

**ISLAMIC WEB PAGES FILTERING AND CATEGORIZATION**

**NURFAZRINA BINTI MOHD ZAMRY**

**UNIVERSITI TEKNOLOGI MALAYSIA**

# ISLAMIC WEB PAGES FILTERING AND CATEGORIZATION

NURFAZRINA BINTI MOHD ZAMRY

A project report submitted in partial fulfillment of the  
requirements for the award of the degree of  
Master of Computer Science (Information Security)

Faculty of Computing  
Universiti Teknologi Malaysia

JUNE 2013

*“To my Beloved Big Happy Family and Wonderful Best- Friends”*

## ACKNOWLEDGMENT

*Bismillahiramanirrahim.* First and foremost, praise to Allah s.w.t to give me strength and His blessing to complete this Project. Special thanks to my supportive supervisor, Dr. Anazida Zainal for her guidance, advice and knowledge. Without her encouragement and patience this research would not have been the same as presented here. Not to forget, thanks to all the lecturers for the advice.

Thank you very much Malaysia Government for sponsoring my Master study. A bunch of appreciation to the Postgraduate staff for their helps. A lot of loves to my bestfriends, who are always there for me when I'm down. Thanks to my fellow classmate for the support and information sharing.

Last but not least, a deepest love to my beloved family for the courage and endless prayers for me. May Allah bless you.

## ABSTRACT

The Internet creates the world without boundaries where people can get lots of information just by surfing the Internet. But still some of the information is not genuine and correct. Because of that, some of the practitioners of deviant teachings can take this opportunity to attract followers just using the Internet especially to distort beliefs of Muslim in Malaysia. Web filtering can be used as protection against inappropriate and prevention of misuse of the network, hence, it can be used to filter the content of suspicious websites and alleviate the dissemination of such website. Currently, process for blocking the deviate teaching website is done manually and in addition there are limited web filtering product offered to filter religion content and very limited for Malay language. This project is aim to classify deviant teachings Website into three categories which is deviate, suspicious and clean. Pre-processing, feature selection and classification are process involved in Web filtering process. In pre-processing three processes are involved: HTML parsing, stemming and stopping to produce the deviant teaching keyword. Three existing term weighting scheme namely TF, TFIDF and Modified Entropy are used as feature selection process in filtering deviant teaching website while Support Vector Machine (SVM) will be used for classification process. Classification is validated by accuracy, precision, recall and F1. 300 Web pages were collected from Internet based on three categories: deviant teaching, suspicious and clean Web pages. As a result, M.Entropy shows the most suitable term weighting scheme to use in Islamic web pages filtering rather than TFIDF and Entropy.

## ABSTRAK

“Internet mewujudkan dunia tanpa sempadan di mana orang ramai boleh mendapatkan banyak maklumat hanya dengan melayari Internet. Namun, beberapa infomasi dalam talian ini belum pasti yang asli dan benar.. Oleh itu, beberapa pengamal ajaran sesat mengambil peluang ini untuk menarik pengikut hanya menggunakan Internet terutamanya untuk memutarbelitkan kepercayaan Islam di Malaysia. Oleh kerana Penapisan Web digunakan untuk perlindungan terhadap laman Web yang tidak berpatutan dan pencegahan penyalahgunaan rangkaian, maka, sistem penapisan Web boleh digunakan untuk menapis kandungan laman-laman web yang mencurigakan dan mengurangkan penyebaran laman web itu. Pada masa ini, proses untuk menapis laman Web ajaran sesat dilakukan secara manual serta produk penapisan web begitu terhad untuk menapis kandungan agama dan sangat terhad dalam Bahasa Melayu. Projek ini bertujuan untuk mengelaskan laman Web ajaran sesat kepada tiga kategori iaitu sesat, mencurigakan dan bersih. Pra-pemprosesan, *feature selection* dan klasifikasi adalah proses yang terlibat dalam proses penapisan Web. Dalam pra-pemprosesan tiga proses terlibat: *HTML parsing*, *stemming* dan *stopping* bertujuan untuk menghasilkan kata kunci ajaran sesat. Tiga *term weighting scheme* yang sedia ada iaitu TFIDF, *Entropy* dan *Modified Entropy* digunakan sebagai *feature selection* dalam penapisan laman Web ajaran sesat manakala *Support Vector Machine* (SVM) digunakan untuk proses pengelasan. Proses pengelasan disahkan oleh ketepatan (*Accuracy*), ketepatan (*Precision*), ingat (*Recall*) dan F1. 300 laman web telah dikumpulkan dari Internet berdasarkan tiga kategori: laman web ajaran sesat, mencurigakan dan bersih. Hasilnya, *M.Entropy* menunjukkan skim pemberat istilah yang paling sesuai untuk digunakan dalam laman web Islam penapisan berbanding TFIDF dan *Entropy*.

