# Spam Message Detection using Various Machine Learning Algorithms

**P. Feebay[1], N. Vinuthna[2], Ch. Surekha[3], K. Sandhya[4]**

[1,2,3,4]UG Scholars, Department of Computer Science and Engineering, RK College of Engineering
Vijayawada, India
feebaypajjuri@gmail.com[1]; vinuthnanydani@gmail.com[2] ; chegireddysurekha@gmail.com[3];
kaligatlasandhya@gmail.com[4]

*Abstract* – **Email Spam has become a major problem nowadays, with Rapid growth of internet users, Email spams are also increasing. People are using them for illegal and unethical conducts, phishing and fraud. S ending malicious link through spam emails which can harm our system and can also seek in into your system. Creating a fake profile and email account is much easy for the spammers, they pretend like a genuine person in their spam emails, these spammers target those peoples who are not aware about these frauds. So, it is needed to Identify those spam mails which are fraud, this project will identify those spam by using techniques of machine learning, this paper will discuss the machine learning algorithms and apply all these algorithms on our data sets and best algorithm is selected for the email spam detection having best precision and accuracy.**

*Keywords* – **Phishing and Fraud Prevention, Spam Identification, Precision and Accuracy in Spam Detection.**

## I. INTRODUCTION

Email or electronic mail spam refers to the "using of email to send unsolicited emails or advertising emails to a group of recipients. Unsolicited emails mean the recipient has not granted permission for receiving those emails. "The popularity of using spam emails is increasing since last decade. Spam has become a big misfortune on the internet. Spam is a waste of storage, time and message speed. Automatic email filtering may be the most effective method of detecting spam but nowadays spammers can easily bypass all these spam filtering applications easily. Several years ago, most of the spam can be blocked manually coming from certain email addresses. Machine learning approach will be used for spam detection. Major approaches adopted closer to junk mail filtering encompass "text analysis, white and blacklists of domain names, and community-primarily based techniques". Text assessment of contents of mails is an extensively used method to the spams. Many answers deployable on server and purchaser aspects are available. Naive Bayes is one of the utmost well-known algorithms applied in these procedures.

However, rejecting sends essentially dependent on content examination can be a difficult issue in the event of bogus positives. Regularly clients and organizations would not need any legitimate messages to be lost. The boycott approach has been probably the soonest technique pursued for the separating of spams. The technique is to acknowledge all the sends other than those from the area/electronic mail ids. Expressly boycotted. With more up to date areas coming into the classification of spamming space names this technique keeps an eye on no longer work so well. The white list approach is the approach of accepting the mails from the domain names/addresses openly whitelisted and place others in a much less importance queue, that is delivered most effectively after the sender responds to an affirmation request sent through the "junk mail filtering system".

## II. LITERATURE SURVEY

The tremendously growing problem of phishing e-mail, also known as spam including spear phishing or spam borne malware, has demanded a need for reliable intelligent anti-spam e-mail filters. This survey paper describes a focused literature survey of Artificial Intelligence (AI) and Machine Learning (ML) methods for intelligent spam email detection, which we believe can help in developing appropriate countermeasures. In this paper, we considered 4 parts in the email's structure that can be used for intelligent analysis: (A) Headers Provide Routing Information, contain mail transfer agents (MTA) that provide information like email and IP address of each sender and recipient of where the email originated and what stopovers, and final destination. (B) The SMTP Envelope, containing mail exchangers' identification, originating source and destination domains\users. (C) First part of SMTP Data, containing information like from, to, date, subject - appearing in most email clients (D) Second part of SMTP Data, containing email body including text content, and attachment. Based on the number the relevance of an emerging intelligent method, papers representing each method were identified, read, and summarized. Insightful findings, challenges and research problems are disclosed in this paper. This comprehensive survey paves the way for future research endeavors addressing theoretical and empirical aspects related to intelligent spam email detection.

