

Comparative Study of Efficient Neural Network Methodology for Text & Image Based Spam Email Filteration

^[1] Dipalee Patil, ^[2] S. R. Ghungrad

^[1] Student, Department of Computer science & Engineering, MSS's College of Engineering & Technology, Jalna, Maharashtra, India

^[2] Faculty, Department of Computer science & Engineering, MSS's College of Engineering & Technology, Jalna, Maharashtra, India

Abstract: Internet users frequently use e-mail for fast data communication of audio, video and textual data but at the same time, they are facing problem due to unwanted e-mail known as spam e-mail. In order to filter this unwanted e-mail, a classifier must be placed in the network or in the computer. Spam e-mail with advertisement text embedded in images presents a great challenge to anti-spam filters. In this paper, we present a fast method to detect image-based spam e-mail. To achieve the objective, Artificial Neural Network is applied for the classification of spam and ham emails. OCR-based modules can be used against image spam, to tolerate the analysis of the semantic content embedded into images.

Keywords: - Spam, Ham, Artificial Neural Network, Image Spam, OCR.

I. INTRODUCTION

Spam E-mail is accentually a junk E-mail send by spammers for their own bonfire intension. Spammers send the e-mail to attract e-mail users, which belongs to various categories like making friend, lottery, adult news, advertisement etc. They collect e-mail addresses from the various sources like website, newsgroup, social website etc. and sends spam E-mail in bulk. Unwanted email can come from anywhere. Current trends say that 95 percent of Internet transmission will be Spam [14] Spam increase the load on the servers and the bandwidth of the ISPs and the added cost to handle this load must be compensated by the customers. In addition, the time spent by people in reading and deleting the spam emails is a waste. Due to this problem, it becomes necessary to distinguish between legitimate and spam emails [8]. Although current anti-spam technologies are quite successful in filtering text, based spam emails [18]. A new trend in email spam is the emergence of image spam. The image spam is substantially more difficult to detect, as they employ a variety of image creation and randomization algorithms. In Image spam the text message is embedded into attached images to defeat the anti-spam filters. Contemporary spam filtering program indulgence spam detection as a text classification problem utilizing machine-learning algorithms such as neural networks and naive Bayesian classifiers to learn spam characteristics. OCR-based modules can be used against image spam, to tolerate the analysis of the semantic content embedded into images. Designing of a software system which can detect the text as well as image based spam and protect user and computer

from the damages like extensive consumption of bandwidth, overload of mail server and wastage of time in detecting spam mails by implementing a decision making model using Artificial Neural Network. A neural network is a set of connected input/output units in which each connection has a weight associated with it. Back propagation is a neural network-learning algorithm. The Back propagation algorithm performs learning on a multilayer feed-forward neural network. During The learning phase, the network learns by adjusting the weights to be able to predict the correct class label of the input tuples.

1.1 The conception and characteristics of spam messages

The so-called junk messages: refers to the messages that the receiver received without their agreement or those that the content is illegal or advertisement which invaded there cleavers' legitimate rights. The three key points are illegality, without the receivers' agreement, invasion of the recipient legitimate rights and interests. Spam messages mainly include the following: advertising information, unhealthy information, false information, winning trick and fraud and so on. Spam messages have the following three characteristics commonly:

- a) Illegal information content
- b) Advertisement
- c) Harassment.

1.2 The causes of the formation of spam messages

There are many causes of the formation of spam messages: Legal defects: So far, there is no a special law or regulations concerning the management of cell phone in our country.

Though we have Regulation of the People's and Internet Information Service Management, in practice, there are some difficulties in dealing with the garbage SMS.

Ineffective supervision: Some unprincipled fellows sell the short message sender and software in the market publicly, which can connect to the computers and cell phones, so the short messages can be sent in a large quantity automatically. There is no effective supervision to such goods. Secondly, short messages interferes with the personification and privacy, the operator can't supervise them effectively, so it becomes the convenient way for the unprincipled men to send a lot of rubbish short messages.

