# SMS Spam Detection System Using NLP

A Project Report

submitted in partial fulfillment of the requirements

of

AICTE Internship on AI: Transformative Learning
with
TechSaksham – A joint CSR initiative of Microsoft & SAP

by

**Rohan Y. Panchal, rp285518@gmail.com**

Under the Guidance of

**Rathod Jay**

# ACKNOWLEDGEMENT

I would like to extend my heartfelt gratitude to all those who contributed to the successful completion of this project.

First and foremost, I am deeply thankful to **Mr. Jay Rathod Sir** for his expert guidance, constructive feedback, and steadfast support throughout the project. His valuable advice and encouragement were pivotal in navigating challenges and accomplishing the goals of this work.

I also wish to express my appreciation to the **TechSaksham initiative by Microsoft & SAP** for offering a remarkable learning opportunity that significantly enriched my skills and knowledge.

Finally, I am truly grateful to my family, friends, and peers for their constant encouragement, understanding, and belief in my abilities, which kept me motivated every step of the way.

Rohan Y. Panchal

# ABSTRACT

SMS Spam Detection Using NLP is a project designed to develop a reliable and efficient system for classifying SMS messages as spam or ham (non-spam). With the increasing volume of unwanted and malicious messages, the need for accurate spam detection systems has become essential to improve communication and ensure cybersecurity.

This project leverages Natural Language Processing (NLP) techniques to preprocess SMS data, extract meaningful features, and analyse message content. The system utilizes the SMS Spam.csv dataset, a widely recognized collection of labelled SMS messages, to train and validate the machine learning models. Preprocessing steps such as tokenization, stemming, and vectorization are applied to prepare the textual data for analysis.

Various machine learning algorithms are explored and evaluated for their ability to differentiate between spam and legitimate messages. Performance metrics such as accuracy, precision, recall, and F1-score are used to assess the effectiveness of these models. A comparative analysis highlights the most suitable algorithm for real-world applications.

By integrating NLP techniques with machine learning, this project provides an effective solution to combat SMS spam. The use of the SMS Spam.csv dataset ensures a robust and reliable system that minimizes disruptions caused by spam while enhancing communication security. This work not only addresses a critical issue but also lays the groundwork for future research and innovation in text-based classification systems.

# TABLE OF CONTENT

## LIST OF FIGURES

# CHAPTER 1

# Introduction

## 1.1 Problem Statement:

The exponential growth of mobile communication has brought convenience to users but has also led to a significant increase in unsolicited and malicious messages, commonly referred to as spam. Spam messages not only disrupt communication but can also pose serious threats, such as phishing scams and malware distribution. Users often struggle to distinguish between legitimate and spam messages, leading to compromised user experience and potential security vulnerabilities.

Traditional rule-based spam detection systems are often inefficient in handling the dynamic nature of spam messages, which frequently evolve to bypass existing filters. These systems also struggle with large-scale data processing and adapting to new patterns in spam content.

This project, **SMS Spam Detection System Using NLP**, addresses these challenges by utilizing Natural Language Processing (NLP) techniques and machine learning algorithms to create a robust solution for classifying SMS messages as spam or ham (non-spam). By analyzing the textual content of messages and identifying linguistic patterns, the system aims to improve detection accuracy and adaptability.

The proposed system is designed to enhance user communication by filtering out unwanted messages, thereby reducing distractions and safeguarding users from potential threats. This solution also aims to provide a scalable and efficient method for processing large volumes of SMS data, making it suitable for real-world applications.

## 1.2 Motivation:

In an era dominated by digital communication, mobile phones have become an essential tool for personal and professional interactions. However, the growing prevalence of spam messages poses a significant challenge. These unsolicited SMS messages not only waste time but also carry potential security risks, such as phishing attacks, fraud schemes, and malware dissemination. Addressing this issue has become a pressing need to safeguard users and improve the efficiency of communication systems.

**Key Motivations Behind the Project**:

- **Increasing Cybersecurity Threats:** Spam messages often serve as a gateway for cybercrimes. Phishing links and malicious attachments can compromise **user**

data and cause severe financial or personal harm. Developing a robust detection system can mitigate these risks.

- **User Experience Enhancement:** Unsolicited messages clutter inboxes and disrupt user workflows. By filtering out spam effectively, users can focus on relevant and important messages, enhancing their overall experience.

- **Limitations of Existing Solutions:** Current rule-based spam filters lack adaptability and struggle to keep up with the evolving tactics of spammers. This motivates the use of advanced technologies like Natural Language Processing (NLP) and machine learning to create a smarter, more adaptable system.