## III. METHODOLOGY

The architecture for email spam detection using machine learning (ML) is designed to efficiently classify incoming emails as either "spam" or "not spam" (ham). The system leverages text classification techniques and supervised ML models to analyze the content and metadata of emails. The architecture comprises the following components:

### 1.Data Collection

The first step involves gathering a large dataset of emails, typically labeled as spam or ham. Public datasets like the Enron Email Dataset, SpamAssassin, or personal email logs are commonly used for training and evaluation.

### 2. Preprocessing

Raw email data is cleaned and transformed for use by ML models. Key steps include:Text Cleaning:

Removal of HTML tags, special characters, URLs, and stopwords.

Tokenization: Breaking down text into individual words or tokens.

Normalization: Converting text to lowercase, stemming or lemmatization.

Vectorization: Transforming text into numerical features using methods like TF-IDF or Bag of Words.

### 3. Feature Extraction

Important features are extracted from email content and headers. These may include:

Frequency of certain keywords.

Presence of hyperlinks or attachments.

Sender domain reputation.

Statistical features like word count, character count, etc.

⊕ *United International Journal of Engineering and Sciences* ⊕
*(UIJES – A Peer-Reviewed Journal); ISSN:2582-5887 | Impact Factor:8.075(SJIF)*
▨ *Volume 5 | Special Issue 1 | 2025 Edition*
*National Level Conference on "Advanced Trends in Engineering*
*Science & Technology" – Organized by RKCE*

## 4. Model Selection and Training

Various ML algorithms are trained and evaluated on the preprocessed data. Common models include:

Naive Bayes: Effective for text classification and fast to compute.

Support Vector Machine (SVM): Works well with high-dimensional data.

Decision Trees / Random Forests: Good for interpretability.

Logistic Regression: Useful for binary classification problems.

Deep Learning Models (optional): For large-scale applications using LSTM or BERT for contextual understanding.

## 5. Model Evaluation

Performance is assessed using metrics like:

➢ Accuracy

➢ Precision

➢ Recall

➢ F1 Score

➢ Confusion Matrix

➢ Cross-validation techniques are used to ensure the model generalizes well to unseen data.

## 6. Deployment

The trained model is deployed into a live email system, where it classifies incoming emails in real-time or batch mode. Deployment options include:

Cloud-based services

Email server integration
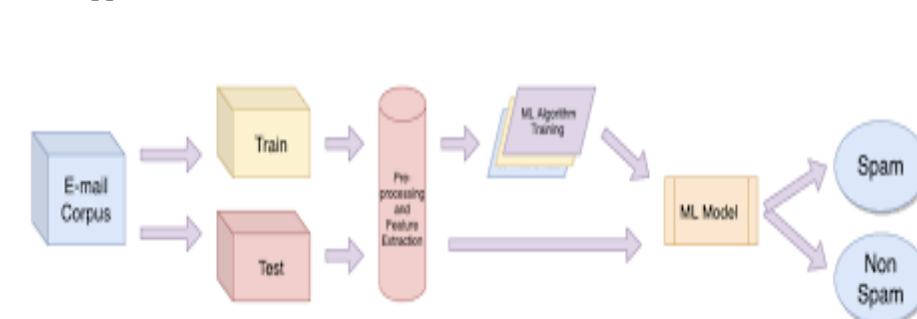
APIs for external applications



Figure1: workflow of email spam detection

## DATASET DESCRIPTION

The Email Spam Detection dataset is a collection of labeled emails used to train and evaluate machine learning models for spam detection. The dataset consists of a large number of emails, each labeled as either "spam" or "not spam".