## TABLE OF CONTENT

CHAPTER	TITLE	PAGE
	<b>DECLARATION</b>	i
	<b>DEDICATION</b>	ii
	<b>ACKNOWLEDGMENT</b>	iii
	<b>ABSTRACT</b>	iv
	<b>ABSTRAK</b>	v
	<b>TABLE OF CONTENTS</b>	vi
	<b>LIST OF TABLES</b>	ix
	<b>LIST OF FIGURES</b>	xi
	<b>LIST OF ABBREVIATION</b>	xiii
<b>1</b>	<b>INTRODUCTION</b>	
	1.1 Introduction	1
	1.2 Problem Background	2
	1.3 Problem Statement	4
	1.4 Project Aim	4
	1.5 Project Objectives	5
	1.6 Project Scope	5
	1.7 Significant of Project	5
	1.8 Project Organization	6
<b>2</b>	<b>LITERATURE REVIEW</b>	
	2.1 Introduction	7
	2.2 Internet	7
	2.2.1 Web Content Issues	8
	2.2.1.1 Deviant Teaching over the Internet	9
	2.2.1.2 Role of Malaysian Department of Islamic Development (JAKIM)	11
	2.3 Overview of Web Filtering System	13

2.4	Web Filtering Approaches	14
2.4.1	Web Filtering Technology Trends	16
2.5	Web Content Filtering with Text Classification	17
2.5.1	Data Collection	19
2.5.2	Pre-Processing	19
2.5.3	Text Representation	21
2.5.4	Feature Selection Method in Web Filtering	22
2.5.4.1	Term Weighting Schemes	22
2.5.4.2	Term Frequency Inverse Document Frequency (TFIDF)	23
2.5.5.3	Entropy	24
2.5.4.2	Modified Entropy	25
2.5.5	Classification Method	26
2.5.5.1	Support Vector Machine (SVM)	26
2.6	Waikato Environment for Knowledge Analysis (WEKA)	28
2.6.1	LibSVM	29
2.7	Summary	30
<b>3</b>	<b>METHODOLOGY</b>	
3.1	Introduction	31
3.2	An Overview of Research Framework	31
3.2.1	Phase 1: Term Identification	34
3.2.2	Phase 2: Features Selection & Classification	37
3.2.3	Phase 3: Test & Validate	37
3.3	Data Set	38
3.4	Evaluation	40
3.5	Summary	41
<b>4</b>	<b>IMPLEMENTATION OF PRE-PROCESSING</b>	
4.1	Introduction	42
4.2	Data Collection	42
4.3	Pre-processing	47
4.3.1	HTML Parsing	47
4.3.2	Stemming	51
4.3.3	Stopping	53
4.4	Summary	55



