

DEFENDING AGAINST PHONE SPAM

¹S. Kaleeswari, ²C. Muruganandam

¹Research Scholar, ²Research Guide & Assistant Professor
PG and Research Department of Computer Science,
Rajah Serfoji Government college, Thanjavur 613 005.
Tamil Nadu - India.

ABSTRACT: Unwanted unsolicited communications ("spam") first became a problem for email because of the negligible cost of sending email. It is increasingly becoming a problem for phone services due to falling costs of phone calls themselves and increasing international phone access in developing countries where labour costs are relatively low. Caller ID allows individuals to ignore calls from unknown sources, but business often revolves around acquiring new customers, and caller IDs can be spoofed. In this project you will investigate technical mechanisms that could help defend against such phone spam. Examples may include applying collaborative filtering techniques (as used to filter email spam) to filter phone calls and in extending phone system software (e.g. the Asterix PBX) to check (e.g. apply ingress filtering to) caller IDs. You will implement prototypes of such mechanisms in software, possibly in the form of a smartphone app or Asterix module.

Keywords: Spam, phone

1. INTRODUCTION

Spam is unsolicited and unwanted messages sent electronically. Email spam is sent/receive over the Internet while SMS spam is typically transmitted over a mobile network. Traditional email spammers are moving to the mobile networks as the return from the email channel is diminishing due to effective filtering, industry collaboration and user awareness. The Short Messaging Service (SMS) mobile communication system is attractive for criminal gangs for a number of reasons. It is becoming cost effective to target SMS because of the availability of unlimited pre-pay SMS packages in countries such as India, Pakistan, China, and increasingly the US. In addition SMS can result in higher response rates than email spam as SMS is a trusted service with subscribers comfortable with using it for confidential information exchange.

SMS spam forms 20 to 30% of all SMS traffic in some parts of Asia such as China and India. Some methods are used to combat SMS spam such as black-and-white listing, traffic analysis and content-based filtering. According to Delany et al., 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.

Problem Statement

Mobile spam calls have been a nuisance for years, but over the last few months, it's felt to me like there's been a surge of them. I get between four and six calls daily, and a quick survey of friends shows that I'm not alone. Every waking day brings with it a new barrage. Robocallers have upped their game by masking their spam with local, genuine-looking phone numbers. Sometimes their nonsense is amusing — like when you get a threatening voicemail about your impending arrest over owed back taxes — but the vast majority of the time, it's an unwelcome distraction. It's all too easy for these scammers to wield the power of the internet and fire off countless calls with ease. And once even just a few people fall for a scam, they've made enough profit to cover their trivial expenses.

Significance of Research

Spamming remains economically viable because advertisers have no operating costs beyond the management of their mailing lists, servers, infrastructures, IP ranges, and domain names, and it is difficult to hold senders accountable for their mass mailings. Because the barrier to entry is so low, spammers are numerous, and the volume of unwanted mail has become very high. In the year 2011, the approximate figure for spam messages is around seven trillion. The costs, such as lost productivity and fraud, are borne by the public and by Internet service providers, which have been forced to add extra capacity to cope with the deluge. Spamming has been the subject of legislation in many jurisdiction.

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.

2. METHODOLOGY

Spam refers to the use of electronic messaging systems to send out unrequested or unwanted messages in bulk. The most common form of spam is email spam, but the term also applies to any message sent electronically that is unsolicited and bulk. This includes instant message spam, search engine spam, blog spam, Usenet newsgroup spam, wiki spam, classified ads spam, Internet forum spam, social media spam, junk fax spam, and so on.

Some experts estimated spam deliveries at nearly seven trillion in 2011. Unfortunately, spammers can be hard to catch, and the numbers will undoubtedly expand. As countries have passed laws outlawing spam, the technology and techniques have evolved. Whereas in the early 90s you'd see spam originating in the United States, most spam now originates overseas. As well, more spam is being sent not from a single location, but from botnets. This opens up the door to even greater security threats, as spam is used for malicious attacks such as phishing.

