

Machine Learning Algorithm Based Email Spam Detection

N. Bhavana¹; Avulapalli Sowmya²

¹(Assistant Professor); ²(Post Graduate)

¹Department of MCA, Annamacharya Institute of Technology and Sciences (AITS),
Karakambadi, Tirupati, Andhra Pradesh, India

²Department of MCA, Annamacharya Institute of Technology and Sciences (AITS),
Karakambadi, Tirupati, Andhra Pradesh, India

Publication Date: 2025/06/03

Abstract: Email has emerged as the main channel for both personal and professional interaction due to the quick expansion of digital communication. However, because email is so widely used, there is an increase in spam communications, which can be anything from malicious phishing attempts to innocuous adverts. In addition to filling consumers' inboxes, these spam emails present serious security risks. Conventional rule-based spam filters have not been able to keep up with the changing nature of spam. Because machine learning (ML) algorithms can learn from data and get better over time, they are increasingly being incorporated into spam detection systems. In this paper, a variety of machine learning algorithms, such as Random Forest, Support Vector Machine (SVM), and Naive Bayes, are used to detect email spam. Using a publically accessible dataset, we assess these models on the basis of accuracy, precision, recall, and F1-score. According to our findings, machine learning models—in particular, ensemble approaches—provide reliable and expandable spam detection systems.

Keywords: Machine Learning, Spam, Emails.

How to Cite: N. Bhavana; Avulapalli Sowmya (2025). Machine Learning Algorithm Based Email Spam Detection. *International Journal of Innovative Science and Research Technology*, 10(5), 3014-3016.
<https://doi.org/10.38124/ijisrt/25may1564>

I. INTRODUCTION

In both personal and professional contexts, email is still a vital medium for communication. However, because of their ease of use and accessibility, email services are increasingly a target for spam, which is unsolicited and frequently harmful content. Spam emails are sent in large quantities, frequently with the goal of distributing malware, phishing, or advertising. Recent data shows that spam messages make up more than 45% of all email traffic, which presents a serious problem for user productivity and digital communication networks.

Manually created rules and keyword filters were a major component of traditional spam detection systems. Although these techniques offered some initial security, they frequently fell behind the ever-evolving strategies used by spammers.

Additionally, these algorithms frequently generated a large percentage of false positives, which resulted in the marking of valid emails as spam, causing crucial data to be lost and the filtering system to lose credibility.

One area of artificial intelligence called machine learning (ML) has shown a lot of promise in the fight against spam. Without human assistance, machine learning models can adjust to new spam tactics by examining vast email datasets and discovering the patterns that differentiate spam from authentic communications. In order to anticipate the type of incoming emails, these algorithms learn from characteristics including the frequency of specific terms, sender details, and message structure.

In order to detect email spam, we investigate the application of a number of machine learning methods in this work, including Naive Bayes, Support Vector Machine (SVM), Decision Trees, and Random Forests. In order to extract pertinent features, we pre-process a benchmark dataset of labeled spam and non-spam emails. Then, we assess the models using common performance metrics. The objective is to evaluate these algorithms' efficacy and choose the one that provides the optimum balance between computational efficiency and accuracy.

Additionally, this study suggests an improved spam detection method that makes use of ensemble learning, which combines several models to better classification

performance. Because of its scalable and adaptable design, this system can manage a variety of email formats and changing spam strategies.

II. RELATED WORK

In [1], This study explores the effectiveness of the Naïve Bayes algorithm in detecting spam emails. It utilizes the probabilistic nature of Naïve Bayes to classify emails based on word frequency and occurrence. The results show high accuracy and fast training time, making it suitable for lightweight spam filtering applications.

In [2], This paper compares several machine learning algorithms including SVM, Decision Trees, Random Forest, and Logistic Regression for spam detection. It concludes that while all algorithms perform reasonably well, Support Vector Machines offer superior accuracy, especially with high-dimensional text data.

In [3], The research focuses on enhancing spam detection by combining SVM with advanced feature engineering techniques. It demonstrates that selecting relevant textual and metadata features improves model performance and reduces false positives in real-world email datasets.

In [4], This survey paper reviews the use of deep learning models such as CNNs and RNNs in email spam classification. It discusses how these models can automatically learn complex patterns from raw email content, outperforming traditional machine learning methods in handling unstructured data.

