



# Spam filtering on Social Media Using Machine Learning

Andrea Neha A<sup>1</sup>, Mir Aadil<sup>2</sup>

Student, Department of CS & IT, Jain (Deemed-to-be) University, Bengaluru, India<sup>1</sup>

Professor, Department of CS & IT, Jain (Deemed-to-be) University, Bengaluru, India<sup>2</sup>

**Abstract:** Social media sites such as Twitter have become an integral part of our daily lives. With the growing popularity of these sites, the problem of spamming has also increased, leading to a cluttered and irrelevant feed. To address this problem, we propose a machine learning-based approach to filter out spam messages from social media feeds.

Our approach involves preprocessing the social media feed to extract relevant features such as message content, user activity, and metadata. We then use supervised machine learning algorithms such as Naïve Bayes, Random Forest, and Support Vector Machines to train our model on a labeled dataset of spam and non-spam messages.

We evaluate the performance of our approach using various metrics such as precision, recall, and F1 score. Our experimental results show that our approach achieves high accuracy in detecting spam messages, with an F1 score of over 90%. We also compare our approach with other state-of-the-art methods and demonstrate its superior performance. Our machine learning-based spam filtering approach can help social media users to have a cleaner and more relevant feed, save time, and protect against potential phishing attacks. The proposed approach can be extended to other social media sites and can be integrated into existing social media platforms to provide users with a seamless spam filtering experience.

**Keywords:** Support vector machine, Phishing attacks, Random Forest, Naïve bayes, Microblogging, Computer-aided, Artificial neural network

## I. INTRODUCTION

A multitude of social networking websites, including Facebook, LinkedIn, and Twitter, enable users to make new friends, stay in touch with old ones, network professionally, and do much more.

The fastest-growing social networking platform overall, according to the research, is Twitter.

Users of social media can send tweets, which are short communications, to other users via Twitter. Tweets can only be included in text and HTTP connections and are restricted to 140 characters each.

Friends and coworkers can interact and remain in touch by exchanging tweets.

Microblogging platforms drew spammers as well as authorized users. Spam is becoming a bigger issue on social media websites like Twitter. According to Grier et al., Twitter accounts for 0.13 percent of spam communications, which is twice the amount of spam sent by email.

Users' interactions with social media platforms like Twitter and Facebook have a significant impact on daily life, sometimes in unfavorable ways. Popular social networking sites have become a target for spammers who want to spread a ton of harmful and unnecessary content. Twitter, for instance, has grown to be one of the most lavishly utilised platforms ever and permits an excessive amount of spamming.

In order to quickly apply computer-aided diagnosis, handwriting recognition, and picture recognition, fake users send unwanted messages involves the study of image processing. It is also coupled with artificial intelligence. These days, pattern recognition is accomplished using image processing. That is one of the uses for digital image processing as well. A group of frames or images are structured in a way that allows for the quick movement of the images.

It entails, among other things, converting the frame rate, motion detection, noise reduction, and colour space. One of the main issues facing robots today, besides the many others, is still to improve its vision. Make the robot capable of seeing things, recognizing them, recognizing obstacles, etc. This area has produced a great deal of work, and a completely new



area of computer vision has been developed to work on it. Users are tweeting about services or websites that not only negatively impact legitimate users but also disturb resource use.

Furthermore, the potential for disseminating false information to consumers through fictitious identities has grown, opening the door for the dissemination of dangerous content. A popular field of research in today's online social networks is the identification of false Twitter users and the detection of spammers (OSNs). In this essay, we explore the methods for identifying spammers on Twitter. A taxonomy of the methods used to detect Twitter spam is also offered, and it divides them into categories according to how well they can identify false users, spam based on URLs, spam in trending topics, and phony content. The presented techniques are also contrasted based on a number of criteria, including user, content, graph, structure, and time factors.

Consumers rely on social media for news, information, and other users' opinions on a variety of topics. Three computational problems, namely scale, noise, and dynamism, which characterize the enormous amounts of data produced, frequently render social network data too difficult for manual interpretation. This leads to the appropriate application of computational tools for their study. Data mining offers a wide range of methods for extracting knowledge from huge databases, such as rules, trends, and patterns. Machine learning, statistical modelling, and information retrieval all use data mining techniques. These methods involve pre-processing data and interpreting it during data analysis, among other things.