- **Scalability and Real-World Applications:** With the rapid increase in SMS usage across sectors like banking, e-commerce, and healthcare, a scalable and efficient spam detection system is critical to maintain trust and reliability in communication channels.

- **Leveraging Technology for Innovation:** The integration of NLP techniques and machine learning offers immense potential for analyzing textual data in ways that were previously unattainable. By tapping into this capability, the project aims to deliver an innovative and effective solution.



**(Fig 1: Workflow of an AI-driven system for email classification and spam detection.)**

This project aspires to contribute meaningfully to the field of spam detection by providing a scalable, accurate, and future-proof system. It reflects the broader mission of combining technology and user-centric design to solve pressing problems in communication and cybersecurity.

## 1.3 Objective:

The objectives are designed to ensure the system meets practical needs while maintaining adaptability for future challenges.

Primary Objectives:

- **Develop an Accurate Classification System:** Build a robust model capable of classifying SMS messages into spam and ham (non-spam) with high accuracy, minimizing false positives and false negatives.
- **Leverage NLP for Text Analysis:** Utilize NLP techniques such as tokenization, stemming, and vectorization to preprocess and analyse SMS data effectively. These techniques will help extract meaningful patterns and linguistic features for spam detection.
- **Utilize a Standard Dataset for Model Training:** Employ the SMS Spam.csv dataset to train, test, and validate the model. This ensures the system is built using diverse and representative data, improving its real-world applicability.
- **Optimize Machine Learning Algorithms:** Experiment with multiple machine learning algorithms (e.g., Naïve Bayes, Support Vector Machines, Random Forest) to identify the most suitable approach for detecting spam messages efficiently.
- **Ensure User-Friendliness:** Develop a user-centric system that is easy to integrate and operate, catering to diverse user needs across various platforms.

## 1.4 Scope of the Project:

**Scope**:

- **Automated Spam Filtering:** The project focuses on building an efficient system to automatically classify SMS messages as either "Spam" or "Ham" (legitimate messages). This will help users reduce the time spent identifying unwanted messages manually.
- **Natural Language Processing (NLP) Integration:** By leveraging NLP techniques, the system will analyze the text content of SMS messages to identify patterns, keywords, and linguistic structures commonly associated with spam messages.
- **Real-Time Detection**: The system aims to operate in real time, ensuring users are notified immediately when a message is classified as spam.
- **Data Security and Privacy**: Since SMS messages often contain personal information, the system will prioritize the secure handling of data and adhere to privacy regulations.

**Limitations:**

- **Dependency on Training Data**: The accuracy of the system heavily relies on the quality and diversity of the training dataset. Poor or biased datasets may lead to incorrect classifications.
- **Contextual Understanding**: NLP models can sometimes misinterpret messages with ambiguous or sarcastic tones, leading to false positives or negatives.
- **Edge Cases**: Messages that do not fit typical spam or legitimate patterns

# CHAPTER 2

# Literature Survey

## 2.1 Literature Review:

The detection of SMS spam has emerged as a significant area of research, driven by the widespread use of mobile communication and the associated challenges posed by unsolicited or malicious messages. With the rapid growth of mobile device usage, spam messages have become a persistent issue, leading to potential risks such as phishing, fraud, and unwanted interruptions. To address these concerns, researchers have increasingly turned to Natural Language Processing (NLP) techniques. By applying NLP methods, studies have sought to develop effective solutions for identifying and filtering spam messages, ensuring a safer and more efficient communication environment for users:

- ➢ Researchers have extensively investigated the application of machine learning algorithms such as Naive Bayes, Support Vector Machines (SVM), and Random Forest for detecting spam messages. These algorithms work by analyzing specific text-based characteristics of SMS messages, including word frequencies, the occurrence of specific terms, and patterns in the arrangement of characters. Based on this analysis, the models classify messages as either spam or legitimate (often referred to as "ham").
- ➢ Deep learning approaches, including Recurrent Neural Networks (RNN) and Transformers, have been applied in recent works to enhance spam detection accuracy by capturing contextual and sequential information in SMS text.
- ➢ Additionally, widely recognized datasets, such as the "SMS Spam Collection Dataset" compiled by Almeida et al., play a crucial role in this domain. These datasets provide researchers with a standardized and reliable benchmark for testing and evaluating the performance of various models, facilitating consistent comparisons across different studies and approaches.
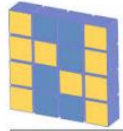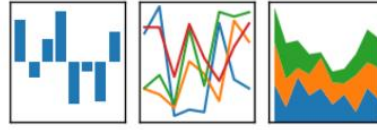
## 2.2 Existing Models, Techniques, or Methodologies:

