

EMAIL SPAM DETECTION USING NLP

I JOSEPH KISHORE¹, SIVA RAMA KRISHNA T.²

^{1,2}Department of CSE, JNTUK UCEV, Vizianagaram, India

¹joekishore999@gmail.com, ²t_srkishna@yahoo.com

ABSTRACT: *Email has become the most widely used and economic form of communication in this digital era. Email users generally get bombarded with unsolicited messages regarding direct marketing often sent to multiple users using bots. Users have to spend a considerable amount of time on clearing such messages. Study shows that there is a sharp increase in spam emails, it is estimated that they are almost 89% of the total email traffic. Spam emails can create havoc by causing financial loss or identity theft of users. Spammers use many techniques to bypass manual filters such as misspelled words by adding extra letters to words (eg: amazing, amaze-on etc.), synonyms of generally used words etc.. Use of Machine learning models can handle such data. Creating text classifiers that precisely filter such emails from the user's mailbox to spam folder is more efficient than manual filters.*

Keywords: Spam, Email, NLP, spammers, ML

1. INTRODUCTION

Spammers began to use several tricky methods to overcome the filtering methods like using random sender addresses or appending random characters to the beginning or the end of the message subject line. Knowledge engineering and machine learning are the two general approaches used in e-mail filtering. In the knowledge engineering approach a set of rules has to be specified according to which emails are categorized as spam or ham. A set of such rules should be created either by the user of the filter, or by the software company that provides a particular rule-based spam-filtering tool. By applying this method, no real promising result shows because the rules must be constantly updated and maintained, which is a waste of time and it is not convenient for most users. Machine learning approach is more efficient than knowledge engineering approach; it does not require specifying any rules. Instead, we use a set of training samples, these samples are a set of pre-classified e-mail messages. A specific algorithm is then used to learn the classification rules from these e-mail messages. Machine learning approaches have been widely studied and there are lots of algorithms that can be used in e-mail filtering. They include Naive Bayes, support vector machines, Neural Networks and K-nearest neighbours.

1.1 MOTIVATION

Spam refers to unsolicited business email, otherwise called junk mail or spam floods the Internet client's electronic mailboxes. These junk emails can contain different sorts of messages, for example, commercial advertising, pornography, business promoting, doubtful product, infections or quasi-legal services

1.2 TYPES OF SPAM DETECTION

Fundamentally, spam can be classified into the accompanying four types:

Usernet Spam: User Network is an open get to arrange on the Internet that gives group talks and group email informing. All the data that goes over the Web is called "Net News" and a running accumulation of messages about a specific topic is known as a "newsgroup". Usernet spam is presenting some commercials on the newsgroups. Spammers focus on the clients that read news from these newsgroups. Spammers present promotion on a substantial measure of newsgroups at once. Usernet spam rob clients of the utility of the newsgroups by overwhelming them with a barrage of promoting or other unrelated posts.

Instant Messaging Spam: Instant informing frameworks, for example, Yahoo Messenger, AOL Instant Messenger (AIM), Windows Live Messenger, Facebook Messenger, XMPP, Tencent QQ, Instant Messaging Client (ICQ), and MySpace talk rooms are all objectives for spammers. A few IM frameworks give a registry of clients, including statistical data, for example, date of birth and gender. Advertisers can gather this data, sign on to the framework, and send undesirable messages, which could incorporate business malware, viruses, and associates to paid destinations [8]. As texting has a tendency to not be stuck by firewalls; Shradhanjali, Verma Toran; International Journal of Advance Research, Ideas and Innovations in Technology. subsequently, it is a particularly helpful route for spammers. It focuses on the clients when they join any visiting space to discover new friends. It ruins the appreciation of individuals and wastes their time moreover.

Mobile Phone Spam: Mobile phone spam is focused on the content informing administration of a cell phone. This can be particularly irritating to clients not just for the bother additionally in light of

the cost they might be charged per instant message gotten in a few markets. This sort of spam more often than not contains a few plans and offers on different items. In some cases, service providers likewise make utilization of this to trap the client for activation of some paid services.