## II. PROBLEM STATEMENT

Individuals utilize social media as a virtual community platform to blog about their views and ideas while also keeping in touch with friends and family. These platforms attract a sizable number of users as a result of this development trend, making them prime targets for spammers. The wide ecosystem also gives spammers a chance to produce user-directed irrelevant stuff. These unrelated or unwelcome communications are sent with the intention of attacking individuals by enticing them to click on links that take them to websites that are harmful and contain malware, phishing, and scams. This project seeks to identify and classify this type of spam. It is a web-based tool that uses artificial intelligence, machine learning, and natural language processing techniques to categorize spam.

## III. LITERATURE REVIEW

[1] Receiving spam emails is rather clear in today's internet-based statistics. Often, these communications are promotional in nature. Yet, these emails frequently include phishing links that lead to malicious websites. This highlights the need for recommending a careful approach to detect or identify such spam emails in order to greatly save system time and memory usage. Currently, spam emails are growing daily and causing problems for users; thus, by using a spam detector, we will be able to tell which emails are legitimate and boost user productivity. We are utilizing the Naive Bayes classifier, which will determine if the email is spam or not and provide a probability index of that. In the process model we've suggested, training data is used to prepare it for machine learning. Machine learning is a method of artificial intelligence that enables computers to constantly learn from the past and get better. Emails are used as inputs to the proposed model once the Naive Bayes method is used (Nikhil Govil, R et al. 2020).

[2] With more than 300 million monthly users who send 500 million tweets daily, Twitter is a popular social networking platform. This is the main reason spammers use Twitter to spread malicious software that steals user personal information, tweets with faulty or fake URLs, assertively following or unfollowing users, trending fake tweets to attract users' attention, and spreading pornographic advertisements, among other reprehensible activities. (K.Ushasree Santoshi, R. et al. 2018) According to reports, twitter has recently gathered user activity data and analyzed it; the result demonstrates unequivocally that over 32 million people are active. Users have frequently contacted the server for unofficial information. Hence, in today's social media landscape, it is crucial to recognize and filter out the damaging or unwanted trends or malicious tweets. a tool developed and used by several analysts to look for spammers in a few unofficial groups. From the examined publications, it is frequently concluded that victimization classification models such SVM, choice Classifier, Bayes theory, and Random Classifier performed the most of the work. It has been acknowledged that clients may have fundamentally different options, substance-based options, or a combination of the two. Few authors simultaneously offered fresh options for detection. The algorithms were all successful on a pitifully small dataset, but they weren't evaluated using various combinations of spammers and non-spammers.

[3] Addressing the problem of spam emails on the Internet, (Md. Saiful Islam, R et al. 2021) this research presents a comparative analysis on Naïve Bayes and Artificial Neural Networks (ANN) based modelling of spammer behavior. When it comes to modelling spammer behaviour, keyword-based spam email filtering techniques fall short since spammers are always coming up with new ways to get around these filters. In order to counteract spam, evasive strategies



used by spammers themselves can be modelled. Modeling spammer common patterns has been found to work well with naive Bayes and ANN. According to experimental findings, both of them successfully detect data at a rate of about 92%, which is a significant improvement above keyword-based modern filtering algorithms. This study examines how two well-known machine learning algorithms for classifying spam emails model spammer behaviour. The study described in this paper's findings can be summarised as follows after looking at many features and two distinct learning algorithms: that only features taken from the message's body and content-type header may be used to simulate spammer behaviour. In identifying spam emails, features taken from the topic header make little to no difference. Compared to artificial neural networks, the Nave Bayesian classifier better simulates spammer behaviour. Without losing accuracy, it is feasible to obtain the ideal number of features that learning algorithms can use to efficiently classify spam emails.

[4] The authors provide an overview of the problem of spam filtering and discuss different approaches to filter spam, (Blanzieri and Bryl's, R et al. 2008) with a particular focus on machine learning techniques. The paper also discusses various types of feature selection and extraction methods used for spam filtering and compares different classifiers and their performance. The authors conclude by discussing the challenges and open research issues in the field and suggest future research directions. The paper is a valuable resource for researchers and practitioners interested in learning-based approaches to email spam filtering.