Profit motivation: According to the report of the «Focus» of CCTV, there are usually three ways of short messages sending: internet, cell phone and short messages sender. At least the designer, sending company and operator should be involved in the short messages sending. The companies buy the short messages with the wholesale price from the operator, then sold them with higher price, even if there is only one cent's profit in every short message, sending one million short messages every time means the great profit to the company. According to the second searching report, the public thinks the operator and rubbish short messages sending company are the users who earn more, they occupy about 67.38% and 61.52% respectively, and next profits are the users who profit from the content of the short messages and the persons who sell the cell phone numbers and short messages senders, they occupy about 58.41% and 39.97%.

No identification for the cell phone number: All the numbers are sold by the operator. One method is to buy it with identification or officer identification. The other one is to grant the number if the customers deposit enough money in phone bills. Therefore, some unprincipled fellows use this method to buy many phone cards with fake identification or no identification. Phone cards can be changed or throw away freely at any time, so the real owners of the cards cannot be found by the law enforcement agency and they can easily escape the punishment of the law.

Ineffective supervision to the content of the short Messages: It is impossible to supervise the content of all the messages because of the large quantity of sending and privacy. Therefore, it also provides the good chance for the unprincipled fellows Loopholes in technique and equipment although the experts are devoting on the research of how to recognize and remove the rubbish messages, and some fulfillment has been gotten, this technique is not ripe enough to recognize and remove them and some loopholes still exist in the equipment.

1.3 The harm caused by spam messages

With the rapid development of Email services, spam messages are increasingly rampant, and the contents of spam messages

are related to all aspects of life. The numbers are very large. Variety of phenomena show that the problem of spam messages has been seriously disrupted people's normal work and life, and caused a great response from the community. Spam messages brought a very bad influence on social harmony.

To grow criminal behavior and affect social stability

Some criminals use mobile phone short message to spread rumors, sow discord, incite the people, provoke ethnic hatred and ethnic discrimination, distribute the concept of evil cults and feudal superstition, resulting in ethnic tensions, which cause mass events and influence the stability of the society.

To interfere with normal communication

Spam messages can be send in mass, so the transmission time will take up network bandwidth, causing congestion, affecting the performance of the network and people's normal communication.

II. LITERATURE SURVEY

Harisinghney A. ; Dixit A. ; Gupta S. ; Arora A.[1] has proposed to detect spam based on the textual content of the email, many text-based anti-spam approaches have been proposed, such as Bayesian filters and Support Vector Machine (SVM) filters. However, these approaches soon lost their effectiveness because spammers have introduced a trick to embed junk information into images [2].

Sujeet More and Dr S A Kulkarni [10] focused mainly on spam words for classification. These deceptive mails have already caused many problems such as filling mailboxes, overwhelming important personal mail, wasting bandwidth on network, consuming users' time and energy to sort through it, not to mention all the other problems associated with spam. This method, can be easily implemented, compares amiably with respect to popular algorithms, like Logistic Regression, Neural Network, Naive Bayes and Random Forest using polynomial kernel as filter [9]. Anti-spam filters that are freely or commercial available rely mostly on manually constructed pattern matching rules that need to be tuned based on each user's incoming messages this is a task requiring time and expertise[21].

According to Jian Zhong, YiLu Zhou [12], Wei Deng, spammers can easily employ a variety of image creation and randomization algorithms to make the message fully legible by the human eye and indiscernible by the most anti-spam engines. However, there are always some regions remaining same or similar, otherwise, human eye can't comprehend the message fully. Namely, image-based spams which sent by the same spammer are always self-similarity for human vision at a given period. This method is a plugged-in for

Spam Assassin and combined with detecting rules of Spam Assassin to identify whether an input email is ham or spam. There has been a resurgence of interest in optical character recognition (OCR) in recent years, driven by a number of factors [4]. OCRopus [6] is a new, open source OCR system emphasizing modularity, easy extensibility, and reuse, aimed at both the research community and large scale commercial document conversions. Breuel, Thomas [6] this paper describes the status of the system, its general architecture, as well as the major algorithms currently being used for layout analysis and text line recognition. The author M. Soranamageswari, Dr. C. Mena [7] has concentrated more on measures of statistical feature i.e. color histogram and mean. A classifier is trained on color histogram and mean value of a block of an image, trying to classify spam images from legitimate ones with minimal effort. In this paper, MeghaRathi [8] analyzed the performance of various classifiers with feature selection algorithm and without feature selection algorithm. Initially this experiment with the entire dataset without selecting the features and apply classifiers one by one and check the results.