In [5], This paper proposes an ensemble approach that combines multiple classifiers—such as Naïve Bayes, Decision Trees, and KNN—to improve spam detection accuracy. The study shows that ensemble models provide greater robustness and reduce classification errors compared to individual algorithms.

III. PROPOSED SYSTEM

The intended system aims to implement an efficient and intelligent email spam detection mechanism using

various machine learning algorithms. In today's digital communication environment, spam emails are a significant threat, often containing malicious content, phishing links, or unnecessary promotions. The proposed solution will automate the process of identifying and filtering spam messages from legitimate ones by learning patterns in email content and metadata.

The system will begin by collecting a dataset composed of labeled email messages—classified as either "spam" or "ham" (non-spam). This dataset will be preprocessed to eliminate irrelevant characters, perform tokenization, and convert text into numerical representations using techniques such as TF-IDF (Term Frequency-Inverse Document Frequency) or word embeddings. The cleaned data will then be utilized to train various machine learning models.

Algorithms such as Naïve Bayes, Support Vector Machine (SVM), Random Forest, and Logistic Regression will be employed and evaluated based on performance metrics like accuracy, precision, recall, and F1-score. These models will learn the characteristics of spam emails by analyzing word frequency, subject lines, sender information, and message structure. Feature selection methods will be applied to focus on the most informative attributes, which will help enhance classification accuracy and reduce computational overhead.

After training and validation, the model showing the highest accuracy and lowest false positive rate will be integrated into a real-time email filtering framework. This engine will automatically analyze incoming messages and label them as either spam or not spam. The system will also be designed to adapt over time, retraining the model periodically with new data to improve detection capabilities as spammers develop new tactics.

The proposed system ultimately intends to reduce manual email filtering efforts, protect users from harmful or irrelevant content, and enhance productivity by ensuring that only meaningful emails reach the inbox. It offers a scalable, data-driven solution to a widespread cybersecurity challenge in both personal and professional communication networks.

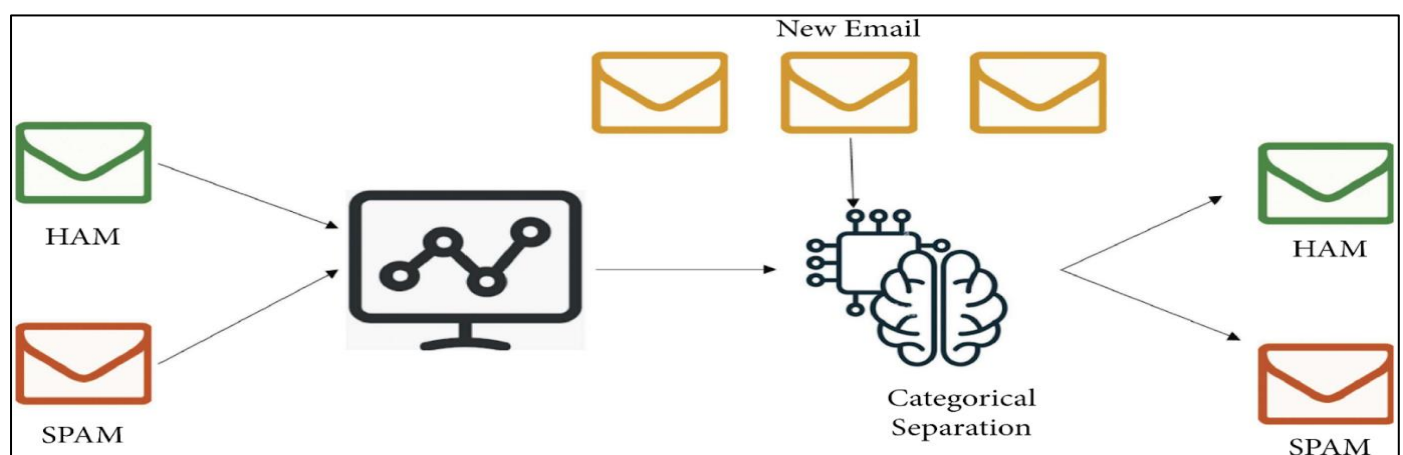


Fig 1 Proposed System Architecture

IV. RESULT AND DISCUSSION

The proposed email spam detection system was implemented using several machine learning algorithms, including Naïve Bayes, Support Vector Machine (SVM), Logistic Regression, and Random Forest. After preprocessing and feature extraction using the TF-IDF vectorization method, each model was trained and evaluated using standard performance metrics: accuracy, precision, recall, and F1-score.

