# A Content-Based Spam E-Mail Filtering Approach Using Multilayer Percepton Neural Networks

**3 authors**, including:

Nandita Chaudhuri
Indian Institute of Technology Guwahati
**6** PUBLICATIONS    **1** CITATION

SEE PROFILE

Some of the authors of this publication are also working on these related projects:

Content based e-mail classification View project

ICoRD 2019 View project

# A Content-Based Spam E-Mail Filtering Approach Using Multilayer Percepton Neural Networks

A.Sesha Rao[#1], P.S.Avadhani [#2], Nandita Bhanja Chaudhuri[#3]

[#1]*Dept. of Computer Science & Systems Engineering, AU College of Engineering, Visakhapatnam, India.*
[#2]*Principal, AU College of Engineering, Andhra University, Visakhapatnam, India*
[#3]*Dept. of Information Technology, Vignan's Institute of Engineering for Women, Visakhapatnam, India*

*Abstract— Nowadays increased spam e-mails are causing inconvenience to internet users and organizations and are considered as a serious wastage of resources, time, memory, space and efforts. Therefore, it is crucial to have an automatic e-mail classification system for the identification of spam e-mails. Spam mails need to be classified and separated from ham (non-spam) mails as they are the source of financial loss and annoyance for the recipients. The spam e-mail classifier performance can be greatly enhanced with the use of Artificial Neural Network classification. It has capability of learning huge amount of data with high dimensionality in a better way. In this paper, Multilayer Perceptron and Back Propagation Training algorithm is explored where 'generalized delta' rule is used for weight adjustments for hidden layers. The Perceptron uses Back Propagation Learning model for calculating its gradient. For fast convergence the learning rate η is changed for every iteration which is proportional to the negative gradient of the instantaneous error with respect to η. To avoid the local minima problem the weights are initialized to small random numbers which are uniformly distributed in the range [ -α/$\sqrt{Ni}$, + α/$\sqrt{Ni}$ ], where Ni is the number of inputs, and α takes value in (1, 3). In this paper, four Multilayer Perceptron (MLP) Network models are constructed. For testing our model bench mark data drawn from UCI, Machine learning Repository is employed for training the neural network. The results of our MLP model are reasonable in terms of TP rate, FP Rate, Accuracy, Precision, Recall, F-measure, MCC, ROC Area, PRC Area.*

**Keywords** —back propagation, delta rule, F-measure, hidden layers, learning rate, local minima, MCC, perceptron, precision, PRC Area, recall, ROC.

## I. INTRODUCTION

Along with the growth of the Internet and e-mail, there has been a dramatic growth in spam in recent years. E-mail spam, known as Unsolicited Bulk Email (UBE), junk mail, or Unsolicited Commercial Email (UCE), is the practice of sending unwanted e-mail messages, frequently with commercial content, in large quantities to an indiscriminate set of recipients [1].

These junk mails can contain various types of messages such as pornography, commercial advertising, doubtful product, viruses or quasi legal services [1]. The inverse of "spam" email is called "ham" which needs by recipients. [2] Links on social networks that take us to free downloads, easy income, pornography and unsolicited text messages offering loans, low priced products etc. Why we need spam detection? Spam causes annoyance and wastes user's time to regularly check and delete this large number of unwanted messages. Flooding of mailboxes with spam e-mails waste storage space and overload the server; thus it may lead to losing legitimate e-mails, degrading the server performance, or even make it totally unavailable. Hence, spam consumes network bandwidth and server storage space.[3] Hence it is crucial to have automatic spam filtration system for every individual user.

[4] The intent of spam content is to catch eyeballs and conduct users to take some actions. Spam email usually contains links to entice the users to click on those links and visit particular Web sites, which we call spam-advertised sites.

Spam has grown steadily since the early 1990s. Recent reports show that spam accounts for approximately 70-90% of email traffic sent on the Internet-about 70 billion messages a day [5,6,7]. Although several spamming botnets were taken down, such as Rustock, Kelihos, and Grum [7,8], spam remains a serious security threat and attackers keep establishing new spam campaigns.

The information and analysis of E-mail statistics 2015-2019 – Executive Summary report, based on primary research was conducted by The Radicati Group, Inc. It reported that over next four years, the average number of email accounts per user ratio will grow from an average of 1.7 to 1.9 accounts. In the year 2015, the number of emails sent and received per day is around 205 billion [ 9].

The task of spam filtering is to rule out deceptive messages automatically from a user's inboxes. There are many techniques evolved which seek, to identify whether a message is spam or ham

based on the content and other characteristics of the message which improves the accuracy of spam filters. These machine learning algorithms are important in the continuous fight against spam. In spite there are large number of methods and techniques available to filter spam, but the volume of spam on the internet is still rising. [10]

Due to security of information, a spam filter must be placed in computer network. There are various techniques which have been used by the different authors among these data mining is one of the popular techniques to develop classifier to classify spam and non-spam data. In order to tackle problems faced by users due to spam e-mail, it is necessary to classify them with the help of intelligent and robust classifier. These classifiers should have the capability to classify spam e-mail against non-spam e-mail. The spam e-mail classifier performance can be greatly enhanced with the use of Artificial Neural Network (ANN) classification algorithms. Artificial Neural Network is a powerful tool used for classification of data [11], it has capability of learning huge amount of data with high dimensionality in better way. There are various parameters of ANN to be set to tune for the better performance of neural network model, these are learning rate, architecture of ANN and momentum term, these all parameters play a very important role in improving the accuracy of ANN model.

In this paper we have used Multilayer Perceptron Network model with BPL. Multilayer Perceptrons are network of linear classifiers. The Back Propagation Learning (BPL) used generalised delta rule for weight adjustments in hidden layers. BPL is trained on spam e-mail data set with different partitions and different learning rate. Different performance measures like precision, recall, F-measure and accuracy of the best model with feature selection have also been calculated.

Rest of the paper is organised as follows. The related work is presented in section II. The proposed design and implementation of the mathematical model for spam mail filtering using Multilayer Perceptron network is described in section III. In section IV, we present the experiential results and evaluation process. Conclusion and scope of the future work are given in section V.

