# A Study of Bayesian Classifiers Detecting Gratuitous Email Spamming

Garima Jain
Student, Department of M. Tech (CSE)
DAV Institute of Engg. &Technology, Jalandhar
Punjab, India 144008

## ABSTRACT

Email turns into the real wellspring of correspondence nowadays. Most people on the earth utilize email for their own or expert utilize. Email is a successful, quicker and less expensive method for correspondence. The significance and use for the email is developing step by step. It gives an approach to effortlessly exchange data universally with the assistance of web. Because of it the email spamming is expanding step by step. As indicated by the examination, it is accounted for that a client gets more spam or insignificant sends than ham or pertinent sends. Spam is undesirable, garbage, spontaneous mass message which is accustomed to spreading infection, Trojans, noxious code, notice or to pick up benefit on irrelevant cost. Spam is a noteworthy issue that assaults the presence of electronic sends. Along these lines, it is vital to recognize ham messages from spam messages, numerous techniques have been proposed for arrangement of email as spam or ham messages. Spam channels are the projects which recognize undesirable, spontaneous, garbage messages, for example, spam messages, and counteract them to getting to the clients inbox. The channel grouping procedures are arranged into two either in view of machine learning method or in view of non-machine learning systems. Machine learning methods, for example, Naïve Bayes, Support Vector Machine, Ad boost, and choice tree and so forth though non-machine learning procedures, for example, dark/white rundown, marks, mail header checking and so on. in this paper we survey these procedures for arranging messages into spam or ham ,non- machine learning techniques, such as black/white list, signatures, mail header checking etc. in this paper we review these techniques for classifying emails into spam or ham.

## Keywords
Ham, Spam, Spamming, Spam Filter, Email Spam, Classifier

## 1. INTRODUCTION
Data mining is the process of mining or extracting knowledge from large databases. Data mining is also known as "Knowledge Discovery Process" or "Knowledge mining". There are many other terms which define data mining such as knowledge extraction, knowledge mining from large amounts of data, data analysis. Data mining is applicable on various kinds of data repositories such as data warehouses, relational databases, transactional databases, data streams, flat files and World Wide Web. Data mining is an essential step in the process of discovery of relevant knowledge. The process of knowledge discovery or knowledge extraction is an iterative process. [1]
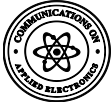
Data are the steps which are used for pre-processing the data, where the data is processed prior to the mining so that and inconsistency, irrelevant or noisy data is removed from the database. This pre-processed data is passed to the data mining algorithms and techniques which produces an output in some forms of patterns. Data mining step interact with the user or a knowledge base. The patterns which are interesting and true are presented to the database user and can be stored as the new knowledge in the knowledge base. Data mining is the essential and most important step in knowledge discovery process because it mines the hidden patterns from the database which is important for the data evaluation and various data analysis task.

Email becomes the major source of communication these days. Most humans on the earth use email for their personal or professional use. Email is an effective, faster and cheaper way of communication. It is expected that the total number of worldwide email accounts is increased from 3.3 billion email accounts in 2012 to over 4.3 billion by the end of year 2016[email statistic report 2012] . Now days, almost every second user in the earth has an email account. The importance and usage for the email is growing day by day. It provides a way to easily transfer information globally with the help of internet.

Spam is an unwanted, junk, unsolicited bulk message which is used to spreading virus, Trojans, malicious code, advertisement or to gain profit on negligible cost. Spams are of many types based on the way of transmission i.e. email spam, social networking spam, web spam, blog or review platform spam, instant message spam, text message spam and comment spam. Spam message can contain text, image, video and also voice data. Spam can be sent via web, fax, telephonic (text messages) [3].

The email spamming is increasing day by day because of effective, fast and cheap way of exchanging information with each other. According to the investigation, it is reported that a user receives more spam or irrelevant mails than ham or relevant mails. About 120 billion of spam mails are sent per day and the cost of sending is approximately zero. According to a spam report of Symantec, the spam rate for December, 2015 [8] was 53.1 percent. Spam not only wastes user time, energy, consumes resources, storage, computation power, bandwidth but also irritates the user with unwanted messages. For example, if you received 100 emails today. Then about approximately 70 emails are spam and only about 30 emails are ham. So, it takes time to identify the ham or important emails from it, which irritated the user. Email user receives hundreds of spam emails per day with a new address or identity and new content which are automatically generated by robot software.

1. **Unsolicited email: -** The email which is not asked for by beneficiary.

