

THE PERFORMANCE OF SOFT COMPUTING TECHNIQUES  
ON CONTENT-BASED SMS SPAM FILTERING

WADDAH WAHEEB HASSAN SAEED

A thesis submitted in partial  
fulfillment of the requirement for the award of the  
Degree of Master of Computer Science (Soft Computing)

Faculty of Computer Science and Information Technology  
Universiti Tun Hussein Onn Malaysia

FEBRUARY 2015

## ACKNOWLEDGEMENT

In the name of Allah, the Most Gracious and Most Merciful, I would like to thank Allah S.W.T. for helping me in the compilation until the completion of this thesis. I would like to express my sincere appreciation to my supervisor, Associate Professor Dr. Rozaida binti Ghazali for the constant support throughout the duration of my studies at FSKTM-UTHM. I am fortunate to have her as my supervisor and I will always be grateful for her guidance and encouragement.

I would also like to thanks to my father, mother and beloved wife and daughter for their infinite patience, love, motivation and prayers. Finally, my appreciation goes to those who have contributed directly or indirectly towards the compilation of this thesis.



PTTA UTHM  
PERPUSTAKAAN TUNKU TUKAMINAH

## ABSTRACT

Content-based filtering is one of the most widely used methods to combat SMS (Short Message Service) spam. This method represents SMS text messages by a set of selected features which are extracted from data sets. Most of the available data sets have imbalanced class distribution problem. However, not much attention has been paid to handle this problem which affect the characteristics and size of selected features and cause undesired performance. Soft computing approaches have been applied successfully in content-based spam filtering. In order to enhance soft computing performance, suitable feature subset should be selected. Therefore, this research investigates how well suited three soft computing techniques: Fuzzy Similarity, Artificial Neural Network and Support Vector Machines (SVM) are for content-based SMS spam filtering using an appropriate size of features which are selected by the Gini Index metric as it has the ability to extract suitable features from imbalanced data sets. The data sets used in this research were taken from three sources: UCI repository, Dublin Institute of Technology (DIT) and British English SMS. The performance of each of the technique was compared in terms of True Positive Rate against False Positive Rate,  $F_1$  score and Matthews Correlation Coefficient. The results showed that SVM with 150 features outperformed the other techniques in all the comparison measures. The average time needed to classify an SMS text message is a fraction of a millisecond. Another test using NUS SMS corpus was conducted in order to validate the SVM classifier with 150 features. The results again proved the efficiency of the SVM classifier with 150 features for SMS spam filtering with an accuracy of about 99.2%.

## ABSTRAK

Penapisan berasaskan kandungan merupakan salah satu kaedah yang paling banyak digunakan untuk mengatasi spam SMS (Short Message Service). Kaedah ini mewakili mesej teks SMS dengan satu set ciri terpilih yang diekstrak daripada set-set data. Kebanyakan daripada set-set data sedia ada mempunyai permasalahan pengagihan kelas yang tidak seimbang. Walau bagaimanapun, tidak banyak perhatian diberi dalam menangani permasalahan ini yang mana ia memberi kesan pada ciri-ciri dan saiz ciri yang dipilih dan menyebabkan prestasi yang tidak diinginkan. Pendekatan pengkomputeran lembut telah digunakan dengan jayanya dalam penapisan spam berasaskan kandungan. Bagi meningkatkan kecekapan pengkomputeran lembut, subset ciri yang bersesuaian perlu dipilih. Oleh itu, kajian ini mengkaji bagaimana tiga teknik pengkomputeran lembut: Fuzzy Similarity, Artificial Neural Network dan Support Vector Machines (SVM) sesuai bagi penapisan spam berasaskan kandungan menggunakan saiz ciri yang bersesuaian yang dipilih menggunakan pengukuran Indeks Gini yang mempunyai kemampuan untuk mengekstrak ciri yang bersesuaian daripada set-set data yang tidak seimbang. Set-set data yang digunakan dalam kajian ini telah diambil dari tiga sumber: repositori UCI, Dublin Institute of Technology (DIT) dan British English SMS. Prestasi teknik-teknik ini telah dibandingkan dari segi True Positive Rate against False Positive Rate,  $F_1$  score dan Matthews Correlation Coefficient. Hasil dapatan menunjukkan bahawa SVM dengan 150 ciri lebih baik daripada kedua-dua teknik bandingan dalam kesemua pengukuran perbandingan. Purata masa yang diperlukan untuk mengelaskan mesej teks SMS adalah pecahan milisaat. Bagi mengesahkan pengelas SVM dengan 150 ciri, pengujian lain menggunakan NUS SMS corpus dijalankan. Hasil dapatan membuktikan bahawa kecekapan pengelas SVM dengan 150 ciri bagi menapis spam SMS dengan ketepatan sekitar 99.2%.

