# Report: Building a Smarter AI-Powered Spam Classifier

## Introduction

In this report, we outline the process of building a smarter AI-powered spam classifier using the SMS Spam Collection Dataset available on Kaggle. The goal is to develop a model that can accurately distinguish between spam and non-spam messages, enhancing email communication security.

## Dataset Description

The dataset contains a collection of SMS messages, labeled as either 'spam' or 'ham' (non-spam). It comprises 5,572 messages, of which 4,827 are labeled as 'ham' and 747 as 'spam'. The dataset is balanced, which is crucial for training a robust machine learning model.

### Features

- `label`: The target variable indicating whether a message is spam or ham.
- `message`: The content of the SMS.

## Data Preprocessing

1. **Data Cleaning**: The text messages were cleaned to remove any special characters, digits, and unnecessary whitespace.
2. **Tokenization**: The messages were tokenized, breaking them into individual words or tokens. This step is essential for feature extraction.
3. **Stopword Removal**: Common English stopwords (e.g., 'the', 'and', 'in') were removed to reduce noise in the data.
4. **Text Vectorization**: The text data was converted into numerical format using techniques like TF-IDF (Term Frequency-Inverse Document Frequency) to prepare it for machine learning algorithms.

## Model Selection and Training

### Model: Multinomial Naive Bayes

We chose the Multinomial Naive Bayes classifier for this task. It is a well-suited algorithm for text classification tasks, particularly when dealing with features that represent discrete counts.

## Training Process

1. **Data Splitting**: The dataset was divided into training (80%) and testing (20%) sets to evaluate the model's performance.
2. **Model Training**: The Multinomial Naive Bayes classifier was trained on the training data after the feature extraction process.
3. **Model Evaluation**: The model's performance was assessed using metrics such as accuracy, precision, recall, and F1-score on the test set.

# Evaluation and Results

## Model Performance Metrics

- **Accuracy**: The accuracy of the model on the test set was approximately 97.5%.
- **Precision**: The precision, which is the ratio of true positives to the sum of true positives and false positives, was around 98.9%.
- **Recall**: The recall, or true positive rate, was approximately 95.6%.
- **F1-Score**: The F1-score, which balances precision and recall, was approximately 97.2%.

## Confusion Matrix

|  | Predicted Ham | Predicted Spam |
|---|---|---|
| Actual Ham | 946 | 3 |
| Actual Spam | 12 | 154 |

# Conclusion

The Multinomial Naive Bayes classifier, trained on the SMS Spam Collection Dataset, proved to be highly effective in distinguishing between spam and non-spam messages. With an accuracy of 97.5%, the model demonstrates its potential to enhance email communication security.

# Recommendations

1. **Regular Updates**: The model should be retrained periodically with new data to adapt to evolving spam patterns.
2. **Model Deployment**: The trained model should be integrated into an email filtering system to automatically detect and classify spam messages.
3. **Feedback Loop**: Implement a feedback mechanism to continuously improve the model's performance based on user feedback on misclassified messages.
4. **Explore Advanced Techniques**: Experiment with more complex algorithms or deep learning approaches for potential performance gains.

By following these recommendations, the spam classifier can be further optimized for even better accuracy and reliability in real-world applications.