-Jaccard.....	43
Figure 5.9 Graph for T2 with S- Cosine.....	43
Figure 5.10. Graph for T2 with Hybrid Jaccard.....	44
Figure 5.11. Graph for T3 with Maximum Value.....	45
Figure 5.12. Graph for T3 with S-Jaccard	45
Figure 5.13. Graph for T3 with S-Cosine	46
Figure 5.14. Graph for T3 with Hybrid Jaccard.....	46

List of Tables

Table 3.1. Soundex Table	25
Table 5.1. Experimental Result For T1	39
Table 5.2. Experimental Result For T2.....	42
Table 5.3. Experimental Result For T3.....	44

List of Equations

Equation 3.1. Jaccard's Similarity.....	26
Equation 3.2. Jaccard's Dissimilarity.....	26
Equation 3.3. Cosine Formula	28
Equation 3.4. Cosine Similarity	28
Equation 4.1. Hybrid S-Jaccard	33
Equation 4.2. Hybrid S-Cosine	34
Equation 4.3. Hybrid Jaccard.....	35