# Spam Detection In Sms Using Machine Learning Through Text Mining

**M.Rubin Julis, S.Alagesan**

**Abstract:**  The development of the cell phone clients has prompted a sensational increment in SMS spam messages. Despite the fact that in many parts of the world, versatile informing channel is right now viewed as "spotless" and trusted, on the complexity ongoing reports obviously show that the volume of cell phone spam is drastically expanding step by step. It is a developing mishap particularly in the Middle East and Asia. SMS spam separating is a similarly late errand to arrangement such an issue. It acquires numerous worries and convenient solutions from SMS spam separating. Anyway it fronts its own specific issues and issues. This paper moves to deal with the undertaking of sifting versatile messages as Ham or Spam for the Indian Users by adding Indian messages to the overall accessible SMS dataset. The paper examinations distinctive machine learning classifiers on vast corpus of SMS messages for individuals.

**Keywords:** Machine Learning, SMS, Spam Detection, Text Mining

————————————◆————————————

## I. INTRODUCTION

 In the most recent decade the consistent development of the spam marvel, to be specific the mass conveyance of spontaneous messages, essentially of business nature, yet in addition with hostile substance, has turned into a principle issue of the SMS benefit for Internet specialist co-ops (ISP), corporate and private clients. Late reviews revealed that more than 60% of all SMS movement is spam. Spam causes SMS frameworks to encounter over-burdens in transfer speed and server stockpiling limit, with an expansion in yearly cost for partnerships of more than several billions of dollars. Furthermore, phishing spam messages are a genuine danger for the security of end clients, since they attempt to persuade them to surrender individual data like passwords and record numbers, using parody messages which are taken on the appearance of originating from trustworthy on-line organizations, for example, budgetary establishments. Despite the fact that it is generally trusted that an adjustment in Internet conventions can be the main successful answer for the spam issue, it is recognized this cannot be accomplished in a brief timeframe. Various types of arrangements have in this manner been proposed up until this point, of conservative, authoritative (for instance the CAN-SPAM act in the U.S.) and innovative nature. The last specifically comprises of the utilization of programming channels introduced at ISP email servers or on the customer side, whose point is to distinguish and naturally erase, or to fittingly deal with, spam messages. Server-side spam channels are considered to be important to lighten the spam issue (Geer, 2004; Holmes, 2005), notwithstanding their disadvantages: for example they can prompt erase true blue messages mistakenly named as spam, and don't take out transfer speed over-burden since they work at the recipient side. At first, hostile to spam channels were basically in  view of catchphrase discovery in email's subject and body. Be that as it may, spammers efficiently acquaint changes with the qualities of their messages to dodge  channels,

which thus push the development of spam channels towards more mind boggling techniques. Traps utilized by spammers can be subdivided into two classes. At the vehicle level, they misuse vulnerabilities of mail servers (like open transfers) to stay away from sender distinguishing proof, and include counterfeit data or  blunders in headers. At the substance level, spammers utilize content darkening procedures to stay away from programmed discovery of average spam catchphrases, for instance by incorrect spelling words and inserting HTML labels inside words. At present, spam channels are comprised of various modules which dissect diverse highlights of messages (to be specific sender address, header, content, and so on.

## II. RELATED WORK

Spam refer to the term, which is related to undesired content with low quality information, [1].Spam referred to the major drawback of mobile business. When comes to the spam detection in campus network[3]they done the analysis using Incremental Learning. For Collecting Spam detection on web pages [4].Moreover Sending out an Spam messages was also analyzed under [5].Data Collection was done privately by a limited company. From the data Collection. There also antispam filter system was evolved. Many parallel and distributed computing system has also processed this spam system. Machine learning algorithm provides accurate result. Text Mining analysis done separates ham and spam separately.