2. **Bulk mailing/mass mailing: -** The email which is sent to huge gathering of individuals.

3. **Nameless messages: -** The email in which the address and personality of the sender are covered up.

Spam messages cost billions of dollars every year to the web access supplier on account of the loss of data transmission. Spam messages causes difficult issue for planned client, web access supplier and a whole web spine arrange. One of the cases to clarify it, might be foreswearing of administration where the spammers send mass messages to the server in this manner deferring pertinent email to achieve the proposed beneficiary.[3] Spam is a noteworthy issue that assaults the presence of electronic sends. In this way, it is vital to recognize ham messages from spam messages, numerous techniques have been proposed for arrangement of email as spam or ham messages. Spam channels are the projects which identify undesirable, spontaneous, garbage messages, for example, spam messages, and avoid them to getting to the clients inbox. The channel arrangement strategies are sorted into two sections:

1. **Based on machine learning system.**

2. **Based on non-machine learning systems**.

Machine learning systems, for example, innocent Bayes, bolster vector machine, neural system, and choice tree and so on though non-machine learning procedures, for example, heuristics, dark/white rundown, marks, [12]Mail heading checking and so on. It is observed that characterization in light of machine learning achievement proportion is high when contrasted with order in view of non-machine learning.

The email is grouped into spam or ham by removing highlights from an email. In this manner the email orders depend on two element determination.

1. **Header based components**

2. **Content based elements**

Both the arrangement of components to recognize spam messages have their own particular advantages and disadvantages. Header components can without much of a stretch avoided by the spammers.

## 2. RELATED WORK

Bo Yu a,*, Zong-ben Xu b*(2008) performed "A near study for substance based element spam characterization utilizing four machine learning algorithms". This paper utilizes the accompanying methods Naıve Bayesian; Neural system; Support vector machine; Relevance vector machine it expresses that NN classifier is more delicate to the change of preparing set in light of the fact that the parameters of NN model must be settled on system size and preparing calculation. The exactness of SVM and RVM classifier is higher than NB classifier. Hence, the RVM characterization is more reasonable to the SVM order as far as applications that require low complexity[1].

TiagoA.Almeidaan (2010) performed" Content-Based Spam Filtering" , utilizing Support Vector Machines.However, there are a few types of Naive Bayes channels. They have directed observational trials utilizing understood, huge and open databases. The outcomes expresses that direct SVM, Boolean NB and Basic NB are the best decision for programmed sifting spams. In any case, SVM obtained the best normal execution for all broke down databases exhibiting a precision rate higher than 90% for all tried corpus [2].

Lording Firte Camelia Lemnaru Rodica Potolea(2010)" Spam Detection Filter using KNN Algorithm and Resampling". It approaches for a spam detection filter.The Messages that are classified with the kNN algorithm based on a set of features extracted from the email's properties and content.[3]

RasimM. Alguliev, Ramiz M. Aliguliyev, and Saadat A. Nazirova(2011)" Characterization of Textual E-Mail Spam Using DataMining Techniques" In this paper, the issue of grouping of spam messages gathering is formalized. The basis capacity is a maximum of likeness between messages in type of bunches, which is characterized by k-closest neighbour calculation.

"Order spam messages utilizing content and intelligibility highlights". They reported a novel spam characterization technique that utilizations highlights, in view of email substance dialect and lucidness consolidated with the beforehand utilized substance based assignment highlights. The components are removed from four benchmark datasets, for example, CSDMC2010, Spam Assassin, Ling Spam, and Enron-spam. They clarify every one of these elements. Elements are isolated three classes i.e. conventional components, test elements, and decipherability highlights. The proposed technique can arrange messages in any dialect on the grounds that the elements are dialect free. They utilize five surely understood machine learning calculations to present spam classifier: Random Forest (RF), Bagging, bolster vector machine (SVM), Naïve Bayes (NB). They assess the classifier exhibitions and inferred that Bagging plays out the best out of five. Finally they contrast their proposed strategy with that of numerous state-to-workmanship hostile to spam channels and presumed that their proposed technique can be a good means to classify spam emails. [5]

Anirudh Harisinghaney, Aman Dixit, Saurabh Gupta, and Anuja Arora (2014) performed a work "Content and Image Based Spam Email Classification Using KNN, Naïve Bayes and Reverse DBSCAN Algorithm" The target of their work is to identify message and spam messages. For this reason they utilize Naïve Bayes, K-Nearest Neighbor and another proposed strategy Reverse DBSCAN (Density-based spatial bunching of utilization with commotion). They utilize enron cropus dataset of content and in addition picture. They separate words from picture by utilizing Google's open source library called, Tasseract.They utilize pre-preparing of information. They demonstrate that preprocessing gives 50 percent better precision comes about with all the three calculations than without utilizing pre-preparing. They presumed that credulous bayes with pre-handling gives the best precision among different calculations.