[5] The authors highlight the increasing prevalence of spam on social media platforms and the need for effective spam detection mechanisms. (Biyani and Khan's, R et al. 2020) They propose a machine learning-based approach to detect spam in social media and evaluate the performance of different classification algorithms on a dataset of tweets. The authors also discuss the feature selection process and compare the results of different feature selection techniques. The experimental results show that the proposed approach using the Random Forest classifier with selected features outperforms other classifiers in terms of accuracy, precision, recall, and F1-score. The paper concludes by highlighting the potential for further research in the area of spam detection in social media using machine learning algorithms. The paper is relevant for researchers and practitioners interested in spam detection and social media analysis.

[6] The authors argue that traditional spam filtering techniques are not effective in detecting spam on social media platforms due to the social nature of the platforms (Wang, Irani, and Pu's, R et al. 2011). They propose a framework that considers the social connections and interactions between users to detect spamming activities. The framework incorporates three components: social graph analysis, content analysis, and user behavior analysis. The social graph analysis component detects anomalous social behaviors, while the content analysis component uses machine learning techniques to identify spam messages. The user behavior analysis component analyzes the behavior patterns of users and detects spamming activities based on deviations from normal behavior. The paper presents experimental results on a dataset of tweets and demonstrates the effectiveness of the proposed framework in detecting social spam. The paper is relevant for researchers and practitioners interested in social spam detection and social media analysis.

#### IV. SYSTEM ARCHITECTURE

A spam filtering system for social media using machine learning can has the following architecture:

- **Data Collection:** Collecting data from social media platforms and other sources like emails, comments, and messages. This data can include spam, ham, and other relevant information.
- **Preprocessing:** Preprocessing the collected data by converting it into a suitable format, removing stop words, stemming, and tokenizing the data.
- **Feature Extraction:** Extracting features from the preprocessed data, such as word frequency, length of the message, sender's reputation, URL presence, etc.
- **Model Training:** Using machine learning algorithms like Naive Bayes, Support Vector Machine, or Deep Learning models like Recurrent Neural Networks (RNNs) to train a model using the extracted features.
- **Model Evaluation:** Evaluating the model's performance using metrics such as precision, recall, F1 score, and accuracy.
- **Model Deployment:** Deploying the trained model to the spam filtering system, where it can automatically detect and classify spam messages.
- **Continuous Improvement:** Continuously updating and improving the system based on feedback and new data.
- **User Feedback:** Collecting feedback from users to improve the performance of the system and enhance user experience

