

100/100/2010

ACKNOWLEDGEMENTS

My special thanks go to Dr. Ruwan Weerasinghe of the University of Colombo, the principal supervisor of my project for all the support and guidance extended to me. I found his expertise and suggestions extremely valuable. My sincere gratitude is also extended to my supervisors Ms. G.S. Makalande and Mr. P. Dias for all their support and guidance. In addition, I would like to thank our coordinator of the Master Degree programme, Dr. Sunethra Weerakoon for her encouragement and support. I also would like to thank LAcNet in helping me to obtain access to news items from the SLNET archives.

Finally, I am very grateful to Mr. Piyasiri Vithanage, the Head of the Department of Economics, University of Ruhuna for his support and encouragement.

Library - USJP



195839

195839

ABSTRACT

The advancement of information technology in the modern world has contributed towards enhancing the quality of life of people around the world. To keep pace with this rapid development, it is important to have links with the enormous network called the Internet. This enables people to have access to information resources, keep abreast of news, send timely E-mail, and have interactive remote conferences. However, a tremendous task facing information consumers today is to identify relevant news items speedily. Hence, designing an information filter for users of Internet news bulletins, is a dire need of the day. The main thrust of this study is, therefore, focused on designing an algorithm to identify news items available on the Internet and categorizing them according to their degree of similarity to each other.

The main concept exploited in obtaining a metric for computing the degree of similarity of two news items is based on calculating and comparing the percentage of proper nouns common to both news items. In order to extract proper nouns from a news item, a filtering process is employed to eliminate pronouns, articles, prepositions, "Be verbs",



determiners and other function words. Subsequently the frequencies in which the extracted words have occurred in the news item are calculated and analyzed. Statistical methods are used to confirm the above results before they are presented to the user. The proposed algorithm was tested and favourable results were obtained by using the news items downloaded from the LAcNet Sri lankan news archives available on the Internet.

Values of degree of similarity obtained for the test data was compared with human classification. Based on these results, it is demonstrated that the proposed algorithm is able to categorize news items into two classes, 'similar' and 'different', successfully. This achievement makes a significant contribution towards achieving the automatic categorization of news items available on the Internet into various topics.

CONTENTS

CHAPTER 1	- INTRODUCTION	
1.1	STATEMENT OF THE PROBLEM	1
1.2	PHYSICAL MOTIVATION	
1.3	OBJECTIVE OF THE STUDY	3
1.4	OVERVIEW OF DISSERTATION	
CHAPTER 2	- THEORETICAL BACKGROUND	6
2.1	INTRODUCTION	6
2.2	SELECTING A STORAGE STRUCTURE	6
2.3	SORTING ALGORITHM	
2.3.1	INTERNAL SORTING	9
2.4	SEARCHING METHODOLOGY	15
2.5	THE HASH TABLE DATA STRUCTURES	17
2.5.1	OPEN HASHING	19
2.5.2	CLOSED HASHING	20
2.6	SOUNDEX CODE ALGORITHM	24
2.7	STATISTICAL METHODS APPLIED	27
CHAPTER 3	METHODOLOGY	30
3.1	SCAN THE TWO NEWS ITEMS AND EXTRACT ALL THE WORDS STARTING WITH CAPITAL LETTERS.	32
3.2	PREPARE A HASH TABLE FOR SUBJECT PRONOUNS, ARTICLES, BE VERBS . . ETC	34
3.3	CALCULATE THE FREQUENCIES IN WHICH THESE EXTRACTED WORDS HAVE OCCURRED IN THE NEWS ITEM	37
3.4	REMOVE EMPTY CELLS IN THE HASH TABLE	41
3.5	GROUP ALL THE SIMILAR WORDS AND ADJUST THEIR FREQUENCIES IN THE HASH TABLE	42



3.6	USE MATCHING ALGORITHM TO OBTAIN THE DEGREE OF SIMILARITY	52
3.7	CONFIRM THE RESULTS OBTAINED IN SEC 3.6 BY USING STATISTICAL THEORIES	54
3.8	MAIN ALGORITHM (STRUCTURE CHART)	55
3.9	USER INTERFACE	56
CHAPTER 4	RESULTS	58
4.1	RESULTS	58
4.1.1	DESCRIPTION OF TEST(INPUT) DATA	58
4.1.2	RESULTS	
4.1.2.1	OUTPUT OF MATCH_NEWS ALGORITHM	
4.1.2.2	OUTPUT OF KOLMÖGOROV TEST	
4.2	DISCUSSION	64
CHAPTER 5	CONCLUSIONS	67
5.1	CONCLUSION	67
5.2	LIMITATIONS AND FUTURE DEVELOPMENTS	
5.1.1	LIMITATIONS	70
5.1.2	FUTURE WORK	72
	BIBLIOGRAPHY	75
	APPENDIX	
A	- TEST DATA (INPUT DATA)	
B	- LIST OF WORDS EXTRACTED FROM NEWS ITEMS "TEST7.TXT" AND "TEST6.TXT"	
C	- CONTENTS IN THE TEXT FILE - "DISCA_LI.TXT"	
D	- HASH TABLE OUTPUTS	
E	- SOUNDEX CODES OF WORDS IN HASH TABLE	
F	- OUTPUT OF MATCH_WORDS ALGORITHM	
G	- PROGRAM LISTING	

