# Unsupervised Learning for Spam Email Filtering

**Sambhangi Chandrahasa[1] , Sai Durga Gade[2] , Sujith[3] , Sai Vineela[4] , Dr.M.S.R Prasad[5]**
[1]Koneru Lakshmaiah Education Foundation, India,  chandrahasasrinivas@gmail.com
[2]Koneru Lakshmaiah Education Foundation, India,  gadesaidurga@gmail.com
[3]Koneru Lakshmaiah Education Foundation, India,  sujithchowdary007@gmail.com
[4]Koneru Lakshmaiah Education Foundation, India,  saivineelachundu@gmail.com
[5]Koneru Lakshmaiah Education Foundation, India,  msrprasad@kluniversity.in

## ABSTRACT

Astounding number of features will have a deleterious effect on some learning classifier's performance, furthermore, the operational time frame during the training process for evaluating the content may be enhanced. A circa-processing period that mostly includes extraction of components and elimination of features in the sector of machine learning consequently plays a significant role in expediting or boosting processing precise classification. The concern addressed throughout this thesis is relevant to the integration of results, prior characterizing machine learning. Feature representation that restores class separability to less dimensionality to detect content. The key benefit with accordance to the envisaged feature representation has always been its rigidity, that mostly empowers data types like those of Random Forest, Assistance Vector Machines, and constraint solver C4.5 to characterize an inbound text as fraud or pseudo-spam at which the component dimension seems to be very limited under a solid truism unaware of the reference source.

**Key words :** Autoencoder, Cosine Similarity, Feature Representation, Machine Learning, Spam Detection, Spam Filtering.

## 1. INTRODUCTION

There are lots of spam emails around the world we can identify some of them are spam and we can't able to identify the rest of spam emails because they look like they have been sent by some professionals or  they are sent by some big organization so spammers are getting advanced and they are not using some words which are mostly found in spam emails. Many phishing mails via emails were advertising of nature. Whether  professional or not, most are not only irritating but also hazardous, as they may comprise references that lead to phishing websites or links that host ransomware or include vulnerabilities as attachments to documents. Spammers acquire email accounts from offices, forums, user files, newsgroups, and viruses that capture position books from users. Spam is also a genre for scammers to scam users into entering personal information on fake websites using falsified emails which really show up to be those of banks or other organizations, such as PayPal. It's named phishing. Targeted phishing is known as scimitar- phishing, where renowned claimant material is being used to generate forged mails. Free mail transfers and accessible intermediary databases, bots typically check and use defenseless outsider systems. SMTP moves mail through one network towards the next— mail servers controlled by ISPs tend to require certain forms of authentication to ensure that somehow the email is that ISP's device. In order to tackle the problems raised by botnets, transparent transactions, and intermediate networks, most email server administrators pre-emptively utilize square efficient IP to impose strict needs on the various servers attempting to send mail. The forward- affirmed reversal in DNS for the active mail server must be effectively set and vast tracts of IP addresses blocked These measures can present issues for those needing to  run a little email server off a residential association.

## 2. LITERATURE SURVEY

### 2.1 Title: An analysis on Spam Mail Filtering Procedure.

This research journal proposes that junk mail  filters can be enforced on all tiers, security software reside at the front of the personal server or at Email Transfer Agent , Email system to include an integrated Extremely pro- Spam and Extremist-Virus framework which offers excellent system edge email security till unwelcome or possibly hazardous messages hit the server. Spam scanners could also be mounted at the email delivery agent level as a privilege to all its subscribers. As shown at email system customers have the customized filters in place which dynamically filter the mail as per the parameters are chosen. So email spam filtering can be done by using this. We found from the analysis that several of the processing approaches are centered on procedures of content identification as there is no strategy that really can promise to include an acceptable solution of Zero cent misdiagnosis and Zero cent adverse reactions.

## 2.2 Title: A Research of Mail Spam Detection Learning Strategies.

Through this publication we are aware of the tactics of spam Some of the potential ways to combat spam would be to upgrade or perhaps even replace the current email propagation regulations with modern, malware-proof versions. The main limitation of the prominently utilized Basic Mail Protocol stack is that it does not offer a good secure framework for ensuring the identity of that same email source. A serious impediment for these kind of measures would be that a substantial number of users voluntarily have to accept a new standard in order to be truly useful. At least one such approach, Source ID, has acquired appropriate prominence and has thus managed to dominate the situation. So there is a query, Why does the email boxes do often full of spam? Scammer reactivity inevitably comes into play, and so is the diverse nature of content reports. Yet one other issue not just to be undervalued it's that we do not customarily guard against unwanted messages in several contexts possible. Throughout other utterances, one implication that server managers and end users should always note is that the extremely pro-spam technology should not only be built and operated, but also decommissioned and would use.