III. PROPOSED WORK

ANN is comparatively based on the neural structure of the brain. The brain basically learns from experience. It is natural proof that some problems that are beyond the scope of current computers are indeed solvable by small energy efficient packages. This brain modeling also promises a less technical way to develop machine solutions. This new approach to computing also provides a more graceful degradation during system overload than its more traditional counterparts.

A typical neural network is an adaptive system made of four main sections:

- A node as a unit that activates upon receiving incoming signals (inputs);
- Interconnections between nodes;
- An activation function (rule) which transforms inside a node, input into output;
- An optional learning function for managing weights of input-output pairs.

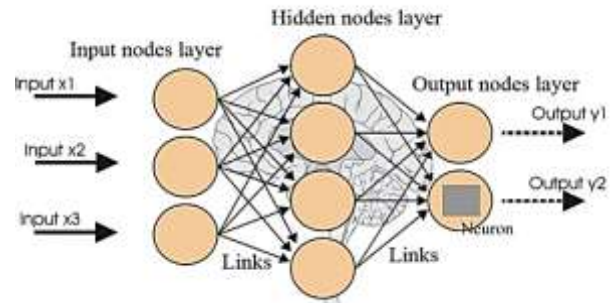


Fig-1: Artificial neural network

The practical execution of the proposed model is as per below flow diagram; where classification takes place based on spam & ham.

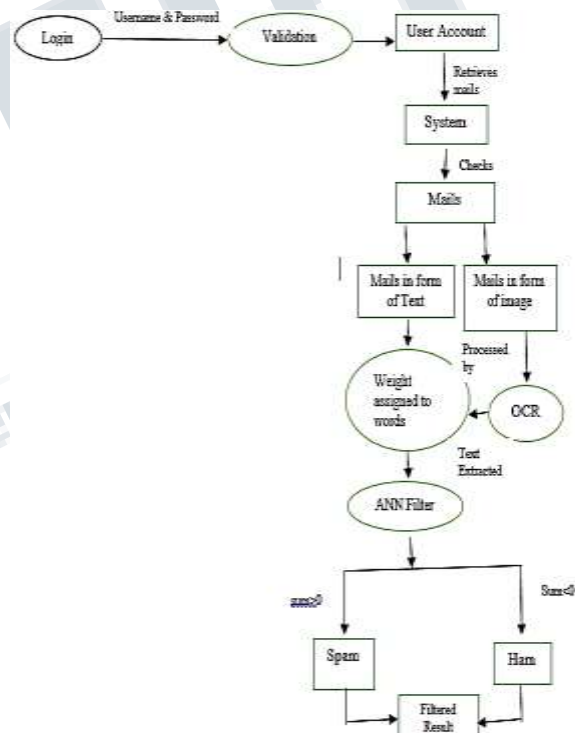


Fig-2: Flow diagram of Proposed System

The proposed system will have the following phases:

3.1 Black Listing and White listing

We need to create such a system where basic filtration done on the basis of domain match

Black Listing: Blacklisting is creating a list of domain names that are used by the spammers, when a mail comes

from that specific domain which is black listed it is considered spam. No further processing is done.

E.g.:

- Marathamarriage.com
- Maxinsurance.com
- Flipkart.in
- Hdfchomeloan.com
- Timesjob.com
- 99acres.com
- Magicbricks.com
- Rishte.com

White Listing: White list is a list of trusted domains and a mail from them is always ham. White listing is a method used to classify user's email addresses as legitimate ones. its listing is not always accurate. Therefore, to counter all these techniques employed by spam filters, spammers now send mails with embedded images containing the spam text.

3.2 Preprocessing of Data

A database of all the words that occur in each mail with the frequency of the word stored in each column will be maintained. Therefore, it is converted to their root form first by applying Porter Stemmer algorithm.

Some steps of this algorithm are:

- Remove the plurals and -ed or -ing suffixes
- Deal with suffixes, -full, -ness etc.
- Take off -ant, -ence etc.

