

# **BUILDING A SMARTER AI-POWERED SPAM CLASSIFIER**

## **Phase 1: Problem Definition and Design Thinking**

**TEAM MEMBERS:** Sai Arjunaa A, Samyukthan A, Shyaam S, Mukesh P, Vikkiran V

**College:** Madras Institute of Technology, Chennai

### **PROBLEM STATEMENT:**

In today's digital age, email communication plays a pivotal role in personal and professional interactions. However, the exponential rise in spam emails has become a significant concern, leading to productivity loss, security risks, and user frustration. The need for an efficient and accurate spam classifier powered by artificial intelligence (AI) has never been more critical.

The objective of this project is to develop an advanced AI-based spam classifier capable of accurately distinguishing between legitimate emails and spam messages. The classifier should exhibit high precision and recall rates while maintaining a low false-positive rate, ensuring that legitimate emails are not erroneously marked as spam.

### **Key Challenges:**

1. **Diverse Types of Spam Content:** Spam emails come in various forms, including phishing attempts, marketing messages, malicious attachments, and more. The classifier must be robust enough to recognize and categorize different types of spam content accurately.
2. **Evolving Spam Techniques:** Spammers continually adapt their tactics to bypass traditional filters. The AI classifier must be able to adapt and learn from new patterns and trends in spam content.
3. **Imbalanced Data:** Obtaining a balanced dataset with an equal distribution of spam and non-spam emails is challenging. The classifier should be able to handle imbalanced data to prevent bias in its predictions.
4. **Real-time Processing:** The classifier should have the capability to process emails in real-time to provide timely protection for users.
5. **Multi-language Support:** As email communication is global, the classifier should be able to effectively classify emails written in different languages.
6. **User Customization and Feedback Loop:** Allowing users to customize the classifier's behavior (e.g., marking false positives/negatives) and incorporating a feedback loop for continuous learning and improvement.
7. **Security and Privacy:** Ensuring that the classifier does not compromise user privacy and that sensitive information within emails is not accessed or stored.
8. **Scalability:** The system should be scalable to handle a large volume of incoming emails, especially in enterprise-level applications.
9. **Integration with Email Platforms:** The classifier should be compatible with popular email platforms (e.g., Gmail, Outlook) and easily integratable into existing email systems.

10. **Model Explainability:** Providing insights into the classifier's decision-making process to build user trust and allow for transparency in its operation.

By addressing these challenges, the goal is to develop an AI spam classifier that enhances email security, reduces the risk of phishing attacks, and improves overall email communication efficiency for users across various domains.

## DESIGN AND THINKING

1. **Data Collection:** We will need a dataset containing labeled examples of spam and nonspam messages. We can use a Kaggle dataset for this purpose.

**Dataset Link:** <https://www.kaggle.com/datasets/uciml/sms-spam-collection-dataset>

2. **Data Preprocessing:** It involves cleaning and preparing the text data to ensure that it is suitable for analysis and modeling. Key steps for data preprocessing in a spam classifier are:

- i. **Lowercasing:** Convert all text to lowercase. This ensures that the classifier doesn't treat "Hello" and "hello" as different words.
- ii. **Removing Punctuation:** Remove any special characters or punctuation marks from the text. These do not typically contribute significantly to the identification of spam.
- iii. **Tokenization:** Split the text into individual words or tokens. This step is essential for building features from the text data.
- iv. **Removing Stop Words:** Remove common words (e.g., "the", "is", "and") known as stop words. These words do not carry much information for the classifier and can be safely removed.
- v. **Stemming or Lemmatization:** Reduce words to their root form. For example, "running", "ran", and "runs" might all be reduced to "run". Stemming and lemmatization help reduce the dimensionality of the feature space.
- vi. **Handling Numerical Values and Symbols:** Decide how to handle numerical values and symbols. In some cases, they may be relevant (e.g., in identifying phone numbers in spam messages), while in others, they can be removed.
- vii. **Handling HTML Tags:** If the data contains HTML tags (common in email data), they should be stripped out to ensure they do not interfere with the analysis.
- viii. **Handling URLs:** Decide whether to keep or remove URLs. They can be a source of noise in the text, but in some cases, they may contain information relevant to spam detection.

- ix. Handling Email Addresses: Similar to URLs, decide whether to keep or remove email addresses from the text data.
- x. Handling Emojis and Special Characters: Decide how to handle emojis and other special characters. They can be informative in some cases, but may need special treatment depending on the context.
- xi. Dealing with Imbalanced Data: If the dataset is imbalanced (i.e., there are significantly more non-spam than spam examples or vice versa), consider techniques like resampling or using different evaluation metrics.
- xii. Vectorization: Convert the preprocessed text data into numerical vectors. Techniques like TF-IDF (Term Frequency-Inverse Document Frequency) or word embeddings (e.g., Word2Vec, GloVe) can be used.