Detecting SMS spam campaigns

Recall that the study relies on URLs embedded in spam content to identify spam campaigns. However, many URLs do not directly point to the destination site. Techniques like URL redirections and URL shortening services are commonly used by spammers possibly to reduce the message length and to avoid content-based detection. In addition, some URLs point to a survey site where manual input (to fill out a survey) is required in order to proceed to the destination site. Due to these reasons, URLs with different forms can point to the same site. There is a need to develop a technique to identify the real site behind each URL and group the spam messages accordingly. Moreover, there is also a need to address the issue that some URLs are expired at the time of the analysis. From the one-year spam reports, 5,249 distinct embedded URLs from the spam reports are identified. 26.7% of these embedded URLs have been shortened through URL shortening services.

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. 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. 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 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. Fuzzy logic has not received much attention for SMS spam filtering. Fuzzy Similarity performs well in email spam filtering. Thus, this research investigates the effectiveness of Fuzzy Similarity in content-based SMS spam filtering.

Fuzzy Similarity

Fuzzy similarity is adapted from the Rocchio algorithm. 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, 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.

Artificial Neural Network

Artificial Neural Network (ANN) is inspired by biological nervous systems, such as the human brain; it can learn and memorize sets of data and adjust its weight matrices to build classifiers that can be used to classify unseen data.

Neuronal Model

According to Haykin, ANN consists of a large number of simple processing units called nodes or neurons. The basic design for ANN is the neuronal model as shown, a set of connecting links associated with weights that link the neurons, adder to sum the weights, an activation function that limit the resulting value from the adder into a specified range and bias which affects the net input of the activation function by increasing or lowering it.

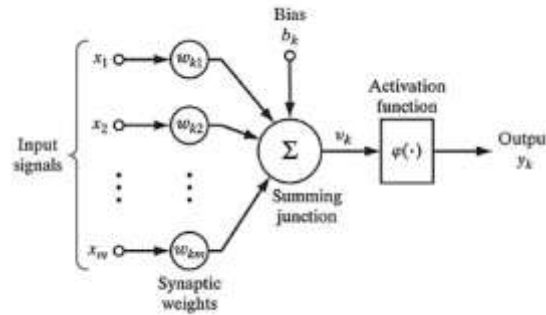


Figure 1: Neuronal model

Activation Functions

As mentioned before, activation functions limit the resulting value from the adder into a specified range. Non-linear and continuous properties are the most important characteristics in the activation functions which allow ANN to work with complex nonlinear domains. The non-linear property allows ANN to map the inputs and outputs in non-linearity, while continuous means that they remain within a specified finite range shows the activation function which is denoted by ϕ and its input v_k .

The most commonly used activation function to construct neural networks is the sigmoid function. It is defined in the range [0-1] and its expression is given.

$$\phi(v) = \frac{1}{1 + e^{-v}}$$

in which v is the neuron's output.

Scaled Conjugate Gradient

Scaled conjugate gradient algorithm (SCG) it uses second order information from ANN but requires on memory usage, where N is the number of the network's weight. Unlike gradient descent algorithm which relies on the user dependent parameters learning rate and momentum, SCG is a fully automated training algorithm which includes no critical user-dependent parameters. There are different types of conjugate gradient algorithms. Each algorithm requires a line search each repetition. This line search is computationally expensive, because it involves several calculations of either the global error function or the derivative of the global error function. SCG avoids line search by using the Levenberg Marquardt approach to scale the step size.

The derivation of this algorithm is found. Details of the SCG algorithm are as follow:

1. Initialize weight vector at the first iteration, w_1 , and set the values of $\sigma > 0$ which represents changes in weight for second derivative approximation, $\lambda_1 > 0$ which is parameter for regulating the indefiniteness of the Hessian, and $\lambda_1 > 0$. Set the initial conjugate solution, p_1 , and the steepest descent direction, r_1 , equal to the error surface gradient, $p_1 = r_1 = -E'(w_1)$. Also, set $k = 1$ and $success = true$.