## III. METHODOLOGY

*A. Text Mining*

 Text mining, conjointly spoken as text data processing, roughly corresponding to text analytics, is that the method of account high- quality data from text.High-quality data is often derived through the fashioning of patterns and trends through suggests that like applied mathematics pattern learning. Text mining sometimes involves the method of structuring the input text (usually parsing, at the side of the addition of some derived linguistic options and also the removal of others, and sequent insertion into a database), derivation patterns among the structured information, and at last analysis and interpretation of the output.'High quality' in text mining sometimes refers to some combination of connection, novelty, and powerfulness.Typical text mining tasks include text categorization, text clustering, concept/entity extraction,

————————————————

- **·** [1]*Assistant Professor, Department of Information Technology, Chennai, Tamilnadu, INDIA*
- **·** *(Email id:rubinjulis@aalimec.ac.in Whatsapp no 6369280529)*
- **·** [2]*Assistant Professor, Department of Information Technology, Chennai,Tamilnadu, INDIA*
- **·** *(Email id:s.alagesan@aalimec.ac.in whatsapp no 9841712894 )*

production of granular taxonomies, sentiment analysis, document summarization, and entity relation modeling (i.e., learning relations between named entities).

Text analysis involves info retrieval, lexical analysis to review word frequency distributions, pattern, info extraction, data processing techniques as well as link and association analysis, visual image, and prognostic analytics.The overarching goal is, basically, to show text into knowledge for analysis, via application of linguistic communication process (NLP) and analytical strategies. In brief Steps concerned in Text Mining

1. Gathering unstructured data from multiple data sources like plain text, web pages, pdf  files,  emails, and blogs, to name a few.

2. Detect and remove anomalies from data by conducting pre-processing and cleansing operations. Data cleansing allows you to   extract and retain the valuable information hidden within the data and to help identify the roots of specific words.

3. Convert all the relevant information extracted from unstructured data into structured formats.

4. Analyse the patterns within the data via Management Information System (MIS).

5. Store all the valuable information into a secure database to drive trend analysis and enhance the decision-making process of the organisation.

*B.*  System Architecture

The main objective of our approach is to classify the spam SMS messages as soon as it received on the mobile phone, regardless of newly created spam message (zero-hour attack). In this, we firstly collected dataset and finalized the features for our experiment. After finalizing features, we extracted the features from the messages (ham and spam) to create a feature vector.

Figure1 shows the system architecture of our proposed approach. In the coaching section, a binary classifier is generated by applying the feature vectors of spam and ham messages. In the testing section, the classifier determines whether or not a replacement message may be a spam or not. At the end we get classification results for different machine learning algorithms and performance is evaluated for each machine learning algorithm such that we can get the best algorithm for our proposed approach. Feature choice may be a vital task for the SMS Spam filtering. Selected options ought to be correlate to the message sort specified accuracy for detection of spam message are often enlarged. There is a length limit for SMS message and it contains solely text (i.e.no file attachments, graphics, etc.) while in the email, there is no text limit and it contains attachments, graphics, etc. SMS message is usually of two types i.e. ham (legitimate) message and spam message.
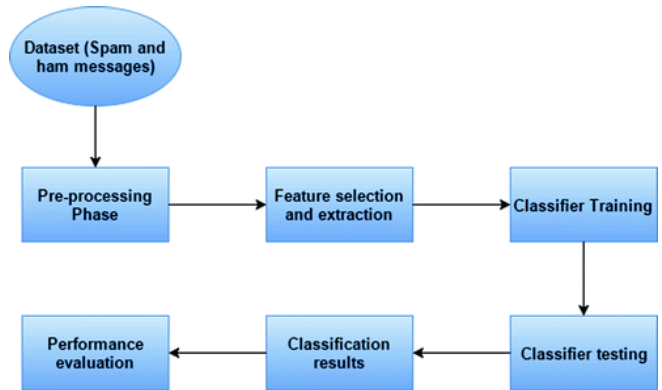


*Fig 1*.System Architecture

Identification of fine feature that may expeditiously filter spam.SMS messages could be a difficult task.Moreover, we have a tendency to study the characteristics of spam messages exhaustive and notice some options, that area   unit helpful within the economical detection of spam SMS.