<b>5</b>	<b>TERM WEIGHTING SCHEMES IN FEATURE SELECTION PROCESS AND CLASSIFICATION IN SVM</b>	
5.1	Introduction	56
5.2	Experiment Setup of Term Weighting Scheme in Feature Selection Process and Classification in SVM	57
5.2.1	Term Feature Ranking	59
5.2.2	Term Weighting Scheme	62
5.2.2.1	Term Frequency Inverse Document Frequency (TFIDF)	62
5.2.2.2	Entropy	63
5.2.2.3	Modified Entropy (M.Entropy)	65
5.3	Implementation of SVM as Classifier	66
5.4	Experiment Result	71
5.4.1	Experiment on Data Set 1	71
5.4.2	Experiment on Data Set 2	75
5.4.3	Experiment on Data Set 3	78
5.5	Discussion and Analysis of Experimental Results	82
5.6	Summary	84
<b>6</b>	<b>CONCLUSION AND RECOMMENDATION</b>	
6.1	Introduction	85
6.2	Research Findings	85
6.3	Research Contributions	86
6.4	Future Works	86
6.5	Conclusion	88
	<b>REFERENCES</b>	89 - 91
	<b>APPENDICES A-D</b>	92 - 120

## LIST OF TABLE

TABLE NO.	TITLE	PAGE
2.1	Overview of the Web Filtering Approaches	15
2.2	The Overview and Briefly Description of the Web Filtering Trends	16
2.3	Differentiate Rules of Sembok's Stemming Algorithm and Enhance Sembok's Stemming Algorithm.	21
3.1	Overall Research Plan	33
3.2	Summarize of the Data Set	38
3.3	Data Set for Experiment Process	39
4.1	The Sample of Raw VSM	55
5.1	Definition of VecA – VecG	58
5.2	Top Terms with Feature Ranking Based On TF, TFIDF and M.Entropy Term Weighting Schemes	61
5.3	Top Terms with Feature Ranking Based On TF, TFIDF and M.Entropy Term Weighting Schemes after Expert Intervention	61
5.4	Calculation Step for TFIDF	63
5.5	Calculation Step for Entropy	64
5.6	Calculation Step for M.Entropy	66
5.7	Purpose of data set used in the study	71
5.8	Accuracy of Classification using Data Set 1	72
5.9	Precision of Classification using Data Set 1	72
5.10	Recall of Classification using Data Set 1	72
5.11	F1 of Classification using Data Set 1	73
5.12	Accuracy of Classification using Data Set 2	75
5.13	Precision of Classification using Data Set 2	76
5.14	Recall of Classification using Data Set 2	76
5.15	F1 of Classification using Data Set 2	76
5.16	Accuracy of Classification using Data Set 3	79

5.17	Precision of Classification using Data Set 3	79
5.18	Recall of Classification using Data Set 3	79
5.19	F1 of Classification using Data Set 3	80
5.20	Summary of Classification results using the TFIDF, Entropy and M.Entropy	82

## LIST OF FIGURES

FIGURE NO.	TITLE	PAGE
2.1	Manual Process of Website Blocking for Deviant Teaching	12
2.2	Basic Technique of Web Filtering	14
2.3	Architecture of Intelligent Content Analysis (Lee et al., 2010)	18
2.4	Pre-processing Process in Web Filtering	19
2.5	Hyper Plane of SVM (Chen & Hsieh, 2006)	27
2.6	WEKA Tool Interface	29
3.1	Research Framework	32
3.2	Summarize of the Data Set	35
3.3	Expected Result after Pre-Processing	36
4.1	Malay Dictionary	44
4.2	Added Dictionary	44
4.3	Flowchart on using Malay and Added Dictionary in Stemming Process	45
4.4	Stop List	46
4.5	Original Website before HTML Parsing Process	48
4.6	Flowchart for HTML Parsing Process	49
4.7	Extracted Text Document after HTML Parsing Process	50
4.8	Flowchart the Stemming Process	52
4.9	Flowchart for Stopping Process	53
4.10	Final Output after Pre-Processing Process	54
5.1	Detailed of term weighting methodology as feature selection method	57
5.2	Detailed of Term Feature Ranking	59
5.3	The Algorithm for TFIDF coding development	62
5.4	The Algorithm for Entropy code development	64
5.5	The Algorithm for M.Entropy code development	65
5.6	Steps Taken To Execute the Data Files in WEKA Tool	67

