



# SMS SPAM CLASSIFICATION PROJECT REPORT

## ABSTRACT

This project aims to develop a SMS spam classifier to automatically identify and filter out unwanted spam messages, improving the messaging experience for users.

Kanchan  
Chowdhury

## Introduction:

The widespread use of mobile phones has led to an increase in unwanted SMS messages, commonly known as spam. These messages can be annoying, intrusive, and potentially harmful. To address this issue, a SMS spam classifier can be developed to automatically identify and filter out spam messages, allowing users to have a more enjoyable and secure messaging experience. In this report, we present a detailed analysis of a SMS spam classifier, including the methodology used and the results obtained.

## Methodology:

### 1. Data Preprocessing:

- **Lowercasing:** Convert all text to lowercase to ensure consistency.
- **Tokenization:** Split the text into individual words or tokens.
- **Special Character Removal:** Remove special characters and punctuation marks.
- **Stopword and Punctuation Removal:** Eliminate common words and punctuation marks that do not contribute to the classification task.
- **Stemming:** Reduce words to their base or root form to handle variations.

### 2. Exploratory Data Analysis:

- **Data Cleaning:** Remove duplicate values and missing entries.
- **Class Distribution:** Analyze the distribution of spam and non-spam messages.
- **Feature Engineering:** Extract features such as the number of characters, words, and sentences in each message.
- **Visualization:** Generate visualizations, including pie charts, histograms, and word clouds, to gain insights into the data.

### 3. Model Building:

- **Feature Extraction:** Use the TF-IDF vectorization technique to convert text data into numerical features.
- **Model Selection:** Evaluate various classification algorithms, including Support Vector Machines (SVM), K-Nearest Neighbors (KNN), Naive Bayes (NB), Decision Trees (DT), Logistic Regression (LR), and Random Forest (RF).
- **Model Training and Evaluation:** Split the data into training and testing sets, train the models on the training set, and evaluate their performance using accuracy and precision metrics.
- **Manual Testing:** Test the trained models on custom messages to assess their effectiveness in classifying spam and non-spam messages.

## **Conclusion:**

A study developed a SMS spam classifier using data preprocessing and machine learning algorithms, achieving high precision with K-Nearest Neighbours, Naive Bayes, and Random Forest.

The dataset analysis reveals spam and non-spam messages, with spam often containing offers and urgent actions. The SMS spam classifier can be improved with additional features and ongoing monitoring. It enhances users' messaging experience by filtering out unwanted messages and blocking spam, promoting a safer mobile communication environment.