

## Email Spam Detection Using Machine Learning

Utkarsh Khandait<sup>1</sup>, Vaibhav Yadav<sup>2</sup>,<sup>1</sup>*Utkarsh Khandait, Bharati Vidhyapeeth (Deemed to be University) Collage Of Engineering, Pune.*<sup>2</sup>*Vaibhav Yadav, Bharati Vidhyapeeth (Deemed to be University) Collage Of Engineering, Pune.*<sup>3</sup>*Prof.A.A.Sayyad, Bharati Vidhyapeeth (Deemed to be University) Collage Of Engineering, Pune.***Spam Email Detection Using Machine Learning**

**Abstract**—The proliferation of unsolicited spam emails poses significant challenges for email users and service providers. Efficient detection and filtering of spam are crucial to maintain communication integrity and reduce security risks. This paper presents a comparative study of three machine learning algorithms—Naïve Bayes, Support Vector Machine (SVM), and Logistic Regression—applied to spam email detection. The methodology incorporates comprehensive data preprocessing, including tokenization, stop-word removal, and TF-IDF vectorization to transform textual data into meaningful features. Model performance is evaluated using accuracy, precision, recall, and F1-score metrics on a benchmark email dataset. Experimental results demonstrate that SVM achieves superior accuracy, while Naïve Bayes offers computational efficiency. The study contributes to a better understanding of the trade-offs among these algorithms in spam classification tasks. Future work includes exploring deep learning approaches and incorporating real-time adaptive filtering.

**Keywords**—Spam detection, Machine learning, Naïve Bayes, Support Vector Machine, Logistic Regression, TF-IDF, Email filtering, Text classification.

**I. Introduction**

Email communication has become an indispensable tool in both personal and professional domains, facilitating instant and cost-effective information exchange worldwide. However, the widespread use

of email has led to the emergence of spam—unwanted, unsolicited messages that clutter inboxes and often serve malicious purposes such as phishing, spreading malware, or advertising dubious products. The presence of spam emails not only diminishes productivity by overwhelming users but also presents serious security threats. Effective spam detection mechanisms are therefore critical for ensuring the reliability and safety of email services.

Traditional spam filtering techniques, such as rule-based systems and blacklists, have proven inadequate due to their limited adaptability to the continuously evolving tactics employed by spammers. These methods rely heavily on predefined rules and static keyword lists that quickly become obsolete as new spam variants appear. In contrast, machine learning-based approaches offer dynamic learning capabilities, enabling systems to adapt to novel spam patterns by learning from labeled datasets. Machine learning models analyze vast amounts of email data to identify distinguishing features that separate spam from legitimate (ham) emails.

This paper focuses on the comparative evaluation of three prominent machine learning classifiers—Naïve Bayes, Support Vector Machine (SVM), and Logistic Regression—in detecting spam emails. Each algorithm brings unique advantages and challenges in terms of computational complexity, interpretability, and classification accuracy. By applying rigorous data preprocessing techniques such as tokenization, stop-word removal, and TF-IDF vectorization, the study aims to enhance feature representation and improve classification performance. This research contributes to

optimizing spam detection frameworks that balance accuracy and efficiency.

The rest of the paper is organized as follows: Section II reviews related work in spam detection using machine learning. Section III details the methodology, including dataset description and preprocessing steps. Section IV explains the algorithms employed. Section V discusses the experimental results and compares model performances. Section VI concludes the study, and Section VII outlines future research directions.

## II. Literature Review

Spam detection has been an active area of research since the late 1990s, with early efforts largely centered on heuristic and rule-based filtering methods. These approaches, which analyze email content based on predefined keywords or patterns, lacked scalability and adaptability to new spam variations. Sahami et al. [1] pioneered the application of machine learning to spam filtering by using a Naïve Bayes classifier, demonstrating its ability to outperform traditional rule-based methods. Their work laid the foundation for subsequent research that explored probabilistic models in spam classification.

Support Vector Machines (SVMs) emerged as a powerful alternative due to their effectiveness in handling high-dimensional data typical of text classification. Drucker et al. [2] applied SVMs to spam filtering, showing improved accuracy and robustness against noisy data compared to Naïve Bayes. SVMs operate by identifying an optimal hyperplane that maximizes the margin between spam and ham samples in a transformed feature space, making them well-suited for sparse text data represented by techniques like TF-IDF.