**CONTENTS**

<b>TITLE</b>	<b>i</b>
<b>DECLARATION</b>	<b>ii</b>
<b>ACKNOWLEDGEMENT</b>	<b>iii</b>
<b>ABSTRACT</b>	<b>iv</b>
<b>ABSTRAK</b>	<b>v</b>
<b>CONTENTS</b>	<b>vi</b>
<b>LIST OF TABLES</b>	<b>x</b>
<b>LIST OF FIGURES</b>	<b>xi</b>
<b>LIST OF SYMBOLS AND ABBREVIATIONS</b>	<b>xii</b>
<b>LIST OF APPENDICES</b>	<b>xiii</b>
<b>CHAPTER 1 INTRODUCTION</b>	<b>1</b>
1.1 Overview	1
1.2 Problem Statement	2
1.3 Aim of Research	5
1.4 Objective of Research	5
1.5 Scope of Research	5
1.6 Significance of Research	6
1.7 Research Outline	6

<b>CHAPTER 2 LITERATURE REVIEW</b>	<b>8</b>
2.1 Introduction	8
2.2 Fuzzy Logic	9
2.2.1 Fuzzy Similarity	10
2.2.2 T-norms and S-norms	11
2.3 Artificial Neural Network	13
2.3.1 Neuronal Model	13
2.3.2 Activation Functions	14
2.3.3 Multilayer Perceptron	15
2.3.4 Back-propagation Algorithm	16
2.3.5 Scaled Conjugate Gradient	16
2.4 Support Vector Machine	19
2.4.1 Support Vector Classification	20
2.4.2 Radial Basis Function Kernel	24
2.5 Content-based Filtering for SMS Spam	24
2.6 Chapter Summary	27
<b>CHAPTER 3 METHODOLOGY</b>	<b>28</b>
3.1 Introduction	28
3.2 Data Collection	29
3.3 Data Preprocessing	30
3.3.1 Removal of Duplicated Messages	31
3.3.2 Removal of non-English Messages	32
3.3.3 Lower case Conversion	32
3.3.4 Feature Abstraction Replacement	32
3.3.5 Tokenization	33



PTTA  
PERPUSTAKAAN TUN TUN AMINAH

3.3.6	Stemming	33
3.4	Dimensionality Reduction	33
3.5	Data Partition	35
3.6	Training and Testing	36
3.6.1	SMS Text Message Representation	36
3.6.2	Training Fuzzy Similarity Model	37
3.6.3	Training Artificial Neural Network Model	38
3.6.3.1	Number of Input and Output Units	39
3.6.3.2	Activation Function	39
3.6.3.3	Inputs and Outputs Scaling	39
3.6.3.4	Number of Hidden Layers and Neurons	40
3.6.3.5	Parameters Setting	40
3.6.4	Training Support Vector Machines Model	41
3.6.4.1	SVM Kernel Parameters Searching	41
3.7	Model Selection	42
3.7.1	True Positive Rate against False Positive Rate	42
3.7.2	$F_1$ Score	43
3.7.3	Matthews Correlation Coefficient	43
3.8	Chapter Summary	44
<b>CHAPTER 4 DATA ANALYSIS AND RESULTS</b>		<b>45</b>
4.1	Introduction	45
4.2	Experimental Design	45
4.3	True Positive Rate against False Positive Rate Analysis	46
4.4	$F_1$ Score Analysis	48
4.5	Matthews Correlation Coefficient Analysis	50

4.6	Feature Characteristics Analysis	51
4.7	Classification Time Analysis	52
4.8	Additional Testing Using NUS SMS Corpus	52
4.9	Chapter Summary	53

## **CHAPTER 5 DISCUSSION AND CONCLUSIONS** **54**

5.1	Introduction	54
5.2	Novelty and Research Contribution	55
5.2.1	Selecting Feature Subsets Using Gini Index Metric	55
5.2.2	Applying Soft Computing Techniques for SMS Spam Filtering with Selected Feature Subsets	55
5.2.3	Evaluating Soft Computing Techniques Performance for SMS Spam Filtering Using Suitable Measures	56
5.3	Recommendations for Future Work	56
5.4	Chapter Summary	57

## **REFERENCES** **58**

## **VITAE** **72**





## LIST OF TABLES

2.1	T-norms and s-norms operators	12
3.1	Class information of the collected data sets	30
3.2	Number of instances before and after removal of message duplicates	31
3.3	Number of instances before and after removal of non-English messages	32
3.4	List of replaced features	33
3.5	Two-way contingency table of a feature $f_i$ and a class $C_j$ for binary classification	34
3.6	Number of hidden neurons with its number of features	40
3.7	C and $\gamma$ values with its number of features	41
4.1	Number of instances in training and testing data set	46
4.2	TPR against FPR classifiers comparison with its number of features [AUC]	47
4.3	Number of instances before and after data preprocessing steps in NUS SMS corpus	53
A.1	Hidden neurons searching with its number of features	65
B.1	C and $\gamma$ values searching with its number of features	67
C.1	ANN simulation results with its number of features	69

## LIST OF FIGURES

2.1	An example of ANN with one hidden layer	14
2.2	Neuronal model	14
2.3	Linearly separable training data	20
2.4	Possible separating hyperplanes with their associated margins	21
2.5	Support vectors	22
2.6	Two class nonlinear separable problem	23
3.1	Research framework	29
3.2	Extract taken from the removal of message duplicates	31
3.3	Representation of vector space model	36
4.1	TPR against FPR classifiers comparison with 150 features	48
4.2	True positive rate classifiers comparison with 150 features	49
4.3	False positive rate classifiers comparison with 150 features	49
4.4	$F_1$ classifiers comparison with its number of features	50
4.5	Precision classifiers comparison with 150 features	50
4.6	MCC classifiers comparison with its number of features	51
4.7	Degree of combination of positive and negative features	52