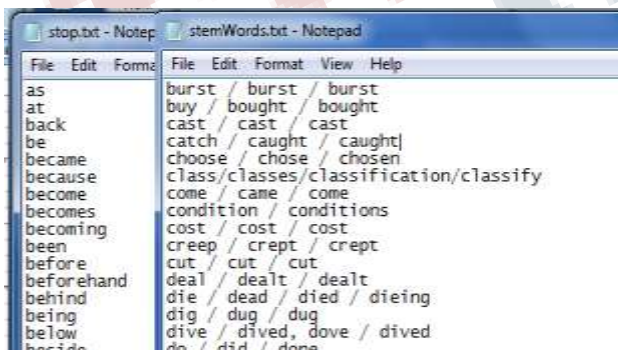


Fig-3: Stop words & Stemmed words

Many words in natural language occur with high frequency but have low information content, such as "a," "an," "the," and most prepositions and conjunctions can be removed with the assumption that no serious loss of information will occur. These so called stop words are specified in a list and can be removed from the token stream.

After the database with the stemmed words, with each mail name in one column and the frequency of occurrence of words in other the system will move on to the next phase.

3.3 Extracting words from Image

To extract the text out of these images is an arduous task. It must be done by sophisticated OCR tools and based on the high level, low level, and combination of both the features of image in a spam mail can be predicted. All those web pages and domains that are notorious for sending spam mails and are not trusted; go on the list of black list. Thus, if a domain that matches from this list, the mail is predicted spam without any further processing. Further, spam is in the eye of the recipient, so a white list is maintained where users can mark those websites they want mails from whether they send "spam" or not. Thus, no processing is done when a white listed domain matches. In OCR processing, the bitmap is analyzed for light and dark areas in order to identify each alphabetic letter or numeric digit. When a character is recognized, it is converted into an ASCII code.

Users have an option of attaching image to their mails. The image is passed through the Google's open source library Tesseract, and words are extracted from it. These words then pass through our different algorithms to predict our mail as spam or ham. Optimum accuracy is achieved for a clear resolution image and more popular fonts like Times New Roman.

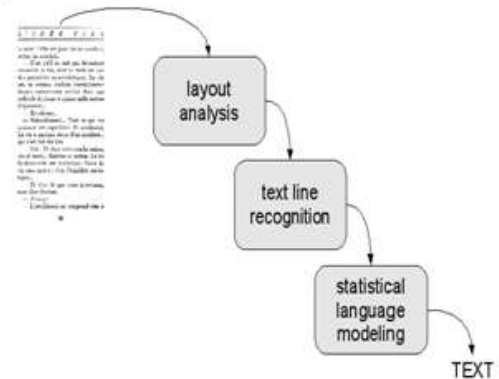


Fig-4: Flow diagram of OCR engine [2]

The overall architecture of the OCRopus OCR system is a strictly feed-forward architecture (no backtracking) with three major components: (physical) layout analysis, text line recognition, and statistical language modeling [6].

3.4 Probability calculation

Term Frequency is the frequent measure of the occurrence of a term in a document. It is possible that a term would appear much more times in long documents than shorter ones, since every document is different in length. Thus, the term frequency is often divided by the document length.

Probability for each word = $\frac{\text{no of occurrence}}{\text{Total No of words}}$

```

Probability of shop = 0.009174311926605505
Probability of save = 0.027522935779816515
Probability of provide = 0.009174311926605505
Probability of wish = 0.009174311926605505
Probability of solicitations = 0.009174311926605505
Probability of e-mail = 0.009174311926605505
Probability of receipt = 0.009174311926605505
Probability of type = 0.009174311926605505
Probability of th = 0.009174311926605505
Probability of cn"please = 0.009174311926605505
Probability of youll = 0.009174311926605505
Probability of free = 0.01834862385321101
Probability of help = 0.009174311926605505
Probability of family = 0.009174311926605505
Probability of kes = 0.009174311926605505
    
```

Fig-5: Probability calculation

3.5 Weight Measure

A weight document is the heart of our system. The fully updated weight document can make this system stronger. The weights are assigned in between 0 to 1 for the spam words and -1 to 0 for the non-spam words.

