

# Detecting spam in emails: "A comparative analysis of machine learning algorithms"

Bharat Ravindra Shinde

*M.L Dahanukar College of commerce*

*Mumbai University*

Email id: [shindebharat461@gmail.com](mailto:shindebharat461@gmail.com)

**Abstract:** *The fast-growing association of spam with scams, phishing, and malware makes them a major security threat to both individuals and organizations. Before machine learning-based spam filters were still efficient and could reach strong performance measures, the main point here is that these filters need to be more robust to issues like dataset shifts and adversarial manipulations. This research digs through the most popular machine learning algorithms (Naive et al. (SVM), Random Forest, and Long Short-Term Memory (LSTM)) in the case of spam detection. Experimental results show the link between accuracy, scalability, and computational efficiency, which implies that the algorithms should be flexible enough for the real world.*

## 1. Introduction

Spam emails, which sometimes take the form of unsolicited advertisements and other times phishing attempts, are a common

cybersecurity problem. Among all the facts, spam accounts for 50–85% of global email traffic, and sometimes technology is used to elude detection mechanisms. However, even though machine learning (ML) algorithms have become more successful at spam detection, the enhancement of algorithms is continually required due to sophisticated spam tactics.

This paper examines ML algorithms for email spam detection, dealing with their adaptability, performance, and computational demands in dynamic environments.

## 2. Literature Review

### Evolution of Spam Detection

Initial spam filters were simple keyword-based systems that often failed, giving false positives. The introduction of ML led to the use of scale-driven approaches based on data for spam detection.

**1. Naive Bayes:** A core algorithm for spam detection

that uses probabilistic models for text classification.

2. **Support Vector Machines (SVM):** Improved precision of methods working in high dimensions of the input space by means of a non-linear kernel.
3. **Random Forest:** An ensemble method that is very good for maintaining the balance between accuracy and computational complexity.
4. **Deep Learning Approaches:** Various types of neural networks, such as LSTM and CNN, have been invented to detect complicated spam patterns.

Challenges in Spam Detection

1. **The Dynamic Nature of Spam:** Spam that evolves with time makes machine models that stay the same drastically ineffective.
2. **Adversarial Manipulations:** On the other hand, spammers may include harmful content in the form of photos with embedded text to be undetected.

3. Methodology

Dataset

A dataset of 50,000 emails, comprising 60% legitimate and 40%

spam emails, was used. Features were extracted using Term Frequency-Inverse Document Frequency (TF-IDF) and Bag of Words (BoW) models.

Algorithms

1. **Naive Bayes:** Quick and resource-efficient but limited in handling complex patterns.
2. **SVM:** Effective in separating linear and non-linear data but computationally intensive.
3. **Random Forest:** Robust in handling diverse datasets with high accuracy.
4. **LSTM:** Captures sequential dependencies in text for unparalleled precision.

Performance was evaluated using metrics such as accuracy, precision, recall, F1-score, and training time.

4. Results

Comparative Analysis

Algorithm	Accuracy (%)	Precision	Recall	F1 - Score	Training Time
Naive Bayes	85.3	0.84	0.82	0.83	Fast
SVM	91.2	0.89	0.88	0.88	Moderate
Random Forest	93.5	0.92	0.91	0.92	Moderate
LSTM	95.2	0.95	0.95	0.95	Slow

## Key Insights

- **Naive Bayes** is the best classifier for resource-constrained environments. However, it cannot deal with complex spam patterns.
- **SVM** delivers a superior level of accuracy, but implementing it entails a substantial amount of processing power.
- **Random Forest** holds a useful compromise between accuracy and efficiency.
- **LSTM** is better at handling problems but requires more power, which can hinder its usage in practical applications.

## 5. Discussion

### Real-World Applicability

- 1. Enterprise Systems:** Random Forest and LSTM enhance the performance of enterprise email systems for spam detection due to their ability to handle large numbers of emails.
- 2. Consumer Applications:** The Naive Bayes algorithm brings about swift and efficient solutions for personal use.

## Limitations and Future Directions

- 1. Dataset Diversity:** A study was carried out on a dataset that might not sufficiently represent global email traffic. Further studies are expected to include multilingual and real-time datasets.
- 2. Scalability:** One of the problems that can occur in such deep learning models like LSTM is optimization for resource efficiency.

## 6. Conclusion

This study has explored ML algorithms for spam detection use, consecutively, LSTM as the most accurate one and Random Forest's balanced performance placed second. Naive Bayes is still in use in a low-resource environment. However, it becomes a question mark with the changing spam strategies that demand the synthesis of more sophisticated algorithms. Their main focus should be on the development of scalable, hybrid models and real-time spam detection systems, which will solve the inevitable problem of spam in the next few years.

## 7. References

1. Androutsopoulos, I., et al. (2000). An evaluation of naive Bayesian anti-spam filtering.
2. Islam, M. R., et al. (2013). Spam detection in email communication: A machine learning perspective.
3. Gangavarapu, T., et al. (2020). Dynamic spam detection using machine learning.