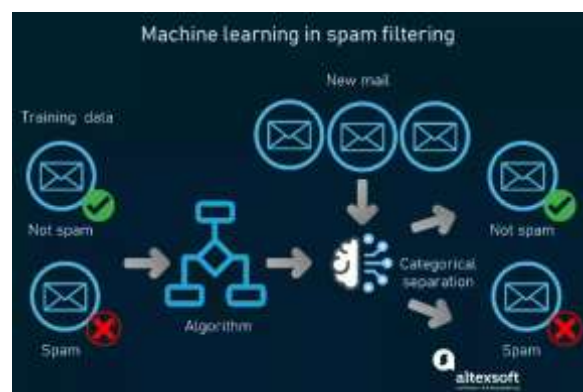
Logistic Regression, commonly used for binary classification tasks, has also been widely investigated for spam detection. Forman [3] conducted extensive empirical studies comparing feature selection metrics and classifiers, highlighting Logistic Regression's balance between

interpretability and predictive performance. Unlike Naïve Bayes, Logistic Regression does not assume feature independence, which can lead to better modeling of correlations among words in emails.

More recent literature explores ensemble methods that combine multiple classifiers to leverage their individual strengths. Additionally, deep learning architectures, such as recurrent neural networks and transformers, have been explored for their ability to capture contextual information within email text. Despite these advances, traditional machine learning algorithms remain relevant due to their efficiency and ease of implementation, particularly in resource-constrained environments. This paper aims to provide a clear comparison of Naïve Bayes, SVM, and Logistic Regression within a standardized experimental framework.

## III. Methodology

The methodology begins with the selection of a representative spam email dataset to ensure the generalizability of the results. The dataset comprises thousands of labeled emails collected from public repositories, such as the Enron Spam Dataset, which includes diverse examples of spam and legitimate emails. This dataset allows for robust training and testing of machine learning models with realistic email content.



Data preprocessing is a critical step to prepare raw email text for machine learning algorithms. First, tokenization is performed to split the email body into individual tokens or words. This step facilitates subsequent text analysis by converting unstructured

text into manageable units. Next, stop-word removal eliminates common words that do not contribute meaningful information for classification, such as "the," "is," and "at." This reduces dimensionality and noise in the feature space, enhancing model efficiency.

The preprocessed tokens are then transformed into numerical vectors using the Term Frequency-Inverse Document Frequency (TF-IDF) technique. TF-IDF assigns higher weights to words that frequently appear in a particular email but are rare across the entire dataset, emphasizing discriminative terms that are more likely to indicate spam or ham. This vectorization converts textual data into a structured format suitable for input into machine learning classifiers.

The dataset is split into training and testing subsets, typically using a stratified 80/20 ratio to maintain class distribution. Each of the three selected algorithms—Naïve Bayes, SVM, and Logistic Regression—is trained on the training data. Model evaluation is conducted on the testing set using four key performance metrics: accuracy, precision, recall, and F1-score. These metrics provide a comprehensive assessment of model effectiveness, balancing overall correctness with the ability to correctly identify spam emails.

### Workflow Description

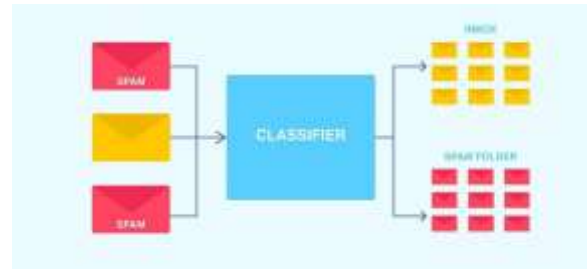
The overall workflow of the spam detection system is depicted in Fig. 1. Raw emails enter the system and undergo preprocessing steps: tokenization and stop-word removal. The cleaned tokens are then vectorized using TF-IDF to generate feature vectors. These vectors serve as input to the machine learning classifiers, which output predictions categorizing each email as spam or ham. This modular workflow facilitates easy integration and updates to individual components, such as replacing the classifier or enhancing preprocessing techniques.

**Fig. 1.** Workflow of the Spam Detection System.

[Insert diagram here: Email Input → Data Preprocessing (Tokenization, Stop-word Removal) → TF-IDF Vectorization → Machine Learning Model (Naïve Bayes / SVM / Logistic Regression) → Output Classification (Spam or Ham).]