5.7	Dataset in .arff format	68
5.8	Overview of Classification Process using WEKA Based on Libsvm	69
5.9	GUI interface for Classification using WEKA	70
5.10	Result Pattern of Different Term Weighting Scheme Based on (a) Precision, (b) Recall, (c) F1 and (d) Accuracy for Data Set 1	74
5.11	Result Pattern of Different Term Weighting Scheme Based on (a) Precision, (b) Recall, (c) F1 and (d) Accuracy for Data Set 2	78
5.12	Result Pattern of Different Term Weighting Scheme Based on (a) Precision, (b) Recall, (c) F1 and (d) Accuracy For Data Set 3	81
5.13	The Summarization of Classification Result for data set 1 to 3	83

## LIST OF ABBREVIATION

<b>ARPANET</b>	Advanced Research Projects Agency Network
<b>CPBF</b>	Class Profile Based Feature
<b>CUDA</b>	Compute Unified Device Architecture
<b>GA</b>	Genetic Algorithm
<b>GNU</b>	GNU's Not Unix
<b>HTML</b>	Hyper Text Markup Language
<b>IE</b>	Internet Explorer
<b>IR</b>	Information Retrieval
<b>ISP</b>	Internet Service Provider
<b>JAKIM</b>	<i>“Jabatan Kemajuan Islam Malaysia”</i> (Department of Islamic Development Malaysia)
<b>JAPAS</b>	<i>“Jawatankuasa Menangani Ajaran Sesat Peringkat Kebangsaan”</i>
<b>JAPPIS</b>	<i>“Jawatankuasa Penyelaras Penyelidikan Islam Peringkat Kebangsaan”</i>
<b>LISP</b>	List Processing
<b>MoHA</b>	Ministry of Home Affair
<b>PCP</b>	Parallel Coordinate Plot
<b>PHP</b>	PHP: Hypertext Preprocessor
<b>PICS</b>	Platform for Internet Content Selection
<b>PSO</b>	Particle Swarm Optimization
<b>SQL</b>	Structured Query Language
<b>SVM</b>	Support Vector Machine
<b>TFIDF</b>	Term Frequency - Inverse Document Frequency
<b>URL</b>	Uniform Resource Locator
<b>VSM</b>	Vector Space Model
<b>WEKA</b>	Waikato Environment for Knowledge Analysis
<b>WWW</b>	World Wide Web

## CHAPTER 1

### INTRODUCTION

#### 1.1 Introduction

Since the Internet was introduced, it brings a very high impact to people around the world. The Internet creates the world without boundaries where people can get much information just by surfing the Internet. Internet comes from the term Internet working, connecting computer network using specific gateway or router. In the beginning, the ideas of the Internet were from the research of development of ARPANET by United State department of Defence. Nowadays, Internet became faster compare to the last twenty years which offers thousands of information whether it was true or false. Through the Internet, people can interact among each other by search for information, sending e-mail, perform electronic transaction and even social networking.

One of the biggest challenge people face from the Internet is the accuracy of content in the Internet. Not all of this online content is accurate, pleasant, or inoffensive (Heins *et al.*, 2006). As the Internet can be accessed from people around the world in faster manner, some of the practitioners of deviant teachings can take this opportunity to attract followers just using the Internet especially to distort beliefs of Muslim in Malaysia. There are some terminologies which this group use to disseminate their doctrine where it is obviously deviate from the Islamic teaching. Some Muslim which does not have a solid foundation of Islam may simply be attracted to this group. Web filtering has provided two major services: protection

against inappropriate and prevention of misuse of the network (Du *et al.*, 2002). Hence, the web filtering system can be used to filter the content of suspicious websites and alleviate the dissemination of such website. Moreover some of the content may threaten the society as there will reflecting in violence, hate as well as undermine national harmony.

## 1.2 Problem Background

Web filtering is commonly used by organizations such as offices and schools to prevent computer users from viewing inappropriate web sites or content, or as a pre-emptive security measure to prevent access of known malware hosts. While web filtering software, a term for software designed and optimized for controlling what content is permitted to a reader, especially when it is used to restrict material delivered over the Web. Early filtering was based on either “self-rating” by online publishers or “third-party rating” by filter manufacturers (Heins *et al.*, 2006).

