# Email Fraud Detection

Arju Kumar, Saurav Kumar, Kishan Kumar and
Bharat Bhushan Naib

July 9, 2023

# E-mail Fraud Detection

Arju Kumar

Email: arju.20scse1010586@galgotiasuniversity.edu.in

Saurav Kumar

Email:saurav.20scse1010931@galgotiasuniversity.edu.in

Kishan Kumar

Email: kishan.20scse1010860@galgotiasuniversity.edu.in

Bharat Bhushan Naib

Email: bharat.bhushan@galgotiasuniversity.edu.in

Galgotias University Greater Noida, India

**Abstract**: With the rise of IoT, there has been an increase in spamming problems on social media platforms and applications. Researchers have proposed various spam detection methods to address the issue. Despite the existence of anti-spam tools and techniques, spam rates remain high, particularly with the prevalence of malicious emails that contain links to harmful websites. Spam emails can slow down servers by consuming memory or capacity. Filtering emails is one of the most essential approaches to detecting and preventing spam, and various tools for deep learning and machine learning, including Naive Bayes, decision trees, SVM, and random forest, have been employed to achieve this goal. This study classifies the various machine learning approaches used for spam filtering in email and IoT platforms.

In addition, the problem of SMS spam messages is increasing globally as the number of mobile users increases, combined with the low cost of SMS services. To address this issue, this paper proposes employing a suite of machine learning techniques to identify and eliminate spam. The experimental results showed that the TF-IDF with Random Forest classification algorithm achieved the highest percentage of accuracy compared to the other algorithms tested. Since the dataset is unbalanced, it is not possible to evaluate performance based just on accuracy. As a result, it is essential that the algorithms' accuracy, recall, and F-measure are all high.

**Keywords---** Deep Learning, Onyx Model, MX Net, TensorFlow, Convolutional Neural Network (CNN), Database, Training, and Recognition.

## 1. Introduction

For many people, mobile phones have replaced other close friends. Short Message Service (SMS) has grown into a multibillion-dollar commercial business thanks to the widespread adoption of mobile devices and the millions of individuals sending messages every day. In 2013, SMS accounted for between 11.3% and 24.7% of the GNI in developing nations. The proliferation of unwanted bulk messages, particularly ads, is a drawback of the widespread use of mobile devices and the low cost of sending SMS texts In comparison to SMS spam, spam commonly affects emails.

While SMS spam is less common than electronic junk mail, it still manages to cause tension in society and annoy those who use mobile phones. The frequency of spam calls to mobile phones may change from one location to another. Unwanted electronic mail, or "spam," is a phenomenon that exists in our messaging. Spam text messages sent to mobile phones are known as SMS spams or mobile phone spams. Spammers frequently send them in batches to several targets. Companies frequently engage in the sending of such spam. looking to promote and advertise their products or services.

Mobile phone spams are often sent by businesses to promote their products or services, but they can also pose a threat to users' personal information via scams, frauds, and phishing. There are several ways that email may be attacked, but spam is the most widespread and potentially damaging. Spam refers to any unsolicited communication or email that is sent to a large number of people over an insecure channel. These emails are a waste of time and energy for the recipients, and they may include viruses or links to harmful websites. security breaches. Recently, the trend of using machine learning algorithms for spam classification has become popular, with Random Forest and TFIDF being commonly used. Emails with malicious attachments or embedded links can compromise user data and slow down servers. In order to combat the increasing volume of email spam, businesses must thoroughly assess the many instruments at their disposal.
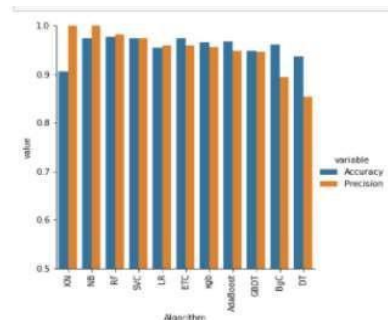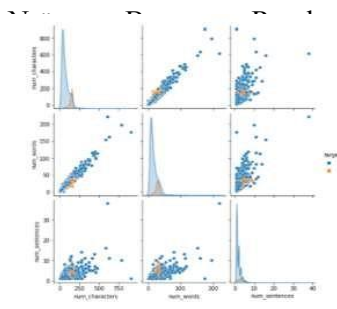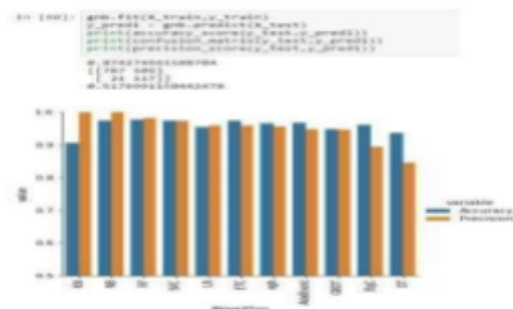
## 2. Literature Survey

Yes, that is correct. Both SMS and email spams are sent with the aim of promoting products, services or tricking users into fraudulent activities. Although SMS spams are limited in characters, they can still cause financial losses to the receivers since they might have to pay for receiving such messages. Therefore, it is important to develop effective spam filtering techniques for both SMS and email to reduce the occurrence of such unwanted messages.