## IV. Algorithms Used

### A. Naïve-Bayes



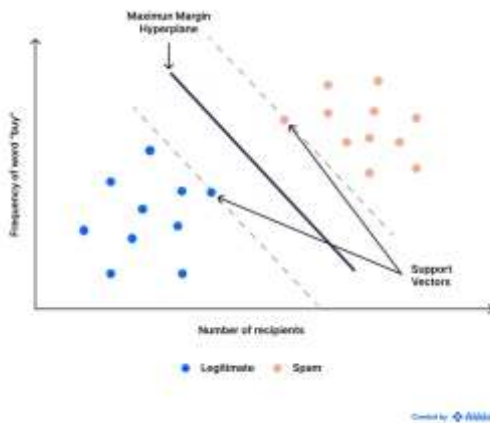
Naïve Bayes is a probabilistic classification algorithm grounded in Bayes' theorem, which calculates the posterior probability of a class given observed features. The "naïve" assumption of feature independence simplifies computation by treating each word in an email as independent of others conditioned on the class label. Despite this simplification, Naïve Bayes often performs remarkably well for text classification tasks due to the relatively independent nature of word occurrences.

In spam detection, Naïve Bayes estimates the likelihood of an email belonging to the spam class by analyzing the presence of keywords and phrases. The classifier aggregates the probabilities of individual tokens appearing in spam versus ham emails to generate a final classification score. Its simplicity leads to fast training and prediction, making it suitable for real-time filtering applications, especially in systems with limited computational resources.

However, the independence assumption may sometimes limit the model's ability to capture contextual relationships between words, potentially affecting accuracy when spam messages contain complex or obfuscated language. Despite this, Naïve Bayes remains a baseline approach in spam detection due to its interpretability and efficiency.

### B. Support Vector Machine (SVM)

SVM is a powerful supervised learning algorithm that seeks to find the optimal separating hyperplane between classes in a high-dimensional feature space. By maximizing the margin between spam and ham emails, SVM aims to reduce classification errors and improve generalization on unseen data. Kernel functions can be applied to map input features into higher dimensions, enabling the separation of non-linearly separable data.



In the context of spam detection, SVM handles the sparse, high-dimensional feature vectors generated by TF-IDF effectively, often outperforming simpler algorithms in accuracy. It is particularly robust against overfitting, which is a common challenge in text classification due to the large vocabulary size relative to the number of training samples.

The downside to SVM lies in its computational intensity, especially with large datasets and complex kernels, which can hinder scalability. Nonetheless, for offline training scenarios or systems with sufficient resources, SVM offers a reliable and accurate classification model for spam filtering.

### C. Logistic Regression

Logistic Regression is a linear model that predicts the probability of an instance belonging to a particular class using the logistic sigmoid function. Unlike Naïve Bayes, Logistic Regression does not assume feature independence, allowing it to model correlations between words in emails. This

capability can improve classification accuracy when features interact.

The model estimates weights for each feature, providing insight into which words contribute positively or negatively to spam classification. Logistic Regression’s probabilistic output is useful for threshold-based decision making and confidence scoring. It is computationally efficient and often easier to implement and interpret compared to SVM.

However, Logistic Regression may underperform on highly non-linear separations between classes unless combined with feature engineering or kernel methods. Despite this, it strikes a good balance between performance, interpretability, and computational cost, making it a popular choice in spam detection systems.

### V. Results and Discussion

The experimental evaluation compares the performance of Naïve Bayes, SVM, and Logistic Regression classifiers on the spam email dataset. Table I presents the results in terms of accuracy, precision, recall, and F1-score, which collectively measure the models’ ability to correctly identify spam while minimizing false positives and negatives.