### 1.1 Artificial Neural Network (ANN)

ANN is a computational model, which is based on Biological Neural Network. To build ANN, artificial neurons, also called nodes are interconnected. An artificial neuron is an abstraction of biological neuron and the basic unit in ANN. Arranging neurons in layers mimics the layered structure of certain portions of brain. The power of neural computation comes from connecting neurons into networks. There are several classes of neural networks, classified according to their learning mechanisms, but structurally they are classified into three fundamentally different classes.
i) Single layer artificial neural networks.
ii) Multilayer artificial neural networks.
iii) Recurrent Networks.
(i) and (ii) are called feed forward neural networks.

In single layer artificial neural networks, the inputs are directly supplied to the outputs through a series of weights. Multilayer artificial neural networks map the input to suitable output via hidden layers. In recurrent neural network, the connection between units forms a directed cycle.

Multilayer artificial neural network is a feed forward neural network which maps a set of inputs x onto a set of outputs y using multi weights connections. It consists of an input layer, an output layer, and one or more hidden layers.

## II. RELATED WORK

Several Machine Language techniques are brought together to explore the best methods [12,13]. The most commonly used technique is ANN as it gives better classification results. Authors Ali Rodan et al., [14] developed an MLP neural network model with Biogeography Based Optimization for identifying e-mail spam. In this work, MLP is trained using Biogeography Based Optimization based on two different spam data sets and compared with other MLPs trainer with back propagation algorithm and common meta heuristic algorithms: Genetic Algorithm, Distributed System Operators, Differential Evolution and Ant Colony Optimization [15]. An anti-spam filtering was presented by his techniques are centred on artificial neural network (ANN) and Bayesian networks.

Rodan et. al., [15] proposed the application of Miltilayer Perceptron for the purpose of e-mail classification where the weights of this network model are found using a new nature-inspired meta heuristic algorithm called Biogeography Based Optimization (BBO). Experiments are conducted using two different datasets revealed that MLP model trained by BBO algorithm gets high generalization performance compared to any other optimization methods.

WangI et. al., [16] proposed Hill Climbing, Simulated Annealing, and Threshold Accepting optimization techniques as feature selection algorithms. Performance of the above three techniques are compared with the Linear Discriminate Analysis. The experimental results show that the proposed strategies not only reduce the dimensions of the e-mail, but also improve the performance of the classification filter. The classification accuracies are 90.0% for LDA, 93.6% for HC, 94.6% for TA and 95.5% or SA as compared to 88.1% for the system without feature selection.

Qiuming et. al. [17], projected a method to filter text spam e-mails using Back Propagation Neural

Network model. In the experiment, they utilized spamassassin corpus as data set, and implemented cross validation method to inspect network simulation result. Experimental result showed that by taking different model parameters, the network will gives different results. To get more accurate classification results, a large number of data sets need to be trained and shall learn how to optimize the network model to achieve better filter results.

Ismailia Idris [18] has proposed neural network model for spam e-mail classification, he has compared the result of both neural network and SVM models in terms of accuracy.

EL-Sayed M., EL-Alfy et al, [19] have introduced abduct network model for spam email classification and compared the results of abduct network model with other Group Method of Data Handling (GMDH) based network.

Harshal Deshmukh et al, [20] used ANN to detect text as well as image based spam to achieve the objective and they have used Enron corpus's data set of spam emails. Various spam detection steps such as pre-processing step, or data clearing step, representation of data and classification of spam or non spam e-mail messages are discussed.

Mohamad and Selamat [21], worked on, a spam filtering technique, which implements a combination of two types of feature selection methods before its classification task. The authors attempted to conduct an experiment to test the efficiency of the proposed Hybrid Feature Selection, which is a combination of Term Frequency Inverse Document Frequency (TFIDF) with the rough set theory in spam email classification problem. The result shows that the proposed Hybrid Feature Selection returns a good result.

Loredana Firte et al,[22] presented a new approach for spam detection filter. The solution proposed is an offline application that uses the K-Nearest Neighbour algorithm and pre- classified email data set for the learning process. K-Nearest Neighbor algorithm classified the messages which are based on features extraction from the emails properties and content .

Nosseir, Khaled Nagati et al, [23] performed a work, Intelligent Word-Based Spam Filter Detection Using Multi-Neural Networks‖. They proposed a character-based technique. A multi-neural networks classifier is used by this approach. A normalized weight values derived from the ASCII value of the word characters are used to train the neural network.

Reena Sharma et al, proposed [24] an efficient spam filtering technique based on neural network. The technique used is RBF a neural network technique in which neuron are trained. The results obtained by using this technique are compared with SVM. The parameter for comparison is precision and accuracy.

Thomas et. al., [25] their study identifies five feature selection techniques namely, namely

TFDF,MI, WMI, CDM and Chi-square which are used in the general text classification for spam filtering. Also, the classification and prediction is performed using different entities of email such as header, body and subject that can be used for effective identification of spam mails. They presented a comparative study of different feature selection methods. Through extensive experiments we demonstrated that Weighted Mutual Information feature selection with header and body of the emails is efficient in email classification based on SVM**.**

Mohammad et. al., [26], utilized Fisher filtering feature selection algorithm. The best error rate is 0.0089. They implemented CAV algorithm which reduced the error rate by 39% compared with other Rnd Tree and SVM.

Moreover, several ML techniques are combined together to produce a more accurate and robust detectionvmethods. For instance, Artificial Neural Network (ANN) is a commonly used technique as it gives accuratevclassification results [27]. ANNs are inspired by the biological neural systems. The most popular and applied type of ANNs is the Feedforward Multilayer Perceptron (MLP).

## III. DESIGN AND IMPLEMENTATION

The model of a Multilayer Perceptron is designed and implemented in this paper to classify spam e-mails. It is flexible structure and non-linearity transformation to accommodate latest spam mails. The experiments are done on publicly available UCI Machine Learning Repository.

Perceptron is a general feed-forward single layer net work with one output neuron. It learns to separate the input pattern. It uses supervised learning algorithm. Hence, the training set consists of a set of input vectors, along with its desired target vector. Input vector components take on a continuous range of values; and target vector components are binary valued. After training, the network accepts a set of continuous inputs and produces the desired binary valued outputs. The perceptron separates linearly separable input patterns. It cannot separate the input pattern sets which are linearly inseparable.