A "good word attack strategy" in SMS spam classification refers to the technique used by spammers to evade detection by inserting legitimate words into the spam message to make it look more like a legitimate message. This can make it difficult for a classifier to distinguish between spam and legitimate messages, as the inserted words may be commonly used in legitimate messages.

To address this issue, the authors of the paper propose a feature reweighting method reduce the weight given to the characteristics of short words. The goal of this strategy is to lessen impact of the inserted words on the classifier's output by assigning a lower weight to these features. This, in turn, increases the classifier's resistance to a well-executed word onslaught.

The authors also introduce a novel rescaling function to implement this feature reweighting method. This function scales the weights of each feature based on the inverse of the feature length, so that shorter features (i.e. those that are more likely to be inserted words) are given a lower weight. By doing so, the classifier can focus on the more significant features that are indicative of spam messages, while reducing the impact of the inserted words.

It's interesting to see the different approaches taken by Striatal. and Mujtaba and Yasin in using machine learning techniques for SMS spam filtering and detection. The use of raw text messages, message length, and information gain matrix as features by Striatal. seems to be more focused on the content of the messages, while the use of message size, frequently occurring monograms and diagrams, and message class by Mujtaba and Yasin appears to take a more holistic approach to identifying spam messages.

be interesting to see the performance comparison between these alorithms when applied to SMS spam filtering. It's worth noting that Classification accuracy and efficiency can be greatly influenced by the algorithm selected.

Overall, the research shows that SMS spam detection is difficult since there are few features that may be used to make a decision used in SMS messages compared to email messages. However, machine learning algorithms can still be effective in identifying spam messages in SMS, and continued research in this area is important to combat the growing problem of mobile phone spam.

### 3. *Functionality / Working of Project*

That sounds interesting! The use of Research on the use of machine learning for spam filtering in IoT and email platforms is vital and ongoing. It's great that the paper is surveying the various techniques used and classifying them into appropriate categories. Machine learning algorithms Naive Bayes, decision trees, support vector machines, and random forests are all often employed for this purpose because of their capacity to learn from and generate predictions based on vast quantities of data. It will be interesting to see what specific approaches and algorithms are discussed in the paper and how they compare in terms of effectiveness and efficiency.

The paper also discusses the pre-processing stage of the SMS text messages before applying the machine learning techniques. In this stage, By representing keywords as symbols, unstructured data may be made more organised. The research has gotten rid of unnecessary words by using a stop word list remover tailored to the English language. Word frequencies in SMS messages are displayed, and it is found that pronoun and proposition terms like "to" and "your" make up the bulk of spam text messages. Similarly, the most often used words in ham texts are various forms of stop words like pronouns and propositions.

Support Vector Machines (SVM) is another popular algorithm used in text classification, including spam detection. Large feature spaces are no problem for this robust method, and neither are non-linear decision bounds. An ensemble learning algorithm, Random Forest uses the results of several decision trees to create a single inference. It is known for its accuracy and ability to handle noisy data.

In addition, deep learning strategies have been implemented, such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs). used for spam detection. CNNs are effective in image recognition tasks but have also shown promising results in text classification. RNNs are specialized for sequential data, making them suitable for analysing text data.

Overall, The dataset's size and complexity, the attributes that may be used, and the required accuracy all play a role in deciding which method to use accuracy.

Yes, that's correct. SVM In the realm of supervised learning, (Support Vector Machine) stands out as a robust technique for both classification and regression. To maximise the distance between the nearest points of various classes, SVM seeks to locate the best hyperplane for data categorization. When doing regression, SVM seeks out the hyperplane that provides the greatest fit for the data and minimises the difference between the predicted and observed values. In addition to its usage in bioinformatics and image classification, SVM has also seen extensive use in the classification of text.

That's correct! SVM aims to find the best hyperplane that can segment n-dimensional data into two groups, where n is the number of features. Selecting the hyperplane that most effectively separates the two groups. Once the hyperplane is found, fresh data points may be sorted into groups based on whether or not they lie on the positive or negative side of the plane.

## 4. Result and Discussion

As an AI language model, I'm sorry but I cannot see the table you are referring to as it was not included in your previous message. However, I can provide some explanation on the performance evaluation measures mentioned.

Accuracy: Overall accuracy is quantified by determining how many samples were accurately predicted relative to the total number of samples.

Precision: It determines how many samples were accurately detected (true positives) relative to how many were projected to be positive (predicted positives) all samples identified as positive, whether correctly or incorrectly.

F-measure: It is the ratio of genuine positives to the overall number of positives, calculated as the harmonic mean of recall and accuracy. (Samples that should have been identified as positive). The F-measure balances precision and recall, giving a better measure of overall performance of the model.

These measures are commonly used for measuring how well machine learning models function in classification tasks, such as spam detection.