**3. Feature Extraction:** After tokenization, the tokens are converted to numerical features. The following techniques can be used for this purpose:

- i. TF-IDF (Term Frequency-Inverse Document Frequency): TF-IDF is a widely used technique in text classification tasks, including spam classification. It assigns weights to terms based on their frequency in a document relative to the entire corpus. Terms that are common in spam but rare in legitimate emails receive higher weights.
- ii. Binary Encoding: Binary encoding can be effective for spam classification because it captures the presence or absence of specific terms. In spam classification, the presence of certain words or phrases can be indicative of spam.
- iii. Count Vectorization: Count vectorization simply counts the frequency of each term in a document. This can be effective for identifying spam based on the occurrence of specific keywords or patterns commonly found in spam messages.
- iv. N-grams: N-grams capture sequences of words, which can be useful in identifying certain patterns or phrases commonly associated with spam. For example, phrases like "limited time offer" or "click here now" are often indicative of spam.
- v. Word Embeddings (Word2Vec, GloVe): Word embeddings capture semantic relationships between words. While not as directly interpretable as some other methods, they can be effective at capturing context and meaning, which can aid in distinguishing spam from legitimate messages.
- vi. Doc2Vec: Doc2Vec provides a way to represent entire documents as fixed-length vectors. This can be useful for capturing the overall context and style of spam messages.

Among these methods, TF-IDF, Count Vectorization, and N-grams are effective for spam classification. They allow the classifier to focus on specific terms and patterns that are indicative of spam. However, it's not uncommon to experiment with multiple methods to see which one performs best on a specific dataset.

- 4. Model Selection:** We can experiment with various machine learning algorithms such as Naive Bayes, Support Vector Machines, and more advanced techniques like deep learning using neural networks. Commonly used machine learning algorithms for spam classifiers are:

Naive Bayes:

Strengths:

- Simple and easy to implement.
- Efficient in terms of memory and computation.
- Can handle a large number of features.

Reasoning: Naive Bayes is a popular choice for text classification tasks like spam detection. It works by calculating probabilities based on the frequency of words in spam and non-spam emails.

Support Vector Machines (SVM):

Strengths:

- Effective in high-dimensional spaces.
- Can handle non-linear decision boundaries with kernel tricks.

Reasoning: SVM aims to find the optimal hyperplane that separates spam and non-spam emails while maximizing the margin between the two classes.

Logistic Regression:

Strengths:

- Simple and interpretable.
- Can be regularized to prevent overfitting.

Reasoning: Although logistic regression is a binary classification algorithm, it can be used effectively for spam detection when extended with techniques like one-vs-rest.

Random Forest:

Strengths:

- Ensemble method that combines multiple decision trees for improved accuracy.
- Can handle a large number of features and complex relationships.

Reasoning: Random Forest builds multiple decision trees and combines their outputs to make a final prediction. It is robust and can handle noisy data.

Gradient Boosting (e.g., XGBoost, LightGBM):

Strengths:

- Can achieve high accuracy and is robust to overfitting.
- Handles complex relationships and interactions between features.

Reasoning: Gradient boosting builds an ensemble of weak learners (usually decision trees) in a sequential manner, with each tree aiming to correct the errors of the previous ones.

Neural Networks:

Strengths:

- Can capture complex relationships in data.
- Deep learning models can automatically learn feature hierarchies.

Reasoning: Neural networks, especially deep learning models, can be used for spam classification. Techniques like recurrent neural networks (RNNs) or convolutional neural networks (CNNs) can be particularly effective.

K-Nearest Neighbors (KNN):

Strengths:

- Simple and easy to implement.
- Can handle non-linear decision boundaries.

Reasoning: KNN classifies data points based on the majority class of their k nearest neighbors. It can be effective for spam classification, especially with appropriate feature engineering.

Ensemble Methods (e.g., AdaBoost, Bagging):

Strengths:

- Combine multiple base models to improve overall performance.
- Can handle noisy data and reduce overfitting.

Reasoning: Ensemble methods like AdaBoost and Bagging combine multiple weak learners (usually simple models) to create a strong classifier.

**5. Evaluation:** We will measure the performance of the model used using the following metrics:

- Accuracy: The ratio of correctly classified instances to the total number of instances. While accuracy is important, it may not be the most informative metric, especially if the dataset is imbalanced.

- Precision: The ratio of true positives to the sum of true positives and false positives. It measures the accuracy of positive predictions.
- Recall (Sensitivity or True Positive Rate): The ratio of true positives to the sum of true positives and false negatives. It measures the ability of the classifier to identify all positive instances.
- F1-Score: The harmonic mean of precision and recall. It provides a balanced evaluation of precision and recall.
- Specificity (True Negative Rate): The ratio of true negatives to the sum of true negatives and false positives. It measures the ability of the classifier to identify negative instances.
- ROC-AUC: The Area Under the Receiver Operating Characteristic Curve. It provides an aggregate measure of the classifier's ability to distinguish between the two classes.
- Confusion Matrix: A table that summarizes the performance of a classification algorithm. It shows the counts of true positives, true negatives, false positives, and false negatives.
- False Positive Rate (FPR): The ratio of false positives to the sum of false positives and true negatives. It measures the rate at which non-spam emails are incorrectly classified as spam.
- False Negative Rate (FNR): The ratio of false negatives to the sum of false negatives and true positives. It measures the rate at which spam emails are incorrectly classified as non-spam.

**6. Iterative Improvement:** We will fine-tune the model and experiment with hyperparameters to improve its accuracy. Hyperparameters are configuration settings that are set before training a machine learning model. They control the learning process but are not learned from the data. These hyperparameters need to be finetuned in order to optimize model performance. This could involve adjusting parameters like learning rates, regularization strengths, or tree depths (depending on the chosen algorithm).