**LIST OF SYMBOLS AND ABBREVIATIONS**

ANN	-	Artificial Neural Network
AUC	-	Area Under Curve
FN	-	False Negative
FP	-	False Positive
FPR	-	False Positive Rate
MCC	-	Matthews Correlation Coefficient
MLP	-	Multilayer Perceptron
MSE	-	Mean Squared Error
RBF	-	Radial Basis Function
ROC	-	Receiver Operating Characteristic
SCG	-	Scaled conjugate gradient
SMS	-	Short Message Service
SVM	-	Support Vector Machine
TN	-	True Negative
TP	-	True Positive
TPR	-	True Positive Rate

**LIST OF APPENDICES**

<b>APPENDIX</b>	<b>TITLE</b>	<b>PAGE</b>
A	Performance of ANN classifier using range of hidden neurons with its number of features	65
B	Performance of SVM classifier using range of C and $\gamma$ values with its number of features	67
C	ANN Simulation results with its number of features	69



PTTA UTHM  
PERPUSTAKAAN TUNKU TUN AMINAH

## CHAPTER 1

### INTRODUCTION

#### 1.1 Overview

SMS which stands for “Short Message Service” is a service used to send short text messages from a mobile device or via the web and received by a mobile device. This service is a very popular type of communication between people, for its ease of use, its fast response and its relatively cheap cost as compared to telephone calls. Thus in 2012, 7.5 trillion SMS messages were sent all over the world (GSMA, 2013). However, not all SMS messages are solicited - mobile device users receive legitimate messages as well as unwanted messages which are called spam.

SMS spam forms 20 to 30% of all SMS traffic in some parts of Asia such as China and India (GSMA, 2011). Some methods are used to combat SMS spam such as black-and-white listing, traffic analysis and content-based filtering (Delany, Buckley & Greene, 2012). According to Delany *et al.* (2012), content-based filtering method is required to counteract the increasing threat of SMS spam and to avoid the disadvantages of other filtering methods. Content-based filtering uses some techniques to analyze the contents of SMS text messages to ascertain whether it is legitimate or spam.

Many studies on content-based SMS spam filtering selected some features (lexical or stylistic) to represent SMS text messages and these selected features are ex-

tracted from SMS data sets with imbalanced class distribution problems. However, not much attention has been paid to handle the imbalanced class distribution problem which could produce unsuitable features or a huge number of features in order to filter SMS spam. Therefore, a suitable feature selection metric is required to select proper features from the imbalanced data sets in order to improve filtering performance. Besides a suitable feature selection metric, a suitable technique which has been engaged in spam filtering is essential. Soft computing techniques have been present in almost every domain (e.g. spam filtering) and their ability have been proven (El-Alfy & Al-Qunaieer, 2008; Guzella & Caminhas, 2009).

In this research, the main purpose is to find out how well suited soft computing techniques, namely Fuzzy Similarity, Artificial Neural Network (ANN) and Support Vector Machines (SVMs) are for content-based SMS spam filtering using appropriate features which are selected by the Gini Index metric.

## 1.2 Problem Statement

SMS spam is a growing problem. Mobile device users in the U.S. received 1.1 billion spam messages in 2007 (Hart, 2008) and 4.5 billion in 2011 (Kharif, 2012). SMS spam can be defined as unsolicited bulk electronic messages. Unsolicited means the recipients receive unwanted messages without their consent and bulk because the sender sends many identical messages to different recipients (Bueti, 2005).

Many reasons motivate spammers to use this service which support the growth of this problem such as the attraction to read all received messages by mobile device users, the accessibility of this service from anywhere, lack of laws and regulations to control the purchase of phone numbers and the handling of this problem in some countries (Liu & Yang, 2012). In addition, there is an increasing number of mobile device users who can be targeted (GSMA, 2013), the limited availability of mobile applications for SMS spam filtering (Almeida, Hidalgo & Yamakami, 2011), the higher response rate for this service and the availability of very cheap bulk pre-paid SMS

packages in some countries in Asia with easy solutions to send bulk messages (Delany *et al.*, 2012) as well as mobile network operators who contribute to this problem by sending messages about their offers.

SMS spam has caused mobile device users and mobile network operators a lot of problems. Spam messages irritate mobile users by filling their in-boxes and wasting their time reading and deleting the messages (Uysal *et al.*, 2012). Some types of SMS spam try to bill mobile device users by tricking them to call premium rate numbers or subscribe to services or, trick the users to call certain numbers to collect confidential information from them to use for other purposes — called phishing (GSMA, 2011). Other types of SMS spam attack mobile device users to steal their money (GSMA, 2011), subject smart-phones to viruses (Murynets & Jover, 2012), harm mobile device operating systems, spread viruses to other mobile device users and violate privacy. Furthermore, in some countries mobile device users pay to receive their messages which may include spam messages (Almeida *et al.*, 2011). Mobile network operators also suffer from this problem. They are prone to lose their subscribers because the performance of the network is affected by the load that SMS spam generates which in turn delay the reception of legitimate messages (Yadav *et al.*, 2011). They may also lose some revenue because they cannot bill the sender(s) a termination fee as some types of SMS spam are sent from fraudulent addresses (Cisco, 2005).

