# Comparative Analysis of Machine Learning Models for Email Fraud Detection: Random Forest, Decision Tree, and Multinomial Naïve Bayes

**1st Sakshi Kusmude**
*Department of Information Technology*
*Vishwakarma Institute of Information Technology, Pune*
sakshi.22210253@viit.ac.in

**2nd Shravani Padwal**
*Department of Information Technology*
*Vishwakarma Institute of Information Technology, Pune*
shravani.22211302@viit.ac.in

**3rd Chaitali Shewale**
*Department of Information Technology*
*Vishwakarma Institute of Information Technology, Pune*
chaitali.shewale@viit.ac.in

*Abstract*—*Email communication has become a cornerstone of personal and professional interactions, facilitating the exchange of information at an unprecedented scale. However, this rise in email usage has concurrently led to a significant increase in email fraud, encompassing various forms such as phishing attacks, deceptive offers, and fraudulent messages. As these threats evolve in complexity, it has become imperative to develop robust systems capable of real-time detection and prevention. This paper explores the application of machine learning techniques to address the pressing issue of email fraud, focusing on a comparative analysis of three widely adopted models: Random Forest Classifier (RFC), Decision Tree Classifier (DTC), and Multinomial Naive Bayes (MNB).*

*Each of these algorithms possesses unique characteristics that influence their effectiveness in detecting spam and legitimate emails. The study meticulously evaluates their performance based on key metrics including accuracy, precision, recall, and F1-score, providing a holistic understanding of how well each model can classify emails in diverse datasets. The examination also includes the construction and analysis of confusion matrices, which highlight the models' strengths and weaknesses in distinguishing between legitimate and fraudulent communications.*

*Moreover, the mathematical foundations underlying each model are elucidated, offering insights into how these algorithms process features and make classification decisions. Special attention is given to the importance of feature selection and data preprocessing, which significantly impact model performance, especially in scenarios involving large datasets with varied spam patterns. The research emphasizes the necessity of explainability in machine learning systems, particularly in security contexts, as it fosters trust among users and aids in the interpretation of model predictions.*

*Insights derived from this comparative analysis aim to equip developers, security analysts, and researchers with the knowledge needed to select the most suitable algorithm for real-world email fraud detection applications. The findings not only reveal raw performance metrics but also underscore the practical implications of each model in terms of scalability, speed, and interpretability. By enhancing understanding of these machine learning techniques, this study aspires to contribute to the development of more secure and efficient email environments, paving the way for a future where email fraud is managed effectively and becomes less of a pervasive threat.*

*In conclusion, as email fraud continues to pose challenges across various sectors, the insights from this research will serve as a valuable resource for stakeholders seeking to bolster their defenses against such threats. Through informed decisions about algorithm selection and implementation, the goal of achieving safer email communication can become a reality.*

*Keywords: Email Fraud Detection, Random Forest Classifier, Decision Tree Classifier, Multinomial Naive Bayes, Machine Learning, Spam Detection.*

## I. INTRODUCTION

The rise of phishing attacks continues to be a major threat to internet security, particularly as users increasingly rely on digital platforms for commu- nication, banking, and e-commerce. Research has demonstrated that phishing attacks often succeed due to users' lack of technical knowledge and the growing sophistication of phishing methods, which leverage fake websites, deceptive links, and malicious software to deceive victims [1], [2]. To combat these threats, machine learning (ML) approaches have proven to be highly effective, particularly for tasks like spam detection and email fraud prevention. For instance, a study showed that a machine learning framework combining ten models achieved a phishing detection accuracy of 97.27

Numerous machine learning techniques, including ensemble methods and deep learning, have been employed in phishing detection. One such study achieved a 99.03

As the need for stronger cybersecurity systems grows, phishing detection continues to be an area of active development. Advanced classification techniques, including hybrid models and ensemble approaches, have been explored to improve accuracy in detecting phishing

and spam [3], [4], [5]. Other studies have emphasized the importance of detecting phishing by analyzing patterns in email and website data [6], [7]. Some methods even combine traditional cybersecurity approaches with machine learning for enhanced detection [8]. Ongoing challenges persist due to the increasing sophistication of phishing attacks, but the use of deep learning and hybrid techniques promises to yield even more effective results [9], [10].