Filtering can be divided into five types which include Browser based filters, Client-side filters, Content-limited ISPs, Network-based filtering and Search-engine filters. Each of this filtering has their own characteristic which can be used by different level of users. Browser based filters usually are implemented by third party browser extension while Content-limited ISPs mostly used by government regulator and parental control as it have offered by ISP went subscribe their services. In Client-site filters, specific software are installed in the client computer which administrator can managed the software itself. Network-based filtering is enforced at the transport layer or application layer where the entire user within the network is bound to the organization policy. Lastly, Search-engine filters may be installed in the search engine like Google and Yahoo where it can be activated to block the certain website when safety filter are activated.



There are four types of web filtering approaches mentioned by Lee *et al.*, (2006). There is Platform for Internet Content Selection (PICS), Uniform Resource Locator (URL) blocking, keyword filtering and intelligent content analysis. From those four approaches intelligent content analysis is the best to classify the web content since it can be used to categorize Web pages into different groups for instance the pornographic and non-pornographic. Moreover keyword filtering sometimes will over-block the Web while URL blocking keeps maintaining the reference list.

The self-regulating nature of the Web community, coupled with the easy making information available on the Web has led some individuals to abuse their freedom of expression by putting up harmful materials on the Web (Lee *et al.*, 2003) and this includes practitioners of deviant teachings to disseminate their doctrine. In Malaysia, Malaysian Department of Islamic Development (JAKIM) is responsible to monitor all Islamic oriented publications in various fields where their objective are to develop the Islamic community mind and spreading Islamic *da'wah* through Islamic oriented printed publications in various fields.

Currently, process for blocking the deviate teaching website is done manually. The processes start when JAKIM receive a report from people either by telephone or via email. Special unit in JAKIM then investigate the reported website to filter the website whether it may contain the deviant teaching issues where the investigation normally last for some months. Since the process will take more time, hence by the mean time some of Internet user may accept doctrine. One of the most popular groups of deviant teachings is spiritual and mystical theosophy or in Malay it was called "*Batiniah dan Tasawuf Teosofi*".

There are many classification method used to classify the Website into respective group. For instance, Naïve Bayes, Decision Tree, Artificial Neurel Network as well as Support Vector Machine are methods used for Website classification. In this project, Support Vector Machine also known as SVM is chosen as classification method to classify the Website into specific group. SVM is the

supervised classification method composed of training data and testing data which used for learning and classifying process respectively. Moreover, SVM is suitable to use in text classification and more effective than other methods.

### **1.3 Problem Statement**

The current passive filtering approach used to block the deviant teachings Website impact on the effectiveness for blocking the website. However, in Malaysia there is no web filtering product offered to filter this kind of religion content especially in Malay language. Furthermore, blocking the Web sometime may cause of misunderstanding when there are some of the Web are intent to educate the readers. The effect of this uncontrolled disseminate of deviant teachings through the website may also result of distort beliefs to Muslim society especially in Malaysia.

### **1.4 Project Aim**

The aim of this project is to classify the deviant teachings Website into three categories which is deviate, suspicious and clean. Nevertheless, by implementing this Malay language Web filtering system may help specific organizations to control the disseminated of deviant teachings over the Internet. In the meantime this project can help to minimize spreading of the deviant teachings through the website and hence improve the distort beliefs to Muslim society.

## **1.5 Project Objective**

The main objectives of the project have been identified in order to complete this research, there are;

- i. To pre-process and identify the deviant teaching keywords from the Web pages.
- ii. To implement the features selection method for deviant teaching terms and classifies the term using SVM.
- iii. To test and validate the accuracy of the classification of deviant teaching.

## **1.6 Project Scope**

In accomplishing this project, a number of scopes have been determined, which are:

- i. The samples of Web pages will be extracted from the Internet.
- ii. Classification will be implemented into three categories which include of deviate, suspicious and clean.
- iii. The Support Vector Machine (SVM) algorithm will be used to classify Web content.
- iv. The Web pages used during the project is only in Malay language.