Many methods have been used to prevent SMS spam due to these problems, such as black-and-white listing which is used by mobile applications such as android applications (GooglePlay, n.d.), traffic analysis (GSMA, 2011), content-based filtering (Hidalgo, Bringas & Sanz, 2006; Almeida *et al.*, 2011; Sohn *et al.*, 2012) and a combination of black-and-white listing and content-based filtering (Deng & Peng, 2006; Mahmoud & Mahfouz, 2012). With black-and-white listing, the mobile device user saves the phone numbers of legitimate and spam message senders into two groups: legitimate group (white list) and spam group (black list). The disadvantages of the black-and-white listing method, is that if the phone numbers are not in the black list, the recipient will receive the spam message(s). In addition, this method will dis-

card legitimate messages that may be sent from a black-listed phone number(s) (Uysal *et al.*, 2012). Another anti-spam method uses traffic analysis to compare the subscriber's volume of sent messages to volume limits, but spammers avoid this method by sending low volumes of messages to observe the operator system response and then determine the operator's volume limit policies (Delany *et al.*, 2012). Content-based filtering method uses some techniques to analyze SMS text message content in order to decide whether it is legitimate or spam. The spammer tries to avoid these filters by making sophisticated message modifications (GSMA, 2011), however, content-based filtering still needs to avoid spammers' traffic analysis tricks (Delany *et al.*, 2012) as well as the black-and-white listing.

Many studies in the literature on content-based SMS spam filtering selected some features to represent SMS text messages and these selected features are extracted from SMS data sets with imbalanced class distribution problem. However, not much attention has been paid to handle the imbalanced class distribution problem which affect the characteristics and the size of the selected features and cause undesired performance. Therefore, in order to select suitable features from the imbalanced data sets, a suitable feature selection scheme is needed. The Gini Index (Shang *et al.*, 2007) is a feature selection metric which has the ability to handle class imbalance problem by selecting proper features (Ogura, Amano & Kondo, 2011) which will improve the performance of filtering. Besides a suitable feature selection metric, a suitable technique which has been engaged in spam filtering is required. Soft computing techniques have been present in almost every domain (e.g. spam filtering) and their ability has been proven (El-Alfy & Al-Qunaieer, 2008; Guzella & Caminhas, 2009).

Therefore, this research investigates the performance of three selected soft computing techniques: Fuzzy Similarity, Artificial Neural Network and Support Vector Machines and whether they are suitable for content-based SMS spam filtering using appropriate size of features selected by the Gini Index metric.



### 1.3 Aim of Research

The aim of this research is to filter SMS spam based on its contents using soft computing techniques, namely Fuzzy Similarity, Artificial Neural Network and Support Vector Machine with appropriate features selected by the Gini Index metric.

### 1.4 Objective of Research

In order to achieve the above mentioned aim of the research, the following are three research objectives:

- i To select feature subsets using the Gini Index metric to represent SMS text messages.
- ii To apply soft computing techniques: Fuzzy Similarity, Artificial Neural Network and Support Vector Machine for SMS spam filtering with feature subsets selected in (i).
- iii To compare the performance of (ii) in terms of True Positive Rate (TPR) against False Positive Rate (FPR),  $F_1$  score and Matthews Correlation Coefficient (MCC).

### 1.5 Scope of Research

This research was to filter English SMS text message into two classes either legitimate or spam based on its contents. The data was taken from three sources: UCI machine learning repository (Bache & Lichman, 2013), Dublin Institute of Technology (DIT) (Delany *et al.*, 2012) and British English SMS (Nuruzzaman, Lee & Choi, 2011). Feature subsets were selected using the Gini Index metric (Shang *et al.*, 2007). Three soft computing techniques: Fuzzy Similarity (Widyantoro & Yen, 2000), Artificial Neural Network which trained using Scaled Conjugate Gradient algorithm (SCG) (Møller, 1993) and Support Vector Machine with Radial Basis Function (RBF) kernel (Chang

& Lin, 2011) were used to filter SMS spam. Results were compared in terms of True Positive Rate (TPR) against False Positive Rate (FPR),  $F_1$  score and Matthews Correlation Coefficient (MCC).

## 1.6 Significance of Research

The efficiency of soft computing techniques for SMS spam filtering with feature subsets selected by the Gini Index metric was examined in this research. Therefore, this research was conducted to establish a comparison in performance between Fuzzy Similarity, Artificial Neural Network and Support Vector Machine to investigate whether they can provide better results based on the selected feature subsets. The outcome of this research could contribute to verifying the best performance with small size features for SMS spam filtering and also contribute to future work in exploring the possibility of other feature selection metrics with soft computing techniques in SMS spam filtering.

## 1.7 Research Outline

The remaining part of this research is arranged in the following chapters. Chapter 2 is concerned with the relevant background in using content-based filtering technique for SMS spam filtering. Likewise, the chapter also highlights soft computing techniques, namely Fuzzy Similarity, Artificial Neural Network and Support Vector Machine.

Chapter 3 describes briefly steps on how to use soft computing techniques for SMS spam filtering, starting from data collection, data preprocessing, dimensionality reduction, data partition, training and testing, and selecting the best soft computing technique based on specified measures.

Simulations results with analysis which evaluate the soft computing techniques are presented in Chapter 4. Feature subset characteristics and classification time are also analyzed. The best soft computing technique with the best feature subsets are tested using another SMS corpus. In order to simplify the discussions, graphs that

summarize the results are provided. Chapter 5 concludes the work done and provides several recommendations to improve and validate the performance of the soft computing techniques for SMS spam filtering.



