# Mobile SMS Spam Detection using Machine Learning Techniques

Samadhan Nagre
Dept of Computer Science & IT
Dr. B.A.M. University Aurangabad

**Abstract—** Spam SMS be unwanted messages to users, which be worrying and from time to time damaging. present be a group of survey papers available on SMS spam detection techniques. study and reviewed their used techniques, approaches and algorithms, their advantages and disadvantages, evaluation measures, discussion on datasets as well as lastly end result judgment of the studies. even though, the SMS spam detection techniques are additional demanding than SMS spam detection techniques since of the local contents, use of shortened words, unluckily not any of the existing research addresses these challenges. There is a enormous scope of upcoming research in this region and this survey can act as a reference point for the upcoming direction of research.

**Keywords** : Mobile SMS spam Detection

## Introduction.

small Message once-over (SMS) is the the majority often and extensively used message medium. The term "SMS" is used intended for together the user activity and all types of small text messaging inside a lot of parts of the world. It has develop into a medium of announcement as well as endorsement of products, banking updates, agricultural information, flight updates and internet offers. SMS be as well working in direct marketing recognized as SMS marketing.

from time to time SMS marketing is a matter of trouble to users. These kinds of SMS

are called spam SMS. Spam is one or additional unwanted messages, which is unwanted to the users, sent or posted as part

of a better collection of messages, every one having considerably matching content. The purposes of SMS spam be announcement and marketing of a variety of products, sending political issues, dispersal unsuitable

adult content and internet offers. so as to is why spam SMS flooding has become a serious problem. All over world.SMS spamming gain reputation over additional spamming approaches

like electronic mail with twitter owing to the rising popularity of SMS Communication.

## Background and related work

SMS spam detection be moderately a fresh research area after that Text SMS electronic mail social tags, and twitter and web spam detection. a number of of the researches of spam detection include [1][2] etc. these researchers are typically conducted following 2011. There be some recognized SMS spam detection technique have some challenge more than SMS  spam detection such as limited message size use of local and shortcut words and incomplete slogan information. These challenges require to be solved. present is scope of research in this field and some research works contain be conducted on it  present be different category of SMS spam filtering such as pallid record and black record . content- based non content-based be two-way approaches and challenge response technique.[4],[5],[8] the techniques using customer surface server in these several machine learning Algorithms such as

## Naïve Bayes.

Bayesian is a probabilistic move toward that starts among a previous faith, observes some data and then updates that faith The probabilistic life form spam and not spam of a word can be intended through the incidence of that word in ham and spam messages with the Bayesian algorithm [12].

## Support vector Machine.  (SVM)

Support vector machines be supervised learning by way of linked algorithms that analyses data used intended for the categorization as well as regression analysis. If a put of teaching example containing spam and rightful SMS is known after that SVM teaching algorithm build a model that can assigns new example keen on spam and rightful group An SVM model is a demonstration of the example because a point in space, mapped so that example of the divide category are

separated by a clear gap so as to is wide as achievable.[9]

### Decision Trees .

A decision hierarchy be a decision support instrument that use a hierarchy similar to or model of decisions and their likely penalty, counting possibility of event outcomes. A decision tree can be used to make choice to whether a fresh message is spam or ham [11]

### Logistic Regression

A logistic regression be a prognostic analysis. Logistic regression be used to explain data and to explain the association flanked by single reliant binary changeable and single or additional supposed ordinal, interval or percentage-level independent variables. from time to time logistic regression be hard to understand, the intellects statistics instrument without difficulty allows you to conduct the analysis, after that in simple English interprets the production.

### Random Forest

Random Forest is a trademark term for an ensemble of decision trees. In random forest we collection of decision trees. To classify a new object based on attributes each tree votes for that class. The forest chooses the classification having the most votes over all the trees in the forest.

```
┌─────────────────────────────┐
│   Spam Filtering process    │
└─────────────────────────────┘
              │
              ▼
┌─────────────────────────────┐
│        Preprocessing        │
└─────────────────────────────┘
              │
              ▼
┌─────────────────────────────┐
│        Tokenization         │
└─────────────────────────────┘
              │
              ▼
┌─────────────────────────────┐
│       Representation        │
└─────────────────────────────┘
              │
              ▼
┌─────────────────────────────┐
│         Selection           │
└─────────────────────────────┘
              │
              ▼
┌─────────────────────────────┐
│          Training           │
└─────────────────────────────┘
              │
              ▼
┌─────────────────────────────┐
│          Testing            │
└─────────────────────────────┘
```

| Study ID | Available | Total No of Messages |
|---|---|---|
| [1] | [13] | 5574 |

### Spam filtering process

A physically confidential spam and ham messages be input or teaching position for a spam filtering

algorithm. The algorithm consists of the following steps.

**Preprocessing.** Removing irrelevant contents like stop word are the part of the data preprocessing

**Tokenization.** Segmenting the message according to words character or symbols called tokens. There are different tokenization approaches such as word tokenization, sentence, word or character N- grams and orthogonal sparse bigrams.