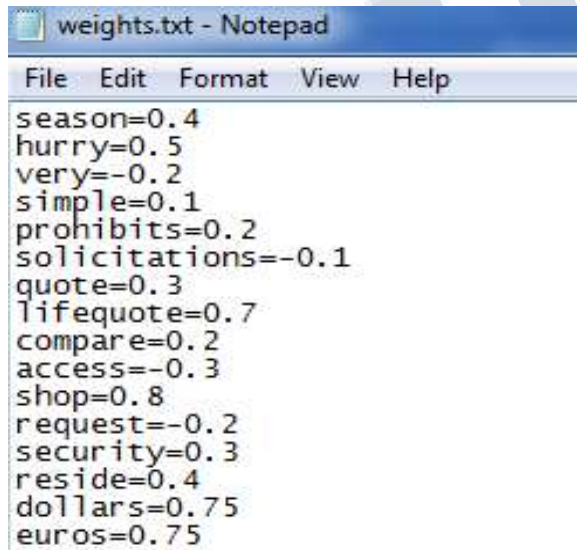


Fig-6: Assigned Weights

Thus, the actual weight of each term will be calculated as:

$$\text{ACTUAL_WEIGHT (T)} = \text{WEIGHT (T)} * \text{PROBABILITY (T)} \tag{1}$$

Where T is the term considered

Sum = \sum actual weight

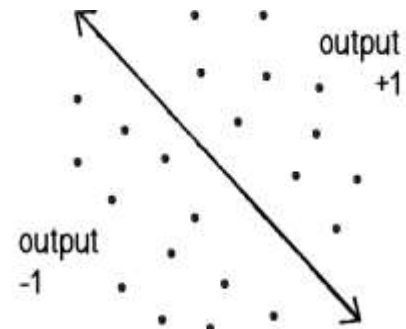


Fig-7: Result output

IV. ALGORITHM APPLIED FOR SPAM MAIL CLASSIFICATION

4.1 K-Nearest neighbor algorithm

The k-nearest neighbor (K-NN) classifier is considered an example-based classifier that means that the training documents are used for comparison rather than an explicit category representation, such as the category profiles used by other classifiers. As such, there is no real training phase. When a new document needs to be categorized, the k most similar documents (neighbors) are found and if a large enough proportion of them have been assigned to a certain category, the new document is also assigned to this category, otherwise not. Additionally, finding the nearest neighbors' can be quickened using traditional indexing methods. To decide whether a message is spam or ham, we look at the class of the messages that are closest to it. The comparison between vectors is a real time process. This is the idea of the k nearest neighbor algorithm:

Stage1. Training

Store the training messages.

Stage2. Filtering

Given a message x, determine its k nearest neighbors' among the messages in the Training set. If there are more spams among these neighbors', classify given Message as spam. Otherwise, classify it as ham.

The K-Nearest Neighbors algorithm is similar to the Nearest Neighbors algorithm, except that it looks at the closest K instances to the unclassified instance. The class then gives the class of then new instance with the highest frequency of those K instances [1]. We are choosing K by trial and error method, for which we obtain the optimal result. The proximity is calculated by finding the Euclidean distance i.e.

$$\text{Euclidean } \sqrt{\sum_{i=1}^k (x_i - y_i)^2} \tag{2}$$

We calculate the proximity of the users mail from our database of mails where $k=20$ [1]. Thus from the majority of the 20 mails, we predict a mail spam or ham. KNN gives a better accuracy than many algorithms, but it has a higher complexity as proximity from each mail is calculated.

4.2 Artificial neural network algorithm

ANN are nonlinear statistical data modeling tool that ties to simulate the functions of biological neural network. It consist of artificial neuron and processes information in the connectionist approach to computation.

The output of a node in an artificial neural network is computed by the equation given below:

$$y_j = \sum_{i=1}^n w_{ij}x_i + \theta_j \quad (3)$$

Where y_j is the value that will be passed to the next layer from node j , n is the number of incoming edges to node j , x is are the input coming from previous layer to node j and θ_j is the bias for node j .