Masurah Mohamad and Ali Selamat (2015) performed a work "An Evaluation on the Efficiency of Hybrid Feature Selection in Spam Email Classification". They exhibit a half and half element choice strategy, in particular The Hybrid Feature Selection, in which they coordinate the unpleasant set hypothesis and term recurrence backwards archive recurrence (TF-IDF) to expand the proficiency result in email channels. They clarify Feature Selection Methods, for example, Information Gain (IG), Gini Index, X2-Statistic, Fuzzy Adaptive Particle Swarm Optimization (FAPSO) and Term

Frequency Inverse Document Frequency (TF-IDF) and Machine Learning Approaches, for example, Naïve Bayes and Rough set hypothesis. They utilize header area and spam practices which are non-content based watchwords. They use dataset comprises of text messages and images. Then they explain their proposed spam filtering framework. In their experimental work they show that rough set theory and TF-IDF were able to work together in order to generate concise and more accurate results. [7]

Izzat Alsmadi and Ikdam Alhami (2015)"Clustering and Classification of Email Contents". In this they clarify different research papers in light of spam location, philosophy characterization on email content and other research objectives. They utilize the information set of general measurement about the email from Google report accommodated Gmail account client. They group the dataset in view of two techniques.1) Classification in light of WordNet class 2) Clustering and Classification assessment. For bunching they utilize K-Means calculation and for grouping they utilize bolster vector machine. Three SVM models are assessed, for example, 1 Best 100 words-VS-email before evacuating stop words, 2. Beat 100 words-VS-emails in the wake of expelling stop words, 3. NGram terms-VS-email. They presumed that the True Positive(TP) rate is appeared to be high for every situation except the False Positive (FP) rate is appeared to be best if there should be an occurrence of NGrambased grouping and order and classification .[8]

Ms.D.Karthika Renuka, Dr.T.Hamsapriya, Mr.M.Raja Chakkaravarthi, Ms.P.Lakshmisurya (2011) performed a work "Spam Classification based on Supervised Learning using Machine Learning Techniques". [9]

Megha Rathi and Vikas Pareek (2013) performed a work "Spam Email Detection through Data Mining-A Comparative Performance Analysis". [10]

Savita Pundalik Teli and Santosh Kumar Biradar (2014) performed a work "Effective Email Classification for Spam and Non-spam" [11]

Rekha and Sandeep Negi (2014) performed a work "A Review on Different Spam Detection Approaches" [12]

# 3. SPAM DETECTION TECHNIQUES
- **Machine Learning Techniques**
- **Non machine Learning Techniques**

There are various spam detection techniques. Out of which some are machine learning. Some of them are defined below:

## 3.1 Machine Learning Techniques
### 3.1.1 Naïve Bayes
A machine learning algorithm, Naive Bayes classifier is based on Baye's theorem of conditioned probability. It is used to recognize an email to be spam or ham. Conditioned Probability is given as

**P (H/X) =P (X/H) P (H) / (P (X).**

Where H denotes hypothesis, X is some evidences, P (H/X) is the probability of given evidence (X) holds by the hypothesis (H). P (X/H) is probability of X conditioned on H. P (H) – prior probability of H, independent on X. There are particularly significant words used in spam emails and ham emails. These words have probability of occurring in both

emails. In advance the filters don't know these probabilities; we must train the filter to build them up. After training the word probabilities are used to compute the probability that an email have that belong to either spam or ham emails. Each particular word or only the most interesting words contribute to email's spam probability. Then, the emails spam probability is computed for every word in the emails. If this total probability exceed over certain threshold then the filters will mark that emails as spam. The Naive Bayesian classifier depends on Bayes hypothesis with autonomy suppositions between indicators. A Naive Bayesian model is anything but difficult to work, with no confused iterative parameter estimation which makes it especially valuable for extensive datasets. Regardless of its straightforwardness, the Naive Bayesian classifier frequently does shockingly well and is broadly utilized in light of the fact that it regularly beats more advanced characterization strategies.