Multilayer neural networks may be formed by simply cascading a group of single layer networks. Such networks provide increase in computational power over a single layer networks by having non linear activation function between layers. Hence multilayer networks consist of one or more hidden layers apart from input layer and output layer. The input layer neurons receive the input signals, perform no computation, and the output layer neurons produce simply the output signals. The function of hidden neurons is to intervene between the external inputs and the network outputs. Connection between neurons is represented as synaptic weights. A two layered neural network require two synaptic weight matrices $w_{i,j}$ and $k_{i,j}$ as shown in our model Fig. 1.
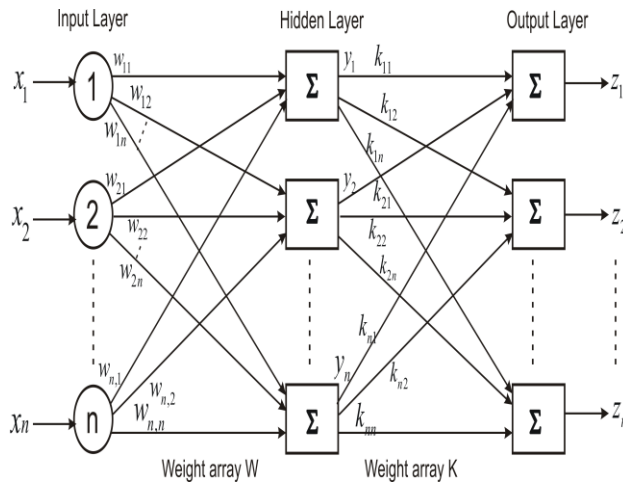
Fig. 1 Architecture of Multilayer Perceptron

Calculation of output for the first (hidden layer) layer consists of multiplying the input vector by the first weight matrix, $w_{ij}$. Output of hidden layer is given by

$$y_i = \sum_{i,j=1}^{n} x_i . w_{i,j} \qquad \rightarrow (1)$$

The output of output layer is given by

$$z_i = \sum_{i,j=1}^{n} y_i . k_{i,j} \qquad \rightarrow (2)$$

Where

$$x_i = (x_1, x_2, \dots \dots .x_n)$$

$$y_i = (y_1, y_2, \dots \dots .y_n)$$

$$z_i = (z_1, z_2, \dots \dots .z_n) \text{ and}$$

$w_{i,j}$ & $k_{i,j}$ are synaptic weight matrices of hidden

layer and output layer.

Each neuron consists of a summation function and activation function as shown below. A single artificial neuron with activation function is shown in Fig. 2.
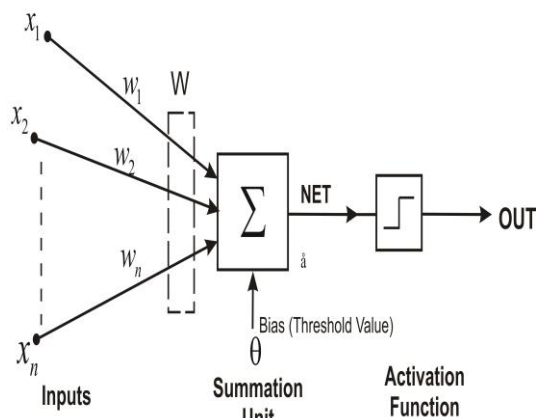


Fig, 2. A Single Node with Summation Unit and Activation Function

Fig. 2 shows the neuron used as fundamental building block for ANN, A set of inputs

$X(x_1, x_2, \dots, x_n)$ are applied either from external or previous layer. Each of these inputs is multiplied by a connecting weight $w_i$, and the weighted sum is given by

$$NET = x_1 w_1 + x_2 w_2 + \dots x_n w_n$$

$$= \sum_{i=1}^{n} x_i \ w_i \qquad \rightarrow (3)$$

After the output 'NET' is calculated activation function F is applied to it there by producing the final output signal 'OUT'.

$$OUT = F(NET) = \frac{1}{1+e^{-NET}} \qquad \rightarrow (4)$$

Where, F is sigmoid function shown in fig. 3.

$$F'(NET) = OUT (1-OUT) \qquad \rightarrow (5)$$

Equation (4) shows that sigmoid function is differentiable and provides the required non linearity and automatic gain control.
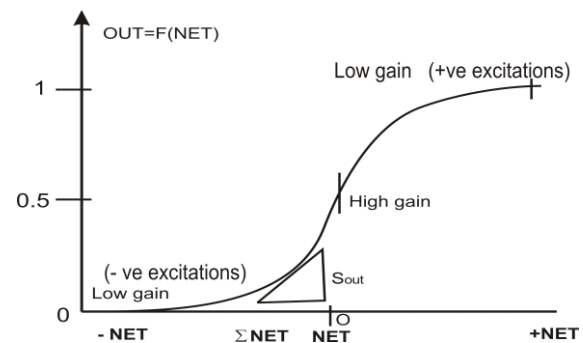


Fig. 3 Sigmoid Function

### 3.1Back Propagation Learning Algorithm for MLP

The difficulty with the multilayer perceptron is adjusting the weights of hidden layers during training process. . More number of hidden layers increases the difficulties. So, back propagation learning is a better substitute to train Multi layer perceptron network.

It adopts supervised learning. It is one of the ways to train artificial neural networks used along with an optimization method such as steepest descent. The method calculates the derivative of a loss function with respect to all the weights in the network, so that the derivative is fed to the optimization method which in turn uses it to update the weights in hidden layers. Hence, it attempts to minimize the loss function.

Back propagation requires a known, desired output (target vector) for each input vector in order to calculate the loss function gradient. It is therefore usually considered to be a supervised learning method; although it is also used in some unsupervised networks such as auto encoders etc. It is a generalization of the delta rule to multi-layered feed-forward networks, made possible by

using the chain rule to iteratively compute gradients for each layer. Back propagation requires that the activation function used by the nodes be differentiable.

Training is accomplished by sequentially applying input vectors while adjusting network weights according to predetermined procedure. Here we adopt supervised training. In supervised training we present both input and desired (target) output vectors to the network.