Email Spam: Email spam is the most well-known type of spam. Email spam focuses on the individual clients with direct emails. Spammers make a rundown of email clients by inspecting Usenet postings, stealing lists of web mail, and searching the web for email addresses. Email spam costs cash to a client of email in light of the fact that while the client is perusing the messages meter is running. Email spam additionally costs the ISPs on the grounds that when a majority of spam sends are sent to the email clients waste the bandwidth of the service providers these expenses are transmitted to clients. All undesirable emails are not spammed messages.

1.3 APPROACHES OF SPAM DETECTION

There are currently different approaches to spam detection. These approaches include blacklisting, detecting bulk emails, scanning message headings, greylisting, and content-based filtering. Blacklisting is a technique that identifies IP addresses that send large amounts of spam. These IP addresses are added to a Domain Name System-Based Blackhole List and future email from IP addresses on the list are rejected. However, spammers are circumventing these lists by using larger numbers of IP addresses. Detecting bulk emails is another way to filter spam. This method uses the number of recipients to determine if an email is spam or not. However, many legitimate emails can have high traffic volumes. Scanning message headings is a fairly reliable way to detect spam. Programs written by spammers generate headings of emails. Sometimes, these headings have errors that cause them to not fit standard heading regulations. When these headings have errors, it is a sign that the email is probably spam. However, spammers are learning from their errors and making these mistakes less often. Greylisting is a method that involves rejecting the email and sending an error message back to the sender. Spam programs will ignore this and not resend the email, while humans are more likely to resend the email. However, this process is annoying to humans and is not an ideal solution.

1.4 APPLICATIONS OF SPAM DETECTION

Spam Detection is now widely used in many fields such as Hospitals, Banks, Online Messaging platforms and in many business applications. Every place where there is potential of spam being carried out is using Spam

Detection Algorithms to provide security and better facility to their customers.

2. LITERATURE SURVEY

2.1. A Study on Email Spam Filtering Techniques [1]: In their paper, they presented our study on various problems associated with spam and spam filtering methods, techniques. They mentioned that Not only is spam frustrating for most email users, it strains the IT infrastructure of organizations and costs businesses billions of dollars in lost productivity. The necessity of effective spam filters increases.

2.2. Machine Learning Methods For Spam E-MAIL Classification [2]: In their paper they review some of the most popular machine learning methods (Bayesian classification, k-NN, ANNs, SVMs, Artificial immune system and Rough sets) and of their applicability to the problem of spam Email classification. Descriptions of the algorithms are presented, and the comparison of their performance on the SpamAssassin spam corpus is presented.

2.3. Detecting Online Spams through Supervised Learning Techniques [3]: With more customers utilizing the online review surveys to educate their administration basic leadership, assessment of reviews which economically affect the reality of organizations. Obviously, crafty people or gatherings have endeavored to manhandle or control online review spam to make benefits, etc, and that tricky recognition and counterfeit sentiment surveys is a subject of continuous research intrigue. In their paper, they clarify how supervised learning strategies are utilized to recognize online spam review surveys, preceding showing its utility utilizing an informational index of lodging reviews. Keywords- online review surveys, supervised learning, unlabeled data, Naive bayes algorithm, classifiers, EM algorithm, Bag of Words, Stop word Filtering, Support Vector Machine Classifier.

2.4. A Case for Unsupervised-learning-based Spam Filtering [4]: Spam filtering has traditionally relied on extracting spam signatures via supervised learning, i.e., using emails explicitly manually labeled as spam or ham. Such supervised learning is labor-intensive and costly, more importantly cannot adapt to new spamming behavior quickly enough. The fundamental reason for needing labeled training corpus is that the learning, e.g., the process of extracting signatures, is carried out by examining individual emails. In this paper, they discussed the feasibility of unsupervised learning-based spam filtering that can more effectively identify new spamming behavior.