Several models and techniques have been developed for SMS spam detection:

- ➢ Libraries and Frameworks in NLP for SMS Spam Detection:
    - o In spam detection projects, Python libraries like NumPy and Pandas are extensively utilized for data manipulation and preprocessing tasks. NumPy provides powerful tools for handling numerical data, enabling efficient array operations and mathematical computations. This is particularly useful for transforming raw text data into numerical formats that can be processed by machine learning algorithms. Similarly, Pandas offers robust data structures, such as Data Frames, that simplify the organization, exploration, and manipulation of datasets. It facilitates tasks like cleaning, filtering, and restructuring data, which are critical steps in preparing SMS datasets for analysis. Together, these libraries form the backbone of the preprocessing pipeline, ensuring that the data is in a format suitable for feature extraction and subsequent model training in spam detection workflows.
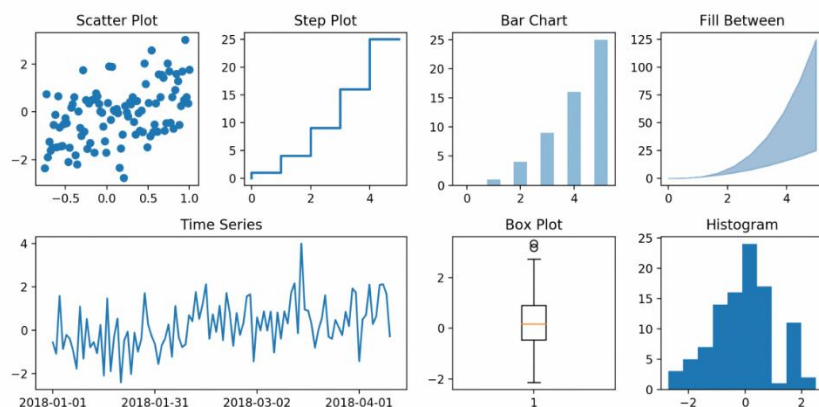
$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$

| Pandas dtype | Python type | NumPy type | Usage |
|---|---|---|---|
| object | str or mixed | string_, unicode_, mixed types | Text or mixed numeric and non-numeric values |
| int64 | int | int_, int8, int16, int32, int64, uint8, uint16, uint32, uint64 | Integer numbers |
| float64 | float | float_, float16, float32, float64 | Floating point numbers |
| bool | bool | bool_ | True/False values |
| datetime64 | NA | datetime64[ns] | Date and time values |
| timedelta[ns] | NA | NA | Differences between two datetimes |
| category | NA | NA | Finite list of text values |

**(Fig: Overview of Pandas and NumPy library)**

o Matplotlib and Seaborn are indispensable tools for visualizing data distributions in spam detection projects. These libraries help researchers and developers better understand patterns and trends in the data, such as differences in character counts or word frequencies between spam and non-spam messages. Matplotlib provides a versatile framework for creating a wide range of plots, including histograms, bar charts, and scatter plots, making it easier to explore the underlying characteristics of the dataset. Seaborn, built on top of Matplotlib, enhances visualization capabilities by offering aesthetically pleasing and high-level statistical plots. It simplifies the creation of detailed and informative graphs, such as box plots or heatmaps, that reveal important insights into the dataset. Together, these libraries play a crucial role in identifying key features and guiding feature selection for building effective spam detection models.