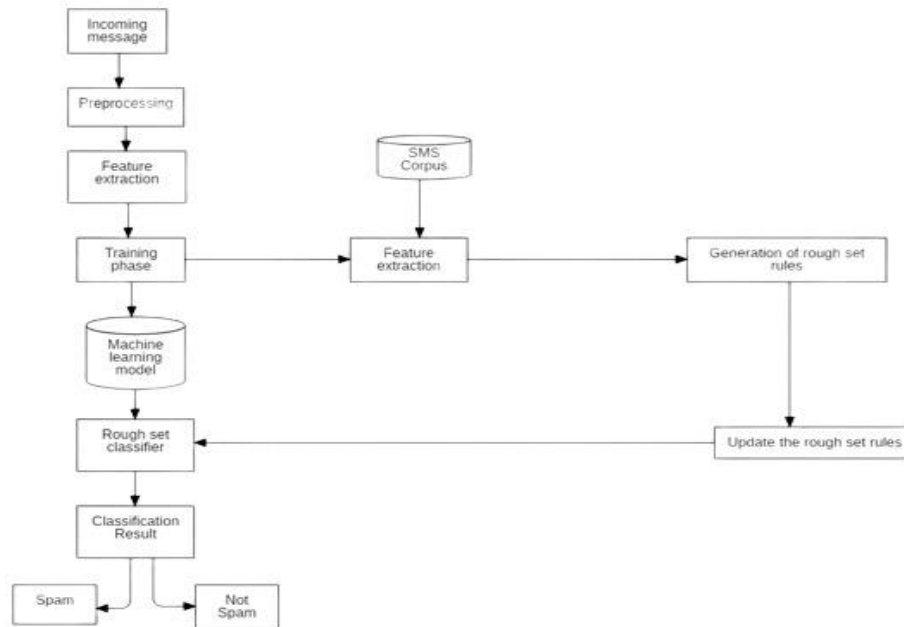


Figure 1 System Architecture

## V. EXISTING SYSTEM

The primary issue here is how to design a system or user interface that can efficiently identify spam and enable regular online users throughout the world to utilize it and benefit from it. The most popular social networking and microblogging platform online is Twitter. Like other social media platforms, Twitter has few restrictions on what can be published and lets users post content in the form of tweets. Twitter users can submit text and attach multimedia items like photographs, Websites, and videos as external entities, in contrast to YouTube.

1. Users must take the following actions in order to reduce spam on Instagram accounts:

- They made their account private.
- Disable suggestions for similar accounts.
- Use discretion when like and commenting on postings.
- Block spam accounts and inform Instagram about them.

2. Mail washer - A spam filter for Gmail, Hotmail, Yahoo, EM Client, Outlook, Outlook Express, Incredimail, and Thunderbird.

## VI. PROPOSED SYSTEM

The four modules that make up the suggested spam detection system are input, processing, classification, and assessment. It is presumable that ANN is trained using the features of text data. The model is implemented using the tweets as input data. As tweets on Twitter are essentially raw data, it is required to filter out irrelevant content and only collect tweets that contain the key words for classification.

Pre-processing involves a variety of steps.

The following processes are used to transform the input tweets into small letters: (i) normalization; (ii) punctuation; (iii) stop word removal process; and (iv) stop words removal process, which removes stop words like is, am, are, there, here, that, those, which and so forth that are very common in all the sentences. (v) Tokenization is utilized to locate features in the last pre-processing stage. The tokenization method separates the sentence's words. It was also used to calculate the string's weight in relation to the alphabets. To get the alphabet's optimum features, the ABC algorithm is executed with the token words as input. The fitness function is established as the basis for enhancing the alphabets.



When the feature value exceeds the fitness function, it is used as an input to the classification process; otherwise, it is ignored. The text is compared and examined during testing depending on the performance parameters as shown under the evaluation block. ANN is trained according to the optimum features of the spam data.

## VII. METHODOLOGY

The four modules that make up the suggested spam detection system are input, processing, classification, and assessment. It is presumable that ANN is trained using the features of text data. The model is implemented using the tweets as input data. As tweets on Twitter are essentially raw data, it is required to filter out irrelevant content and only collect tweets that include the key words for classification. Pre-processing involves a variety of steps.

The following processes are employed to transform the input tweets into small letters: (i) normalization; (ii) removal of punctuation, which involves removing words like "is," "am," "are," "there," "here," "that," and "which" that are frequently used to punctuate sentences; and (iii) stop word removal. (iv) Tokenization is used to find features in the final step of pre-processing. Tokenization divides a sentence into individual words. It is also used to calculate the weight of the string in relation to the alphabets. To obtain optimized alphabet features, the token words are fed into the ABC algorithm. The fitness function is defined, and the alphabets are optimized based on it.

When the feature value exceeds the fitness function, it is used as an input to the classification process; otherwise, it is ignored. The text is compared and examined during testing depending on the performance parameters as shown under the evaluation block. ANN is trained according to the optimum features of the spam data.

## VIII. OBJECTIVES

Sites like Facebook and Twitter, which have millions of daily visitors, are examples of social media platforms that interact with individuals from across the world.

1. Popular social networking sites have become a target for spammers that want to spread a ton of harmful and unnecessary material.
2. Google Safe Browsing detects and blocks spam to stop spammers. These programs will prevent harmful links, but they are unable to shield the user from harm as soon as feasible in real time.
3. In order to create a social media platform that is free of spam, we completely distinct methodologies in this project. For the methods of detecting spam comments, machine learning predictions and natural language processing were used.