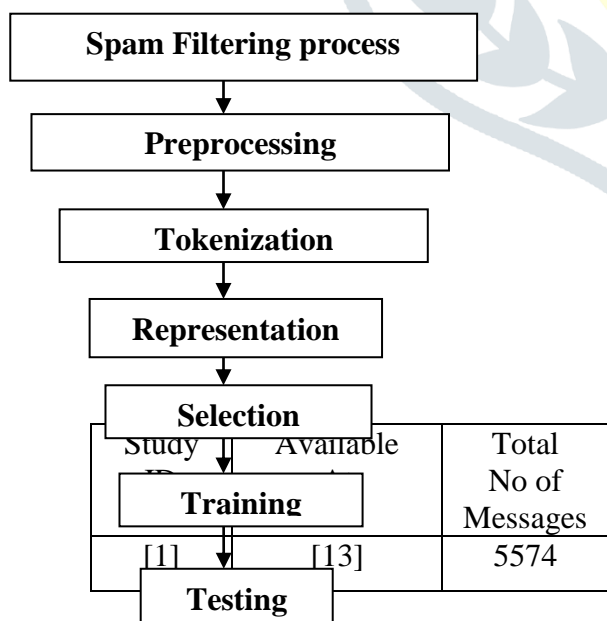**Representation.** Conversion to attribute value pairs

**Selection.** Selecting important attribute values which contain crash on top of categorization quite than choosing every one pairs of attribute value.

**Training.** teach the algorithm by way of the chosen quality values.

**Testing.** experiment the recently inwards data by way of the teaching model.

### Study Selection Procedure :

near select related studies we mostly searched on Google academic. we  have together a number of papers from it present are a number of additional conference and journals such when IEEEExplore, IJCSI ITJ ACM, be create from side to side Google academic instrument The file of journals and conference as of anywhere we containchosen paper contains a lot of references. We also searched used for the referenced papers and have in use a number of of them as our applicable paper . we used the Google scholar`s connected articles and cited feature four our searching procedure.

Table 1.SMS Spam Detection dataset description

### Validation of the study :

Our SLR be conducted in the direction of examine every one the used approaches and techniques inside SMS spam detection. The intimidation  to the strength of our review be that present might be selection bias and be short of  enough resources. We try to get to every one possible and pertinent information resources . some resources strength not contain been published in a straight line. one

more danger be a number of resources are not available intended for community use.

## Result Analysis:

A initial we physically searched on top of Google using the topic spam detection to increase an impression in spam detection field, It resulted inside a lot of SMS spam detection connected papers. then we modified our search using merely Mobile SMS spam detection Through our study selection procedure we have chosen 13 Paper paper  published in different conferences and journals connecting merely on the way to Mobile SMS spam detection  in the middle of the 13 studies.

## Dataset Description:

A preparation dataset be wanted intended for a number type of machine learning classification algorithms. Result of the machine knowledge algorithms depend on the dataset. because a effect spam detection algorithms be able to run without a dataset. within we established dissimilar openly available dataset apply in different studies. link of the dataset and a quantity of statistics such as total number of SMS number of spam and ham messages are shown in table [13]

## Conclusion.

This paper present the result of the systematic literature review on SMS spam detection. We chose a total of 13 research paper on this field and reviewed their proposed techniques. Advantages and disadvantages. And challenges they addressed. we also examined their evaluation procedures. We demonstrated the publicly available dataset information  which is a prior need for a spam filtering algorithm.

## REFERENCES

[1] K. Yadav, S. K. Saha, P. Kumaraguru, and R. Kumra,'' take control of your smses: designing an usable spam  sms filtering system," in 2012 IEEE 13th International Conference on Mobile Data Management. IEEE, 2012, pp. 352–355.

[2] S. J. Warade, P. A. Tijare, and S. N. Sawalkar, "An approach for sms spam detection."

[3] A. Narayan and P. Saxena, "The curse of 140 characters: evaluating the efficacy of sms spam detection on android," in Proceedings of the Third ACM workshop on Security and privacy in smartphones & mobile devices ACM, 2013, pp. 33–42.

[4] A. S. Onashoga, O. O. Abayomi-Alli, A. S. Sodiya, and D. A.Ojo "An adaptive and collaborative server side sms spam filtering  scheme using artificial immune system," Information Security Journal: A Global Perspective, vol. 24, no. 4-6, pp. 133–145, 2015.

[5]  J. W. Yoon, H. Kim, and J. H. Huh, "Hybrid spam filtering for mobile communication," computers & security, vol. 29, no. 4, pp. 446–459, 2010.

[6] Q. Xu, E. W. Xiang, Q. Yang, J. Du, and J. Zhong, "Sms spam detection using noncontent features," IEEE Intelligent Systems, vol. 27, no. 6, pp. 44–51, 2012.

[7] I. Ahmed, D. Guan, and T. C. Chung, "Sms classification based on naïve bayes classifier and apriori algorithm frequent itemset," International Journal of machineLearning and computing, vol. 4, no. 2, p. 183, 2014.

[8] J. M. G´omez Hidalgo, G. C. Bringas, E. P. S´anz, and F. C. Garc´ıa, "Content based sms spam filtering," in Proceedings of the 2006 ACM symposium on Document engineering. ACM, 2006, pp. 107–114.

[9] https://en.wikipedia.org/wiki/Support_ vector machine [last Accessed: 05-11-2016]

[10]https://en.wikipedia.org/wiki/Decision tree[Last Accessed:05-11-2016]

[11] http://en.wikipedia.org/wiki/K-

[12]http://fastml.com/bayesian-machine-learning/ [Last Accessed: 05-11-2016

[13]SMS Spam Collection data set from UCI Machine learning Repository,"http://archive.ics,uci.edu/ml/data set/SMS+Spam+Collection