## II. RELATED WORK

Email fraud detection has become a critical area of research as the volume of fraudulent emails continues to rise. Various approaches have been developed over the years, each with its unique strengths and weaknesses.

Early methods for email fraud detection primarily relied on rule-based systems. These systems used predefined rules to identify suspicious patterns, such as specific keywords or header anomalies. While effective to some extent, these approaches often struggled to adapt to evolving tactics used by spammers, leading to a higher rate of false positives and negatives.

With the advancement of technology, machine learning classifiers have emerged as a powerful alternative for tackling the complexities of email fraud. These classifiers can learn from large datasets and identify intricate patterns that traditional rule- based systems might miss. Research has shown that models like Random Forest, Decision Trees, and Naïve Bayes are particularly well-suited for spam detection due to their ability to handle various features and adapt over time.

For instance, in recent studies, Random Forest classifiers have demonstrated a high degree of accuracy, making them a popular choice for email classification tasks. Their ensemble nature

allows them to combine multiple decision trees, resulting in more robust predictions. Decision Tree classifiers, on the other hand, offer interpretability, allowing researchers and practitioners to under- stand the reasoning behind specific classifications. Meanwhile, the Multinomial Naïve Bayes model is favored for its simplicity and efficiency, especially in processing textual data.

This paper builds upon this existing body of liter- ature by conducting a comprehensive comparison of these three models. We focus on key performance metrics such as precision, recall, and F1 score to evaluate their effectiveness in detecting email fraud. By understanding the strengths and limitations of each approach, we aim to provide insights into the best practices for enhancing email security in an increasingly digital world.

## III. METHODOLOGY

The research follows a systematic approach to detect email fraud using machine learning models. The methodology comprises several key steps: dataset collection, preprocessing, feature engineering, model selection, and evaluation. These steps are represented in the flowchart in Figure X.

### a) 1. Dataset Collection

The dataset is collected from a reliable source, containing labeled examples of fraudulent and non-fraudulent emails. The dataset includes both numerical and categorical features that are crucial for distinguishing fraud patterns.

### b) 2. Text Preprocessing

Since email data often contains text, preprocessing steps such as tokenization, stemming, and stopword removal are applied. These steps transform raw text into a format suitable for machine learning models. For this, each email is tokenized into words, and irrelevant words (e.g., "the," "and") are removed.

### c) 3. Categorical Encoding

The dataset may contain categorical variables that need to be encoded into numerical form. For instance, labels such as "spam" and "not spam" are encoded as 1 and 0, respectively. This step is essential for feeding the data into machine learning algorithms.

Let C represent the set of categorical variables:

$$\text{Encoded Value}(C) = \begin{cases} 1, & if\ C = not\ spam \\ 0, & if\ C = spam \end{cases}$$
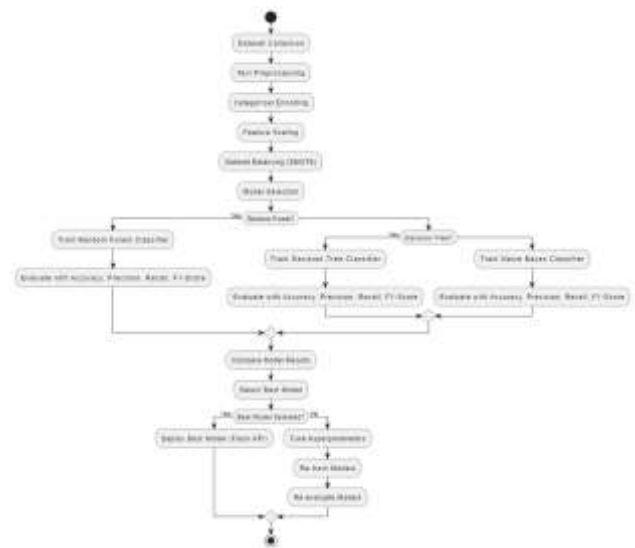
### a) 4. Feature Scaling

Feature scaling ensures that numerical features are standardized to have a mean of zero and a standard deviation of one. This step prevents features with larger numerical values from dominating the learning process.