## **1.7 Significance of Project**

This project is prepared to accomplish the objectives to filter and categorize the Islamic Web specific on deviant teachings pages by implementing the Support Vector Machine algorithm. Thus, this project is important to categorize the level of

deviant teachings to aware the readers especially Muslim since the content may lead to distorted of their beliefs. Hence it can minimize the spread of the deviant teaching Furthermore the study can help organization like JAKIM, school, university community to filter their web content within the networks.

## **1.8 Project Organization**

This chapter has covered about the introduction of the project including the problem background, problem statement, project aims, project objectives, project scopes, and significance of the project. This project will be continue with Chapter 2 that discuss more on literature review of the project collected from some resources such as Internet, books, journals and collecting primary data using different research methods. Chapter 3 will covered the project methodology where it highlights the methodology used in this study and discuss about the data used during this study. Chapter 4, the implementation of pre-processing which highlight the second objective is discussed. Next, Term Weighting Schemes in feature selection process and classification of Web pages in SVM will be discussed in detail. Lastly, chapter 6 will concludes the overall of this project.

## REFERENCES

- Ahmad, F., Yusoff, M. and Sembok T. M. T. 1996. *Experiments with a Stemming Algorithm for Malay Words*, Journal of the American Society for Information Science. 47(12). pp. 909-918.
- Barry, M. L., Vinton, G. C. and David, D. C. 1997. *The Past and Future History of the Internet*, Communications Of The ACM, Vol. 40, No. 2, pp. 102-108.
- Barry, M. L., Vinton, G. C., David, D. C. 2009, *A Brief History of the Internet*, ACM SIGCOMM Computer Communication Review. Vol.39, pp.22-31.
- Blum, A. and Langley, P.1997. *Selection of Relevant Features and Examples in Machine Learning*. Artificial Intelligence. 97(1-2). pp. 245–271.
- Chang, C. and Lin, C.J. 2012. *LIBSVM: A Library for Support Vector Machines*. Department of Computer Science National Taiwan University, Taipei, Taiwan
- Chang, C. and Lin, C.J. 2012. *LIBSVM -- A Library for Support Vector Machines Website* [Online]. Available: <http://www.csie.ntu.edu.tw/~cjlin/libsvm/> [Accessed December 2012].
- Chen, R. C. and Hsieh, C. H. 2006. *Web page classification based on a support vector machine using a weighted vote schema*, Systems with Applications. Vol. 31. pp. 427–435
- Cummins, R., and O’Riordan, R. 2007. *Evolved Term-Weighting Schemes in Information Retrieval: An Analysis of the Solution Space*, Dept. of Information Technology, National University of Ireland, Galway, Ireland
- Du, R., Safavi-Naini, R. and Susilon W. 2003. Web filtering using text classification, The 11th IEEE International Conference on Networks. pp.325-330.
- Gupta, S., Kaiser, G., Neistadt, D. and Grimm P. 2003. *DOM-based content extraction of HTML document*, In Proceedings of the 12nd international conference in World Wide Web, ACM New York, NY, USA, 207-214