2. If $success = true$ then calculate second order information sec_k :

$$\sigma = \frac{\sigma}{|p_k|}$$

$$sec_k = \frac{E'(w_k + \sigma_k p_k) - E'(w_k)}{\sigma_k}$$

$$\delta_k = P_k^T sec_k.$$

3. Scale sec_k and δ_k

$$sec_k = sec_k + (\lambda_k - \lambda_{k-1}) p_k,$$

$$\delta_k = \delta_k + (\lambda_k - \lambda_{k-1}) |p_k|^2.$$

4. If $\delta_k \leq 0$ then make the Hessian matrix positive definite:

$$sec_k = sec_k + (\lambda_k - 2 \frac{\delta_k}{|p_k|^2}) p_k$$

$$\lambda_k = 2 (\lambda_k - \frac{\delta_k}{|p_k|^2}).$$

$$\delta_k = -\delta_k + \lambda_k |p_k|^2.$$

$$\lambda_k = \lambda_{k-1}$$

5. Calculate step size α_k :

$$\mu_k = P_k^T r_k.$$

$$\alpha_k = \frac{\mu_k}{\delta_k}.$$

6. Calculate the comparison parameter Δ_k :

$$\Delta_k = \frac{2\delta_k [E(w_k)] - E(w_k + \alpha_k p_k)}{\mu_k^2}$$

7. If $\Delta_k \geq 0$ then a successful reduction in error can be made:

$$\begin{aligned} w_{k+1} &= w_k + \alpha_k p_k \\ r_{k+1} &= -E'(W_{k+1}) \\ \lambda_k &= 0. \end{aligned}$$

Success = true.

7a. If $k \bmod N = 0$ then restart algorithm by: $p_{k+1} = r_{k+1}$, else create new conjugate direction:

$$\beta_k = \frac{|r_{k+1}|^2 - r_{k+1} r_k}{\mu_k}$$

$$p_{k+1} = r_{k+1} - \beta_k p_k.$$

7b. If $\Delta_k \geq 0.75$ then reduce the scale parameter: $\lambda_k = 0.5 \lambda_k$ else a reduction in the error is not possible: $\lambda_k = \lambda_k$, success = false.

8. If $\Delta_k < 0.25$ then increase the scale parameter: $\lambda_k = 4\lambda_k$

9. If the steepest descent direction $r_k \neq 0$ then set $k = k + 1$ and go to 2 else terminate and return w_{k+1} as the desired weights.

Generally, there are many algorithms which can be used in back-propagation learning. The difference between these training algorithms is the way how they adjust the weights of the network. It is difficult to know which training algorithm will produce the best accuracy for a given problem. The most widely used algorithm is the gradient descent algorithm. The standard gradient descent algorithm is generally very slow because it requires small learning rates for stable learning; momentum variation is introduced to improve the convergence of the standard algorithm by increasing the learning rates while maintaining stability, but it is still too slow for many practical applications. For these reasons, the SCG algorithm is used to train multilayer perceptron in this research.

Support Vector Machine

Support Vector Machine (SVM) is a classification method for linear and nonlinear data. It is presented by based on early work on statistical learning theory. SVMs have been applied successfully in many real world problems such as handwritten characters and digit recognition, image classification, object recognition, speaker identification and text categorization. The general idea of SVM is separating two data sets with maximum distance between them. The following subsections relate the SVM concepts.