**(Fig: Plot types in Matplotlib)**

## Seaborn Plots



**(Fig: Plot type in seaborn)**

➢ Text Preprocessing Techniques:
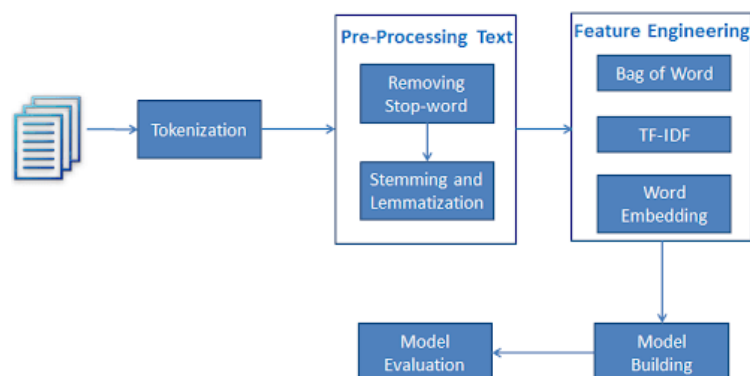  o The Natural Language Toolkit (NLTK) is a powerful Python library widely used in SMS spam detection projects for essential preprocessing tasks. It provides tools for tokenization, which involves breaking down text into smaller units such as words or phrases. Additionally, NLTK includes functionalities for stop word removal, where common but insignificant words (e.g., "the", "and," "is") are filtered out to reduce noise in the dataset. Stemming, performed using tools like Porter Stemmer, is another critical feature of NLTK. It simplifies words to their root form (e.g., "running" becomes "run"), helping to normalize the text and further improve the model's ability to recognize patterns.
  o The Python string library complements these preprocessing efforts by enabling the removal of punctuation and other non-textual characters from SMS messages. This step ensures that the data is clean and focused solely on meaningful textual content, improving the accuracy and reliability of subsequent analysis and classification tasks. Together, these libraries play a vital role in preparing SMS datasets for effective spam detection.



**(Fig: Text pre-processing and feature engineering in NLP.)**

- ➢ Machine Learning Models:
  - o Naive Bayes classifiers, such as GaussianNB, MultinomialNB, and BernoulliNB, are frequently utilized in SMS spam detection tasks because of their simplicity and effectiveness in text classification. These classifiers are based on the Bayes theorem and are particularly well-suited for handling high-dimensional data, such as text.
  - o The GaussianNB classifier is ideal for datasets where features are continuous and follow a normal distribution. MultinomialNB is commonly applied to count-based data, making it well-suited for text data represented as word frequencies or term frequency-inverse document frequency (TF-IDF) values. BernoulliNB, on the other hand, is designed for binary/Boolean features and is particularly effective when text data is represented in terms of binary word occurrence (e.g., whether a specific word appears in a message).
  - o The straightforward implementation and relatively low computational requirements of these classifiers make them a popular choice for spam detection projects. Despite their simplicity, Naive Bayes classifiers often deliver competitive performance, especially when combined with well-preprocessed datasets and appropriately selected features.
- ➢ Metrics for Evaluation:
  - o Metrics like **accuracy**, **precision**, and **confusion matrix** are standard for evaluating the performance of spam detection models.

## 2.3 Gaps or Limitations in Existing Solutions:

- ➢ Lack of Generalization: Many models are trained on specific datasets and fail to generalize well to new or diverse datasets.
- ➢ Dynamic Spam Tactics: Spammers constantly evolve their strategies, making static models ineffective over time.
- ➢ Multilingual Challenges: Most studies focus on English messages, leaving a gap in the detection of spam in regional or less-common languages.
- ➢ High False Positives: Models sometimes misclassify legitimate promotional or transactional messages as spam, affecting user experience.
- ➢ Limited Privacy Preservation: Few solutions address data privacy concerns, which are crucial for handling SMS messages.

Addressing These Gaps:

- • Using a hybrid approach that combines traditional and deep learning techniques to improve accuracy and generalization.
- • Incorporating adaptive learning mechanisms to handle evolving spam tactics.
- • Expanding the dataset to include multilingual messages for broader applicability.
- • Emphasizing privacy-preserving techniques, such as Federated Learning, to ensure data security.

# CHAPTER 3

# Proposed Methodology

## 3.1 System Design (Workflow Diagram)



SMS Spam Detection Workflow

START

Upload Input Data

Text Preprocessing

Feature Extraction

Develop Classification Model

Evaluate the Model

Develop Command-Line Application

User Input SMS

Classify as Spam — Yes — Is it Spam? — No — Classify as Non-Spam

END

**System Design for Spam Email Classification**

The Spam Email Classifier is a system designed to classify emails as spam or ham (non-spam) using machine learning techniques. This system integrates the model into a user-friendly web application, providing an efficient tool for email management. The design is divided into two main phases: Model Development and Application Deployment. Each phase is critical in ensuring the system's functionality, accuracy, and usability. This detailed design document explains the steps involved in each phase.

**Phase 1: Model Development**

The Model Development phase focuses on building and training a machine learning model capable of identifying spam and ham emails with high accuracy. The primary steps in this phase are outlined below.