- Guyon, I. and Elisseeff, A. 2003. *An Introduction to Variable and Feature Selection*, Journal of Machine Learning Research. Vol. 3. pp.1157-1182
- Finkelstine, S. 2002. *BESS's Secret LOOPHOLE censorware vs.privacy & anonymity*. [Online]. Available: <http://sethf.com/anticensorware/bess/loophole.php>. [Accessed 18/10/12].
- Heins, M., Cho, C. and Feldman, A. 2006. *Internet Filters, A Public Policy Report*. Brennan Center For Justice, NYU School Of Law, New York, United State of America.
- Ismail, S.Z. 2010, "Menangani Ajaran Sesat Di Kalangan Umat Islam: Perspektif Undang-Undang Dan Pentadbiran". Shariah Journal, Vol. 18, No. 2 (2010) pp. 247-276.
- Kohavi, R. and John G.H. 1997. *Wrappers for feature subset selection*. Artificial Intelligence 97. pp. 273-324.
- Lee, P.Y., Hui, S.C. and Fong A.C.M 2002. *Neural Networks for Web Content Filtering*, IEEE Intelligent Systems. pp. 48 - 57
- Lee, P.Y., Hui, S.C., Fong, A.C.M. 2003. *A structural and content-based analysis for Web filtering*, Internet Research. Vol. 13 Issue: 1 pp. 27 - 37
- Lee, Z. S., Maarof, M. A., Selamat, A., and Shamsuddin S. M. 2008. *Enhance TermWeighting Algorithm as Feature Selection Technique for Illicit Web Content Classification*, Eighth International Conference on Intelligent Systems Design and Applications ISDA. Nov. 26-28, Kaohsiang City, Taiwan, pp. 145-150
- Lee, Z. S. 2010. *Enhanced Feature Selection Method For Illicit Web Content Filtering*. Ph.D Thesis, Universiti Teknologi Malaysia
- Liu, Y., Loh, H. and Sun, A. 2009. *Imbalanced text classification: A term weighting approach*. Expert Systems with Applications. 361, 670-701.
- Lin, P. C., Liu, M.D., Lin, Y.D., And Lai, Y. C. 2008. *Accelerating Web Content Filtering by the Early Decision Algorithm*. IEICE TRANS. INF. & SYST. Vol.E91-D, No.2, pp-251-257.
- Mazlam, N.F. 2012. *Enhancement of Stemming Process For Malay Illicit Web Content*. Master Thesis, Universiti Teknologi Malaysia.
- Othman, A. 1993. *Pengantar perkataan Melayu untuk sistem capaian dokumen Introduction to Melayu Words for Document Retrieval System*. M.S. thesis, University Kebangsaan Malaysia.

- Lee, P.Y., Lee, S.C. and Fong, A.C.M. 2003. *A Structural And Content-Based Analysis For Web Filtering*, Internet Research. Vol. 13 Issue: 1, pp. 27 – 37,
- Sebastian, F. 2002. Machine learning in automated text categorization. ACM Computing Surveys, Vol 34, 1–47.
- Raghavan, V. and Wong, S. 1986. *A critical analysis of Vector Space Model for Information Retrieval*. Journal of the American Society for Information Science. 35(5)
- Salleh, S.F. 2012. *Comparative Study On Term Weighting Schemes As Feature Selection Method For Malay Illicit Web Content Filtering Content*. Master Thesis, Universiti Teknologi Malaysia.
- Salton, G., Wong, A. and Yang, G. 1975. *A Vector Space Model for automatic Indexing*
- Selamat, A. and Omatu, S. 2004. *Web Page Feature Selection and Classification Using Neural Networks*, Information Sciences. 158 , pp. 69–88
- Sembok, Tengku M.T., Uman, I.I.H. and Rocessing, I.N. 2005. *Word Stemming Algorithms and Retrieval Effectiveness in Malay and Arabic Documents Retrieval Systems*, Conference on Research and Development in Information Retrieval. pp.95–97.
- Verikas, A. and Bacauskiene M. 2002. *Feature Selection With Neural Networks*, Pattern Recognition Letters , 23,1323–1335
- Wikipedia, “*Contentcontrolsoftware*”,[http://en.wikipedia.org/wiki/Contentcontrol\\_software](http://en.wikipedia.org/wiki/Contentcontrol_software) visited 7/8/12.
- Wikipedia, “*Web Crawler*” [http://en..org/wiki/Web\\_crawler](http://en..org/wiki/Web_crawler) visited 18/10/12.
- Wikipedia, “*tf-idf*” [http:// http://en..org/wiki/Tf%E2%80%93idf](http://http://en..org/wiki/Tf%E2%80%93idf) visited 18/10/12.
- Yang, Y. and Liu, X. 1999. *A Re-Examination Of Text Categorization Methods*, Proceedings Of The 22nd International ACM SIGIR Conference On Research And Development In Information Retrieval. pp. 42–49.