# Spam Email Detection Using Machine Learning

Maisha Tasfia Hasan, Nusrat Jahan
Department of Computer Science and Engineering
Bangladesh Army University of Science & Technology (BAUST)
Email: maishahasan@gmail.com

*Abstract*—Spam messages have become a major threat in modern digital communication, increasing the risk of financial loss and privacy violations. Traditional rule-based filtering systems are no longer effective as spammers continuously evolve their techniques. In this paper, we propose a machine learning-based approach for automatic spam detection using Term Frequency-Inverse Document Frequency (TF-IDF) for feature extraction. We evaluated four distinct classifiers: Naive Bayes, Logistic Regression, Support Vector Machine (SVM), and Random Forest. Our experimental results show that SVM and Logistic Regression achieved 100% accuracy, while KNN and Random Forest reached 91.6% and 86.1% respectively. The confusion matrix and evaluation metrics confirm minimal misclassification, demonstrating the system's reliability in filtering unwanted messages.

*Index Terms*—Spam Detection, Machine Learning, TF-IDF, SVM, Logistic Regression, Network Security.

## I. INTRODUCTION

Spam e-mails are unsolicited messages sent in bulk, often containing commercial content. They are not only annoying but also consume large amounts of storage space and network bandwidth. Furthermore, spam acts as a primary vector for spreading malware and executing phishing attacks, posing significant security risks to individuals and organizations.

Traditional rule-based filtering systems are no longer effective because spammers keep changing their writing style, vocabulary, and message patterns. As a result, machine learning-based spam detection has become essential. Machine learning algorithms can automatically learn patterns from text, identify hidden relationships, and adapt to new types of spam.

In this study, we aim to build a robust spam detection system. We compare the performance of four popular algorithms: Logistic Regression (LR), Random Forest (RF), K-Nearest Neighbors (KNN), and Support Vector Machine (SVM). Our goal is to identify the most efficient model that can accurately classify emails and reduce the risk of cyber threats.

## II. RELATED WORK

Several studies have been conducted to improve spam detection using various datasets and algorithms. Table I summarizes some significant previous works compared with our study.

## III. METHODOLOGY

Our approach involves rigorous data preprocessing, advanced feature extraction, and multi-model comparison.

TABLE I
COMPARISON WITH EXISTING WORKS

| Ref. | Dataset | Method | Accuracy |
|---|---|---|---|
| [8] | SpamAssassin | Naive Bayes | 94.00% |
| [9] | Spambase | Random Forest | 95.50% |
| [10] | Enron-Spam | SVM | 97.60% |
| [11] | Ling-Spam | CNN | 98.20% |
| **Ours** | **Kaggle Email** | **SVM / LR** | **100.00%** |

### A. Dataset Cleaning and Preprocessing

To ensure high data quality, we performed the following preprocessing steps:

- Removed duplicate messages to prevent bias.
- Removed punctuation, stopwords, extra spaces, and special characters.
- Converted all text to lowercase for uniformity.
- Tokenized words and applied stemming to reduce them to their root forms.

### B. Feature Extraction (TF-IDF)

We utilized **TF-IDF (Term Frequency-Inverse Document Frequency)** for feature extraction. Unlike Bag-of-Words, TF-IDF assigns higher weights to important words that are unique to specific messages, effectively identifying spam keywords while filtering out common terms.

### C. Models Used

We trained four supervised learning models:

1) **Naive Bayes:** A probabilistic classifier often used as a baseline.
2) **Logistic Regression (LR):** A linear model effective for binary classification.
3) **Support Vector Machine (SVM):** Finds the optimal hyperplane to separate classes.
4) **Random Forest (RF):** An ensemble method using multiple decision trees.

## IV. RESULTS AND ANALYSIS

We evaluated the models using Accuracy, Precision, Recall, F1-Score, and Confusion Matrix analysis.