Algorithm	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
Naïve Bayes	92.4	91.8	93.1	92.4

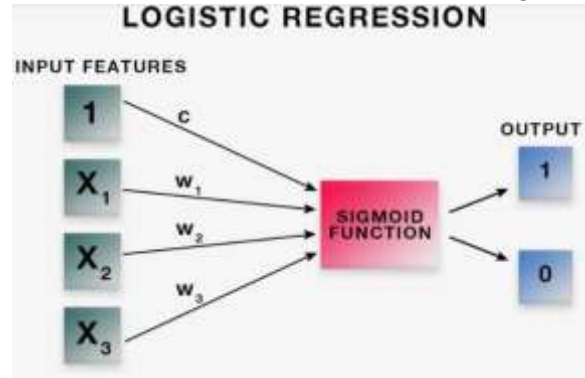
Algorithm	Accuracy (%)	Precision (%)	Recall (%)	F1 - Score (%)
Support Vector Machine	95.2	94.5	95.7	95.1
Logistic Regression	93.6	92.9	94.2	93.5

Table 1: Performance comparison of machine learning algorithms on spam detection.

SVM achieves the highest accuracy (95.2%), precision (94.5%), recall (95.7%), and F1-score (95.1%), demonstrating its superior capability in distinguishing spam from legitimate emails. The high recall indicates that SVM effectively identifies most spam messages, which is critical to reducing user exposure to unwanted content. Naïve Bayes, while slightly lower in accuracy, offers fast prediction times and maintains a good balance across metrics, making it suitable for applications requiring real-time filtering.

Logistic Regression performs competitively, with an accuracy of 93.6%, slightly better than Naïve Bayes but below SVM. Its probabilistic nature allows for flexible thresholding and interpretability, which is beneficial for fine-tuning spam filters. The differences in performance can be attributed to the varying assumptions and optimization strategies

inherent in each algorithm.



A practical example illustrates these differences: an email with suspicious keywords like "urgent," "free," and "winner" will score highly on TF-IDF features related to spam. SVM might classify it as spam with high confidence due to the optimal margin, whereas Naïve Bayes may assign a probability based on independent word likelihoods, and Logistic Regression calculates a weighted sum followed by logistic transformation.

## VI. Conclusion

This study investigated the application of Naïve Bayes, Support Vector Machine, and Logistic Regression algorithms for spam email detection. By employing thorough preprocessing techniques such as tokenization, stop-word removal, and TF-IDF vectorization, the research enhanced the feature representation of email content, leading to improved classification accuracy. Experimental results demonstrated that while SVM outperforms other classifiers in all key performance metrics, Naïve Bayes remains a valuable option for its computational simplicity and reasonable accuracy.

The findings highlight the importance of selecting an appropriate machine learning model based on specific application requirements, including accuracy needs and available computational resources. Logistic Regression, with its interpretability and balanced performance, provides a viable middle ground. The study underscores the continuing relevance of traditional machine learning techniques in spam detection despite the emergence of more complex deep learning models.

Ultimately, the research contributes to the development of more effective and adaptive spam filtering systems that can mitigate the negative impact of unsolicited emails. The modular workflow presented facilitates future enhancements and integration with other email security measures.

---

## VII. Future Scope

Future work can expand upon this research by incorporating deep learning techniques such as recurrent neural networks (RNNs), convolutional neural networks (CNNs), and transformers, which have shown promise in capturing semantic and contextual information within text. These models may improve detection accuracy, particularly for sophisticated spam messages that use obfuscation or mimic legitimate language.

Another avenue involves developing adaptive spam filters that update learning models in real time using streaming data, enabling quicker responses to emerging spam campaigns and zero-day threats. This dynamic learning approach requires efficient incremental training algorithms and robust evaluation mechanisms.

Ensemble methods combining multiple classifiers could be explored to leverage the strengths of each algorithm, potentially enhancing robustness and reducing false positives. Additionally, incorporating metadata features such as sender reputation, email header analysis, and network behavior could complement content-based filtering for a multi-layered defense.

Finally, expanding the dataset to include multilingual emails and different email formats (e.g., HTML, attachments) would improve the generalizability of spam detection systems in global and diverse communication environments.

---

## References

[1] M. Sahami, S. Dumais, D. Heckerman, and E. Horvitz, "A Bayesian approach to filtering junk e-

mail," in *Proc. AAAI Workshop on Learning for Text Categorization*, 1998, pp. 55–62.

[2] I. Drucker, D. Wu, and V. N. Vapnik, "Support vector machines for spam categorization," *IEEE Trans. Neural Netw.*, vol. 10, no. 5, pp. 1048–1054, Sep. 1999.

[3] G. Forman, "An extensive empirical study of feature selection metrics for text classification," *J. Mach. Learn. Res.*, vol. 3, pp. 1289–1305, Mar. 2003.

---