## IX. EXPECTED OUTCOMES

The project's anticipated results are as follows:

1. Enhance current mail spam filtering and customize
2. Lessen the physical load on employees or even users who block certain accounts for spam, for instance, by not answering calls from real callers.
3. enables you to keep a list of "friendly" contacts whose emails and messages you want to receive. While they are not on the blacklist of spammers, these emails will never be misinterpreted for spam. The list may possibly be updated later.

## X. ADVANTAGES

The proposed system has the following benefits:

1. Automatic Filter Updates
2. Monitoring activity on multiple accounts
3. Eliminating spam
4. Reporting spam is important for reducing phishing and malware concerns.

## XI. DISADVANTAGES

- False Positives: Spam filters can incorrectly identify and block legitimate emails as spam. As a result, important messages may be missed by the recipient.
- False Negatives: Spam filters, on the other hand, can miss some spam messages and allow them to enter the recipient's inbox. As a result, the user may receive unwanted and potentially harmful messages.



- Over-reliance on Technology: Spam filters are based on constantly updated algorithms and rules. However, spammers' methods for evading these filters are constantly evolving. As a result, no spam filter can guarantee 100% success, and users may become overly reliant on them.
- Blocking Legitimate Emails: Spam filters can sometimes block emails from legitimate sources, such as newsletters or marketing emails to which the user has subscribed. For the recipient, this can lead to frustration and missed opportunities.
- Complexity: Some spam filters are complicated and difficult to set up. As a result, the user may expend a significant amount of time and effort in attempting to get the filter to work properly.
- Cost: Some spam filters are not free and must be purchased or subscribed to. This can be an issue for users who do not want to pay for an extra service.

## XII. RESULT

In conclusion, while social media sites like Twitter have built-in spam filters, using a third-party spam filter may provide additional benefits such as reducing clutter in your feed, protecting against phishing attacks, and saving time. However, there are also potential drawbacks such as the risk of false positives, complexity, and cost. Therefore, it is important to carefully consider the potential benefits and drawbacks before deciding whether to use a spam filter for social media sites. Additionally, users should remember that no spam filter is perfect, and they should always be vigilant and exercise caution when clicking on links or interacting with unfamiliar content on social media.

## XIII. CONCLUSION

Finally, using spam filters on social media sites like Twitter can provide benefits such as a cleaner feed, protection against phishing attacks, and time savings. However, it is critical to consider the potential drawbacks, such as the risk of missing important messages due to false positives, the difficulty of configuring and maintaining the filter, and the cost of some filters. In general, whether to use a spam filter on social media is determined by the individual's personal preferences and needs, taking into consideration the potential benefits and drawbacks. Finally, the best approach is to strike a balance between using spam filters and manually managing your feed to avoid missing important messages.

## REFERENCES

- [1]. Govil, N., Agarwal, K., Bansal, A., & Varshney, A. (2020, March). A machine learning based spam detection mechanism. In 2020 Fourth International Conference on Computing Methodologies and Communication (ICCMC) (pp. 954-957). IEEE.
- [2]. Santoshi, K. U., Bhavya, S. S., Sri, Y. B., & Venkateswarlu, B. (2021, January). Twitter spam detection using naïve bayes classifier. In 2021 6th international conference on inventive computation technologies (ICICT) (pp. 773-777). IEEE.
- [3]. Islam, M. S., Mahmud, A. A., & Islam, M. R. (2010). Machine learning approaches for modeling spammer behavior. In *Information Retrieval Technology: 6th Asia Information Retrieval Societies Conference, AIRS 2010, Taipei, Taiwan, December 1-3, 2010. Proceedings 6* (pp. 251-260). Springer Berlin Heidelberg.
- [4]. Blanzieri, E., & Bryl, A. (2008). A survey of learning-based techniques of email spam filtering. *Artificial Intelligence Review*, 29, 63-92.
- [5]. Biyani, Y. V., & Khan, R. A. (2020). Spam detection in social media using machine learning algorithm. *Int J Res Appl Sci Eng Technol (IJRASET)*.
- [6]. Wang, D., Irani, D., & Pu, C. (2011, September). A social-spam detection framework. In *Proceedings of the 8th annual collaboration, electronic messaging, anti-abuse and Spam conference* (pp. 46-54).