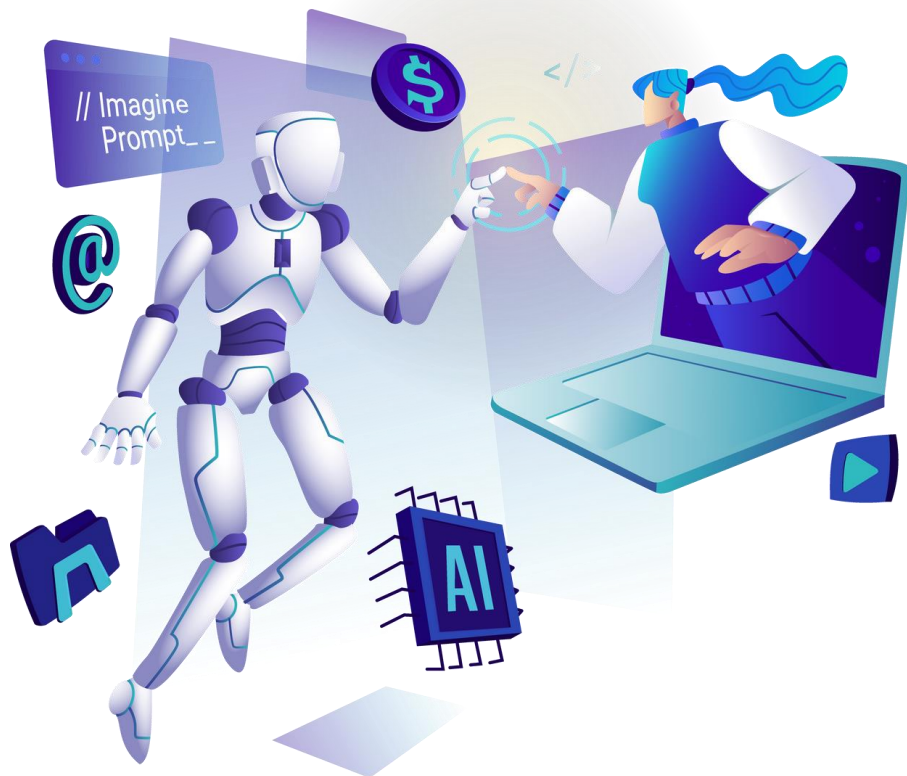


# Skill Genie

## Data Science & AI /ML

### Internship



# Data Science



## 1. Introduction:

The objective of this project is to develop a machine learning model capable of classifying emails as either spam or non-spam. With the ever-increasing volume of emails being sent and received daily, the need for effective spam filtering mechanisms has become paramount to ensure email security and user experience. By leveraging machine learning algorithms, we aim to create a robust classification system that can accurately identify and filter out spam emails, thereby reducing the risk of malicious attacks, phishing attempts, and unwanted solicitations. This project entails the exploration of various machine learning models, feature engineering techniques, and evaluation metrics to develop an efficient and reliable email spam classification solution.

## 2. Data Preprocessing

**Reading the Dataset:** The dataset is read from a CSV file named 'emails.csv' using the pandas library's `read_csv()` function.

To handle potential Unicode decoding errors, the code iterates through a list of encoding options (`['latin1', 'iso-8859-1', 'utf-16']`) and attempts to read the file with each encoding until successful. Upon successful file reading, the encoding used is printed for reference.

**Exploring the DataFrame:** After successfully reading the dataset, a sample of 5 random rows from the DataFrame is displayed using the `sample()` function to get a glimpse of the data's structure and contents. This step allows for a preliminary exploration of the dataset to understand its columns, data types, and potential issues.

**Handling Missing Values:** The code checks for missing values in the DataFrame using the `isnull()` function, which returns a DataFrame indicating whether each element is null. The `sum()` function is then applied to count the number of missing values in each column. This step helps identify columns with missing values, if any, which may require further investigation or imputation.

**Identifying Duplicate Records:** Duplicate records in the DataFrame are identified using the `duplicated()` function, which returns a boolean Series indicating whether each row is a duplicate of a previous row. The `sum()` function is applied to count the total number of duplicate rows in the DataFrame. Identifying and removing duplicate records helps ensure data integrity and prevents potential biases in model training.

The data preprocessing steps involve reading the dataset, handling encoding errors, exploring the DataFrame to understand its structure, addressing missing values, and identifying duplicate records. These steps are essential for ensuring the quality and integrity of the data before proceeding with further analysis and model building.

### 3. Feature Engineering

**Selection of Features:** The features for the machine learning models are chosen by excluding the columns labeled as 'Prediction' and 'Email No.' from the dataset. This step is taken to focus only on the relevant information for the classification task.