### 3.2 Preceptron Training

The training set consists of a set of input vectors each with its desired target vector. Input vector components take on a continuous range of values and target vector components are binary valued. Perceptron training cannot be applied for multilayer networks. However, multi layer network provides increase in computational power over a single layer network by having non linear activation function between layers.

To update the weights one must calculate the error, i.e., difference between the actual output and the desired (target) output. This error is easily calculated at the output layer. At the hidden layers, however, there is no idea about the target vector. So, some other technique must be used to calculate an error at the hidden layers that will cause optimization of the output error.

MLP algorithm is a systematic method for training multilayer neural network. It considers only feed forward networks. MLP provides automatic gain control where,
(i) Large signals can be accommodated by the network without saturation.
(ii) While the small signals are allowed to pass through without excessive attenuation.
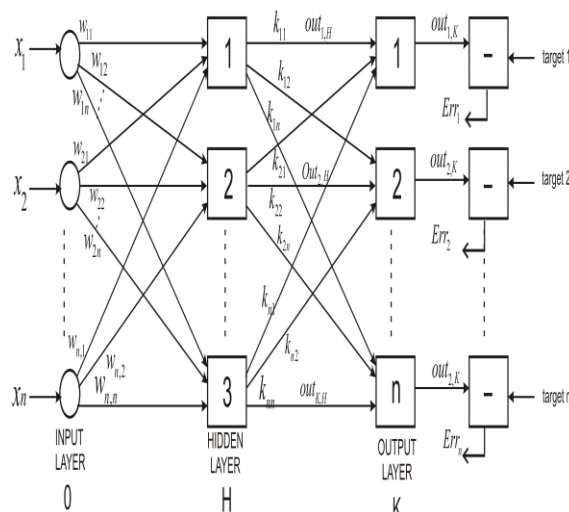


Fig. 4 Architecture of Multilayer Perceptron

A typical MLP starts as a network of nodes arranged in three layers, input, hidden, output layers. The architecture of MLP is shown in Fig. 4. There

is no theoretical limit on number of hidden layers but 2 or 3 layers are preferable. First layer neurons serve as distribution points. They perform no computation. Each neuron in the subsequent layers produces 'NET' and output signals after applying sigmoid function. The input layer is designated as 0. The weights in layer 0 terminate on the neurons of subsequent layer H and so on.

### 3.3 MLP Algorithm
### 3.3.1 Two Passes of Computation

The algorithm constitute a 'forward pass', in that signal propagates from the network input layer to its output layer and a 'reverse pass', here, calculates error signal propagates backward through the network where it adjusts the weights. Thus, the MLP learning algorithm cycles through two distinct phases a 'forward pass', followed by 'backward pass' through the layers of network. The algorithm alternates between these phases several times as it scans the training data.

Calculations are performed on a layer by layer basis. First the outputs of the neuron in hidden layer H are calculated. These are then used as inputs to subsequent hidden layers followed by output layer K. Then layer K neurons outputs are calculated. Each of the network output is subtracted from its corresponding component of the target vector. This error is used to adjust the weights of the network in all the layers. After enough repetitions, the error between actual output and target outputs should be reduced to an acceptable value. At this point, the network is used for recognition and weights are not changed further.

**i) Forward Pass:** The first pass is referred to as forward pass. In this pass the signal propagates from the network input layer to its output layer. During this the synaptic weights remain unaltered throughout the network and the function signals are computed on a node by node basis. The 'NET' value of each node in the first layer is calculated as weighted sum of its node's output. The final output of each node is computed using activation function. The error is calculated at the output layer as follows.

$$NET_{i,H} = \sum_{i,j=1}^{n} x_i . w_{i,j} \qquad \rightarrow (6)$$

$$OUT_{i,H} = F\left(NET_{i,H}\right) = \frac{1}{1 + e^{-NET_{I,H}}} \quad \rightarrow (7)$$

Where, $w_{i,j}$ is the synaptic weight connecting neuron i to neuron j in the hidden layer. $x_i$'s are the input signal of node i.

The output vector of one layer is the input vector to the next layer. So calculating the output of final layer requires the application of eqn. 7 in each layer starting from the network input layer to its output layer.

**ii) Reverse Pass:** This phase begins with the computation of error at each node in the output layer. The error in the output layer is calculated by subtracting the computed output from the desired output. The reverse pass on the other hand, starts at the output layer by pushing the computed error signals backward through the network layer by layer. Because a target value is available for each neuron in the output layer only adjusting the associated weights is usually accomplished using a modified delta rule. So, we compute local gradient δ for each node.

### ii. a) Weight adjustment for output layer

For a node located in the output layer, the δ is simply equal to the error signal of that node multiplied by first derivative of sigmoid function . If K is the output layer, error of the $nq^{th}$ node in the $K^{th}$ layer is given by

$$e_{q,K} = (tar\,get\;output\; - \;camputed\;\;output\;) \quad \rightarrow (8)$$

$$Local\;gradient\;,\;\delta_{q,K} = F'\left(NET_{q,K}\right)e_{q,K}$$

$$= out(1 - out).e_{q,K} \quad \rightarrow (9)$$

$$\Delta w_{pq,K} = \eta\,\delta_{q,k}.out_{P,H} \qquad \rightarrow (10)$$

$$w_{pq,K}(t+1) = w_{pq,K}(t) + \Delta w_{pq,K} \qquad \rightarrow (11)$$

$\Delta w_{pq,K}$ : Weight factor from node p in the hidden layer to node q , in the output layer

$w_{pq,K}(t)$ :The value of weight from neuron p in the hidden layer to node q in the output layer at step t.

$w_{pq,K}(t+1)$: Value of the weight at step t+1 (after adjustment)

$\delta_{q,K}$ :Local gradient from node q in the output layer
Out $_{p,H}$: the value of out from node p in the hidden layer H.
Out : output of node in output layer K.
$\eta$ : learning rate.

### ii. b) Weight adjustment for hidden Layer

MLP trains the hidden layers by propagating the output error back through the network layer by layer adjusting the weights at each layer. Since hidden layers have no target vector, the above training process may not work. MLP trains hidden layers by propagating the output layer 'error', $\delta$, back through the network layer by layer, there by adjusting weights in each layer.