Multi-layer networks uses a variety of learning techniques, the most popular being back-propagation. This technique is also sometimes called backward propagation of errors, because the error is calculated at the output and distributed back through the network layers. Here, the output values are compared with the correct answer to compute the value of some predefined error-function. By various techniques, the error is then fed back through the network. Using this information, the algorithm adjusts the weights of each connection in order to reduce the value of the error function by some small amount. The inputs are provided to the ANN for which there is a prominent answer. Hence, it is required to check whether the network has correct guess or not. If the incorrect guess is made, the network learns from its mistake and adjusts its weights. With the help of sign activation function, the output will be either -1 or 1. The input data will be classified according to the sign of the output. If the sum is, positive it is classified as +1 and the negative sum is classified as -1. The Bias value is taken between 0 and 1. It is fixed to 0.5.

Error = Desired value – Guessed value

Here to adjust the perceptron's weight error is the determining factor. The error is 0 if the perceptron's guessed answer equals the desired answer. The error is -2 if the desired answer is -1 and guessed is +1. Similarly, the error becomes +2 if the desired and guessed answer is +1 and -1 respectively.

V. RESULT & ANALYSIS

Below are the details of result measures, which gives better idea in comparison with each other.

Experiment	Precision	Recall	Accuracy	Time taken (sec)
1	95.33%	69.36%	96.20%	15
2	94.02%	77.86%	95.80%	14
3	93.40%	65.40%	95.73%	13
4	96.97%	70.38%	97.00%	15
5	93.00%	71.55%	96.33%	14
6	94.00%	67.01%	96.22%	15
7	96.57%	70.69%	94.50%	13
8	94.55%	68.33%	95.20%	12
9	96.77%	70.38%	93.88%	12

Table -1 Performance measure calculated for KNN

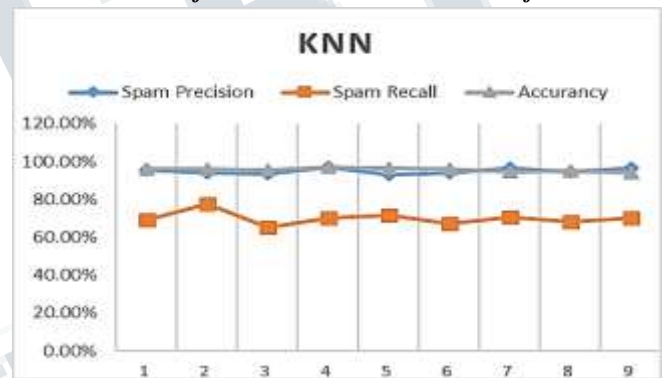


Chart-1 Performance parameter of KNN

Thus above experiments on KNN gives results as per the above chart , in the same parameter we have conducted the experiment using ANN ,where readings shows results as below :

Experiment	Precision	Recall	Accuracy	Time taken(sec)
1	94.24%	78.22%	97.11%	8
2	95.76%	73.44%	98.22%	7
3	93.56%	73.89%	97.78%	6
4	94.44%	77.21%	96.82%	9
5	94.64%	73.78%	97.55%	10
6	98.00%	74.00%	98.20%	9
7	96.20%	76.25%	98.35%	8
8	94.08%	78.22%	98.00%	9
9	98.33%	78.00%	98.80%	9