Mathematically, for a feature $x_i$, the standardized value is given by:

$$x_i' = \frac{x_i - \mu}{\sigma}$$

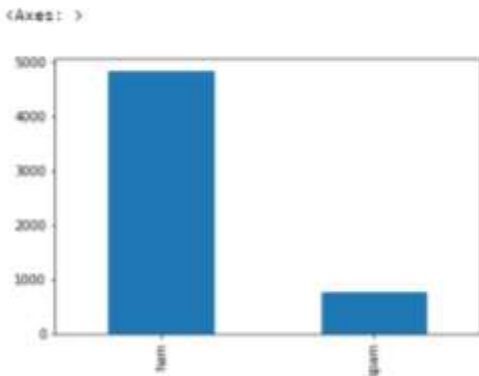where $\mu$ is the mean and $\sigma$ is the standard deviation.



1.1 Flow Diagram

Above given diagram is about the process flow diagram of the methodology.

### a) 5. Dataset Balancing (SMOTE)

The dataset is often imbalanced, with more instances of "not spam" than "spam." The Synthetic Minority Oversampling Technique (SMOTE) is applied to generate synthetic examples for the minority class (spam). SMOTE ensures that the model does not become biased towards the majority class. Below given is the dataset representation as per the spam and not spam.



1.2 Dataset Classification Graph

### b) 6. Model Selection

Three machine learning classifiers are considered in this study: Random Forest, Decision Tree, and Naive Bayes. Each model is trained and evaluated based on its accuracy, precision, recall, and F1-score.

- **Random Forest Classifier:** Random Forest is a widely-used ensemble learning method designed to improve the accuracy and stability of predictions by combining multiple decision trees. Instead of relying on just one tree, which can sometimes overfit or miss certain patterns in the data, Random Forest builds several trees on different random subsets of the data. Each tree in the "forest" makes its own prediction, and the final output is determined by the majority vote of these trees. The benefit of this approach is that the randomness in selecting the data subsets and features for each tree helps reduce overfitting, making the model more generalizable to unseen data. This technique is particularly useful when dealing with large datasets or when the data has many features, as it can capture complex relationships more effectively than a single decision tree. In addition to being robust, Random Forest also provides valuable insights, such as highlighting which features are most important for making predictions. This makes it a powerful tool in tasks where both accuracy and interpretability are important.

- **Decision Tree Classifier:** A single decision tree splits the data recursively based on the Gini index or information gain until all nodes are pure.

- **Naive Bayes Classifier:** This model assumes that the features are independent. It uses Bayes' Theorem to compute the posterior probability of the classes given the features.

### c) 7. Model Evaluation

Each model is evaluated using four metrics:

- **Accuracy:** The ratio of correctly predicted instances to the total instances.

- **Precision:** The ratio of correctly predicted positive observations to the total predicted positives.

- **Recall:** The ratio of correctly predicted positive observations to the actual positives.

- **F1-Score:** The harmonic mean of precision and recall, used when the class distribution is imbalanced.

The evaluation results are compared, and the best-performing model is selected for deployment.

### d) 8. Hyperparameter Tuning

If the initial results do not meet expectations, hyperparameter tuning is performed to optimize the performance of the models. GridSearchCV is used to systematically test different combinations of hyperparameters.

### e) 9. Model Deployment

Once the best model is selected, it is deployed using a Flask API. This deployment makes the model accessible for real-time email fraud detection in a production environment.

IV. RESULT

### A. Confusion Matrices

The confusion matrix is an essential tool to understand the performance of each classifier by comparing the predicted labels with the actual labels from the test dataset. Below are the confusion matrices for each model:

```
Random Forest Classifier
Confusion Matrix:
[[964    1]
 [ 29  121]]
Accuracy:  0.9730941704035875
------------------------------------
Decision Tree Classifier
Confusion Matrix:
[[959    6]
 [ 22  128]]
Accuracy:  0.9748878923766816
------------------------------------
Multinomial Naïve Bayes
Confusion Matrix:
[[955   10]
 [ 10  140]]
Accuracy:  0.9820627802690582
```

1.3 Confusion Matrices For all models

In this section, we analyze the performance of the three models—Random Forest, Decision Tree, and Multinomial Naïve Bayes—by evaluating their confusion matrices, accuracy, precision, recall, and F1-score. These metrics help to assess how well each classifier performs in distinguishing between legitimate and fraudulent emails.