## 2.3 Title: Email spam detection using algorithm for value and graph mining.

The whole report presents a blended strategy for the characterization of fake email leveraging perspective- based email identification paradigm as key methodology critiqued by intelligence benefit computation to boost the effectiveness of spam diagnosis. Extraction function as well as mail identification. Studies show that word filter is remarkably effective in distinguishing scam emails from the monolithic work email cluster. This thesis has shown itself to be plausible to enforce the filtering system throughout the framework – driven system of email labeling. Research operation has affirmed the potential outcome of upgrading the email archiving concept of perspective-based categorisation model. The analysis of the entirely new document theoretical constructs the renowned as well as the ' information-based email evaluation mechanism ' and ' Linger ' will be used by the experiment to tackle the unfulfilled demand. That is to guarantee that content about the kinship of email messages has still been retained in the area of the smaller spatial function This feature reduction technique assured a strengthened detection quality.

## 3. METHODOLOGY

The issue addressed in this thesis is linked to the processing of content, prior to categorizing machine learning. The extraction function and feature collection for component depiction will be seen as an evaluation process for either the possible configurations of the existing feature range, that retains classification of the distinguish ability in memory even more than practicable, with the purpose of trying to identify email documentation as being either junk mail or pseudo-spam with those of the minimum dimensional space.

### 3.1 Distributed memory

Distributed memory is being used to acquire a specified length of a binary variable of each email address. This models retrieve sentence structure and the semantic meaning from the word file.

### 3.2 Bag of words

A bag-of-words paradigm seems to be an interpreting description used throughout the rendering of dynamic languages and in the retention of data. The context is actually represented even in this model as the bag of its words, completely disregarding grammar and even sometimes word order but retaining dizzying array. The layout of the bag-of-words had been used for machine vision too. The bag-of-words template is frequently used in data analysis methods. Where the incidence of every term is used as a training element for a classifier.
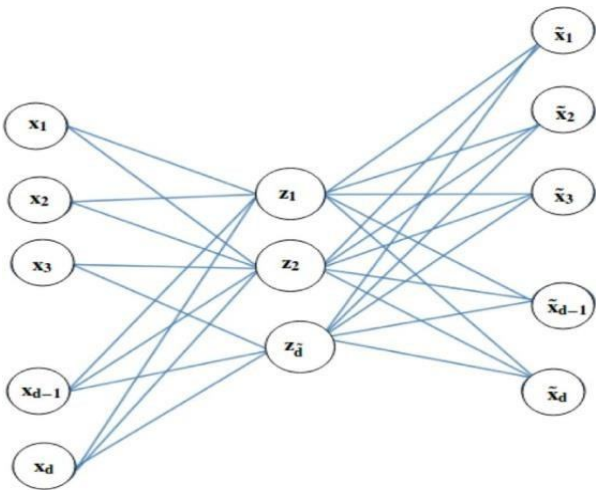
### 3.3 Autoencoder



**Figure 1:** Autoencoder

Autoencoder is an unsupervised training method which condenses a huge amount of space of content into the accompanying low space of functions. Autoencoder converts information into a lower functional area in such a manner that we can recreate data into the initial information. This needs mislabeled information to understand the pattern, because it is unmonitored. The triple of this sequential levels of a plain Autoencoder consists of the inside layer, hidden layer, outside

layer. The inside level operates from the leaving-hand side, the inside level is accompanied by the secret layer to the correct-hand side, and inevitably we have the result layer. Subsequently these forms inside level, secret level and output level obey one another. Each cell is consisting of neurons. The left half is alluded to this as the encoding phase and the other half is referred to as the decoding phase. The center node is known as the bottle neck.

### 3.4 Cosine Similarity

Cosine similarity determines the dimensional variance through several vectors of records. That method may be used in terms of functionality to eliminate the resemblance between two reports. Considering two function parameters, say a(x) and a(y), then the measure of Cosine Similarity is given by:

$$\text{Cos } \theta = \frac{a(x)a(y)}{\|a(x)\|\ \|a(y)\|}$$

From this review, after obtaining its quantitative depiction utilizing Distributed Memory + Bag Of Words for each email a(x) ∈ A, we measure it's own Cosine Similarity gauge against some of A emails. Some steps are used as characteristics. The Cosine Similarity calculation could then be used to generate functionality that guarantee that content-related information on and around the proximity of identical email messages is recorded.