**1. Dataset Loading**

The system begins by importing a dataset of labeled emails, such as **spam.csv**. The dataset contains two categories: "spam" and "ham". These labels are essential for teaching the model to identify patterns in the data and predict the correct category for new emails. The dataset consists of various features such as the email content, subject line, and sender information, which help in recognizing spam characteristics.

A carefully selected dataset is essential for training a reliable and accurate model. Inaccurate or poorly labeled data can lead to a suboptimal model performance. The dataset used should be diverse, containing various types of spam and ham emails to ensure the model generalizes well.

**2. Preprocessing the Email Content**

Emails often contain unstructured data with varying formats, which can introduce noise into the system. Therefore, the email text needs to undergo preprocessing before it can be used for feature extraction. Key preprocessing steps include:

- **Eliminating Unnecessary Characters**: Special characters, punctuation, and other irrelevant symbols are removed, as they do not contribute to the classification process.
- **Standardizing Text Case**: Converting all text to lowercase ensures consistency, preventing words with different cases (e.g., "Free" and "free") from being treated as separate entities.
- **Tokenizing Text**: The email content is broken down into smaller units (tokens), such as individual words or terms. This step is necessary for feature extraction, as each token will be analyzed separately.

By standardizing and cleaning the data, the system ensures that only meaningful information is retained for further analysis. This reduces noise and improves the model's performance.

## 3. Feature Extraction from Text

Once the email text has been cleaned, it needs to be transformed into a numerical format to be processed by the machine learning model. The system employs **CountVectorizer** for feature extraction.

CountVectorizer converts the email content into a matrix where each row represents an email, and each column corresponds to a unique word found in the entire dataset. The values in the matrix indicate the frequency of words appearing in each email. This transformation allows the machine learning model to process the textual data as numerical features, which are crucial for training the model.

For example, if the word "free" appears five times in an email, the corresponding cell in the matrix would contain the number 5. By converting the text into a numerical format, the model can analyze the distribution of words and learn patterns associated with spam or ham.

## 4. Model Construction

The next stage in building the spam email classifier is constructing the model that will perform the classification task. For this system, the **Naive Bayes Classifier** is selected due to its efficiency and ease of implementation for text classification problems.

**Naive Bayes Methodology**

The Naive Bayes algorithm is a probabilistic classifier that predicts the category (spam or ham) of an email by calculating the likelihood of it belonging to each category based on its features. In this case, the features are typically the words or terms within the email content.

During the training process, the model analyzes the relationships between the words present in the email and their corresponding categories (spam or ham). The Naive Bayes classifier calculates the conditional probabilities of words given a category and uses these probabilities to classify new emails. For a given email, the model calculates the probability of the email being spam and the probability of it being ham. The category with the higher probability is chosen as the predicted label for the email.

The model is trained using labeled data, which means that each email in the dataset is already classified as either spam or ham. This allows the algorithm to learn the typical features associated with both categories. Once the model is trained, it can then be used to classify new emails based on the relationships it has learned from the training data.

**5. Model Evaluation**

After the model has been trained, its performance must be thoroughly assessed to ensure it can accurately classify new, unseen emails. To achieve this, the model is tested on a separate dataset that was not used during training. This allows us to evaluate the model's ability to generalize and perform well on fresh data.

Several key metrics are used to evaluate the performance of the model:

**Accuracy**

Accuracy is a key metric that represents the overall correctness of the model's predictions. It is calculated as the ratio of correct classifications (both spam and ham) to the total number of classifications made. A higher accuracy indicates that the model is correctly identifying more emails, providing an indication of its overall performance

**Precision**

Precision measures the proportion of emails that the model classified as spam and were indeed spam. This metric is critical in understanding how well the model avoids incorrectly labeling non-spam emails as spam. A high precision value means fewer legitimate (ham) emails are wrongly classified as spam, reducing the occurrence of false positives.

**Recall**

It is the proportion of actual spam emails that the model successfully classifies as spam. A high recall value means the model is good at detecting most of the spam emails, ensuring that the spam is caught.

**F1 Score**

It is particularly helpful when dealing with imbalanced datasets, where there are significantly more ham emails than spam emails. By combining precision and recall, the F1 score provides a balanced evaluation of the model's ability to correctly identify spam while minimizing both false positives and false negatives.

These performance metrics help to assess the classifier's ability to generalize and perform well in practical scenarios. By analyzing them, one can identify areas of improvement, such as optimizing precision to reduce false positives or boosting recall to improve the identification of spam emails. Together, they provide a clear picture of the model's reliability and effectiveness.