Their study is motivated by three key observations of today's Internet spam: (1) the vast majority of emails are spam, (2) a spam email should always belong to some campaign, (3) spam from the same campaign are generated from some template that obfuscates some parts of the spam, e.g, sensitive terms, leaving other parts unchanged. They presented the design of an online, unsupervised spam learning and detection scheme. The key component of our scheme is a novel text-mining-based campaign identification framework that clusters spam into campaigns and extracts the invariant textual fragments from spam as campaign signatures. While the individual terms in the invariant fragments can also appear in ham, the key insight behind our unsupervised scheme is that their learning algorithm is effective in extracting co-occurrences of terms that are generated by campaign templates and rarely appear in ham. Using large traces containing about 2 million emails from three sources, they show their unsupervised scheme alone achieves a false negative ratio of 3.5% and a false positive ratio of at most 0.4%. These detection accuracies are comparable to those of the de-facto supervised-learning-based filtering systems such as SpamAssassin (SA), suggesting that unsupervised spam filtering holds high promise in battling today's Internet Spam.

2.5. Machine learning for email spam filtering: review, approaches and open research problems

[5]: The upsurge in the volume of unwanted emails called spam has created an intense need for the development of more dependable and robust anti spam filters. Machine learning methods of recent are being used to successfully detect and filter spam emails. We present a systematic review of some of the popular machine learning based email spam filtering approaches. Their review covers surveys of the important concepts, attempts, efficiency, and the research trend in spam filtering. The preliminary discussion in the study background examines the applications of machine learning techniques to the email spam filtering process of the leading internet service providers (ISPs) like Gmail, Yahoo and Outlook emails spam filters. Discussion on the general email spam filtering process, and the various efforts by different researchers in combating spam through the use of machine learning techniques was done. Their review compares the strengths and drawbacks of existing machine learning approaches and the open research problems in spam filtering. they recommended deep learning and deep adversarial learning as the future techniques that can effectively handle the menace of spam emails.

2.6.A Bayesian Approach to Filtering Junk E-Mail [6]:In addressing the growing problem of

junk E-mail on the internet, they examined methods for the automated construction of filters to eliminate such unwanted messages from a user's mail stream. By casting the problem in a decision theoretic framework, they are able to make use of probabilistic learning methods in conjunction with a notion of differential misclassification cost to produce filters which are especially appropriate for the nuances of this task. While this may appear, at first, to be a straight-forward text classification problem, they show that by considering domain-specific features of this problem in addition to the raw text of Email messages, they can produce much more accurate filters. Finally, they show the efficacy of such filters in a real world usage scenario, arguing that this technology is mature enough for deployment.

2.7. Email Spam Detection using Extended KNN algorithm [7]:

E-mail is the cheaper and fast way of communication. E-mail is used in both personal and professional levels of life. Various types of email are lies on social websites. Spam is one of them. Spam is the undesired messages on the internet site which is nothing but wastes the time and resources. Spam messages are sent by the spammer for marketing, promotion, and spreading the virus. Various detection and filtering approaches are used to manage the spam. One of the most useful and simple approaches is the KNN algorithm which is a content based approach. In this paper, the authors are trying to improve the KNN algorithm which can be later used for better Spam Email Detection.

2.8. Spam Detection filter using KNN algorithm and resampling [8]:

Spamming has become a time consuming and expensive problem for which several new directions have been investigated lately. This paper presents a new approach for a spam detection filter. The solution developed is an offline application that uses the k-Nearest Neighbour (KNN) algorithm and a pre-classified email data set for the learning process.

2.9. Improved Bayesian Anti-Spam filter Implementation and Analysis of Independent Spam Corporuses [9]:

Spam emails are causing major resource wastage by unnecessarily flooding the network links. Though many anti-spam solutions have been implemented, the Bayesian spam score approach looks quite promising. A proposal for a spam detection algorithm is presented and its implementation using Java is discussed, along with its performance test results on two independent spam corporuses - Ling-spam and Enron-spam. They used the Bayesian calculation for single keyword sets and multiple keywords sets, along with its keyword contexts to

improve the spam detection and thus to get good accuracy.