Many other malware identification documentation uses word make a difference-based component techniques and materials to the address realistic which semantics inside the documentation. Utilizing word mean a thing-based materials and techniques like those of Word Factor-Inverse Record Regularity and Bag of Words, nevertheless, we will experience the distinct shortcomings which can include computational complexity curse. Parameters that are not quite familiar are commonly eliminated to ensure adequate efficiency. Furthermore, this avoidance of those less insightful contexts is taken into account by first determining those employing grading techniques like those such as chi-square and information gain. Subsequently, the interpretation established by the functionality starts to lose valuable features for pinpointing phishing emails. The whole research suggests using such a neural network type process to produce in a comprehensive, complex feature representation that will lead to greater identification efficiency in a smaller-dimensional space.

Unsupervised training mechanisms identified distributed bag of words and Distributed Memory are often used in this research to determine a specified- length computational parameter with each message. All such frameworks extract grammatical structure as well as semantic significance from the use of a word document. In addition, To that the cosine similarity is utilized with these email applications to produce a new attribute field employing Cosine Similarity techniques of every email communication. Single invisible level auto encoder will then be used mostly for feature elimination. The feasibility of the implemented methodology to feature representation will be researched through determining the effectiveness of certain coordinated identification mechanisms along the characteristics obtained.

The conclusions in comparison to classifier efficiency demonstrate that perhaps the unsupervised development of the paradigm incorporated leads to the increased evaluation of the effectiveness compared to the conventional feature processing methods and the component preference methods for determining email phishing where element dimensions are reduced. Particularly in comparison to the old conventional feature manufacturing methodologies, the substantial benefit of the suggested feature training strategy will be its tendency to quickly identify scam emails with little or no components throughout the evaluation procedure. This process will definitely be more accurate compared to other inefficient processes. This analysis, by comparison, incorporates unsupervised strategies for developing features to assure the term alignment is often not overlooked and catches the semantic connection between phrases. Some produced attributes can be used to train multiple classifiers because then classifier accuracy can be evaluated.

This publication's suggested solution is now being contrasted with several of the frameworks that have been in use previously. But since the complexity of the function seems to be limited, it thus indicates classifiers could still evaluate the stored data for training and data testing. The concern contemplated throughout this research is relevant to that of the interpretation of information, before identifying machine learning. Feature extraction and component discovery for object recognition can indeed be perceived as a selection process amongst all conceivable combinations of the best original set of features, That retains object separability within the available space as much or more than conceivable, with the ultimate objective of recognizing mails transcripts as being either junk mail or non-spam with just the minimum complexity.

### 4. EXPERIMENTAL SETUP

This segment portrays a thorough analysis of the experimental frame-up of the planned technique and the observational configurations of the conventional methods.

## 4.1 Experimental Design of the Suggested Feature Representation Method

The complete process of determining unsolicited emails among diverse receivers as outlined throughout this analysis incorporates seven main procedures to be followed, which include: Data gathering and Pre-processing, Unsupervised Content Analysis, Information Feature Representation, Feature Conversion, Feature Elimination, Hyperparameter Standardization, and Classification. This step of the pre-processing analysis of all other compounds paradigm-information (e.g. response) from that of the message. The whole semi-processing period involves cutting punctuation marks and void letters, transmission to lowered category of all characters including expressions. For unsupervised function processing as well as information presentation Distributed Memory along with Dynamic Box Of Words , create and practice consistently shared parameters that might encompass paradigm- information of the re-processed electronic mail records.

Throughout the Component Conversion procedure, of the Cosine Similarity calculation will be used to construct attributes to make sure that the knowledge is preserved in terms of features towards the connectedness of equivalent email addresses. Throughout the Feature Removal procedure, autoencoder will be used to minimize the feature memory towards a more resilient feature description. This is just to guarantee that within the of smaller spatial component storage, the details regarding the proximity of email traffic will still be retained. Because email documentation can sometimes be full of noise owing to typographic inaccuracies and phrases which have already been purposefully incorporated to bypass junk mail sensors, the information attribute representation paradigm implemented by Le along with Mikolov that contribute classifiers to terrible performance whilst using loud and annoying input to understand parameters of reports.

## 4.2 Experimental implementation of the supervised feature restriction strategies

This whole segment addresses mechanisms that are used for the detection of unsolicited emails in most other previous research. These are the conventional function depictions would be used for compatibility with the suggested methodology on the Trec07 data as well as Enron data. There seems to be a significant difference here between paradigms with which this segment brings into account especially relative to that of the strategy which has been introduced earlier. Comparatively approaches throughout this segment are monitored for feature limitation. The practices used mostly for diminishing functionality are IG as well as chi-square.

Most of these approaches involve subject definition to determine the prominence of a parameter attribute.