**Phase 2: Application Deployment**

Once the model has been successfully developed and evaluated, the focus shifts to integrating the model into a web application for end-user use. This phase aims to make the classifier accessible and easy to use.

1. **Developing the Web Application**

   To deploy the model, a **Streamlit** web application is created. Streamlit is a Python framework that allows for the quick and easy development of interactive web applications. It is ideal for this project because it allows for seamless integration of machine learning models into a web interface.

The application includes an intuitive interface where users can input the content of an email they wish to classify. Streamlit handles all the backend operations, including data preprocessing and model inference, ensuring that users can interact with the system without needing technical knowledge.

2. **User Input**

The web application provides a simple input form where users can paste the content of an email for classification. Once the email text is submitted, the application processes the input and initiates the classification process.

This input form is designed to be user-friendly and easy to navigate. The user can paste any email content, and the system will automatically handle the preprocessing and feature extraction steps.

3. **Email Classification Process**

After the user submits an email, the system undergoes the following steps to classify the email:

- **Text Cleaning**: The email text is cleaned by removing unnecessary characters and symbols.
- **Tokenization**: The email text is broken down into smaller units (tokens).
- **Feature Extraction**: The processed text is converted into numerical features using **CountVectorizer**.

Once the email is transformed into numerical features, it is passed through the trained **Naive Bayes Classifier**, which predicts whether the email is spam or ham based on the patterns it learned during training.

4. **Displaying Results**

After the classifier makes a prediction, the result is displayed on the web interface. The user is informed whether the email is classified as spam or ham. This immediate feedback allows users to quickly assess whether an email is potentially harmful.

- **Spam**: If the email is classified as spam, the system displays a message indicating it is spam.

- **Ham**: If the email is classified as legitimate (ham), the system shows a corresponding message.

The display of results is clear and user-friendly, ensuring that users can easily interpret the classification outcome.

5. **Continuous Workflow**

The **Spam Email Classifier** system effectively combines a robust machine learning model with an easy-to-use web application to provide an efficient solution for detecting spam emails. The **Model Development** phase ensures that the classifier is trained and evaluated accurately, while the **Application Deployment** phase ensures that the model is accessible to users through a simple and intuitive web interface.

The **Spam Email Classifier** system effectively combines a robust machine learning model with an easy-to-use web application to provide an efficient solution for detecting spam emails. The **Model Development** phase ensures that the classifier is trained and evaluated accurately, while the **Application Deployment** phase ensures that the model is accessible to users through a simple and intuitive web interface.

Together, these phases create a powerful tool for reducing the impact of spam in email communications, improving email management, and enhancing user experience. This system is scalable, reliable, and practical, offering users a simple way to classify emails and identify potential spam threats.

## 3.2 Requirement Specification

The successful implementation of the spam email classification system requires a combination of both hardware and software components. These components will ensure that the system operates efficiently and effectively, enabling accurate predictions and a smooth user experience.

### 3.2.1  Hardware Requirements:

For the system to operate effectively, the following hardware specifications must be met:

**Processor:**  A minimum of a 2 GHz dual-core processor is required to handle data preprocessing and model training tasks. Higher processing power is recommended for better performance and faster computation, especially during the training phase**.**

**RAM:** At least 4 GB of RAM is required to ensure smooth performance when handling large datasets and processing complex machine learning models. 8 GB of RAM is recommended for optimized performance, particularly for training the Naive Bayes classifier.

**Storage:** The system needs a minimum of 500 MB of free disk space for storing the dataset, model, and the necessary libraries. Additional storage may be required depending on the size of the dataset and the model.

**Internet Connection**: A stable internet connection is required for accessing online resources, deploying the web application via Streamlit, and testing the model on the cloud or remote servers if needed.

### 3.2.2 Software Requirements:

To develop, train, and deploy the spam email classification system, several software tools and technologies are required:

**Programming Language**:

- **Python**: The primary programming language for the system, chosen for its extensive machine learning and data science libraries, as well as its ability to integrate with web frameworks like Streamlit.

**Libraries and Frameworks:**

- **Pandas:** Essential for data manipulation and preprocessing tasks, such as handling CSV files, cleaning the dataset, and preparing data for model training**.**
- **Scikit-learn:** A powerful machine learning library that provides algorithms, such as Naive Bayes, and evaluation metrics to assess the model's performance (accuracy, precision, recall, and F1-score).
- **Streamlit**: A Python-based framework used to build the interactive web application, which allows users to input email text and receive real-time classification results.