Consider a single node p, in the hidden layer just before the output layer. In the forward pass, this node propagates its output (OUT) to nodes in the output layer through the inter connecting weights. During training these weights operate in reverse

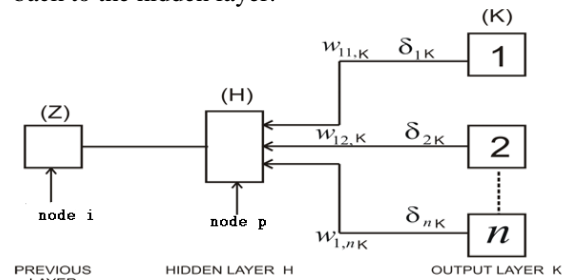direction, passing value of δ from the output layer back to the hidden layer.



Fig. 5 Weight Adjustments in Hidden Layer

Local gradient $\delta_{p,H}$ from node p in hidden layer H is obtained as a product of the associated derivative $F'\left(NET_{P,H}\right)$ and the weighted sum of the local gradients $\left(\delta_{1,K},\delta_{2,K},.........\;.......,\delta_{n,K}\right)$ computed for the neurons in the next hidden layer or output layer that are connected to node p.

$$\delta_{P,H}(t) = F'\left(NET_{P,H}\right)\sum_{q=1}^{n}\delta_{q,K}.w_{pq,K}$$

$$= out_{p,H}\left(1 - out_{p,H}\right)\sum_{q=1}^{n}\left(\delta_{q,K},w_{pq,K}\right) \qquad \rightarrow (12)$$

With $\delta$ in hand, the weights feeding the first hidden layer can be adjusted as

$$\Delta w_{pq,H} = \eta\,\delta_{P,H}.out_{i,H} \qquad \rightarrow (13)$$

Updated weight of hidden layer is given by

$$w_{pq,H}(t+1) = w_{pq,H}(t) + \Delta w_{pq,H} \qquad \rightarrow (14)$$

Where $w_{pq,H}(t).\;w_{pq,H}(t+1),\delta_{p,H},out_{i,H}$ are

defend as above

For each node in a given hidden layer, the local gradient and all weight associated with that layer must be adjusted. This is repeated moving back towards the input layer by layer, until the weights are adjusted.

The performance of MLP law depend on the initial setting of the weights, learning rate parameter, activation formation of the units and presentation of data.

### 3.4 Weight Initialization

Owing to its gradients descent nature, backdrop is very sensitive to initial conditions if $w^0$ is close to local minima then the convergence will be very fast. On the other hand back drop convergence very slow if $w^0$ starts the search in a relatively flat region of the error surface. (for away from local minima). In practice the weight are normally initiated to small random numbers that are uniformly distributed in a small range of values. The range is typically $\left[-\alpha/\sqrt{N_i},+\alpha/\sqrt{N_i}\right]$ where $N_i$ is the number of

inputs to the i[th] unit. The value of $\alpha$ is it typically in the range (1, 3)

### 3.5 Learning Rate

Better convergence in learning can be achieved by adopting learning rate, $\eta$, which are suitable for each iteration. For this the change is made proportional to the negative gradient of the instantaneous error with respect to $\eta$.

$$\Delta\eta_{ji}(t+1) = -\gamma \; \frac{\partial \; e^2(t)}{\partial \; \eta_{ji}(t)} \qquad \rightarrow (15)$$

Where $\gamma$ is proportionality constant

### 3.6 Momentum Term

One way to increase the rate of learning is by using a momentum term in the weight change (Rumel hart, Hinton and Williams 1986). The method involves adding a term to the weight factor that is proportional to the amount of previous weight change. Once an adjustment is made, it is 'remembered' and serves to modify all the subsequent weight adjustments. The adjustment eqns. are

$$\Delta w_{pq,K}(t+1) = \eta\left(\delta_{q,K}\;.out_{p,H}\right) + \alpha\left[\Delta w_{pq,K}(t)\right] \qquad \rightarrow (16)$$

$\alpha$ : is momentum coefficient is commonly used set to around 0.9

### 3.7 Implementation of the Model
#### 3.7.1 Dataset Collection

Data is collected from the UCI Machine Learning Repository. It is a standard dataset which are applied to this model. It contains both samples of spam and non-spam e-mail data. The dataset has 57 continuous attributes and 1 nominal class label. There are 4501 number of instances among that 1813 are spam's which is around 39.4%.
Source:
http://www.ics.uci.edu/~mlearn/MLRepository.html

#### 3.7.2 Preprocessing

The data that is collected is in Comma Separated Version (.CSV) file format. For applying the neural network techniques, a tool called Weka is used. This tool supports Attribute Relation File Format (ARFF). So the data had been converted to ARFF. The original data set contains 58 attributes but only the important attributes needs to be selected. This is called as attribute selection mechanism. To achieve this, an algorithm called Principal Component Analysis (PCA) is used along with Ranker Search method which is shown in Fig.6. The resultant of this is 47 continuous attributes and 1 nominal class label.


Fig.6: Result of Principal Component Analysis

After this, the data is partitioned into two set- 70% for training and 30% for testing.
The experiments are also conducted by collecting variant samples of the datasets which are demonstrated in Table-1, Table-2, Table-3, and Table-4.

**TABLE I**
Sample1: 25% of the Training samples and 805 Instances



**TABLE III**
Sample 2: 50% of the samples and 1610 Instances



**TABLE IIIII**
Sample 3:75% of the samples and 2416 Instances

**TABLE IVV**
Sample 4:100% of the samples and 3221 Instances



### 3.7.3 MLP Training Phase (Classification)

Once pre-processing phase is completed, the spam-based data set is ready to be trained by MLP. The aim of this phase is to train MLP based on the amount of samples of spam and non-spam (ham) messages to be able to distinguish between them. The samples are selected randomly from spam-based data sets. The weights of all input nodes are initialized to small random numbers and given to the network.

### i)   MLP-1 Training Implementation

MLP is a form of supervised learning for multilayer perceptron. It is most often used as training algorithm. Four different neural network structures for MLP learning  are proposed in our work.