Calculation: Bayes hypothesis gives a method for computing the back likelihood, P(c|x), from P(c), P(x), and P(x|c).Bayes classifier expect that the impact of the estimation of an indicator (x) on a given class (c) is autonomous of the estimations of different indicators. This suspicion is called class restrictive freedom.

$$p(x|C_k) = \prod_{k=1} p(t_k|c_i)$$

### 3.1.2 Gaussian naive Bayes
When dealing with continuous data, a typical assumption that the continuous values associated with each class are distributed according to a Gaussian distribution. For example, suppose the training data contain a continuous attribute, x. We first segment the data by the class, and then compute the mean and variance of x in each class Then, the probability distribution of v given a class c, p(x=v|c)} p(x=v|c), can be computed by plugging v into the equation for a Normal distribution parameterized.

Another common technique for handling continuous values is to use binning to discrete the feature values, to obtain a new set of Bernoulli-distributed features; some literature in fact suggests that this is necessary to apply naive Bayes, but it is not, and the discretization may throw away discriminative information.[4]

$$p(x = v \,|\, c) = 1\backslash\sqrt{(2\pi\sigma\_(c\ )^2)}\, e^\wedge((-(v - \mu\_c\ )^\wedge 2)/(2\sigma\_c^2))$$

In likelihood hypothesis, the typical (or Gaussian) dissemination is an exceptionally regular constant likelihood appropriation. Typical conveyances are essential in measurements and are regularly utilized as a part of the normal and sociologies to speak to genuine esteemed arbitrary factors whose appropriations are not known. [1][2] The typical dissemination is helpful as a result of as far as possible hypothesis. In its most broad shape, under a few conditions (which incorporate limited change), it expresses that midpoints of irregular factors freely drawn from autonomous dispersions merge in appropriation to the ordinary, that is, turn out to be regularly conveyed when the quantity of arbitrary factors is adequately huge. Physical amounts that are relied upon to be the whole of numerous free procedures, (for example, estimation mistakes) frequently have conveyances that are almost normal. [3] Moreover, numerous outcomes and techniques, (for example, spread of vulnerability and minimum squares parameter fitting) can be determined

scientifically in express shape when the important factors are regularly disseminated.

The typical conveyance is in some cases casually called the ringer bend. Nonetheless, numerous different conveyances are ringer formed, (for example, the Cauchy, Student's t, and strategic appropriations). The terms Gaussian capacity and Gaussian ringer bend are additionally equivocal in light of the fact that they once in a while allude to products of the ordinary dissemination that can't be straightforwardly deciphered as far as probabilities.

### 3.1.3 Multinomial naive Bayes

With a multinomial occasion show, tests (highlight vectors) speak to the frequencies with which certain occasions have been produced by a multinomial where p {i} is the likelihood that occasion happens (or K such multinomial's in the multiclass case). An element vector {x = (x_ {1} checking the quantity of times occasion was seen in a specific example. This is the occasion demonstrate normally utilized for report order, with occasions speaking to the event of a word in a solitary archive (words supposition) classifier turns into a straight classifier when communicated in log-space[2]

On the off chance that a given class and highlight esteem never happen together in the preparation information, then the recurrence based likelihood gauge will be zero. This is hazardous in light of the fact that it will wipe out all data in alternate probabilities when they are increased. Hence, it is frequently attractive to join a little example adjustment called without a doubt gauges to such an extent that no likelihood is ever set to be precisely zero. Along these lines of regularizing Bayes is called Laplace .

Rennie et al. examine issues with the multinomial suspicion with regards to record characterization and conceivable approaches to lighten those issues, including the utilization of weights rather than crude term frequencies and report length standardization, to create an innocent Bayes classifier that is aggressive with bolster vector machines.[2] in the field of machine taking in, the objective of factual order is to utilize a question's qualities to recognize which class (or gathering) it has a place with. A straight classifier accomplishes this by settling on an arrangement choice in light of the estimation of a direct mix of the attributes. A question's qualities are otherwise called highlight values and are regularly displayed to the machine in a vector called a component vector. Such classifiers function admirably for down to earth issues, for example, report characterization, and all the more by and large for issues with numerous factors (highlights), achieving exactness levels tantamount to non-direct classifiers while setting aside less opportunity to prepare and use.