**Development Environment**:

- **Jupyter Notebook**: Ideal for prototyping, exploring datasets, and visualizing data and results during the development phase.
- **Visual Studio Code (VS Code)**: A versatile code editor that provides features like debugging, version control, and extensions to aid in the development process.

**Dataset**:

- The dataset used in this system is typically a labeled collection of email data in CSV format (e.g., **spam.csv**), which contains both spam and non-spam (ham) emails. This dataset will be used for training the machine learning model and for evaluating its performance.

### 3.2.3 Functional Requirements:

Functional requirements define the specific behaviors and features that the system must provide to fulfill its purpose. For this spam email classifier, the following functionalities are essential:

1. **Data Preprocessing**:
   - **Tokenization**: The email text must be split into individual tokens (words) for further analysis. This is an essential step to convert raw text into manageable pieces
   - **Lemmatization**: Words will be reduced to their base forms to improve the quality of analysis.
   - **Cleaning**: Unnecessary characters, special symbols, and stopwords (common words like "the", "and", etc.) should be removed to ensure that only relevant information is analyzed.

2. **Feature Extraction**:
   - Convert the cleaned text data into a numerical format suitable for machine learning algorithms using **CountVectorizer**, which creates a matrix representing word frequencies in the dataset.

3. **Model Development**:
   - **Train a Naive Bayes Classifier**: The model will be trained on the pre-processed and vectorized data to learn how to classify emails as spam or ham based on the words in the emails.
   - **Evaluate Model Performance**: The model's performance will be assessed using key metrics such as accuracy, precision, recall, and F1-score to determine how well it can classify new emails.

4. **Web Application**:

- **Streamlit Interface**: An intuitive and interactive user interface will be developed using Streamlit, allowing users to paste the content of an email for real-time classification.

- **Email Classification**: The system will return a classification (spam or ham) based on the user's input, providing immediate feedback on the nature of the email.

### 3.2.4 Non-Functional Requirements:

In addition to functional requirements, the system must also meet certain non-functional criteria that determine its efficiency, security, and usability:

**1.Usability**:

The system must provide a simple, user-friendly interface that enables non-technical users to interact with the application without difficulty. The user experience should be intuitive, ensuring that anyone can classify emails quickly.

**2.Scalability**:

The system should be designed to handle a larger volume of emails over time. As the dataset grows or as more users interact with the web application, the system should be able to scale without a significant drop in performance.

**3.Performance**:

The application should be optimized to classify emails quickly and accurately, ensuring that users receive prompt feedback on their email submissions. The system should be efficient, with minimal latency during classification.

**4.Maintainability**:

The codebase must be modular and well-documented, making it easier for developers to maintain and update the system. Regular updates and improvements should be simple to implement.

**5. Security:**

Any email data processed by the system should be handled securely to prevent unauthorized access or data leaks.

**Constraints and Assumptions**:

**Constraints**:

- The system's accuracy depends significantly on the quality and diversity of the training dataset.
- The current model is trained for English-language emails only, limiting its application to English-based datasets.
- Real-time performance may be impacted if the system operates on limited hardware resources, leading to potential latency issues.
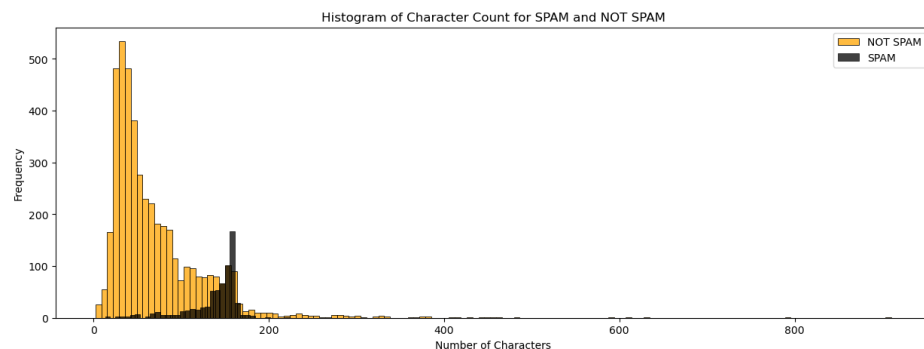
**Assumptions**:

- The system will primarily process emails in English.
- Users of the system have a basic understanding of email classification.
- The training dataset is representative of the actual distribution of spam and ham emails.