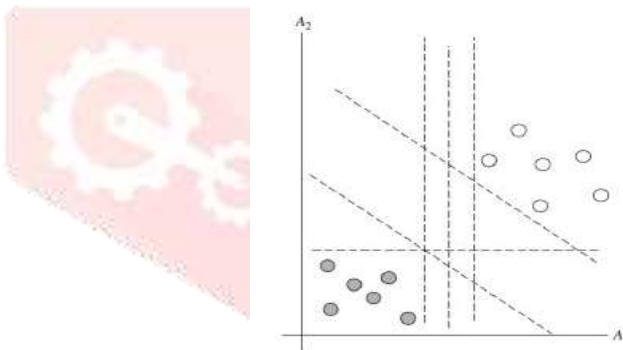


Figure 2: Linearly separable training data

Support Vector Classification

Assume that a two-class classification problem is given as follows. Let D be the data set which is given as $(X_1, y_1), (X_2, y_2), \dots, (X_l, y_l)$, where X_i is the set of training samples with associated class labels, $y_i \in \{-1, +1\}$. It is also assumed that each training sample is represented using two attributes A_1 and A_2 , as shown. It can be seen that the training data are linearly separable, because a straight line can be drawn to separate the two classes.

The shortest distance from the maximum marginal hyperplane to one side of its margin is equal to the shortest distance from the maximum marginal hyperplane to the closest training sample of either class. For linear classification, the two classes and the maximum marginal hyperplane separating them can be identified.

$$\begin{aligned} wX_i + b &\geq 1, y_i = 1 \\ wX_i + b &\leq -1, y_i = -1 \end{aligned}$$

in which w is a weight vector and b is a bias. By combining Equations, $y_i \times (wX_i + b) \geq 1, i = 1, \dots, l$ in which l denotes the number of support vectors. The samples that lie on the margin are called support vectors shows one support vector from each class which is encircled with a thicker border.

$$d = \frac{2}{\|w\|}$$

where $\|w\|$ is the Euclidean norm of w .

Better separation between the two classes can be achieved by maximizing d . In order to maximize d while making sure that all the training samples are on the correct side of the hyperplane, w must be minimized. This problem can be solved using the Lagrange function $L = f - \alpha g$, where f is the function that will be minimized which is equal to $1/2 \|w\|^2$ a simplification of the weight vector calculation, α is the Lagrange multiplier and g is the constraint which is equal to $y_i (wX_i + b) \geq 1$. Therefore by applying the Lagrange function, it yields:

$$L(w,b,\alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^l \alpha_i [y_i(wX_i + b) - 1]$$

can be solved by minimizing according to w and b , maximizing according to $\alpha_i \geq 0$ values, w can be obtained

$$w = \sum_{i=1}^l \alpha_i y_i X_i, \alpha_i \geq 0$$

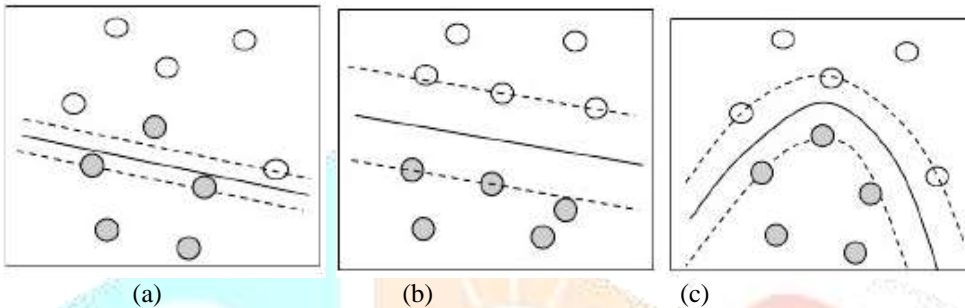


Figure 3: Two class nonlinear separable problem

According to $\sum_{i=1}^l \alpha_i \cdot y_i = 0, \alpha_i \geq 0$ can be replaced into the Lagrangian formula which yields:

$$L(w,b,\alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^l \alpha_i \cdot y_i X_i, \alpha_i \geq 0$$

An upper bound for the Lagrangian multipliers is introduced by C , where $0 \leq \alpha_i \leq C$. C is called a penalty parameter of the error term and is best determined experimentally.

Optimum separation is provided by a nonlinear separating surface in the original space. The dot product computation involved is expensive to compute for the transformed data samples. Thus, kernel function $K(X_i, X_j) = (\phi(X_i) \cdot \phi(X_j))$ is used so that all calculations are made in the input space. Many kernels have been proposed by researchers such as Linear, Polynomial, Sigmoid and Radial Basis Function kernels.