## CHAPTER 2

### LITERATURE REVIEW

#### 2.1 Introduction

Many real world problems cannot be solved using hard computing techniques that deal with precision and certainty due to the fact that either these real-world problems are difficult to model mathematically or computationally expensive or require huge amounts of memory (Shukla, Tiwari & Kala, 2012). However, in some cases, human experts can deal with these problems successfully, e.g. face recognition. According to Zadeh, soft computing is “an emerging approach to computing, which parallels the remarkable ability of the human mind to reason and learn in an environment of uncertainty and imprecision” (Zadeh, 1994). From this definition, it is clear that soft computing is inspired by natural processes — especially the human brain. Therefore, soft computing techniques are needed to offer simple, reliable and low cost solutions to these types of problems with best results.

The development of soft computing techniques has attracted the interest of researchers from different disciplines over the past two decades. Soft computing techniques are applied in various domains such as bioinformatics, biomedical systems, data mining, image processing, machine control, robotics, time series prediction, wireless networks, etc.(Shukla *et al.*, 2012).

Classification problem is one of three main categories of problems for which

soft computing is applied (Shukla *et al.*, 2012). A classification problem relates an object depending on its attributes into a known group or class. If there are many differences among the classes based on their attributes then the classification problem becomes quite simple. However, if the classes are quite similar, it becomes rather difficult. Therefore, soft computing is needed to offer solutions to these problems.

In this research, three soft computing techniques are used, namely Fuzzy Similarity, Artificial Neural Networks and Support Vector Machine to classify SMS text messages into two classes either legitimate or spam. These techniques have been used for email spam filtering (El-Alfy & Al-Qunaieer, 2008; Guzella & Caminhas, 2009). Therefore, in order to be more certain about these techniques, this chapter provides a discussion on them. This chapter also reviews the related works regarding the problem under study; the content-based filtering for SMS spam.

## 2.2 Fuzzy Logic

The concept of fuzzy logic was introduced in 1965 by Zadeh as a new concept to deal with problems in which the imprecision is the absence of precisely defined criteria of class membership (Zadeh, 1965). The acceptance of fuzzy logic started in the second half of the 1970s after the success of the first practical application which is called fuzzy control. Since then, fuzzy logic has been applied in many mathematical and practical areas including clustering, optimization, operations research, control and expert systems, medicine, data mining and pattern recognition (Zimmermann, 2010).

Fuzzy logic deals with fuzzy sets which are an extension of the definition on crisp sets. Unlike the characteristic function for crisp sets, the characteristic function (membership function) of fuzzy sets is represented by a degree of relevance in the range [0,1]. This provides flexibility in dealing with uncertainty in systems such as spam filtering (El-Alfy & Al-Qunaieer, 2008). Fuzzy logic has not received much attention for SMS spam filtering. Fuzzy Similarity (Widyantoro & Yen, 2000) performs well in email spam filtering (El-Alfy & Al-Qunaieer, 2008). Thus, this research investigates

the effectiveness of Fuzzy Similarity in content-based SMS spam filtering.

### 2.2.1 Fuzzy Similarity

Fuzzy similarity is adapted from the Rocchio algorithm (Rocchio, 1971). In this algorithm, a cluster center is created for each category from training samples and the similarity between each test sample and a category is measured using cosine coefficient. In fuzzy similarity which was proposed by (Widyantoro & Yen, 2000), a fuzzy term-category relation is developed, whereby the Rocchio cluster is represented by a set of membership degree of words to a particular category. Based on the fuzzy term-category relation, the similarity between a document and a category's cluster center is calculated using fuzzy conjunction and disjunction operators, and the calculated similarity represents the membership degree of document to the category.

Fuzzy similarity has two finite sets, a set of terms  $T = t_1, t_2, \dots, t_n$  and a set of categories  $C = c_1, c_2, \dots, c_n$ . A fuzzy relation  $R : T \times C \rightarrow [0, 1]$ , whereby the membership value of the relation, which denotes by  $\mu_R(t_i, c_j)$ , specifies the degree of relevance of term  $t_i$  to category  $c_j$ . The membership values of this relation are extracted from a training set.

Every training example in the training set is represented by a set of term-frequency pairs  $d = \{(t_1, o_1), (t_2, o_2), \dots, (t_m, o_m)\}$  where  $o_j$  is the occurrence frequency of term  $t_j$  in the document. Given a set of training documents  $D$ , the membership value of the relation  $R(t_i, c_j)$ , denoted by  $\mu_R(t_i, c_j)$ , is calculated as follows. First, all documents are grouped according to their category. Next, the occurrence frequency of each term for each category is collected by summing up the term frequency of individual documents in that category. Then the value of  $\mu_R(t_i, c_j)$  is calculated from the total number of occurrences of term  $t_i$  in category  $c_j$  divided by the total number of term frequency  $t_j$  in all categories as expressed in Eq. (2.1).

$$\mu_R(t_i, c_j) = \frac{\sum_{\{w_i \in d_k \wedge d_k \in D \wedge c(d_k) = c_j\}} w_i}{\sum_{\{w_i \in d_k \wedge d_k \in D\}} w_i} \quad (2.1)$$

Now, the membership values of fuzzy term-category relation are known, the similarity between a document and the category's membership values of the term is given by Eq. (2.2),

$$Sim(d, c_j) = \frac{\sum_{t \in d} \mu_R(t, c_j) \otimes \mu_d(t)}{\sum_{t \in d} \mu_R(t, c_j) \oplus \mu_d(t)} \quad (2.2)$$

in which  $\mu_d(t)$  is the membership degree that term  $t$  belongs to  $d$  for each term  $t$  in  $d$ ,  $\otimes$  and  $\oplus$  denote fuzzy conjunction (t-norm) and fuzzy disjunction (s-norm) operators, respectively. The category of the document is the category that has the highest similarity measure.