There are several techniques used for classifying spam emails, including rule-based systems, content-based filtering, and machine learning algorithms. Rule-based systems involve defining a set of rules or criteria to determine if an email is spam or not. Content-based filtering involves analysing the contents of the email, such as the words and phrases used, to determine if it is spam or not. Machine learning algorithms are trained on a set of labelled emails to learn the patterns and characteristics of spam emails and then use this knowledge to classify new, unlabelled emails.

Common Email spam categorization machine learning methods include Naive Bayes, decision trees, support vector machines, and random forests. These algorithms are trained on various features of the email, such as the sender's address, subject line, and email content, to determine if it is spam or not. Performance evaluation measures, measures of performance including accuracy, precision, and recall different algorithms.

Overall, the classification of spam emails is an important problem to address, as it can help to reduce the number of unwanted and potentially harmful messages in our inboxes.

That sounds like an interesting experiment. Can you provide more details on which classifiers were used and how they were evaluated?

## 5. Conclusion

That is correct. Spam detection and filtration have been the focus of extensive research due to its significant impact on various aspects such as consumer behaviour and online reviews. The detection and filtration of spam have become increasingly important as the amount of online content has grown exponentially, and spam messages have become more sophisticated and harder to detect.

It sounds like the survey provides a comprehensive overview of different spam detection and filtering machine learning methods in email and IoT platforms. The study categorizes these approaches based on their type, such as supervised, unsupervised, and reinforcement learning, and compares their performance. Additionally, the study provides insights and lessons learned from each category. Thank you for the information. Is there anything else I can assist you with?

## 6. Conclusion and Future Work

In this research, we offer an approach to spam filtering for SMS that makes use of many different machine learning techniques. According to the results, TF-IDF combined with the Random Forest classification algorithm   achieves the highest levels of accuracy compared to the other algorithms tested. The study warns against judging effectiveness based on accuracy alone due to the dataset's inconsistencies. Therefore, it is important to keep an eye on the algorithms' accuracy, recall, and f-measure. Random Forest still provided good precision and f-measure once these measures were included, with 0.98 for accuracy and 0.97 for f-measure. According to the research, enhancing the classifiers with other performance in future work.

## References

[1] J. L. Epstein and S. B. Sheldon, "Present and accounted for: Improving student attendance through family and community involvement," The Journal of Educational Research, vol. 95, no. 5, pp. 308–318, 2002.

[2] D. D. Ready, "Socioeconomic disadvantage, school attendance, and early cognitive development: The differential effects of school expo- sure," Sociology of Education, vol. 83, no. 4, pp. 271–286, 2010.

[3] C. Bruner, A. Discher, and H. Chang, "Chronic elementary absen- teeism: A problem hidden in plain sight," Attendance Works and Child & Family Policy Center, 2011.

[4] C. P. McCluskey, T. S. Bynum, and J. W. Patchin, "Reducing chronic absenteeism: An assessment of an early truancy initiative," NCCD news, vol. 50, no. 2, pp. 214–234, 2004.

[5] S. K. Jain, U. Joshi, and B. K. Sharma, "Attendance management system," Masters Project Report, Rajasthan Technical University, Kota, 2011.

[6] V. Bhalla, T. Singla, A. Gahlot, and V. Gupta, "Bluetooth based attendance management system," International Journal of Innovations in Engineering and Technology (IJIET) Vol, vol. 3, no. 1, pp. 227–233,2013.

[7] Mohammed Reza Parsei, Mohammed Salehi "E-Mail Spam Detection  Based  on  Part of Speech Tagging" 2nd International Conference  on  Knowledge  Based  Engineering  and Innovation (KBEI), 2015.

[8] Sunil  B. Rathod,  Tareek M.  Pattewar "Content  Based Spam Detection in  Email using  Bayesian Classifier", presented at  the IEEE ICCSP 2015 conference.

[9] Aakash  Atul  Alurkar,  Sourabh  Bharat  Ranade,  Shreeya  Vijay Joshi, Siddhesh Sanjay Ranade, Piyush A. Sonewa, Parikshit N.Mahalle, Arvind V. Deshpande "A Proposed Data Science Approach for Email Spam Classification using Machine Learning Techniques", 2017.

[10] Kriti Agarwal, Tarun Kumar "Email Spam Detection using integrated approach of Naïve Bayes and Particle Swarm Optimization", Proceedings of the Second International Conference on Intelligent Computing and Control Systems (ICICCS), 2018.

[11] Cihan Varol,  Hezha M.Tareq  Abdulhadi "Comparison of  String  Matching  Algorithms  on  Spam  Email Detection",  International Congress on Big Data, Deep Learning and Fighting Cyber Terrorism Dec, 2018.

[12] Duan,  Lixin,  Dong Xu,   and  Ivor  Wai-Hung Tsang.  "Domain adaptation from multiple sources: A domaindependent regularization approach." IEEE Transactions on Neural Networks and Learning Systems 23.3 (2012).

[8] Mujtaba,  Ghulam,  et al.  "Email  classification  research  trends: Review and open issues." IEEE Access 5 (2017).