*C.*   Proposed Method of Text Mining

The strategy of the proposed work is introduced as a matter of first importance the SMS messages are recovered from the information sources and gathered together  for  examination.  Content  pre-handling procedures are connected on these gathered SMS information to make it decipherable for the investigation. Subsequent to halting, another pre- preparing system, in particular, stemming, is performed. The reason for stemming is to change over the words to their root shapes with the goal that every one of the events of the words are introduced consistently in the content accumulation for consummate examination. For case the words 'works', 'working' and 'worked' will be changed over to the word 'work' to its root significance to make beyond any doubt that the word is viewed as consistently as a novel word for printed examination. Aside from the ceasing and stemming an exceptional pre-handling strategy to be specific, homoglyphing, is connected on the SMS content accumulation, which is for the most part not required for other printed information. Most SMS informing clients utilize homoglyphs as often as possible in SMS informing. Homoglyphs are framed by supplanting specific characters with comparable appearance characters in a word, for instance, the digit zero and the capital letter 'O' (i.e. '0' and 'O'), or the digit '1', the lowercase letter 'l' and the capitalized 'I' (i.e. '1', 'l' and 'I'), seem comparative when perusing. Homoglyphs are the immense wellspring of disarray since it is troublesome for PC projects to separate the homoglyph words in spite of the fact that they are intended to be a similar word in a similar setting. Since the homoglyphs are much of the time utilized in SMS informing, pre-handling is required with a specific end goal to secret the homoglyphs to their genuine frame for better investigation of SMS content gathering. SMS messages are shaped in an etymological style with the assistance of shortened forms for speedier informing, which is named SMS dialect or text ease.

*D.*    Dataset Collection

The SMS dataset that we've got used for our experiment contains 2608 messages out of that 2408 collected from SMS Spam Corpus in public on the market and two hundred collected manually that consists of 25 spam messages and 175 ham messages. The SMS Spam Corpus v.0.1 consists of following 2 sets of messages:
a. SMS Spam Corpus v.0.1 Small - It consists of 1002 ham messages and 82 spam messages. This corpus is helpful and has been employed in the analysis.
b. SMS Spam Corpus v.0.1 Big - It consists of 1002 ham messages and 322 spam messages. This corpus is helpful and has been employed by the analyzers in their research work.
Kaggle (Dataset Collection Supporters) supports a variety of dataset publication formats, but we strongly encourage dataset publishers to share their data in an accessible, non-proprietary format if possible. Not only are open, accessible data formats better supported on the platform, they are also easier to work with for more people regardless of their tools. Supported File Types:

**CSVs**
The simplest and best-supported file type available on Kaggle is the "Comma-Separated List", or CSV, for tabular data. CSVs uploaded to Kaggle should have a header row consisting of human-readable field names. A CSV representation of a shopping list with a header row, for example, looks like this:
1. id,type,quantity
2. 0,bananas,12
3. 1,apples,7
CSVs are the most common of the file formats available on Kaggle and are the best choice for tabular data. On the Data tab of a dataset, a preview of the file's contents is visible in the data explorer. This makes it significantly easier to understand the contents of a dataset, as it eliminates the need to open the data in a Kernel or download it locally.CSV files will also have associated column descriptions and column metadata. The column descriptions allows you to assign descriptions to individual columns of the dataset, making it easier for users to understand what each column means. Column metrics, meanwhile, present high-level metrics about individual columns in a graphic format.

# IV.  ALGORITHMS
In machine learning, tasks square measure usually classified into broad classes. These classes square measure supported however learning is received or however feedback on the training is given to the system developed. Two of the most widely adopted machine learning methods are supervised learning which trains algorithms based on example input and output data that is labelled by humans, and unsupervised learning which provides the algorithmic rule with no tagged information so as to permit it to search out structure inside its input file.