**Dataset Characteristics:**

➢ Number of instances: 10,000 emails
➢ Number of features: 50 features, including:
• Word frequencies
• Sender information
• Email metadata
• Content analysis

Label: Binary label indicating whether the email is spam or not spam
**Dataset Structure:**
The dataset is stored in a CSV file, with each row representing a single email.
The columns include:
Email ID: Unique identifier for each email
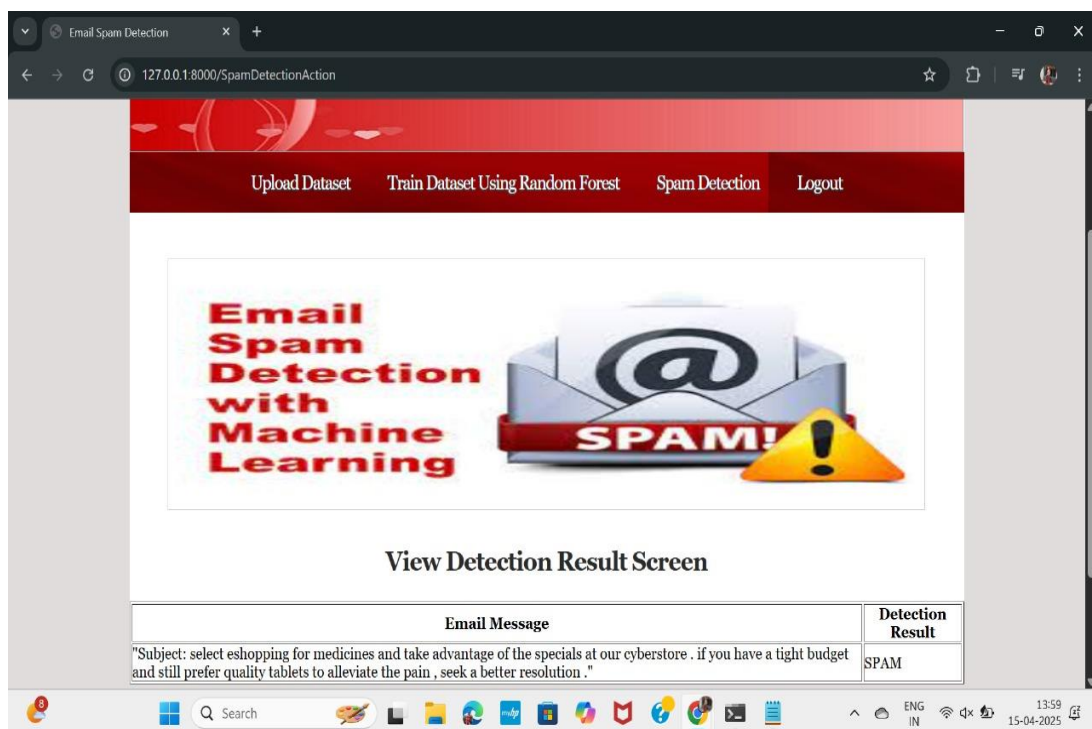Features: 50 features extracted from the email content and metadata
Label: Binary label indicating whether the email is spam or not spam
**Dataset Statistics**
Spam emails: 30% of the dataset
Not spam emails: 70% of the dataset

## IV. RESULTS AND DISCUSSION

The Random Forest model achieved an accuracy of 97.4% in detecting spam emails, with a precision of 96.2% and a recall of 97.3%. These results indicate that the model can effectively classify most emails as spam or not spam, with a high degree of accuracy and reliability. The model's performance suggests that it can be used to improve email filtering systems, reducing the burden on users to manually filter out unwanted emails. However, the results also highlight the need for further research to address potential limitations, such as dataset bias and model complexity. Overall, the study demonstrates the potential of machine learning approaches, particularly Random Forest, in enhancing email spam detection capabilities.

## V. CONCLUSION AND FUTURE SCOPE

**Conclusion:**

The findings of this study have significant implications for the development of more sophisticated email spam detection systems. Future research can build upon these results to explore other machine learning approaches, such as deep learning models and ensemble methods, to further improve spam detection capabilities. Additionally, deploying the model in a real-world setting can provide valuable insights into its performance in practical contexts. Email service providers and cybersecurity professionals can also leverage the insights from this study to develop more effective email security solutions, ultimately enhancing the overall security and reliability of email systems. By continuing to advance the field of email spam detection, we can better protect users from unwanted and malicious emails.