The probability of watching a histogram x is given by:

$$p(x|C_k) = \frac{(\sum_i xi)!}{\prod_i x_i^{!}} \prod_i pki^{x_i}$$

### 3.1.4 Bernoulli naive Bayes

The multivariate Bernoulli occasion demonstrate, elements are autonomous Boolean (double factors) portraying inputs. Like the multinomial model, this model is well known for archive order tasks, [9] where paired term event components are utilized as opposed to term frequencies. In the event that x{i}is a Boolean communicating the event or 0nonappearance of the Ith term from the vocabulary, then the probability of a report given a class C_{k} is given by[9] This event model is especially popular for classifying short texts. It has the benefit of explicitly modelling the absence of terms. Note that a naive Bayes classifier with a Bernoulli event model is not the same asa multinomial NB classifier with frequency counts truncated to one.

$$p(x|C_k) = \prod_{i=1}^{n} p_{ki}^{x_i} \ (1 - p_{ki})^{(1-x_i)}$$

## 4. CONCLUSION

For The impact of ensemble hybrid feature ranking method is analyzed on the benchmark classifier, Naïve Bayes.

As we have noticed that naïve classifier is the best far so on using this with "Swarm" hybrid ensemble feature ranking method, the proposed swarm intelligence algorithm can be used to solve intrusion detection as classification problems.

In spite of the way that the expansive autonomy suppositions are regularly erroneous, the innocent Bayes classifier has a few properties that make it shockingly helpful by and by. Specifically, the decoupling of the class contingent element dispersions implies that every appropriation can be freely assessed as a one-dimensional conveyance. This lightens issues coming from the scourge of dimensionality, for example, the requirement for information sets that scale exponentially with the quantity of elements. While gullible Bayes frequently neglects to create a decent gauge for the right class probabilities [12] this may not be a necessity for some applications. For instance, the innocent Bayes classifier will settle on the right MAP choice run arrangement in as much as the right class is more likely than some other class.

## 5. REFERENCES

[1] Yu, Bo and Xu, Zong-ben,"A comparative study for content-based dynamic spam classification using four machine learning algorithms", Elsevier Knowledge-Based Systems, 2008.

[2] Almeida, Tiago A and Yamakami, Akebo,"Content-based spam filtering", IEEE Neural Networks (IJCNN), International Joint Conference, pp.1-7,jul,2010.

[3] Firte, Loredana and Lemnaru, Camelia and Potolea, Rodica," Spam detection filter using KNN algorithm and resampling ",IEEE, 6th International Conference on Intelligent Computer Communication and Processing ,pp.27—33,Romania, Aug. 26-28, 2010.

[4] Rasim, MA and Ramiz, MA and Saadat, AN," Classification of Textual E-mail spam using Data Mining Techniques", the Journal of Applied Computational Intelligence and Soft Computing, JAN 2011.

[5] Rushdi Shams and Robert E. Mercer, "Classification spam emails using text and readability features," *IEEE 13th International Conference on Data Mining*, 2013.

[6] Anirudh Harisinghaney, Aman Dixit, Saurabh Gupta, and Anuja Arora , "Text and image based spam email classification using KNN, Naïve Bayes and reverse DBSCAN Algorithm, " *ICROIT 2014, India,* Feb 6-8 2014.

[7] Masurah Mohamad and Ali Selamat, "An evaluation on the efficiency of hybrid feature selection in spam email classification," *IEEE International Conference on Computer Communication, and Control Technology (14CT 2015)*, April. 2015.

[8] Izzat Alsmadi and Ikdam Alhami, "Clustering and Classification of email contents," *Journal of King Saud University-Computer and Information Sciences*, vol. 27, pp. 46-57, Jan. 2015.

[9] Ms.D.Karthika Renuka, Dr.T.Hamsapriya, Mr.M.Raja Chakkaravarthi, Ms.P.Lakshmisurya, "Spam Classification based on Supervised Learning using Machine Learning Techniques," *IEEE*, 2011.

[10] Megha Rathi and Vikas Pareek, "Spam Email Detection through Data Mining-A Comparative Performance Analysis," *I.J. Modern Education and Computer Science, vol. 12, pp. 31-39,* 2013.

[11] Savita Pundalik Teli and Santosh Kumar Biradar, "Effective Email Classification for Spam and Non-spam," *International Journal of Advanced Research in Computer and software Engineering, vol. 4,* June 6, 2014.

[12] Rekha and Sandeep Negi, "A Review on Different Spam Detection Approaches," *International Journal of Engineering Trends and Technology (IJETT), Vol. 1*, May 6, 2014.