**Approaches:**
As a field, machine learning is closely associated with machine statistics, so having a background knowledge in statistics is useful for understanding and leveraging machine learning algorithms. For those who may not have studied statistics, it can be helpful to first define correlation and regression, as they are commonly used techniques for investigating the relationship among quantitative variables.
Correlation may be a live of association between 2 variables that aren't selected as either dependent or freelance.
Regression at a basic level is employed to look at the link between one dependent and one variable quantity. Because regression statistics is wont to anticipate the variable quantity once the variable quantity is thought, regression enables prediction capabilities. Approaches to machine learning are continuously being developed. For our functions, we'll bear many of the favored approaches that square measure getting used in machine learning at the time of writing.

*A.*    Logistic Regression

In this machine learning algorithm the dependent variable is categorical and measures the relationship between the independent variable and categorical dependent variable using the logistic function.

B. Knearest Neighbours
In K means that formula, for each test data point, we would be looking at the K nearest training data points and take the most frequently occurring classes and assign thatclass to the test data. Therefore, K represents the quantity of coaching purpose|datum|information is lying in proximity to the take a look at information point that we have a tendency to ar reaching to use to seek out the category.

C. NaiveBayes Classifier
A Naive Thomas Bayes Classifier could be a supervised machine- learning formula that uses the Bayes' Theorem, which assumes that features are statistically independent. The theorem depends on the naive assumption that input variables are freelance of every alternative, i.e. there is no thanks to understand something concerning alternative variables once given a further variable. Regardless of this assumption, it has proven itself to be a classifier with good results.
Naive Bayes Classifiers rely on the Bayes' Theorem, which is based on conditional probability or in simple terms, the likelihood that an event (A) will happen given that another event (B) has already happened. Essentially, the theorem allows a hypothesis to be updated each time new evidence is introduced. The equation below expresses Bayes' Theorem within the language of probability:

$$P(A｜B) = P(B｜A)\ P(A)$$

$$P(B) \hspace{3cm} (1)$$

Let's explain what each of these terms means

1.    "P" is the symbol to denote probability.
2.    P(A | B) = The probability of event A (hypothesis) occurring given that B (evidence) has occurred.
3.    P(B | A) = The probability of the event B (evidence) occurring given that A (hypothesis) has occurred.
4.    P(A) = The probability of event B (hypothesis)

500

occurring.

5.        P(B) = The probability of event A (evidence) occurring.

### B. Support Vector Machines

The objective of the support vector machine algorithmic rule is to search out a hyperplane in associate degree N-dimensional space(N — the quantity of features) that clearly classifies the information points.
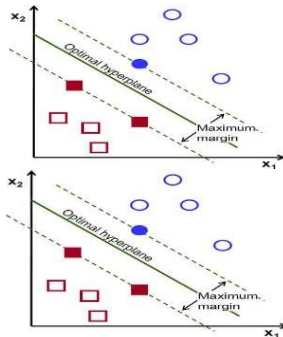


***Fig 2.****Comparison Based on Accuracy*

To separate the 2 categories of knowledge points, there area unit several attainable hyper planes that might be chosen. Our objective is to search out a plane that has the most margin, i.e the most distance between knowledge points of each categories. Maximizing the margin distance provides some reinforcement so future knowledge points will be classified with additional confidence.

### C.Decision Tree Classifier

The classification technique could be a systematic approach to create classification models from Associate in Nursing input audiotape set. as an example, call tree classifiers, rule-based classifiers, neural networks, support vector machines, and naive Bayes classifiers area unit totally different technique to unravel a classification downside. every technique adopts a learning rule to spot a model that most closely fits the connection between the attribute set and sophistication label of the computer file. Therefore, a key objective of the educational rule is to create prophetic model that accurately predict the category labels of antecedently unknown records. Decision Tree Classifier could be a easy and wide used classification technique. It applies a simple plan to unravel the classification downside. call Tree Classifier poses a series of fastidiously crafted questions on the attributes of the take a look at record. whenever time it receive a solution, a follow-up question is asked till a conclusion regarding the category label of the record is reached.

## V. EVALUATION METRICS

In order to judge the effectiveness of our planned approach, we are going to contemplate eight attainable outcomes i.e. true positive rate, false positive rate, true negative rate, false negative rate, f1 score, accuracy, precision, and recall. These area unit the quality metrics to gauge any spam detection system. These analysis metrics area unit represented in short as follows:- True Positive Rate (TP) - It denotes the percentage of spam messages that were accurately classified by the machine learning algorithm. If we denote spam messages as S and spam messages that were accurately categorized as P, then

TP=P/S.