Nevertheless, Auto encoder sometimes doesn't demand much knowledge or understanding regarding the level classification of the learning case, as compared to IG as well as chi-square for feature restriction. Box of Words with differential values including TF-IDF are often used for abstraction of feature throughout this segment of the BoW with frequency parameters. Such approaches contrast from DM+DBOW, because they would not recognize term order. Each phrase is handled fully autonomous.

It is therefore not possible to distinguish these phrases with identical meaning however unique composition of expression. As more of a consequence, we wound up not having certain stop words essential to demarcate certain phishing emails from those of the genuine mails. Almost all of the stop words including zero symbols become eliminated during most of the Pre-processing period. All of the terms as well as phrases in their lowercase letters have been translated. Stemming can also be used to transform most of the phrases to its origin form.
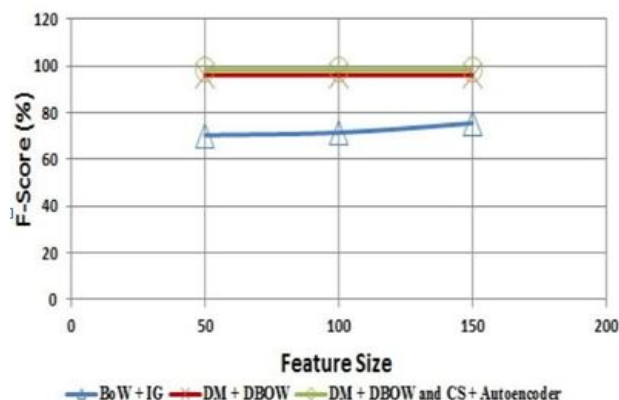
## 5. RESULTS
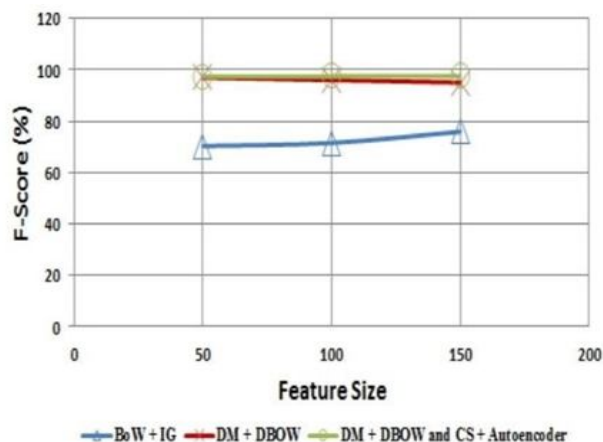


**Figure 2:** SVM Performance on Tree07 dataset.



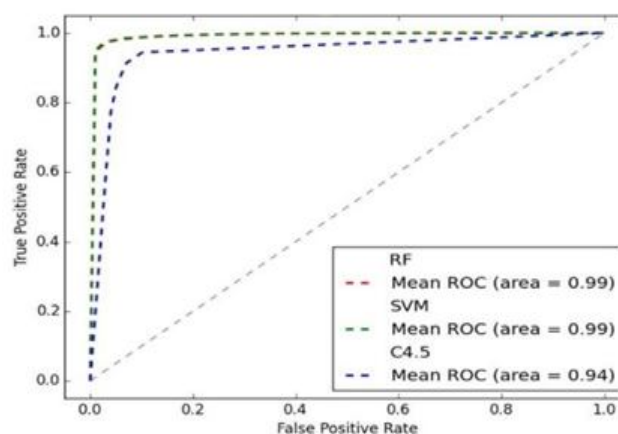**Figure 3:** RF Performance onTree07 dataset.

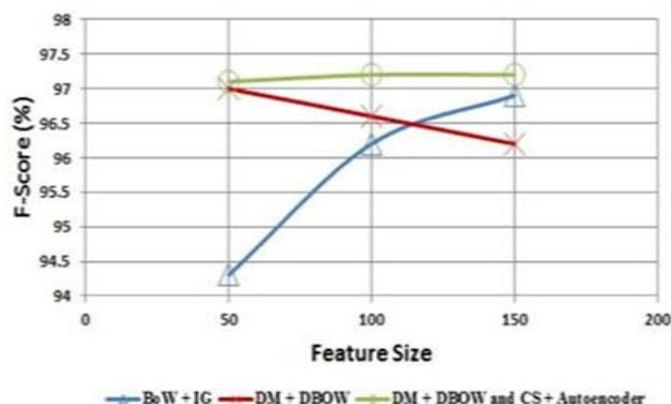**Figure 4:** ROC curve based on enron dataset for performance evolution.