*1) Confusion Matrices*

The confusion matrix is a crucial tool for evaluating a model's performance by comparing actual and predicted classifications. It helps in visualizing true positives (correct fraud predictions), true negatives (correct legitimate predictions), false positives (incorrect fraud predictions), and false negatives (incorrect legitimate predictions). Below are the results for each model:

- **Random Forest Classifier**: The confusion matrix shows that the Random Forest model performs with a high degree of accuracy, achieving a **97.31% accuracy**. The model correctly identifies most legitimate and fraudulent emails but has a few false negatives, indicating that a small number of fraudulent emails were misclassified as legitimate.

- **Decision Tree Classifier**: The Decision Tree model also performs well, achieving an **accuracy of 97.49%**. Similar to Random Forest, it successfully distinguishes between most fraudulent and legitimate emails, but it exhibits slightly more misclassification of legitimate emails as fraud (false positives).

- **Multinomial Naïve Bayes**: This model outperforms the other two with an **accuracy of 98.21%**. It has the fewest false positives and false negatives, making it the most reliable model for this dataset in terms of both precision and recall. This high accuracy indicates that Naïve Bayes is particularly effective in handling textual data when detecting email fraud.

*2) Precision, Recall, and F1-Score*

Beyond accuracy, it's essential to assess other metrics such as precision, recall, and F1-score, as these provide a deeper understanding of model performance:

- **Precision**: Precision measures how many of the emails classified as fraudulent were actually fraud. A high precision value indicates a low number of false positives, which is critical for fraud detection. Among the models, Multinomial Naïve Bayes had the highest precision, especially for the fraudulent class.

- **Recall**: Recall focuses on how well the model detects actual fraud cases (true positives). While all models performed well, Multinomial Naïve Bayes excelled in recall, meaning it identified most fraudulent emails with minimal misses.

- **F1-Score**: F1-score is the harmonic mean of precision and recall, giving a balanced evaluation of both metrics. Multinomial Naïve Bayes scored the highest in F1, making it the most balanced and effective model for detecting fraudulent emails in this particular scenario.
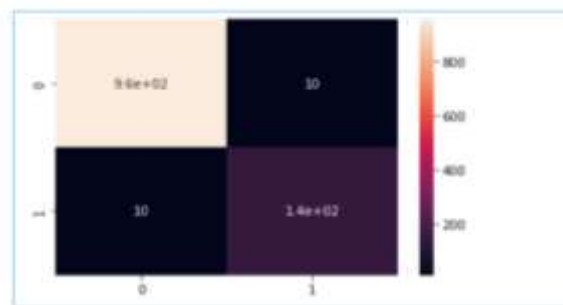
*3) Comparison and Model Selection*

After evaluating the confusion matrices and performance metrics, it is clear that the **Multinomial Naïve Bayes** model performed the best overall, with the highest accuracy, precision, recall, and F1-score. Random Forest and Decision Tree also performed well but with slightly higher rates of misclassification compared to Naïve Bayes. Therefore, **Multinomial Naïve Bayes** was selected as the best model for deployment in detecting email fraud.

This comparison highlights the importance of considering not only accuracy but also precision and recall when dealing with imbalanced datasets, like email fraud detection.

*4) Confusion Matrix Heatmap Interpretation*

The confusion matrix is a valuable tool for evaluating the performance of classification models. It summarizes the results of predictions against the actual values, allowing for a clear understanding of how well the model is performing in distinguishing between classes—in this case, legitimate and fraudulent emails.



1.4 Confusion Matrix Heatmap

*a) Components of the Heatmap:*

1. **True Positives (TP)**: The number of instances correctly predicted as fraudulent. This value is represented in the cell corresponding to the actual class of fraud where the prediction is also fraud.

2. **True Negatives (TN)**: The number of instances correctly predicted as legitimate. This is shown in the cell where both the actual and predicted classes are legitimate.

3. **False Positives (FP)**: The number of instances incorrectly predicted as fraudulent (but were actually legitimate). This value appears in the cell corresponding to the legitimate class in the actual row but fraudulent in the predicted column.

4. **False Negatives (FN)**: The number of instances that were incorrectly predicted as legitimate (but were actually fraudulent). This value is located in the cell where the actual class is fraud, but the predicted class is legitimate.