True Negative Rate (TN) - It denotes the percentage of ham messages that were accurately categorized as ham messages by the machine learning algorithm. If we denote ham message as H and ham messages that were accurately categorized as        ham by Q, then

TN=Q/H.

3.False Positive Rate (FP) - It denotes the percentage of ham messages that were wrongly categorized as spam by the machine learning algorithm. If we denote ham messages as H and ham messages that were wrongly classified as spam by R, then

FP=R/H.

4.False Negative Rate (FN) - It denotes the percentage of spam messages that were incorrectly classified as ham message by the machine learning algorithm.

FN=T/S.

5.Precision- It denotes the percentage of messages that were spam and actually classified as spam by the classification algorithm. It shows the exact correctness. It is given as:-

Precision=TP/TP+FP.

6.Recall - It denotes the percentage of messages that were spam and classified as spam. It shows the completeness. It is given as:-

Recall=TP/TP+FN.

7.F-measure - It is defined as the harmonic mean of Precision and Recall. It is given as:-
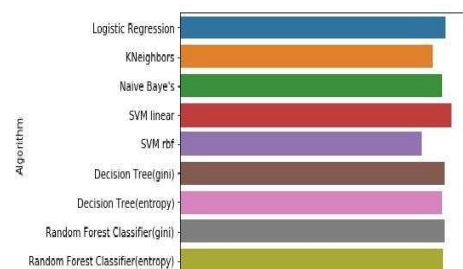F-measure=2*Precision*Recall/Precision+Recall.

8.Receiver Operating Characteristics (ROC) area - In this an area is plotted between True Positive Rate and False Positive Rate for different threshold values.

Here The metrics are as follows:

I.      Accuracy
II.     Time to train the Model
III.    Time To Predict
IV.     Accuracy-Time Ratio( Train)
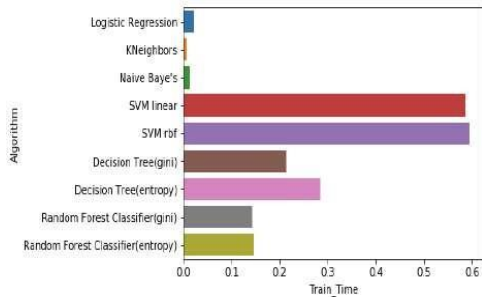V.      Accuracy-Time Ratio( Test)

### A. Accuracy

*B*. Training set



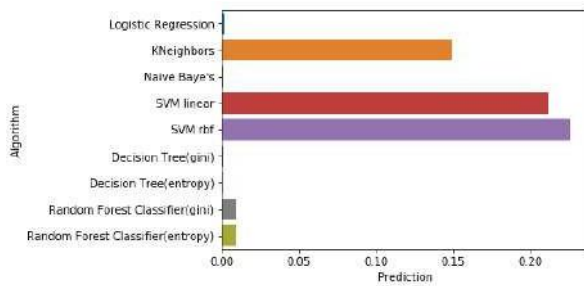**Fig 3.** *Comparison Based on Training Set*

C.Prediction Time



**Fig 4 Comparison** *Based on Prediction Time*

# VI RESULT REPORTS

|  | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Ham | 0.96 | 1.00 | 0.98 | 1355 |
| Spam | 0.98 | 0.70 | 0.82 | 196 |
| Avg/Total | 0.96 | 0.96 | 0.96 | 1551 |

*A.*    Linear Regression

**Table 1.***Classification Report of Linear Regression*

*B*. K Neighbors Classifiers

|  | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Ham | 0.91 | 1.00 | 0.95 | 1355 |
| Spam | 1.00 | 0.31 | 0.47 | 196 |
| Avg/Total | 0.92 | 0.91 | 0.89 | 1551 |