3. RESULT AND DISCUSSION

Show the accuracy in detecting spam SMS

Valid (non-spam) Message (Self)	Invalid (spam) Message (non-Self)	Total error
84%	65%	20%

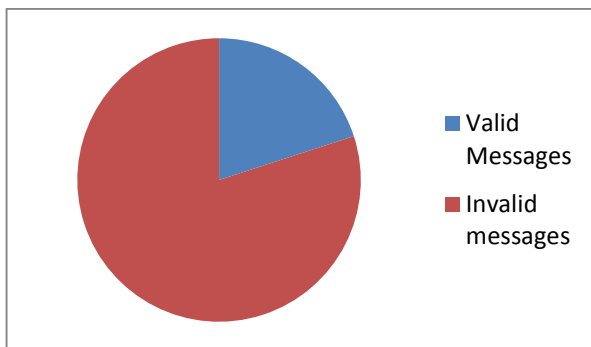


Figure 14: Identifying Valid and Invalid messages

Support Vector Machine (SVM)

SVM: The relationship between database size and accuracy rate to classify spam messages

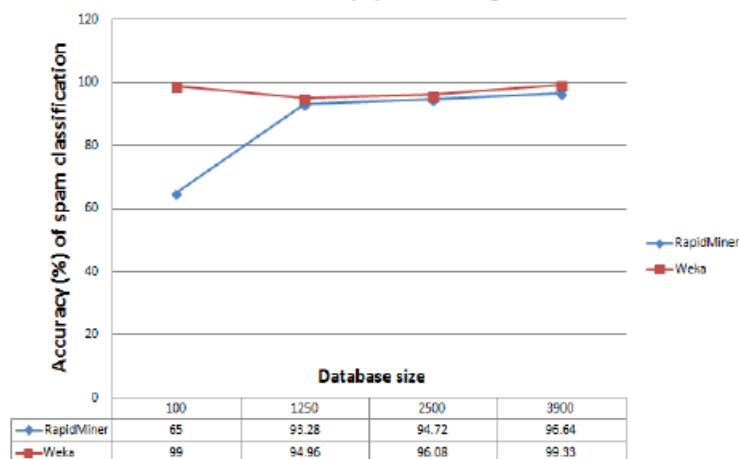


Figure 15: The relationship between database size and accuracy rate of classify spam messages

4. CONCLUSION:

In this study expand our analysis on SMS spam data presented. We investigate the characteristics and classification time are also analyzed, the relevant background in using content-based filtering technique for SMS spam filtering. SMS spam automatically using Text Classification techniques that consist of training, filtering, and updating processes a computer or a large amount of SMS data in advance to filter SMS spam, especially for the training. This increases hardware maintenance and communication costs. The selected articles were Applying Text Classification on an independent mobile phone also ensures security and privacy, as spammers do not have the chance to inspect the filtering system and users do not have to store SMS anywhere.

Future work for this project involves continuing to study SMS spammer behaviour in order to current algorithms based on changes in behaviour, and to develop methodologies as spammer behaviours being evolved.

REFERENCES:

1. E. B. Cleff, "Privacy issues in mobile advertising," *Int. Rev. Law Comput. Technol.*, vol. 21, no. 3, pp. 225_236, 2007.
2. I. Murynets and R. P. Jover, "Analysis of SMS spam in mobility networks," *Int. J. Adv. Comput. Sci.*, vol. 3, no. 1, pp. 1_8, 2013.
3. A. Skudlark, "Characterizing SMS spam in a large cellular network via mining victim spam reports," in *Proc. 20th ITS Biennial Conf.*, Rio de Janeiro, Brazil Nov./Dec. 2014, pp. 1_23.
4. H. Saeed and W. Waheeb, "The performance of soft computing techniques on content based SMS spam filtering," M.S. thesis, Dept. Elect. Eng., Univ. Tun Hussein Onn Malaysia, Johor, Malaysia, 2015.
5. 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.