### 2.2.2 T-norms and S-norms

There are various t-norms and s-norms which are frequently used in the literature. In order to define any t-norms and s-norms operations, there are some axioms that should be satisfied. For t-norms operation, any binary operation  $t$  should satisfy the following axioms in order to be a t-norm operation, given  $x, y, z \in [0, 1]$ :

$$\text{Axiom1.} \quad t(x, 1) = x \quad (\text{boundary condition})$$

$$\text{Axiom2.} \quad y \leq z \text{ implies } t(x, y) \leq t(x, z) \quad (\text{monotonicity})$$

$$\text{Axiom3.} \quad t(x, y) = t(y, x) \quad (\text{commutativity})$$

$$\text{Axiom4.} \quad t(x, t(y, z)) = t(t(x, y), z) \quad (\text{associativity})$$

Almost the same axioms are defined for s-norms operation, given  $x, y, z \in [0, 1]$ :

$$\text{Axiom1.} \quad s(x, 0) = x \quad (\text{boundary condition})$$

$$\text{Axiom2.} \quad y \leq z \text{ implies } s(x, y) \leq s(x, z) \quad (\text{monotonicity})$$

$$\text{Axiom3.} \quad s(x, y) = s(y, x) \quad (\text{commutativity})$$

$$\text{Axiom4. } s(x, s(y, z)) = s(s(x, y), z) \quad (\text{associativity})$$

The boundary condition is to range the results to be in  $[0,1]$ . Monotonicity and commutativity are to ensure that a decrease in the degree of membership in set  $X$  or  $Y$  cannot produce an increase in the degree of membership in the intersection or union. Commutativity ensures that the fuzzy intersection and fuzzy union are symmetric therefore there is no consideration for order. The last axiom, associativity, allows taking the intersection of any number of sets in any order of pairwise grouping desired (Klir & Yuan, 1995).

Among the various t-norms and s-norms as shown in Table 2.1, the standard fuzzy intersection and the standard fuzzy union have special features. One of the desirable features is that the standard fuzzy intersection, min operator, and the standard fuzzy union, max operator, prevent the compounding of errors in the operands which is lacking in most alternative norms (Klir & Yuan, 1995). For example, If any error  $e$  is associated with the membership values  $\mu_A(x)$  and  $\mu_B(x)$ , then the maximum error associated with the membership value of  $x$  in  $\mu_{\bar{A}}(x)$ ,  $\mu_{A \cup B}(x)$  and  $\mu_{A \cap B}(x)$  remains  $e$  (Klir & Yuan, 1995). For that, the standard fuzzy intersection, min operator, and the standard fuzzy union, max operator, are selected in this research.

Table 2.1: T-norms and s-norms operators

t-norms $t(x, y)$	s-norms $s(x, y)$
Standard intersection $t(x, y) = \min(x, y)$	Standard union $s(x, y) = \max(x, y)$
Algebraic product $t(x, y) = x \cdot y$	Algebraic sum $s(x, y) = x + y - x \cdot y$
Bounded difference $t(x, y) = \max(0, x + y - 1)$	Bounded sum $s(x, y) = \min(1, x + y)$
Drastic intersection $t(x, y) = \begin{cases} x & \text{when } y = 1 \\ y & \text{when } x = 1 \\ 0 & \text{otherwise.} \end{cases}$	Drastic union $s(x, y) = \begin{cases} x & \text{when } y = 0 \\ y & \text{when } x = 0 \\ 1 & \text{otherwise.} \end{cases}$



## REFERENCES

- Almeida, T., Hidalgo, J.M.G. & Silva, T.P. (2013). Towards SMS spam filtering: Results under a new dataset. *International Journal of Information Security Science*, 2(1), pp. 1–18.
- Almeida, T.A., Hidalgo, J.M.G. & Yamakami, A. (2011). Contributions to the study of SMS spam filtering: new collection and results. In: *Proceedings of the 11th ACM symposium on Document engineering*, ACM, pp. 259–262.
- Bache, K. & Lichman, M. (2013). UCI Machine Learning Repository [<http://archive.ics.uci.edu/ml>]. Irvine, CA: University of California, School of Information and Computer Science.
- Basheer, I. & Hajmeer, M. (2000). Artificial neural networks: fundamentals, computing, design, and application. *Journal of microbiological methods*, 43(1), pp. 3–31.
- Bishop, C.M. (1995). *Neural Networks for Pattern Recognition*. New York, NY, USA: Oxford University Press, Inc., ISBN 0198538642.
- Boser, B.E., Guyon, I.M. & Vapnik, V.N. (1992). A training algorithm for optimal margin classifiers. In: *Proceedings of the Fifth Annual Workshop on Computational Learning Theory*, COLT '92, New York, NY, USA: ACM, ISBN 0-89791-497-X, pp. 144–152.
- Bueti, S. (2005). International Telecommunication Union: WSIS Thematic Meeting on Cybersecurity, ITU Survey on Anti-Spam Legislation Worldwide.