The first structure coined as *MLP-1,* consists of input layer with 46 nodes and 25% of the training spam-based dataset. There are two hidden layers in this structure - 8 nodes at the first hidden layer, 6 nodes at the second hidden layer. There are two nodes at the output layer i.e., spam and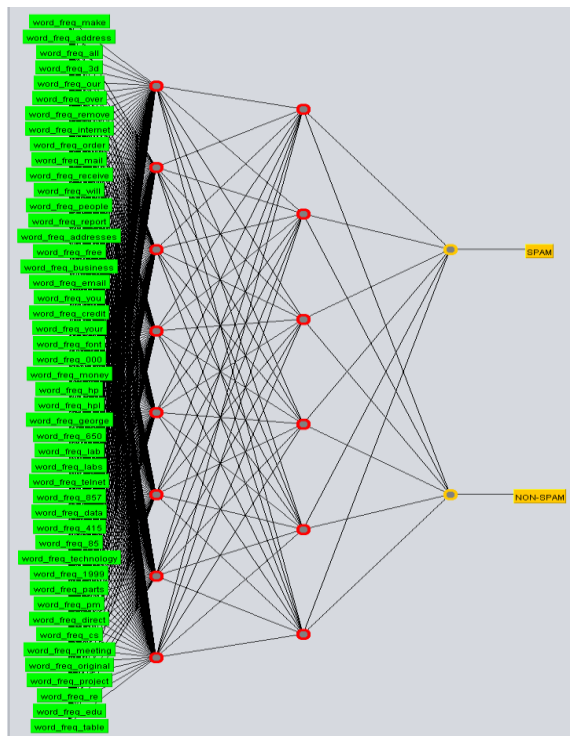 non-spam. Fig.7 demonstrates MLP-1 structure. The algorithm in section  3.3 clarifies how MLP-1 operates for  all the below mentioned   structure MLP-1, MLP-2, MLP-3, and MLP-4.shown in Fig.7, Fig.8, Fig.9, Fig.10.



Fig. 7 Input Layer, Hidden Layers, and Output Layer of MLP-1

### ii)   MLP-2 Training Implementation

The second structure named  *MLP-2,* consists of input layer with 46 nodes and 50% of the training spam-based dataset.  We have used two hidden layers for this structure - 8 nodes at the first hidden layer, 6 nodes at the second hidden layer. There are two nodes at the output layer i.e., spam and non-spam.  Fig.8 demonstrates MLP-2 structure.



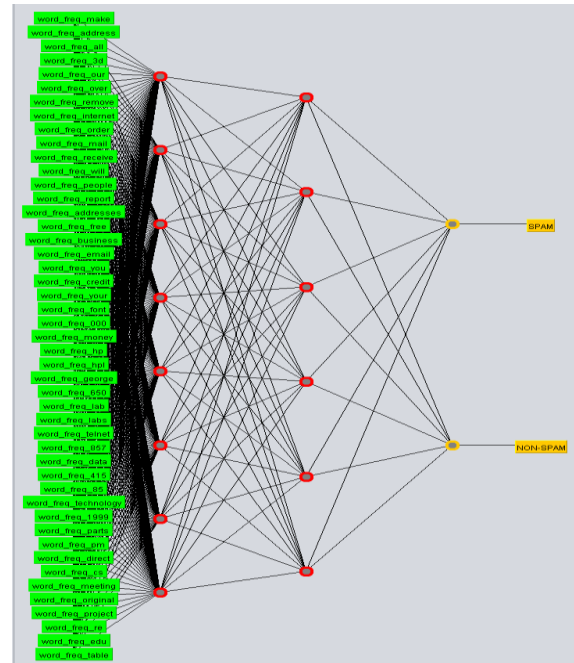Fig. 8 Input Layer, Hidden Layers, and Output Layer of MLP-2

### iii) MLP-3 Training Implementation

The third structure coined as *MLP-3,* consists of input layer with 46 nodes and 75% of the training spam-based dataset. There are two hidden layers in this structure - 8 nodes at the first hidden layer, 6 nodes at the second hidden layer. There are two nodes at the output layer i.e., spam and non-spam. Fig.9 demonstrates MLP-3 structure.



Fig. 9 Input Layer, Hidden Layers, and Output Layer of MLP-3

### iv) MLP-4 Training Implementation

The fourth structure coined as *MLP-4,* consists of input layer with 46 nodes and 100% of the training spam-based dataset. There is only one hidden layer with multiple nodes. However there are two nodes at the output layer i.e., spam and non-spam. Fig.10 demonstrates MLP-4 structure.



Fig. 10 Input Layer, Hidden Layers, and Output Layer of MLP-4

### 3.7.4 Testing Phase

When the training phase for four network structures is completed, the Neural Network (NN) models are ready to identify if the incoming e-mail is legitimate or spam. The testing phase uses these NN models to determine the type of message as a spam or an e-email. Testing is done for checking the accuracies of a model.

In this paper, 30% of the total dataset are selected for testing. The test samples should be pre-processed using PCA algorithm. When the pre-processing is completed, the test sample propagates forward the network on every in the layer. Then the weight sum of the inputs to the node is calculated, and the bias value is added to the sum, finally the activation function to the node (i.e., the desired output) is calculated. The network decides if the test sample is legitimate email or a spam depends on the desired output.

Various samples of the datasets are applied to the above NN models and classification of the spam and non-spam (ham) mails are achieved.