Table-2 Performance measure calculated for ANN

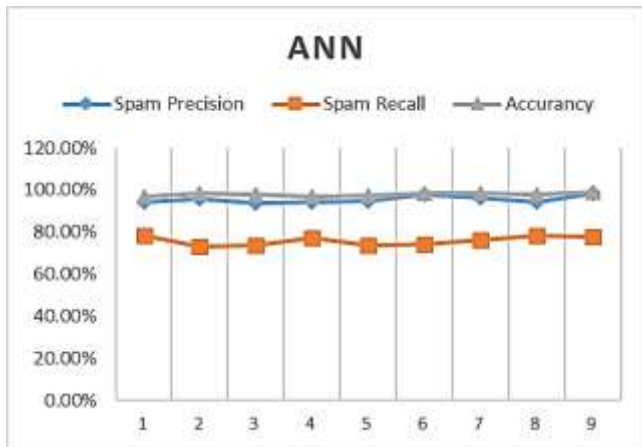


Chart-2 Performance Parameter of ANN

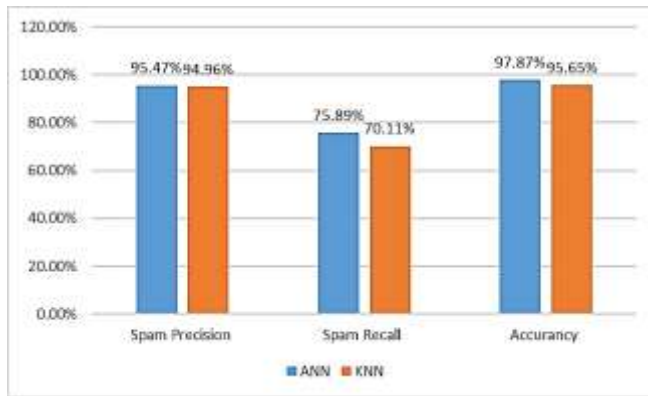


Chart-3 Comparing Average Performance between ANN & KNN

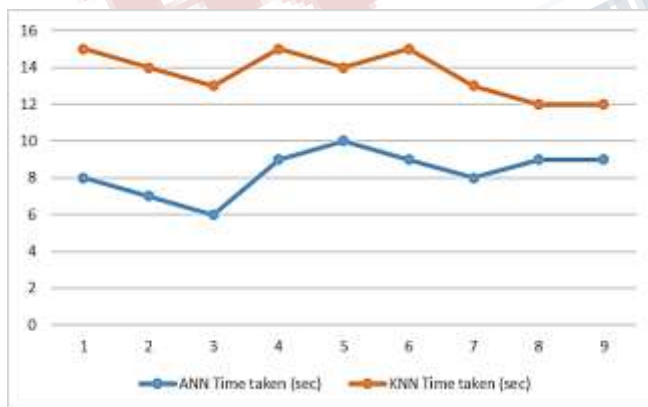


Chart-4 Comparing Average time between KNN & ANN after trying around 100 sets mail on both algorithm we come up with above results where time taken to filter the mail and showing spam & ham is on an average 8 seconds; while with same mail detection in KNN took around average 14 seconds

also number of cases where spam mails detected as ham is on higher side in KNN, i.e. Accuracy is also better in ANN. Via: 95.65% in KNN & 97.87% in ANN.

In this comparison and analysis, we can see that spam precision is almost same in both the algorithms but if we go with more accurate and faster, parameter ANN took the edge. Moreover, since time is the important factor while performing any task we could see the better scope in ANN comparative with previous working with KNN

VI. CONCLUSION

In this internet dependent era users are struggling to fight against spammer so we need to create such a system which will detect / filter/ delete such a spam mails according to threat levels, so self-train mechanism in such a system will definitely going to help. Therefore, we tried to work with self-train ANN mechanism against existing filtering techniques. In this paper we have worked on Artificial Neural Network algorithm for detecting spam mail in comparison with KNN on different parameter via spam precision ,recall , accuracy and time to perform; which work faster with the help of weight age algorithm & self-train mechanism, which also work on text base & image base filtering using different techniques like OCR . There is better scope in identifying such a spam mails more accurately for text as well as multimedia messages;

VII. ACKNOWLEDGEMENT

Author will like to thanks the MSS trust and Principal DR. Biradar , MSSCET for giving me all their support & help. I express my deep sense of gratitude towards Prof. G. P. Chakote, Head of the Department of Computer for his valuable guidance and encouragement. I wish to acknowledge my extreme gratitude to my guide Prof. S. R. Ghungrad for guiding me throughout the work on project. I have been greatly benefited by his valuable suggestion and ideas.

REFERENCES

[1] Harisinghney A. ; Dixit A. ; Gupta S. ; Arora A. "Text and Image based spam email classification using KNN, naïve Bayes and Reverse DBSCAN algorithm" Optimization, Reliability and Information Technology(ICROIT) , 2014 International Conference on DOI:10.1109/ICROIT.2014.6798302, page(s):153-155, 2014

[2] N. Nhung and T. Phuong."An Efficient Method for Filtering Image-Based Spam E-mail". Proc. IEEE