- Chang, C.C. & Lin, C.J. (2011). LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2, pp. 27:1–27:27, software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- Chen, T. & Kan, M.Y. (2013). Creating a live, public short message service corpus: The nus sms corpus. *Language Resources and Evaluation*, 47(2), pp. 299–335.
- Cisco (2005). *SMS SPAM AND FRAUD PREVENTION*.
- Cristianini, N. & Shawe-Taylor, J. (2000). *An introduction to support vector machines and other kernel-based learning methods*. Cambridge university press, ISBN 9780521780193.
- Delany, S.J., Buckley, M. & Greene, D. (2012). SMS spam filtering: methods and data. *Expert Systems with Applications*, 39(10), pp. 9899–9908.
- Demuth, H., Beale, M. & Hagan, M. (2008). Neural network toolbox 6. *User's guide*.
- Deng, W.W. & Peng, H. (2006). Research on a naive bayesian based short message filtering system. In: *Machine learning and cybernetics, 2006 international conference on*, IEEE, pp. 1233–1237.
- El-Alfy, E.S. & Al-Qunaieer, F.S. (2008). A fuzzy similarity approach for automated spam filtering. In: *Computer Systems and Applications, 2008. AICCSA 2008. IEEE/ACS International Conference on*, IEEE, pp. 544–550.
- GooglePlay (n.d.). Retrieved on December 2, 2014 from <https://play.google.com/store>.
- GSMA (2011). *SMS Spam and Mobile Messaging Attacks: Introduction, Trends And Examples*. GSMA: Spam Reporting Service.
- GSMA (2013). *GSMA: The Mobile Economy 2013*.
- Guzella, T.S. & Caminhas, W.M. (2009). A review of machine learning approaches to spam filtering. *Expert Systems with Applications*, 36(7), pp. 10206–10222.

- Han, J., Kamber, M. & Pei, J. (2011). *Data mining: concepts and techniques*. Morgan kaufmann.
- Hart, K. (2008). Advertising Sent To Cellphones Opens New Front In War on Spam. *Washington post*. Retrieved on December 2, 2014 from [http://articles.washingtonpost.com/2008-03-10/news/36788122\\_1\\_text-messages-spam-e-mail-block-texts](http://articles.washingtonpost.com/2008-03-10/news/36788122_1_text-messages-spam-e-mail-block-texts).
- Haykin, S. (2009). *Neural networks and learning machines*, volume 3. Pearson Education Upper Saddle River.
- Hidalgo, J., Bringas, G. & Sáenz, E. (2006). Content based SMS spam filtering. In: *Proceedings of the 2006 ACM symposium on Document engineering*, ACM, pp. 107–114.
- Hsu, C.W., Chang, C.C. & Lin, C.J. (2003). A practical guide to support vector classification.
- Japkowicz, N. & Shah, M. (2011). *Evaluating Learning Algorithms*. Cambridge University Press.
- Joe, I. & Shim, H. (2010). An sms spam filtering system using support vector machine. In: *Future Generation Information Technology*, Springer, pp. 577–584.
- Jurman, G., Riccadonna, S. & Furlanello, C. (2012). A comparison of mcc and cen error measures in multi-class prediction. *PloS one*, 7(8), p. e41882.
- Karami, A. & Zhou, L. (2014). Improving static sms spam detection by using new content-based features.
- Kharif, O. (2012). Mobile Spam Texts Hit 4.5 Billion Raising Consumer Ire. *Bloomberg Business Week News*. Retrieved on December 2, 2014 from <http://www.businessweek.com/news/2012-04-30/mobile-spam-texts-hit-4-dot-5-billion-raising-consumer-ire>.

- Khemapatapan, C. (2010). Thai-english spam sms filtering. In: *Communications (APCC), 2010 16th Asia-Pacific Conference on*, IEEE, pp. 226–230.
- Klir, G.J. & Yuan, B. (1995). *Fuzzy Sets and Fuzzy Logic: Theory and Applications*. Upper Saddle River, NJ, USA: Prentice-Hall, Inc., ISBN 0-13-101171-5.
- Kotsiantis, S., Kanellopoulos, D. & Pintelas, P. (2006). Data preprocessing for supervised learning. *International Journal of Computer Science*, 1(2), pp. 111–117.
- Levenberg, K. (1944). A method for the solution of certain problems in least squares. *Quarterly of applied mathematics*, 2, pp. 164–168.
- Liu, G. & Yang, F. (2012). The application of data mining in the classification of spam messages. In: *Computer Science and Information Processing (CSIP), 2012 International Conference on*, IEEE, pp. 1315–1317.
- Liu, Y., Loh, H.T. & Sun, A. (2009). Imbalanced text classification: A term weighting approach. *Expert systems with Applications*, 36(1), pp. 690–701.
- Mahmoud, T.M. & Mahfouz, A.M. (2012). SMS spam filtering technique based on artificial immune system. *IJCSI International Journal of Computer Science Issues*, 9(1).
- Maier, J. & Ferens, K. (2009). Classification of english phrases and sms text messages using bayes and support vector machine classifiers. In: *Electrical and Computer Engineering, 2009. CCECE'09. Canadian Conference on*, IEEE, pp. 415–418.
- Marquardt, D.W. (1963). An algorithm for least-squares estimation of nonlinear parameters. *Journal of the Society for Industrial & Applied Mathematics*, 11(2), pp. 431–441.
- Mathew, K. & Issac, B. (2011). Intelligent spam classification for mobile text message. In: *Computer Science and Network Technology (ICCSNT), 2011 International Conference on*, volume 1, IEEE, pp. 101–105.