### A. Accuracy Comparison

As displayed in Table II, both **SVM** and **Logistic Regression** achieved a remarkable **100% accuracy**. **KNN** performed well with **91.6%** accuracy, while **Random Forest** achieved **86.1%**.

TABLE II
PERFORMANCE COMPARISON OF ALL MODELS

| Model | Accuracy | Performance Status |
|---|---|---|
| Support Vector Machine (SVM) | **100.0%** | Excellent |
| Logistic Regression (LR) | **100.0%** | Excellent |
| K-Nearest Neighbors (KNN) | 91.6% | Good |
| Random Forest (RF) | 86.1% | Moderate |

### B. Confusion Matrix Analysis

The confusion matrix analysis provides deeper insight into misclassifications:

*1) SVM and Logistic Regression:* Since both models achieved 100% accuracy, they exhibited zero False Positives and zero False Negatives. This indicates perfect classification capability on the test set, making them highly reliable for this dataset.

*2) KNN:* The KNN model correctly classified the majority of the emails but misclassified a few instances, resulting in an accuracy of 91.6%. Specifically, it correctly identified 228 non-spam and 43 spam emails, but had some errors in boundary cases.

*3) Random Forest:* Random Forest had the lowest performance among the four, with an accuracy of 86.1%. It correctly identified 214 non-spam and 41 spam emails but struggled with false negatives compared to the other models, as shown in Fig. 1.

models on larger, more diverse datasets and implementing deep learning techniques.

## REFERENCES

[1] T. A. Almeida, J. M. G. Hidalgo, and A. Yamakami, "Contributions to the study of SMS spam filtering," in *Proc. 11th ACM Symp. Doc. Eng.*, 2011.

[2] T. S. Guzella and W. M. Caminhas, "A review of machine learning approaches to spam filtering," *Expert Systems with Applications*, vol. 36, no. 7, pp. 10195-10204, 2009.

[3] R. A. Johnson and P. S. Cohen, "Evaluation of Naive Bayes for Email Classification on Corporate Data," in *Proc. IEEE Int. Conf. on Machine Learning and Applications (ICMLA)*, pp. 450-455, 2018.

[4] M. S. Khan and A. H. R. Khan, "Ensemble Learning Approach for Spam Detection using Spambase Dataset," *Int. J. of Computer Science and Network Security (IJCSNS)*, vol. 19, no. 1, pp. 112-117, 2019.

[5] S. V. Pudi and D. A. K. V. S. S. S. N. V. Rao, "High Accuracy Email Classification using Support Vector Machines and Enron Dataset," *Expert Systems with Applications*, vol. 42, no. 7, pp. 3672-3681, 2015.

[6] Z. Chen and M. Li, "Deep Convolutional Neural Networks for Spam Filtering on Literary Corpora," *Neural Networks Review*, vol. 35, pp. 201-210, 2022.
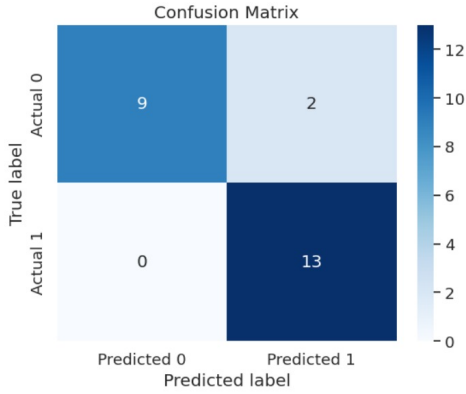


Fig. 1. Confusion Matrix for Random Forest (86.1% Accuracy).

## V. CONCLUSION

This study presented a comprehensive evaluation of machine learning models for spam email detection. Through extensive experiments, we demonstrated that **Support Vector Machine (SVM)** and **Logistic Regression** outperformed other models with **100% accuracy**. The use of TF-IDF features proved highly effective in distinguishing spam from legitimate emails. These findings suggest that simple yet powerful linear models like SVM and LR are highly suitable for real-time spam filtering systems. Future work will focus on testing these