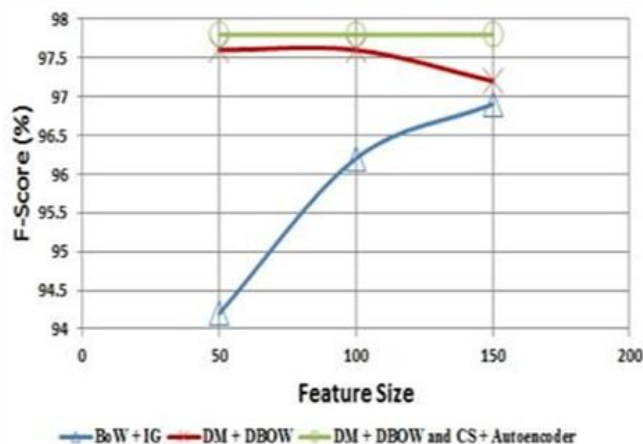**Figure 5:** RF Performance on enron dataset.



**Figure 6:** SVM Performance on enron dataset.

## 6. CONCLUSION AND FUTURE SCOPE

The envisaged solution, distributed memory as well as dispersed phrase bags with cosine similarity and Autoencoder, demonstrated durability based on information

frames that have been properly considered for annual performance review.

Furthermore, additional examination on some other applications is needed to establish that somehow the potential strategy may not be constrained to dealing with problems specific to malware detection solely. Some of the significant obstacles described with respect to that of the potential strategy would be that the estimate of cosine similarity may lead to a massive component range first before using Autoencoder with object elimination, which can often result in storage complication concern. Future experiments would clearly demonstrate how and why the dilemma regarding term parsing and analysis of cosine similarity can only be. addressed. Grouping protocols and arbitrary post-sampling methodology will indeed be investigated about how to use the these for cosine similarity.

## REFERENCES

[1] Bryan Klimt and Yiming Yang. Introducing the enron corpus. In Proceedings of the First Conference on Email and Anti-Spam, CEAS'2004.

[2] Chih-Chin Lai and the Ming-Chi Tsai. An empirical performance for comparison of machine learning methods for spam e-mail categorization. Hybrid Intelligent Systems, pages 44–48, 2004.

[3] Larry Seltzer. Should senders pay for the mess we calle-mail?eWeek, http://www.eweek.com/article2/0,4149,1273186,00.asp, Accessed: 31.05.06, 2003.

[4] Wisaeng, K. "A comparison of different classification techniques for bank direct marketing." International Journal of Soft Computing and Engineering (IJSCE) 3, no. 4 (2013): 116-119.

[5] Chae, M. K., Abeer Alsadoon, P. W. C. Prasad, and Sasikumaran Sreedharan. "Spam filtering email classification using gain and graph mining algorithm." In Anti-Cyber Crimes, 2017 2nd International Conference on, pp. 217-222. IEEE, 2017. https://doi.org/10.1109/Anti-Cybercrime.2017.7905294

[6] Tarek Hassan, Peter Cole and Chun Che Fung "Towards Eradication of SPAM: A Study on Intelligent Adaptive SPAM Filters" Proceedings of the 5th PEECS Symposium, Perth, Western Australia, pp 203-206, September 2004.

[7] Cournane, A., & Hunt, R. 2004. "An Analysis of the tools used for the generation and prevention of spam." Computer & Security, 23 (2), 154166.

[8] David Mertz, "Comparing a Half-Dozen Approaches to the Eliminating Unwanted Email", August.

[9] D. Patil and Y. Dongre, "A Clustering Technique for Email Content Mining," Int. J. Comput. Sci. Inf. Technol., vol. 7, no. 3, pp. 73–79, 2015. https://doi.org/10.5121/ijcsit.2015.7306

[10] J S. Wasi, S. Jami and Z. Shaikh, "Context-based email classification model", Expert Systems, vol. 33, no. 2, pp. 129-144, 2015.
https://doi.org/10.1111/exsy.12136

[11] Face detection and classification based on local binary patterns Y. Rama Devi, M. Sandeep Kumar, C. Nagaraju Volume 3, No.6, November - December 2014 International Journal of Advanced Trends in Computer Science and Engineering

[12] Image enhancement based on the fuzzy logic and thresholding techniques N. Janaki Devi, R.V. Kiran Kumar , M. Sandeep Kumar Volume 3, No.6, November - December 2014 International Journal of Advanced Trends in Computer Science and Engineering.

[13] Context aware framework in IoT: a survey Sagar Sukode1, Prof. Shilpa Gite2, Dr. Himanshu Agrawal3 Volume 4, No.1, January – February 2015 International Journal of Advanced Trends in Computer Science and Engineering.