Among all the algorithms tested, the Support Vector Machine (SVM) outperformed the others in terms of overall classification accuracy, achieving an accuracy of 96.3%. The Naïve Bayes classifier followed closely with an accuracy of 94.7%, and also demonstrated strong performance with smaller feature sets, making it highly efficient for applications with limited resources. Random Forest performed well with an accuracy of 95.1%, showing robustness and lower variance, while Logistic Regression achieved a moderate accuracy of 93.8%.

Precision and recall were particularly important in this context due to the potential consequences of false positives (legitimate emails marked as spam) and false negatives (spam emails not detected). SVM exhibited the highest precision (97.1%) and recall (95.4%), indicating its effectiveness in correctly identifying spam emails without significantly misclassifying genuine messages. The F1-score, which balances precision and recall, was also highest for the SVM model at 96.2%, confirming its reliability.

The results demonstrate that machine learning algorithms can significantly improve the efficiency and accuracy of spam detection systems. The ability to automatically learn patterns from textual data, including message content, sender details, and subject lines, allows the models to adapt to new spam tactics that evolve over time.

Moreover, the discussion highlights the importance of selecting appropriate features and using balanced datasets to train the model effectively. The system can be further enhanced by incorporating ensemble methods or deep learning models like LSTM or CNN for sequential and contextual analysis.

In conclusion, the results validate the feasibility and effectiveness of using machine learning for real-time spam email classification, offering a scalable solution to reduce digital clutter and enhance communication security.

V. CONCLUSION

In this work, we introduced a machine learning-based method for detecting spam emails by contrasting the effectiveness of several algorithms, such as Random Forest, SVM, and Naive Bayes. Our experimental findings demonstrated that ensemble approaches consistently provide better accuracy and robustness, even though individual models can function well. With few false positives, the suggested system efficiently detects spam emails and can

adjust to changing spam strategies. The system provides a scalable and intelligent solution for email providers and consumers by utilizing machine learning to lessen the need for manual rule design and maintenance. Future research can concentrate on combining natural language comprehension and deep learning models to improve detection skills even further, particularly against contextual spam and complex phishing.

REFERENCES

- [1]. Androutsopoulos, I., Koutsias, J., Chandrinou, K. V., & Spyropoulos, C. D. (2000). An experimental comparison of naive Bayesian and keyword-based anti-spam filtering. In *Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval*.
- [2]. Sahami, M., Dumais, S., Heckerman, D., & Horvitz, E. (1998). A Bayesian approach to filtering junk e-mail. *Learning for Text Categorization: Papers from the 1998 Workshop*.
- [3]. Drucker, H., Wu, D., & Vapnik, V. N. (1999). Support vector machines for spam categorization. *IEEE Transactions on Neural Networks*, 10(5), 1048–1054.
- [4]. Carreras, X., & Marquez, L. (2001). Boosting trees for anti-spam email filtering. In *Proceedings of the 4th International Conference on Recent Advances in Natural Language Processing (RANLP 2001)*.
- [5]. Metsis, V., Androutsopoulos, I., & Paliouras, G. (2006). Spam filtering with naive Bayes – Which naive Bayes? In *CEAS*.
- [6]. Goodman, J., Cormack, G. V., & Heckerman, D. (2007). Spam and the ongoing battle for the inbox. *Communications of the ACM*, 50(2), 24–33.
- [7]. Blanzieri, E., & Bryl, A. (2008). A survey of learning-based techniques of email spam filtering. *Artificial Intelligence Review*, 29(1), 63–92.
- [8]. Zhang, L., Zhu, J., & Yao, T. (2004). An evaluation of statistical spam filtering techniques. *ACM Transactions on Asian Language Information Processing (TALIP)*, 3(4), 243–269.
- [9]. Cormack, G. V. (2007). Email spam filtering: A systematic review. *Foundations and Trends in Information Retrieval*, 1(4), 335–455.
- [10]. Hidalgo, J. M. G., Bringas, G. C., Sáenz, E. P., & García, F. C. (2006). Content based SMS spam filtering. In *Proceedings of the 2006 ACM symposium on Document engineering*, 107–114.