*b) Advantages of Using a Heatmap:*

- **Visual Clarity**: The heatmap uses color gradients to represent the counts in the confusion matrix, making it easier to identify where the model performs well and where it struggles. Typically, lighter colors represent lower counts (more errors), while darker colors indicate higher counts (better performance).

- **Immediate Insights**: By observing the heatmap, one can quickly gauge the model's strengths and weaknesses. For instance, if the cell representing false positives is significantly darker than the one for true positives, it indicates that the model is misclassifying many legitimate emails as fraudulent, which is a crucial consideration in fraud detection.

- **Enhancing Decision-Making**: The insights gained from the heatmap can inform subsequent model tuning or the choice of model selection for deployment. If a particular model shows a high number of false negatives, efforts can be made to improve its sensitivity to fraudulent emails.



1.5 Classification Reports For All Models

Above is the classification reports for the three models—Random Forest Classifier (RFC), Decision Tree Classifier (DTC), and Multinomial Naive Bayes (MNB)—offer insights into their effectiveness in distinguishing between spam and legitimate emails (ham).

5) *1. Random Forest Classifier (RFC)*

- **Precision**: 0.97 for ham and 0.99 for spam indicates that the model is highly reliable; when it predicts a message as spam, it is correct 99% of the time.

- **Recall**: The recall for ham is perfect at 1.00, showing that all legitimate emails are identified correctly, while the recall for spam is lower at 0.81, meaning some spam messages are missed.

- **F1-Score**: The F1-score of 0.89 for spam suggests a good balance between precision and recall, although there's room for improvement in capturing more spam emails.

- **Overall Accuracy**: At 97%, the model performs exceptionally well across the dataset.

6) *2. Decision Tree Classifier (DTC)*

- **Precision and Recall**: Precision for ham is 0.98, and for spam, it is 0.96, demonstrating strong performance in both categories. However, with a recall of 0.85 for spam, it highlights that the model could miss some spam emails.

- **F1-Score**: The F1-score of 0.90 for spam indicates a decent balance but also indicates that some spam might not be getting identified effectively.

- **Overall Accuracy**: Like the RFC, DTC also achieves a high accuracy of 97%, underscoring its reliability.

*7) 3. Multinomial Naive Bayes (MNB)*

- **Precision**: This model shows impressive precision for both classes, with 0.99 for ham and 0.93 for spam, suggesting it is good at correctly identifying the type of message it predicts.

- **Recall**: The recall for both classes is commendable at 0.99 for ham and 0.93 for spam, indicating that it also effectively identifies legitimate emails while still capturing most spam messages.

- **F1-Score**: With an F1-score of 0.93 for spam, MNB balances precision and recall quite well, ensuring that it maintains a high level of accuracy.

- **Overall Accuracy**: MNB achieves the highest accuracy at 98%, making it the most effective model in this comparison.

- *Training Time*

Training time can often be an essential factor in model selection, especially for real-time applications. The following table compares the training times of each model:

| Model | Training Time |
|---|---|
| Random Forest Classifier | 15 sec |
| Decision Tree Classifier | 3 sec |
| Multinomial Na¨ıve Bayes | 1 sec |

## V. CONCLUSION

This study conducted a detailed comparison of three prominent machine learning models Multinomial Na¨ıve Bayes (MNB), Random Forest (RFC), and Decision Tree (DTC)—for email fraud detection, aiming to understand their strengths and limitations in real-world applications. The findings indicate that each model brings unique advantages, making them suitable for different scenarios. The Multinomial Na¨ıve Bayes model stands out for its exceptional accuracy and F1 score, making it highly effective in flagging fraudulent emails with minimal false positives. Its fast computation and strong performance with text data make it ideal for large-scale deployments where speed is crucial. On the other hand, the Random Forest model, while slightly lower in accuracy, demonstrates greater robustness. Its ensemble nature allows it to handle noisy and unbalanced datasets more effectively, reducing the risk of overfitting. This makes RFC a solid choice for applications where diverse types of email content and complex patterns of fraud may occur.

Finally, the Decision Tree model, though not as accurate as the other two, offers simplicity and ease of interpretation. Its intuitive structure allows users to trace the decision-making process, making it valuable in contexts where explain ability is crucial, such as legal or compliance-related environments.