- Møller, M.F. (1993). A scaled conjugate gradient algorithm for fast supervised learning. *Neural networks*, 6(4), pp. 525–533.
- Mujtaba, G. & Yasin, M. (2014). SMS Spam Detection Using Simple Message Content Features.
- Murynets, I. & Jover, R.P. (2012). Crime Scene Investigation : SMS Spam Data Analysis Categories and Subject Descriptors, pp. 441–452.
- Najadat, H., Abdulla, N., Abooraig, R. & Nawasrah, S. (2014). Mobile SMS Spam Filtering based on Mixing Classifiers.
- Nuruzzaman, M.T., Lee, C., bin Abdullah, M.F.A. & Choi, D. (2012). Simple SMS spam filtering on independent mobile phone. *Security and Communication Networks*, 5(10), pp. 1209–1220.
- Nuruzzaman, M.T., Lee, C. & Choi, D. (2011). Independent and personal SMS spam filtering. In: *Computer and Information Technology (CIT), 2011 IEEE 11th International Conference on*, IEEE, pp. 429–435.
- Ogura, H., Amano, H. & Kondo, M. (2011). Comparison of metrics for feature selection in imbalanced text classification. *Expert Systems with Applications*, 38(5), pp. 4978–4989.
- Özkan, C. & Erbek, F.S. (2003). The comparison of activation functions for multi-spectral landsat tm image classification. *Photogrammetric Engineering & Remote Sensing*, 69(11), pp. 1225–1234.
- Porter, M.F. (1980). An algorithm for suffix stripping. *Program: electronic library and information systems*, 14(3), pp. 130–137.
- Rafique, M.Z. & Abulaish, M. (2012). Graph-based learning model for detection of SMS spam on smart phones. In: *Wireless Communications and Mobile Computing Conference (IWCMC), 2012 8th International*, IEEE, pp. 1046–1051.
- Rocchio, J.J. (1971). Relevance feedback in information retrieval.

- Rumelhart, D.E., McClelland, J.L. & PDP Research Group, C. (1986). *Parallel Distributed Processing: Explorations in the Microstructure of Cognition, Vol. 1: Foundations*. Cambridge, MA, USA: MIT Press, ISBN 0-262-68053-X.
- Salton, G., Wong, A. & Yang, C.S. (1975). A vector space model for automatic indexing. *Communications of the ACM*, 18(11), pp. 613–620.
- Samarasinghe, S. (2006). *Neural networks for applied sciences and engineering: from fundamentals to complex pattern recognition*. CRC Press.
- Shahi, T.B. & Yadav, A. (2013). Mobile SMS Spam Filtering for Nepali Text Using Naïve Bayesian and Support Vector Machine. *International Journal of Intelligence Science*, 4, p. 24.
- Shang, W., Huang, H., Zhu, H., Lin, Y., Qu, Y. & Wang, Z. (2007). A novel feature selection algorithm for text categorization. *Expert Systems with Applications*, 33(1), pp. 1–5.
- Shukla, A., Tiwari, R. & Kala, R. (2012). *Real life applications of soft computing*. CRC Press.
- Sohn, D.N., Lee, J.T., Han, K.S. & Rim, H.C. (2012). Content-based mobile spam classification using stylistically motivated features. *Pattern Recognition Letters*, 33(3), pp. 364–369.
- Sohn, D.N., Lee, J.T. & Rim, H.C. (2009). The contribution of stylistic information to content-based mobile spam filtering. In: *Proceedings of the ACL-IJCNLP 2009 Conference Short Papers*, Association for Computational Linguistics, pp. 321–324.
- Tagg, C. (2009). *A corpus linguistics study of SMS text messaging*. Ph.D. thesis, University of Birmingham.
- Uysal, A.K., Gunal, S., Ergin, S. & Gunal, E.S. (2012). A novel framework for SMS spam filtering. In: *Innovations in Intelligent Systems and Applications (INISTA), 2012 International Symposium on*, IEEE, pp. 1–4.

- Vapnik, V.N. & Chervonenkis, A.Y. (1971). On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability & Its Applications*, 16(2), pp. 264–280.
- Widyantoro, D.H. & Yen, J. (2000). A fuzzy similarity approach in text classification task. In: *IEEE International conference on fuzzy systems*, volume 2, pp. 653–658.
- Yadav, K., Kumaraguru, P., Goyal, A., Gupta, A. & Naik, V. (2011). SMSAssassin: Crowdsourcing driven mobile-based system for SMS spam filtering. In: *Proceedings of the 12th Workshop on Mobile Computing Systems and Applications*, ACM, pp. 1–6.
- Zadeh, L.A. (1965). Fuzzy sets. *Information and control*, 8(3), pp. 338–353.
- Zadeh, L.A. (1994). Fuzzy logic, neural networks, and soft computing. *Communications of the ACM*, 37(3), pp. 77–84.
- Zhang, J., Li, X., Xu, W. & Li, C. (2011). Filtering algorithm of spam short messages based on artificial immune system. In: *Electrical and Control Engineering (ICECE), 2011 International Conference on*, IEEE, pp. 195–198.
- Zimmermann, H.J. (2010). Fuzzy set theory. *Wiley Interdisciplinary Reviews: Computational Statistics*, 2(3), pp. 317–332.
- Zobel, J. & Moffat, A. (1998). Exploring the similarity space. In: *ACM SIGIR Forum*, volume 32, ACM, pp. 18–34.