2.10. An Integrated approach for Malicious Tweets Detection using NLP [10]: Many previous works have focused on detection of malicious user accounts. Detecting spams or spammers on Twitter has become a recent area of research in social networks. However, they present a method based on two new aspects: the identification of spam-tweets without knowing the previous background of the user; and the other based on analysis of language for detecting spam on twitter in such topics that are trending at that time. Trending topics are the topics of discussion that are popular at that time. This growing micro blogging phenomenon therefore benefits spammers. Their work tries to detect spam tweets based on language tools. They first collected the tweets related to many trending topics, labelling them on the basis of their content which is either malicious or safe. After a labelling process we extracted many features based on the language models using language as a tool. They also evaluate the performance and classify tweets as spam or not spam. Thus, their system can be applied for detecting spam on Twitter, focusing mainly on analysing of tweets instead of the user accounts.

2.11. E-Mail Spam Detection and Classification using SVM and Feature Extraction [11]: Today emails have become a standout amongst the most well-known and efficient types of correspondence for Internet clients. Hence because of its fame, the email will be misused. One such misuse is the posting of unwelcome, undesirable messages known as spam or junk messages. Email spam has different consequences. It diminishes productivity, consumes additional space in mailboxes, additional time, expands programming damaging viruses, and materials that contain conceivably destructive data for Internet clients, destroys the stability of mail servers, and subsequently, clients invest lots of time for sorting approaching mail and erasing undesirable correspondence. So there is a need for spam detection so that its outcomes can be reduced. In this paper, they propose a novel method for email spam detection using SVM and feature extraction which achieves an accuracy of 98% with the test datasets.

3 SPAM DETECTION APPROACHES

3.1 UNSUPERVISED LEARNING

Unsupervised learning is a type of machine learning that looks for previously undetected patterns in a data set with no pre-existing labels and with a minimum of human supervision. In contrast to supervised learning that usually makes use of human-labelled data, unsupervised learning, also

known as self-organization allows for modelling of probability densities over inputs. It forms one of the three main categories of machine learning, along with supervised and reinforcement learning. Semi-supervised learning, a related variant, makes use of supervised and unsupervised techniques.

3.1.1 CLUSTERING

“Clustering” is the process of grouping similar entities together. The goal of this unsupervised machine learning technique is to find similarities in the data point and group similar data points together. Grouping similar entities together help profile the attributes of different groups. In other words, this will give us insight into underlying patterns of different groups. There are many applications of grouping unlabeled data, for example, it can be identified that different groups/segments of customers market each group in a different way to maximize the revenue. Another example is grouping documents together which belong to similar topics etc.

3.2 SUPERVISED LEARNING

Supervised learning is the machine learning task of learning a function that maps an input to an output based on example input-output pairs. It infers a function from labeled training data consisting of a set of training examples. In supervised learning, each example is a pair consisting of an input object (typically a vector) and the desired output value (also called the supervisory signal). A supervised learning algorithm analyzes the training data and produces an inferred function, which can be used for mapping new examples. An optimal scenario will allow for the algorithm to correctly determine the class labels for unseen instances.

3.3 REGRESSION ANALYSIS

Regression analysis is a reliable method of identifying which variables have an impact on a topic of interest. The process of performing a regression allows you to confidently determine which factors matter most, which factors can be ignored, and how these factors influence each other.

3.4 CLASSIFICATION:

Classification is a process of categorizing a given set of data into classes, It can be performed on both structured or unstructured data. The process starts with predicting the class of given data points. The classes are often referred to as target, label or categories. The classification predictive modeling is the task of approximating the mapping function from input variables to discrete output variables. The main goal is to identify which class/category the new data will fall into.