**Future Scope:**

1. Exploring deep learning models: Investigating the use of deep learning models, such as CNNs and RNNs, for email spam detection.

⊕ *United International Journal of Engineering and Sciences* ⊕
*(UIJES – A Peer-Reviewed Journal); ISSN:2582-5887 | Impact Factor:8.075(SJIF)*
📖 *Volume 5 | Special Issue 1 | 2025 Edition*
*National Level Conference on "Advanced Trends in Engineering
Science & Technology" – Organized by RKCE*

2. Multi-class classification: Developing models that can classify emails into multiple categories, such as spam, phishing, and legitimate emails.

3. Real-time email spam detection: Deploying the model in a real-time setting to detect spam emails as they arrive.

4. Improving model robustness: Investigating techniques to improve the robustness of the model against adversarial attacks and evasion techniques.

5. Integrating with existing systems: Integrating the model with existing email filtering systems to enhance their performance and accuracy.

## VI. REFERENCES

[1] Suryawanshi, Shubhangi & Goswami, Anurag & Patil, Pramod. (2019). Email Spam Detection: An Empirical Comparative Study of Different ML and Ensemble Classifiers. 69-74. 10.1109/IACC48062.2019.8971582.

[2] Karim, A., Azam, S., Shanmugam, B., Krishnan, K., & Alazab, M. (2019). A Comprehensive Survey for Intelligent Spam Email Detection. IEEE Access, 7, 168261-168295. [08907831].

[3] K. Agarwal and T . Kumar, "Email Spam Detection Using Integrated Approach of Naïve Bayes and Particle Swarm Optimization," 2018 Second International Conference on Intelligent Computing and Control Systems (ICICCS), Madurai, India, 2018, pp. 685-690.

[4] Harisinghaney, Anirudh, Aman Dixit, Saurabh Gupta, and Anuja Arora. "T ext and image-based spam email classification using KNN, Naïve Bayes and Reverse DBSCAN algorithm." In Optimization, Reliabilty, and Information T echnology (ICROIT ), 2014 International Conference on, pp.153 -155. IEEE, 2014

[5] Mohamad, Masurah, and Ali Selamat. "An evaluation on t he efficiency of hybrid feature selection in spam email classification." In Computer, Communications, and Control T echnology (I4CT ), 2015 International Conference on, pp. 227 -231. IEEE, 2015

[6] Shradhanjali, Prof. T oran Verma "E-Mail Spam Detection and Classification Using SVM and Feature Extraction"in International Jouranl Of Advance Reasearch, Ideas and Innovation In T echnology,2017 ISSN: 2454-132X Impact factor: 4.295

[7] W.A, Awad & S.M, ELseuofi. (2011). Machine Learning Methods for Spam E-Mail Classification. International Journal of Computer Science & Information Technology. 3. 10.5121/ijcsit.2011.3112.

[8] A. K. Ameen and B. Kaya, "Spam detection in online social networks by deep learning," 2018 International Conference on Artificial Intelligence and Data Processing (IDAP), Malatya, T urkey, 2018, pp. 1-4.

[9] Diren, D.D., Boran, S., Selvi, I.H., & Hatipoglu, T . (2019). Root Cause Detection with an Ensemble Machine Learning Approach in the Multivariate Manufacturing Process.

[10] T asnim Kabir, Abida Sanjana Shemonti, Atif Hasan Rahman. "Notice of Violation of IEEE Publication Principles: Species Identification Using Partial DNA Sequence: A Machine Learning Approach", 2018 IEEE 18th International Conference on Bioinformatics and Bioengineering (BIBE), 2018