## II. RESULT ANALYSIS

There are four networks in our model. The input for each model is a sample dataset and outputs are the desired classes. For these models, the performance metrics during training phase is demonstrated on Table 5.

| | | | | |
|---|---|---|---|---|
| **Total number of instances** | 344 | 690 | 1034 | 1380 |

**TABLE V**

Performance metrics of the four models during Training Phase.

| Performance Metrics | MLP-25% | MLP-50% | MLP-75% | MLP-100% |
|---|---|---|---|---|
| **Correctly classified instances** | 705-87.58 % | 1460-90.68% | 2226-92.14% | 2971-92.23 % |
| **Incorrectly classified instances** | 100-12.42 % | 150-9.31 % | 190-7.86% | 250-7.77 % |
| **Mean statistic** | 0.7514 | 0.7642 | 0.7701 | 0.6867 |
| **Mean absolute error** | 0.169 | 0.1517 | 0.1249 | 0.1361 |
| **Root mean squared error** | 0.2838 | 0.269 | 0.2425 | 0.2527 |
| **Relative absolute error** | 33.79% | 39.03% | 36.27% | 58.44% |
| **Root Relative squared error** | 56.76% | 61.02% | 58.43% | 74.07% |
| **Total number of instances** | 805 | 1610 | 2416 | 3221 |

In MLP-1, the correctly classified and incorrectly classified instances are 87.58% and 12.42% respectively during the training phase. In MLP-2, the correctly classified and incorrectly classified instances are 90.68% and 9.31% respectively during the training phase. In MLP-3, the correctly classified and incorrectly classified instances are 92.14% and 7.86% respectively during the training phase. In MLP-4, the correctly classified and incorrectly classified instances are 92.23% and 7.77% respectively during the training phase.

The performance metrics of the models during the testing phase is demonstrated on Table 6.

**TABLE VI**

Performance metrics of the four models during Testing Phase.

| Performance Metrics | MLP-25% | MLP-50% | MLP-75% | MLP-100% |
|---|---|---|---|---|
| **Correctly classified instances** | 309-89.83% | 593-85.95% | 913-88.29% | 1127-81.66 % |
| **Incorrectly classified instances** | 35-10.17% | 97-14.05% | 121-11.71% | 253-18.34 % |
| **kappa statistics** | 0.7887 | 0.6445 | 0.7674 | 0.4806 |
| **Mean absolute error** | 0.185 | 0.1933 | 0.1797 | 0.2525 |
| **Root mean squared error** | 0.2943 | 0.3049 | 0.3061 | 0.3626 |

In MLP-1, the correctly classified and incorrectly classified instances are 89.83% and 10.17% respectively during the testing phase. In MLP-2, the correctly classified and incorrectly classified instances are 85.95% and 14.05% respectively during the testing phase. In MLP-3, the correctly classified and incorrectly classified instances are 88.29% and 11.71% respectively during the testing phase. In MLP-4, the correctly classified and incorrectly classified instances are 81.66% and 18.34% respectively during the testing phase.

The performance metrics of the models during the testing phase is demonstrated on Table 6.

The models can be compared in terms of various accuracy measures. The comparison of the models based on precision is shown in Fig. 11
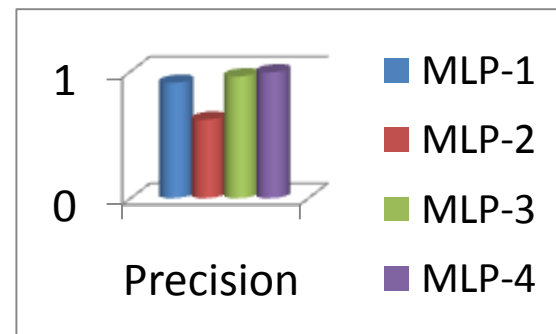


Fig.11: Comparison of MLP-1, MLP-2, MLP-3, and MLP-4 based on Precision

The precision of MLP-1 is 0.910, MLP-2 is 0.619, MLP-3 is 0.960, and MLP-4 is 0.992.

The comparison of the models based on recall is shown in Fig. 12.
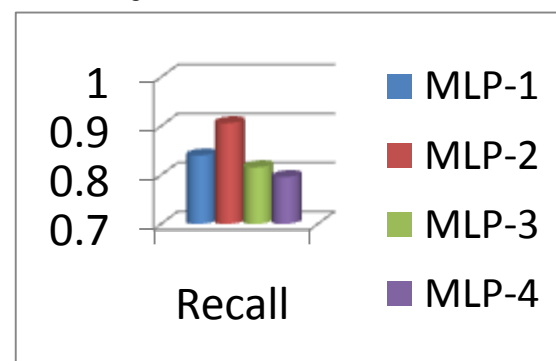


Fig.12: Comparison of MLP-1, MLP-2, MLP-3, and MLP-4 based on Recall

The precision of MLP-1 is 0.840, MLP-2 is 0.906, MLP-3 is 0.816, and MLP-4 is 0.796.

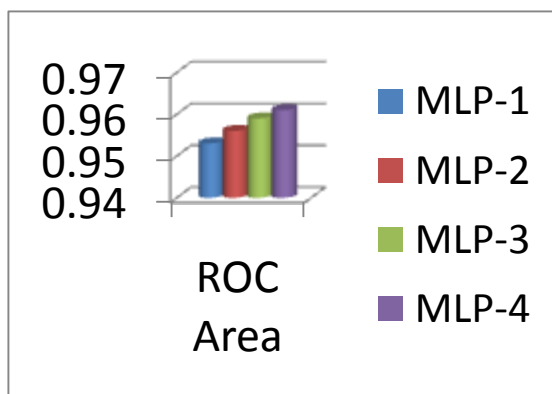The comparison of the models based on ROC Area is shown in Fig. 13.

Fig.13: Comparison of MLP-1, MLP-2, MLP-3, and MLP-4 based on ROC Area

The ROC Area of MLP-1 is 0.953, MLP-2 is 0.956, MLP-3 is 0.959, and MLP-4 is 0.961.

## III. CONCLUSION AND FUTURE WORK

Spam e-mails have become a big problem in the Internet world and it creates several problems for e-mail users. These spam e-mails may come from fake institutions which pretend to be authentic institutions like bank, government organizations or some others financial institutions. In this paper, we have presented several Neural Network techniques to handle the e-mail data, pre-process, and classify them. This model is beneficial to the e-mail data miners. We have tested the model with large number of instances and the outcome of the model is impressive. On an average, all of the four models have 95% of correctly classified instances. The results of the tested spam-based dataset with MLP is comparable in terms of TP Rate, FP rate, Accuracy, Precision, Recall, F-measure, MCC, ROC Area, PRC Area. From the results, it can be determined that all the models outperformed and can be successfully implemented to classify the spam e-mails.

The future work would include large number of e-mail clients. Moreover, one can modify the model to set the parameters like to change the number of hidden layers instead of trial and error settings of the nodes. Additionally, the convergence time of MLP is slow which can be improved by using appropriate optimization techniques for fast convergence. A rigorous research would be carried to detect and prohibit the spam senders.