3.5 DECISION TREE

Decision Tree Analysis is a general, predictive modelling tool that has applications spanning a number of different areas. In general, decision trees are constructed via an algorithmic approach that identifies ways to split a data set based on different conditions. It is one of the most widely used and practical methods for supervised learning. Decision Trees are a non-parametric supervised learning method used for both classification and regression tasks. The goal is to create a model that predicts the value of a target variable by learning simple decision rules inferred from the data features. The decision rules are generally in the form of if-then-else statements. The deeper the tree, the more complex the rules and fitter the model. A decision tree is a tree-like graph with nodes representing the place where we pick an attribute and ask a question; edges represent the answers to the question, and the leaves represent the actual output or class label. They are used in non-linear decision making with simple linear decision surfaces. Decision trees classify the examples by sorting them down the tree from the root to some leaf node, with the leaf node providing the classification to the example. Each node in the tree acts as a test case for some attribute, and each edge descending from that node corresponds to one of the possible answers to the test case. This process is recursive in nature and is repeated for every subtree rooted at the new nodes.

3.6 SUPPORT VECTOR MACHINE

In machine learning, support-vector machines (SVMs, also support-vector networks) are supervised learning models with associated learning algorithms that analyze data used for classification and regression analysis. Given a set of training examples, each marked as belonging to one or the other of two categories, an SVM training algorithm builds a model that assigns new examples to one category or the other, making it a non-probabilistic binary linear classifier.

In addition to performing linear classification, SVMs can efficiently perform a non-linear classification using what is called the kernel trick, implicitly mapping their inputs into high-dimensional feature spaces. The objective of the support vector machine algorithm is to find a hyperplane in an N-dimensional space(N — the number of features) that distinctly classifies the data points. We have to find a plane that has the maximum margin, i.e the maximum distance between data points of both classes. Maximizing the margin distance provides some reinforcement so that future data points can be classified with more confidence.

3.7 NAIVE BAYES ALGORITHM

It is a classification technique based on Bayes' Theorem with an assumption of independence

among predictors. In simple terms, a Naive Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature. For example, a fruit may be considered to be an apple if it is red, round, and about 3 inches in diameter. Even if these features depend on each other or upon the existence of the other features, all of these properties independently contribute to the probability that this fruit is an apple and that is why it is known as 'Naive'. Naive Bayes model is easy to build and particularly useful for very large data sets. Along with simplicity, Naive Bayes is known to outperform even highly sophisticated classification methods.

$$P(c|x) = \frac{p(x_1|c)p(x_2|c)p(x_3|c)p(x_4|c)\dots p(x_n|c)p(c)}{p(x)}$$

$$P(c|x) = p(x_1|c)p(x_2|c)p(x_3|c)p(x_4|c)\dots p(x_n|c)p(c)$$

4. PROPOSED METHODOLOGY

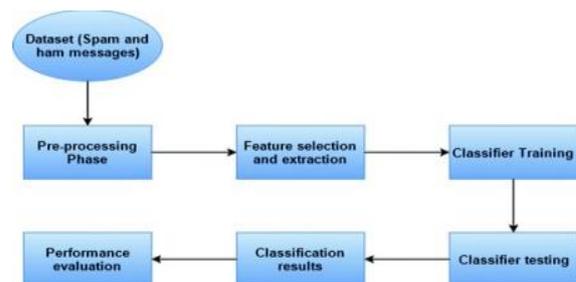


Fig 4.1. Text classification Steps

4.1 DATASET

Two types of datasets are considered for this project. One dataset is a combination of the reported word frequencies that are identified in spam and ham mails in numerical format obtained from UCI Machine Learning Repository and the second dataset is collection of textual mails labelled as spam and ham obtained from kaggle.