- International Conference on Research, Innovation and Vision for the Future (RIVF07), IEEE Press, Mar. 2007 , pp. 96-102. doi: 10.1109/RIVF.2007.369141.
- [3] Ketari, Lamia Mohammed, Munesh Chandra, and Mohammadi Akheela Khanum. "A Study of Image Spam Filtering Techniques." *Computational Intelligence and Communication Networks (CICN)*, 2012 Fourth International Conference on IEEE, 2012.
- [4] R. Smith, "An Overview of the Tesseract OCR Engine," in *Proc. International Conference on Document Analysis and Recognition*, 2007
- [5] Klimt, Bryan, and Yiming Yang. "The Enron corpus: A new dataset for email classification research." *Machine learning: ECML 2004*. Springer Berlin Heidelberg, 2004. 217-226.
- [6] Breuel, Thomas M. "The OCRopus open source OCR system." *DRR 6815* (2008): 68150.
- [7] M. Soranamageswari, Dr. C. Meena "A Novel Approach towards Image Spam Classification" *International Journal of Computer Theory and Engineering*, Vol.3, No.1, February, 2011, 1793-8201.
- [8] Megha Rathi, Vikas Pareek "Spam Mail Detection through Data Mining – A Comparative Performance Analysis" *I.J. Modern Education and Computer Science*, 2013, 12, 31-39.
- [9] Saab, S.A., Miri, N. Awad M. "Ham or Spam? A comparative study for some content-based classification algorithms for email filtering" DOI:10.1109/MELCON.2014.6820574 , 2014, Pages:339-343
- [10] More S. , Kulkarni, S. A. "Data mining with machine learning applied for email deception" *Optical Imaging Sensor and Security (ICOSS)*, @013 International Conference on DOI:10.1109/ICOISS.2013.668403 , 2013, pages:1-4
- [11] Xiao Mang Li, Ung Mo Kim "A hierarchical Framework for content-based image spam filtering" *Information Science and Digital technology (ICIDT)*, 2012 8th International Conference on Volume: 1, 2012, pages: 149-155
- [12] Jian Zhong, Yilu, Wei Deng "Filtering image-based spam using multifractal analysis and active learning feedback-driven semi-supervised support vector machine" Conference Anthology, IEEE DOI:10.1109/ANTHOLOGY.2013.6784950
- [13] Bouzerdoum, Havstaad, Beghdadi, "Image Quality assessment Using a Neural Network Approach", IEEE, 2004, pp-330-333.
- [14] Brain Whit worth, Ellizbet Whit worth "Spam and the social technical gap", IEEE Computer society, 2004, pp 38-45
- [15] L. Zhang, J. Zhu, and T. Yao, "An evaluation of statistical spam filtering techniques," *ACM Transactions on Asian Language Information Processing*, vol. 3, pp. 243–269, 2004.
- [16] C.-C. Lai and M.-C. Tsai, "An empirical performance comparison of machine learning methods for spam e-mail categorization," in *Fourth International Conference on Hybrid Intelligent Systems (HIS '04)*. USA: IEEE Computer Society, 2004, pp. 44–48.
- [17] L. F. Cranor, B. A. LaMacchia, "Spam!" *Communications of the ACM*, vol. 41, pp. 74-83, 1998.
- [18] C. Pu and S. Webb. 2006. Observed trends in spam construction techniques: A case study of spam evolution. In *Proc. of the 3rd Conf. on Email and Anti-Spam*.
- [19] Email Metrics Report First, Second and Third Quarter 2011. <http://www.maawg.org/lemetricsreport>
- [20] G. Fumera, I. Pillai, & F. Roli. Spam filtering based on the analysis of text information embedded into images, *Journal of Machine Learning Research*, 7, 2006, 2699-2720.
- [21] Z. Wang, W. Josephson, Q. Lv, M. Charikar, K. Li. Filtering image spam with near duplicate detection. *Proc. 4th Conf. on Email and Anti-Spam, USA*, 2007.
- [22] Sabri AT, Adel. Mohammads H, Al-Shargabi B, Hamdeh MA. Developing New Continuous Learning Approach for Spam, Detection using Artificial Neural Network (CLA_ANN). *Eur. J. Sci. Res.* 2010; 42(3): 525-535

BIOGRAPHIES



Dipalee Patil, B.E (Computer Science & Engineering), from PESCOE, Aurangabad She has 3 yrs. experience in NIT College of Engineering as assistant professor, Nagpur. Her area of interest is Software Testing & ERP.