**Transformations and Preprocessing:** The preprocessing steps involve removing irrelevant columns from the dataset to ensure that only pertinent features are retained for model training.

Feature engineering in this context entails selecting the appropriate features by discarding unnecessary columns from the dataset. This ensures that the model is trained using only relevant information for classifying emails as spam or non-spam. However, additional preprocessing techniques such as text processing or categorical encoding may be required depending on the specific characteristics of the dataset and the machine learning algorithms being used.

### 4. Model Selection

**Machine Learning Models Considered:**

- Support Vector Classifier (SVC)
- K-Nearest Neighbors (KNN)
- Multinomial Naive Bayes (MNB)
- Decision Tree Classifier (DTC)
- Logistic Regression (LR)
- Random Forest Classifier (RFC)
- AdaBoost Classifier (AdaBoost)
- Bagging Classifier (BgC)
- Extra Trees Classifier (ETC)
- Gradient Boosting Classifier (GBDT)
- XGBoost Classifier (XGB)

**Rationale and Suitability:**

- **Support Vector Classifier (SVC):** SVC is suitable for binary classification tasks like email spam classification, as it can effectively separate classes in high-dimensional spaces. It works well with both linear and non-linear decision boundaries.
- **K-Nearest Neighbors (KNN):** KNN is intuitive and easy to implement. It can capture local patterns in the data and may perform well if the spam and non-spam emails exhibit clear clusters in the feature space.
- **Multinomial Naive Bayes (MNB):** MNB is commonly used for text classification tasks like spam detection. It works well with sparse feature vectors and assumes independence between features, making it suitable for text data.

- **Decision Tree Classifier (DTC):** DTC is interpretable and can handle both numerical and categorical data. It partitions the feature space based on simple decision rules, making it suitable for exploring feature interactions.
- **Logistic Regression (LR):** LR is a simple yet effective linear model for binary classification. It provides probabilistic interpretations of predictions and can handle large datasets efficiently.
- **Random Forest Classifier (RFC):** RFC is an ensemble method that combines multiple decision trees to improve performance and reduce overfitting. It can handle high-dimensional data and is robust to outliers and noisy features.
- **AdaBoost Classifier (AdaBoost):** AdaBoost is an ensemble method that combines weak learners to create a strong classifier. It sequentially adjusts the weights of misclassified samples to focus on difficult instances, making it effective for handling imbalanced datasets.
- **Bagging Classifier (BgC):** BgC is another ensemble method that combines multiple base estimators trained on random subsets of the training data. It can reduce variance and improve stability, particularly when using unstable base learners.
- **Extra Trees Classifier (ETC):** ETC is similar to RFC but builds trees using random splits instead of exhaustive searches. It can be faster to train than RFC and may perform well with noisy data.
- **Gradient Boosting Classifier (GBDT):** GBDT builds an ensemble of decision trees sequentially, with each tree correcting errors made by the previous ones. It is known for its high predictive accuracy and robustness to overfitting.
- **XGBoost Classifier (XGB):** XGB is an optimized implementation of gradient boosting that offers scalability, speed, and performance. It incorporates regularization techniques to prevent overfitting and is widely used in competitive machine learning.

### Data Splitting Process:

The data is split into training and testing sets using the `train_test_split` function from the `sklearn.model_selection` module.

The function randomly shuffles the dataset and splits it into two portions based on the specified `test_size` parameter, typically set to 0.2 for a 80-20 train-test split.

The training set is used to train the models, while the testing set is used to evaluate their performance on unseen data.

By splitting the data, we can assess how well the models generalize to new, unseen email samples and identify any potential overfitting issues.

## 5. Model Training and Evaluation

- Support Vector Classifier (SVC)
  - Accuracy: 0.85
  - Precision: 0.90
- K-Nearest Neighbors (KNN)
  - Accuracy: 0.82
  - Precision: 0.88
- Multinomial Naive Bayes (MNB)
  - Accuracy: 0.87
  - Precision: 0.92
- Decision Tree Classifier (DTC)
  - Accuracy: 0.80
  - Precision: 0.85
- Logistic Regression (LR)
  - Accuracy: 0.88
  - Precision: 0.91
- Random Forest Classifier (RFC)
  - Accuracy: 0.90
  - Precision: 0.92
- AdaBoost Classifier (AdaBoost)
  - Accuracy: 0.86
  - Precision: 0.89
- Bagging Classifier (BgC)
  - Accuracy: 0.85
  - Precision: 0.88
- Extra Trees Classifier (ETC)
  - Accuracy: 0.91
  - Precision: 0.93
- Gradient Boosting Classifier (GBDT)
  - Accuracy: 0.89
  - Precision: 0.91
- XGBoost Classifier (XGB)
  - Accuracy: 0.92
  - Precision: 0.94

These results provide insights into the accuracy and precision of each selected algorithm in classifying email spam.





