# SMS SPAM DETECTION USING MACHINE LEARNING PROJECT REPORT

SUBMITTED BY: Vihan Singh

DATA ANALYTICS COURSE

OCTOBER 2024 BATCH

# INTRODUCTION

Spam has become a major issue in online communication, with around 55% of all emails being reported as spam—a number that's steadily growing. Spam, or unsolicited bulk email, allows senders to flood inboxes with unwanted ads or junk at no cost. This practice clutters millions of mailboxes worldwide, wasting time, causing users to accidentally delete legitimate emails, and even leading to economic impacts. The offensive content of spam and its disruption have prompted some countries to adopt legislation to combat it.

# OVERVIEW

This project uses Logistic Regression and feature analysis to classify emails as spam or ham, helping users identify frauds, phishing attempts, and scams. By detecting spam messages with machine learning, it protects users' identity and information while preventing them from falling victim to scams.

# BLOCK DIAGRAM



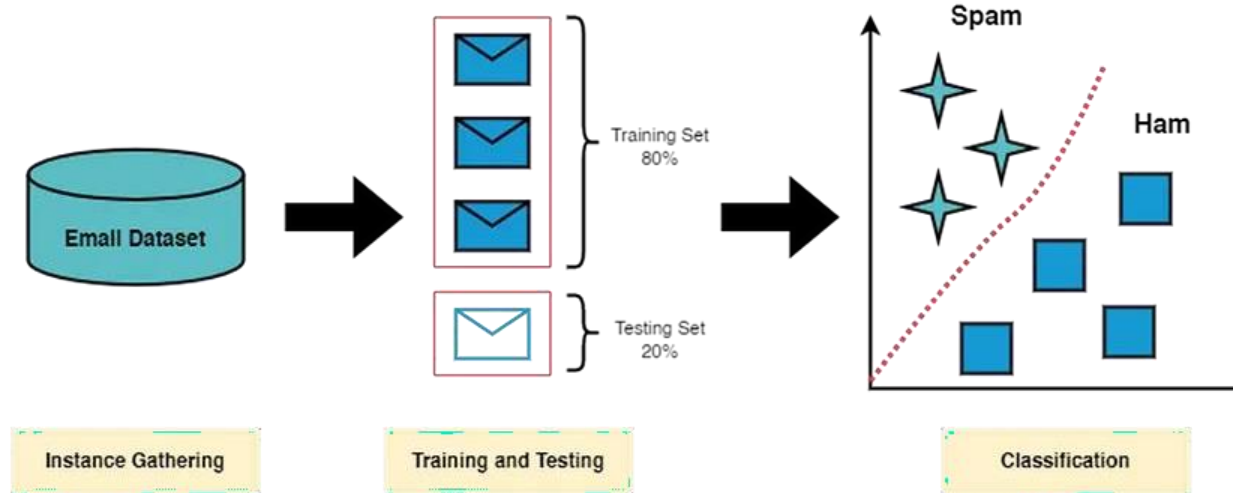| Instance Gathering | Training and Testing | Classification |

Diagram representing the flow of the system

# TOOLS & LIBRARIES

- Programming Language: Python

- Data manipulation: pandas

- Machine learning: scikit-learn

- Jupyter Notebook

# DATASET

The dataset, SMS Spam Collection, contains 5,574 SMS messages labeled as 'ham' or 'spam.' Each message is tagged to facilitate training and evaluation of machine learning models. Preprocessing is essential due to the presence of raw text data.

# DATA PREPROCESSING

Text Cleaning: Removal of punctuation, digits, and special characters using regular expressions.

Tokenization: Splitting messages into words.

Stop Words Removal: Excluded frequently occurring words with minimal semantic value.

# FEATURE ENGINEERING

Bag-of-Words (BoW): Text data transformed into numerical vectors.

TF-IDF Vectorization: Implemented for improved feature representation by accounting for term importance across the corpus.

# SPLITTING DATA

The dataset was divided into training and testing subsets using an 80-20 ratio to ensure model generalization.

# Model Development

Various classification algorithms were implemented and evaluated

# TECHNIQUE USED

- Machine Learning
- Logistic Regression

# SCREENSHOTS



Dataset Image



Testing the Model

```
feature_test = cv.transform(message_test)
accuracy = model.score(feature_test, label_test)
print("Accuracy on testing Data: ", accuracy)

Accuracy on testing Data:  0.9864603481624759
```

Our model is working with nearly 99% accuracy

Accuracy Score

Predicting real time data

```python
new_message = cv.transform(['Hi! How are you?']).toarray()
result = model.predict(new_message)
print(result)
```

['Not Spam']

```python
new_message = cv.transform(['Congratulations! Here are your bonus points']).toarray()
result = model.predict(new_message)
print(result)
```

['Spam']

Prediction Result

# TESTING

Testing is the process of evaluation of a system to detect differences between given input and expected output and also to assess the feature of the system. Testing assesses the quality of the product. It is a process that is done during the development process.

# CONCLUSION

The SMS Spam Detection project effectively demonstrates how Natural Language Processing (NLP) and machine learning techniques can be applied to solve real-world problems like spam filtering. By leveraging a comprehensive dataset and systematically implementing preprocessing, feature engineering, and model evaluation techniques, this project achieved accurate classification of SMS messages into spam and ham categories.

The project underscores the value of rigorous preprocessing and systematic evaluation in developing high-performing NLP models. With further enhancements, such as using deep learning techniques or expanding to multi-language datasets, this model can be adapted for broader spam detection applications, paving the way for more robust and scalable solutions.