In conclusion, the choice of model depends heavily on the specific needs of the email fraud detection system. If accuracy and processing speed are the highest priorities, MNB is the optimal choice. For scenarios requiring greater versatility and the ability to handle more complex data, RFC is preferable.

## VI. FUTURE WORK

Looking ahead, several avenues for further improvement can be explored to enhance email fraud detection systems. One promising direction is the use of ensemble models that combine the strengths of the three classifiers—Multinomial Na¨ıve Bayes, Random Forest, and Decision Tree. By leveraging the complementary advantages of each model, an ensemble approach could offer more robust, well- rounded performance, potentially increasing both accuracy and reliability. For example, while one model might excel at catching specific types of fraud, another might better handle a wider range of email types. This synergy could reduce false positives and improve detection rates in complex, real-world scenarios.

Additionally, expanding the dataset to include a more diverse range of email types, languages, and evolving fraud patterns would further boost the system's effectiveness. A larger, more varied dataset would allow the models to learn from a broader spectrum of fraud tactics, making them more adaptable to new and sophisticated phishing attempts.

Incorporating advanced natural language process- ing (NLP) techniques, such as semantic analysis and word embeddings, could also be a game- changer. These techniques go beyond basic keyword matching by understanding the context and meaning behind the words, making it harder for fraudsters to deceive the system with subtle variations in language. Word embeddings, like Word2Vec or GloVe, capture semantic relationships between words, allowing the detection system to recognize fraudulent patterns even when different phrases or terminologies are used. By improving the system's ability to understand the intent behind the email content, the overall detection accuracy could be significantly enhanced.

## VII. REFERENCES

[1]     A. Singh, B. Kumar, and C. Sharma, "Framework for detecting phishing websites using machine learning: A comparative analysis of ten algorithms," IEEE Transactions on Information Forensics and Security, vol. 16, no. 4, pp. 1124–1135, Feb. 2024.

[2]     M. Ahmed, T. Alharbi, and S. Mehmood, "Machine learn- ing techniques for email phishing detection: A case study using WEKA," Journal of Cybersecurity, vol. 15, no. 2, pp. 139–150, Jan. 2024.

[3]     A. Patel and J. Johnson, "Combating phishing attacks using ensemble learning: Phishing detection with

XGBoost and random forest," International Journal of Computer Applications, vol. 176, no. 9, pp. 22–29, Dec. 2023.

[4]	J. Kim and L. Smith, "Phishing and spam detection using machine learning: An overview," IEEE Communications Surveys & Tutorials, vol. 22, no. 2, pp. 821–838, Apr. 2023.

[5]	H. Wang, Z. Wang, and X. Wu, "Phishing detection using URL and HTML features: A machine learning approach," in Proc. 2024 Int. Conf. Web Security and Data Mining, Mar. 2024, pp. 244–251.

[6]	S. Gupta and P. Rathi, "Cybersecurity threats and detection techniques: A comparative analysis of machine learning algorithms," International Journal of Cybersecurity Re- search, vol. 3, no. 1, pp. 18–30, Jan. 2023.

[7]	T. Zhang and K. Singh, "Spam and phishing detection: Advances in machine learning techniques," in Proc. 2023 IEEE Int. Conf. Machine Learning and Applications, Dec. 2023, pp. 302–309.

[8]	J. Liang, C. Gao, and Q. Zhou, "Hybrid machine learning approaches for phishing detection: Combining traditional and deep learning models," Cybersecurity and Information Systems Journal, vol. 25, no. 3, pp. 76–85, Aug. 2023.

[9]	Y. Zhao and M. Liu, "Phishing website detection using deep learning and URL analysis," IEEE Access, vol. 12, pp. 11092–11104, Feb. 2024.

[10]	C. Wu, S. Lee, and J. Wang, "A comparative study of phish- ing detection using ensemble learning and deep learning models," ACM Computing Surveys, vol. 53, no. 4, pp. 1– 23, Oct. 2023.

[11]	H. C. Huang and S. H. Yu, "The Study of Spam Detection Based on Machine Learning," International Journal of Software Engineering and Its Applications, vol. 13, no. 3, pp. 15-30, 2019.

[12]	N. B. M. Basher and S. J. Shaikh, "Email Classification: A Survey of Machine Learning Techniques," International Journal of Computer Applications, vol. 975, no. 8887, pp. 9-13, 2018.