4.1.1 NUMERICAL DATASET

The dataset contains nearly 4601 mails with each mail containing 58 attributes The dataset contains nearly 4601 mails with each mail containing 58 attributes The dataset contains nearly 4601 mails with each mail containing 58 attributes each, in which the last attribute comes under label, which is only 0 or 1(0 each, in which the last attribute comes under label, which is each, in which the last attribute comes under label, which represents spam and 1 represents ham). represents spam and 1 represents ham).

Each of the 57 attributes is a word frequency that constitute to defining Each of the 57 attributes is a

word frequency that constitute to defining Each of the 57 attributes is a word frequency that constitute to defining whether mail is spam or ham whether mail is spam or ham

4.1.2 TEXTUAL DATASET

The textual dataset is a collection of mails labelled as spam or ham in the first column. This data is usually unstructured since each mail is of different length and some mails may contain unwanted symbols.

4.2 PRE-PROCESSING PHASE

Pre-processing refers to the transformations applied to our data before feeding it to the algorithm. Data Preprocessing is a technique that is used to convert the raw data into a clean data set. Real-world data is often incomplete, inconsistent, and/or lacking in certain behaviors or trends, and is likely to contain many errors. Data preprocessing is a proven method of resolving such issues. Here, data preprocessing is performed on the data set for removing noise in the data set that we are using here.

4.2.1 REMOVING PUNCTUATIONS AND STOP WORDS

In extracting features from the dataset, we may have stop words such as 'a','an','the','is' and so on in huge quantity which doesn't help in classification so Removing punctuations and stop words is very crucial step to reduce our features dimensionality

4.2.2 SPLITTING INTO TRAINING AND TESTING SETS

In order to classify the data, the dataset must be divided into two parts, test data(used to test) and train data(used as reference to test the test data). The given dataset is divided into test data and train data by taking the split percentage. Usually test data has less percent compared to train data, as the data used for reference should be more. Here in this paper, we are taking a split percentage as 0.25 for test data and 0.75 for train data. We split data into lists, X-train, Y-train, X-test, Y- test.

4.3 FEATURE SELECTION AND EXTRACTION

4.3.1 NUMERICAL DATASET

This dataset is already available with required features as 57 attributes showing particular word frequency corresponding to that mail

4.3.2 TEXTUAL DATASET

This dataset is unstructured so we need to extract features for classification purposes. The feature extraction is nothing but separating distinct and useful words with their frequencies in each mail. This is done by word vectorization and collecting them as a bag of words.

4.4 CLASSIFIER TRAINING

The dataset splitted into train sets is used for this Training phase. We calculate required statistical measures such as mean, standard deviation, or word probabilities etc.. on this training set which is later used for testing purposes.

4.5 CLASSIFIER TESTING

The information obtained from training phases such as mean, standard deviation, word frequencies or probabilities are used here to evaluate test set performance.

4.6 CLASSIFICATION RESULTS

The results obtained from the testing phase are taken into consideration to finalize the results showing confusion matrix.

4.7 PERFORMANCE EVALUATION

This phase is for testing new data that the model hasn't seen previously in its training phase or testing phase.

4.8 Classification:

The paper used KNN and Naive Bayes Algorithms to classify the two datasets considered.

4.8.1 KNN: First assume a K-value. K=7 considered best for this dataset.

4.8.1.1 NUMERICAL DATASET

Training phase:In KNN training phase is nothing but storing values in memory for further purpose. No additional calculations are required.

4.8.1.2 TEXT DATASET

The training phase and testing phase is the same as described in the numerical dataset. What differs is that it uses classifier defined in sklearn neighbors.KNeighborsClassifier

4.8.1.3 THE PROBLEM OF CHOOSING BEST K VALUE

Generally K is chosen at random to find the best k

value In K Nearest Neighbor algorithm. The probability of a mail being ham or spam does not change much by increasing the K-value. General strategy is taking the K value as the square root of test set length. The below figure shows different probabilities associated with different k values in this paper.