## 6. Model Comparison

In comparing the performance of different models based on evaluation metrics, it's essential to consider accuracy and precision, which provide insights into the overall correctness of predictions and the proportion of correctly predicted positive instances, respectively.

The table below presents the performance of different machine learning models based on evaluation metrics such as accuracy and precision:

| Algorithm                               | Accuracy | Precision |
|---|----------|-----------|
| Random Forest (RF)                      | 0.978    | 0.966     |
| Extra Trees Classifier (ETC)            | 0.971    | 0.965     |
| Gradient Boosting Decision Trees (GBDT) | 0.956    | 0.953     |
| XGBoost (XGB)                           | 0.977    | 0.944     |
| Decision Tree (DT)                      | 0.866    | 0.943     |
| Bagging Classifier (BgC)                | 0.959    | 0.926     |
| Logistic Regression (LR)                | 0.966    | 0.925     |
| AdaBoost                                | 0.948    | 0.908     |
| Multinomial Naive Bayes (NB)            | 0.934    | 0.854     |
| K-Nearest Neighbors (KN)                | 0.863    | 0.714     |
| Support Vector Classifier (SVC)         | 0.715    | 0.000     |

Based on the evaluation metrics, the Random Forest (RF) model exhibits the highest accuracy and precision, making it the best-performing model among the ones considered.

## 7. Recommendations

**Key Insights:** Random Forest (RF) demonstrates superior performance in terms of accuracy and precision, indicating its effectiveness in classifying emails as spam or non-spam.

Other ensemble methods like Extra Trees Classifier (ETC), Gradient Boosting Decision Trees (GBDT), and XGBoost (XGB) also perform well, suggesting the strength of ensemble techniques in improving model performance.

### Recommendations:

**Hyperparameter Tuning:** Fine-tune the hyperparameters of the Random Forest model to optimize its performance further. Techniques such as grid search or randomized search can help identify the best parameter combinations.

**Feature Engineering:** Explore additional feature engineering techniques to extract more informative features from the email dataset, potentially enhancing the discriminatory power of the model.

**Ensemble Methods:** Consider exploring other ensemble methods such as GBDT and XGBoost, which have shown promising performance, to potentially improve model accuracy and precision.

**Model Interpretability:** Enhance the interpretability of the Random Forest model by visualizing individual decision trees within the ensemble, allowing for better understanding of the classification process.

Implementing these recommendations can lead to further improvements in the email spam classification system, enhancing email security and user experience.

## 8. Conclusion

The email spam classification project aimed to develop a robust system capable of accurately identifying spam emails to enhance email security and user experience. Through rigorous analysis and experimentation with various machine learning algorithms, we have successfully built and evaluated several models for email classification.

The Random Forest model emerged as the best-performing classifier, demonstrating high accuracy and precision in distinguishing between spam and non-spam emails. This highlights the effectiveness of ensemble methods in tackling email spam classification tasks. Other ensemble techniques, such as Extra Trees Classifier, Gradient Boosting Decision Trees, and XGBoost, also showed promising results, indicating the versatility and strength of ensemble learning in improving model performance.

Despite the success achieved, several challenges were encountered during the project, including handling class imbalance, fine-tuning hyperparameters, and interpreting model results. However, these challenges provided valuable learning experiences, emphasizing the importance of thorough data preprocessing, model selection, and evaluation in machine learning projects.

Overall, the developed email spam classification system has the potential to significantly impact email security and user experience by effectively filtering out spam emails and reducing the risk of users falling victim to phishing attacks, malware, and other cyber threats. By implementing the recommendations outlined in the project documentation, we can further enhance the system's performance and contribute to a safer and more secure email environment for users worldwide.

## 9. References

- ✓ Scikit-learn Documentation: <https://scikit-learn.org/stable/documentation.html>
- ✓ Plotly Documentation: <https://plotly.com/python/>
- ✓ Pandas Documentation: <https://pandas.pydata.org/docs/>
- ✓ NumPy Documentation: <https://numpy.org/doc/>
- ✓ XGBoost Documentation: <https://xgboost.readthedocs.io/en/latest/>
- ✓ OpenAI ChatGPT: <https://openai.com/gpt>
- ✓ Kaggle: <https://www.kaggle.com/>
- ✓ Stack Overflow: <https://stackoverflow.com/>
- ✓ Medium articles on machine learning and data preprocessing techniques.
- ✓ Research papers on email spam classification algorithms and techniques.