### REFERENCES

[1] Christina V. et al, "A Study on Email Spam Filtering Techniques", International Journal of Computer Application, Dec 2010, Vol. 12, pp. 7-9.

[2].Priyanka Sao, Pro. Kare Prashanthi, "E-mail Spam Classification Using Naïve Bayesian Classifier", International Journal of Advanced Research in Computer Engineering & Technology (IJARCET), June 201, Vol. 4 Issue 6., pp 2792 - 2796.

[3].Anjali Sharma, Manisha, Dr. Manisha, Dr. Rekha Jain, "Unmasking Spam in Email Messages", International Journal of Advanced Research in Computer and Communication Engineering, Feb. 2015, Vol. 4, Issue 2, pp 35 - 39

[4].Shuang Hao," A Thesis on Early Detection of Spam-related Activity", Georgia Institute of Technology, Dec 2014.

[5]."MAAWG email metrics program: The network operators perspective, report 15." http://www.maawg.org/sites/maawg/files/news/MAAWG_2011_ Q1Q2Q3_Metrics_Report_15.pdf, 2011.

[6]."Email spam record activity" http://www.guardian.co.uk/technology/ 2011/jan/10/email-spam-record-activity, 2011.

[7]."Symantec Internet security threat report" http://www.symantec.com/content/en/us/enterprise/other_resourc es/b istr_main_report_ v18_2012_21291018.en-us.pdf, 2013.

[8].Microsoft releases new threat data on Rustock." http://blogs.microsoft. com/blog/microsoft-releases-new-threat-data-on-rustock/, 2011.

[9] "Email Statistics Report, 2015-2019", The Radicati Group, Inc. A Technology Market Research Firm Palo Alto, CA, USA, http://www.radicati.com.

[10] Sujeet More, Dr S A Kulkarni, " Data Mining With Machine Learning Applied for Email Deception", IEEE Proceedings of International Conference on Optical Imaging Sensor and Security, Coimbatore, Tamil Nadu, India, July 2-3, 2013.

[11] H.S. Hota, et al, "Tuned Artificial Neural Network Model for E-mail Data Classification with Feature Selection", International Journal of Computer Applications, April 2013, Vol. 67, pp.20-25.

[12] Herrero, A., Snasel, et al, "Combined Classifiers with Neural Fuser for Spam Detection". [Advances in Intelligent Systems and Computing] International Joint Conference CISIS12-ICEUTE12-SOCO12 Special Sessions, 2013, Vol.189.

[13] Idris, I, Selamat, A. et al, "A Combined Negative Selection Algorithm-Particle Swarm Optimization for an Email Spam Detection System", Engineering Applications of Artificial Intelligence, 2015, vol.39, pp.33-44.

[14].Simon, D. "Biogeography-Based Optimization",IEEE Transactions on Evolutionary Computation, 2008, vol.12, pp.702-713.

[15] Ali Rodan, et al, "Optimizing Feed forward Neural Networks Using Biogeography Based Optimization for E-Mail Spam Identification", Int. Jl. of Communications, Network and System Sciences, 2016, vol. 9, pp. 19-28.

[16] Ren WangI, Amr M. Youssef , Ahmed K. Elhakeem, "On Some Feature Selection Strategies for Spam Filter Design", The 2006 IEEE Canadian Conference on Electrical and Computer Engineering (CCECE/CCGEI),Ottawa, May 2006, pp. 2186 - 2189.

[17] Qiuming Ma, et al, "Text Spam Neural Network Classification Algorithm", IEEE International Conference on Communications, Circuits and Systems, 2010, vol. 9, pp. 466-469.

[18] Ismaila Idris, "E-mail spam classification with ANN and Negative selection algorithms", International Journal of Computer Science & Communication Networks,2011, Vol.1(3), pp 227-231.

[19] El-Sayed M. El-Alfy et al., "Using GMDH-based networks for improved spam detection and email feature analysis", Applied soft computing, 2011, vol. 11, pp. 477-488.

[20] Harshal Deshmukh, "Spam Mail Detection Using Artificial Neural Network"; Imperial Journal of Interdisciplinary Research (IJIR), 2016, Vol-2, Issue-5, pp.1233-1236.

[21] Masurah Mohamad, Ali Selamat, " An Evaluation on the Efficiency of Hybrid Feature Selection in Spam Email Classification ", IEEE 2015 International Conference on Computer, Communication , and Control Technology(14 ICT 2015), Sarawak, Malaysia, April 21-23 2015, pp. 227 – 231.

[22] L. Firte, C. Lemnaru and R. Potolea, ―Spam Detection Filter Using KNN Algorithm and Resampling‖, in Proceedings of the 2010 IEEE 6th International Conference on Intelligent Computer Communication and Processing, IEEE-2010, pp. 27-33.

[23] A. Nossier, khaled kagati, and Islam Taj-Eddin, ―Intelligent Word- Based Spam Filter Detection Using Multi Neural Networks‖, International Journal of Computer Science, March 2013, Vol. 10, Issues 2, No. 1, 1664- 0784

[24] Reena Sharma, et al, "E-mail Spam Detection Using SVM and RBF", I.J. Modern Education and Computer Science, April 2016, vol.4, pp.57- 63.

[25] Josin Thomas, Nisha S. Raj , Vinod P. "Towards Filtering Spam Mails using Dimensionality Reduction Methods", 5th IEEE International Conference - Confluence The Next Generation Information Technology Summit - 2014, pp. 163 - 168.

[26].Mohammad-Ali Oveis-Gharan, Kaamran Raahemifar, " Multiple Classifications for Detecting Spam email by Novel Consultation Algorithm", The IEEE-2014 Canadian Conference on Electrical and Computer Engineering (IEEE-CCECE), 14-15 May 2014, Toronto, Canada, pp. 1- 5.

[27] Arram, A., Mousa, H. and Zainal, A. "Spam Detection Using Hybrid Artificial Neural Network and Genetic Algorithm", 13th International Conference on Intelligent Systems Design and Applications (ISDA), IEEE- 2013, pp.336-340.