**Table 2.***Classification Report of K Neighbours*

*C.* Classification Report of Naïve Bayes Classsifiers

|  | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Ham | 0.94 | 1.00 | 0.97 | 1355 |
| Spam | 1.00 | 0.59 | 0.74 | 196 |
| Avg/Total | 0.95 | 0.95 | 0.94 | 1551 |

**Table 3**.*Classification Report of NaiveBayes Classifiers*

*D.*  Classification Report of Support Vector Machines

|  | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Ham | 0.98 | 1.00 | 0.99 | 1355 |
| Spam | 0.98 | 0.87 | 0.92 | 196 |
| Avg/Total | 0.98 | 0.98 | 0.98 | 1551 |

**Table 4.***Classification Report of SVM*

*E.*  Classification Report of Decision Tree Classifiers

|  | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Ham | 0.96 | 1.00 | 0.98 | 1355 |
| Spam | 0.99 | 0.72 | 0.83 | 196 |
| Avg/Total | 0.96 | 0.96 | 0.96 | 1551 |

**Table 5.** *Classification Report of Decision Tree*
F.*Classification Report of Random Forest Classifier*

|  | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Ham | 0.96 | 0.98 | 0.97 | 1355 |
| Spam | 0.85 | 0.73 | 0.79 | 196 |
| Avg/Total | 0.95 | 0.95 | 0.95 | 1551 |

**Table 6.** *Classification Report of Random Forest*

502

According to the Comparison of Accuracy result Support Vector Machines provides the best accuracy score is 98% among all other algorithms. The Comparison of Training Set Result shows the maximum Support Vector Machines RBF trained more times among all other algorithms. The Comparison of Prediction time shows the best result is Naive Bayes because it predicts short span of time among all other algorithms. Naive Bayes provides the fair accuracy and moreover short amount of time. Naive Bayes collects the spam and ham messages effectively

## VII. CONCLUSION

The aims and objectives of the project, which achieved throughout the course, defined at the very first stage of the process. To collect all the information, the research work involved a careful study on the different filtering algorithms and existing anti-spam tools. These large scale research papers and existing software programs are one of the sources of inspiration behind this project work. The whole project was divided into several iterations. Each iteration was completed by completing four phases: inception, where the idea of work was identified; elaboration, where architecture of the part of the system is designed; construction, where existing code is implemented; transition, where the developed part of the project is validated. However, there are still some parts that can be improved: for example, adding additional filtering techniques or changing aspects of the existing ones. The changes such as incrementing or decrementing the number of interesting words of the message and reorganizing the formula for calculating interesting rate can be done later.

## VIII. FUTURE SCOPE

In the future, we plan to deal with more challenging problems such as the analysis and management of report in spam SMS filters storing. Solution for this problem is another focus of work in the future.

## REFERENCES

[1] Camponovo G, Cerutti D., "The spam issue in mobile business: A comparative regulatory overview", Proc. 3rd Int. Conf. Mobile Bus., pp. 1-17..

[2] Cleff E.B., "Privacy issues in mobile advertisin"', Int. Rev. Law Comput.Technol., vol. 21, pp. 225-236.

[3] Fu J, Lin P, Lee S. , "Detecting spamming activities in a campus network using incremental learning", J. Netw. Comput. Appl., vol. 43, pp. 56-65.

[4] Hua J, Huaxiang Z., "Analysis on the content features and their correlation of Web pages for spam detection" , China Commun., vol. 12, no. 3, pp. 84-94.

[5] Reaves B, Scaife N, Tian D, Blue L, Traynor P, Butler K.R."Sending out an SMS: Characterizing the security of the SMS ecosystem with public gateways", Proc. IEEE Symp. Secur. Privacy (SP), pp. 339- 356..

[6] Wang et al C.,"A behavior-based SMS antispam system", IBM J. Res. Develop., vol. 54, no. 6, pp. 3:1-3:16.

[7] Yamakami T, "Impact from mobile SPAM mail on mobile internet services' in Parallel and Distributed Processing and Applications", Berlin, Germany:Springer, pp. 179-184