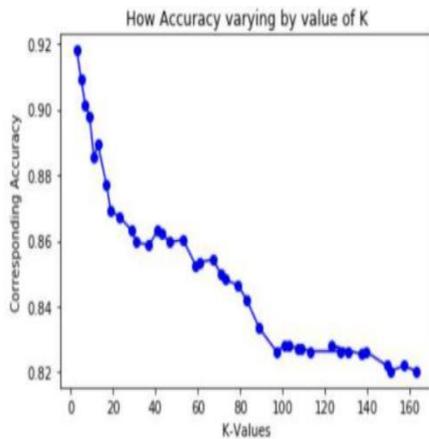


Fig 4.2 How accuracy is varying by changing k value

4.8.2 NAIVE BAYES

4.8.2.1 NUMERICAL DATASET: TRAINING PHASE

In Naive Bayes training phase is summarizing the train set by calculating mean In Naive Bayes training phase is summarizing the train set by calculating mean In Naive Bayes training phase is summarizing the train set by calculating mean and standard deviation of each attribute.

5. CONCLUSIONS

In this paper, the authors gave a new idea to classify the emails as Spam or Ham using the Machine Learning Algorithms. The idea was to use KNN and Naïve Bayes algorithm for E-mail Spam Detection and to improve the terms of parameters like accuracy, precision. This paper,, Email Spam Detection is capable of Text Classification only. Therefore, at this stage this is unable to classify images in the mails and the increase in misspellings in the Email may lead to the decrease in the Accuracy. There is a wide scope of enhancement in this area. Following enhancements can be done: Image Classification can be done on the basis of its contents. Furthermore, Misspellings can be classified on the basis of modules present in Python.

REFERENCES

[1] Christina V, Karpagavalli S and Suganya G “A study on Email Spam Filtering Techniques”.

- [2] W.A. Awad and S.M. ELse ufi “Machine Learning Methods for Spam Email Classification”.
- [3] M.S Minu, Kamagiri Mounika, N.Suhasini and Bezawada Tejaswi “Detecting Online Spams through Supervised Learning Techniques”
- [4] Feng Qian, Y. Charlie Hu and Z. Morley Mao “A Case for Unsupervised-learning-based Spam Filtering”
- [5] Emmanuel Gbenga Dada, Joseph Stephen Bassi , Haruna Chiroma, Shafi'i Muhammad Abdulhamid, Adebayo Olusola Adetunmbi and Opeyemi Emmanuel Ajibuwa “Machine learning for email spam filtering: review, approaches and open research problems”
- [6] Mehran Sahami, Susan Dumais, David Heckerman and Eric Horvitz “Bayesian Approach to Filtering Junk E-Mail”
- [7] Ritu Saini and Er. Geetanjali Chawla “Email Spam Detection using Extended KNN algorithm”
- [8] Loredana Firt, Camelia Lemnar and Rodica Potolea “Spam Detection filter using KNN algorithm and resampling”
- [9] Biju Issac, Wendy Japutra Jap and Jofry Hadi Sutanto “Improved Bayesian Anti-Spam filter Implementation and Analysis on Independent Spam Corporuses”
- [10] Sagar Gharge and Manik Chavan “An Integrated approach for Malicious Tweets Detection using NLP”
- [11] Shradhanjali and Prof. Toran Verma “E-Mail Spam Detection and Classification using SVM and Feature Extraction”
- [12] Hanif Bhuiyan, Akm Ashiquzzaman, Tamanna Islam Juthi , Suzit Biswas and Jinat Ara “A Survey of Existing E-Mail Spam Filtering Methods Considering Machine Learning Techniques”
- [13] Dr. Swapna Borde, Utkarsh M. Agrawal, Viraj S. Bilay and Nilesh M.Dogra “Supervised Machine Learning techniques for Spam Email Detection”
- [14] Konstantin Tretyakov “Machine Learning Techniques in Spam Filtering”
- [15] Harjot Kaur and Er. Prince Verma “Survey on E-MAIL SPAM DETECTION using supervised approach with Feature Selection”