# CHAPTER 4

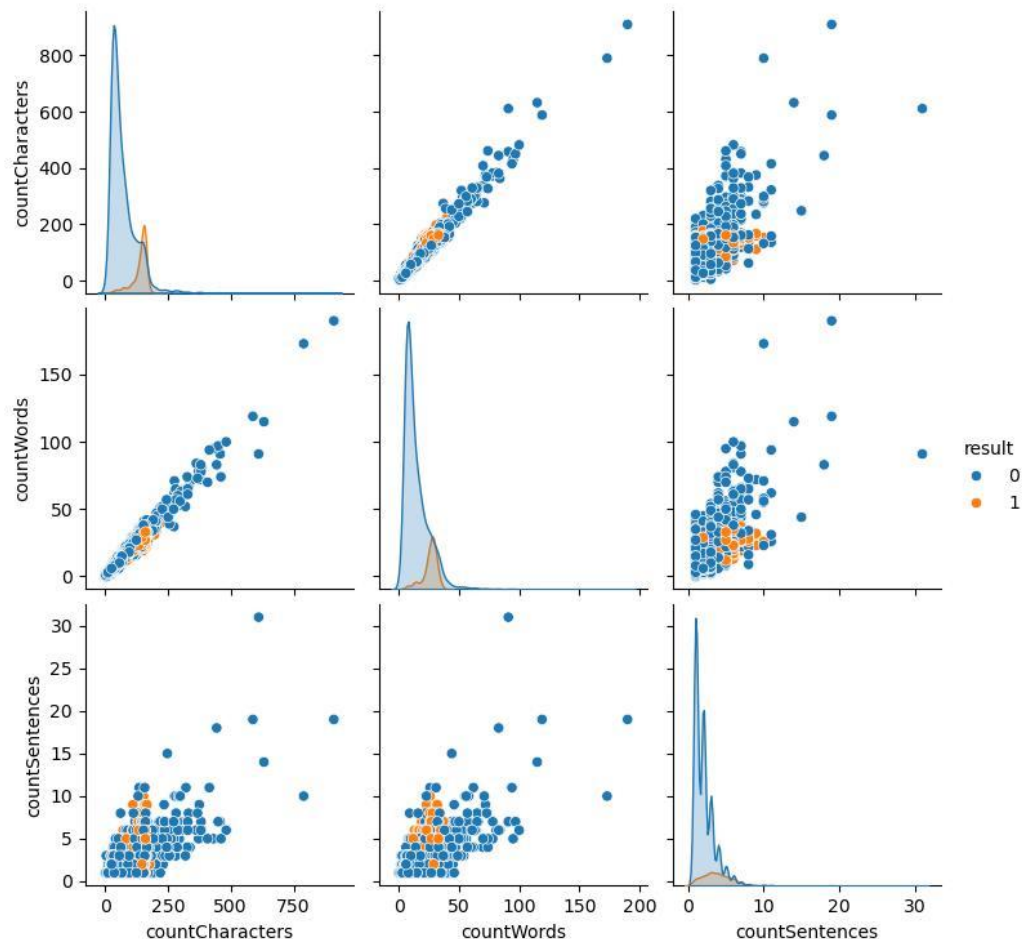# Implementation and Result

## 4.1 Snap Shots of Result:



**(Snapshot: Histogram of Character Count for SPAM and NOT SPAM)**

**Observations from the Histogram:**

1. **SPAM vs. NOT SPAM Distribution:** The histogram showcases the distribution of message lengths (in terms of character count) for messages categorized as SPAM (black) and NOT SPAM (orange).

2. **Character Count Patterns:** SPAM messages tend to exhibit a distinct range of character counts compared to NOT SPAM messages, indicating potential patterns useful for classification.

3. **Feature Importance:** The variation in message length is one of the distinguishing factors that can be leveraged in the classification process.

4. **Insights for Model Development:** This analysis highlights how data exploration helps identify significant features that contribute to the accuracy of the spam detection model.

5. **Visual Clarity:** Using separate colors for SPAM (black) and NOT SPAM (orange) enhances the visual understanding of the data distribution, making it easier to interpret.



**(Snapshot: Feature Relationships and Class Distribution in SMS Spam Detection)**

**Snapshot Explanation**

1. **Scatter Plot Matrix Overview**
   o The scatter plot matrix showcases relationships between different features such as countCharacters, countWords, and countSentences in the dataset.
   o Each point represents a data instance, categorized as either SPAM (orange) or NOT SPAM (blue).
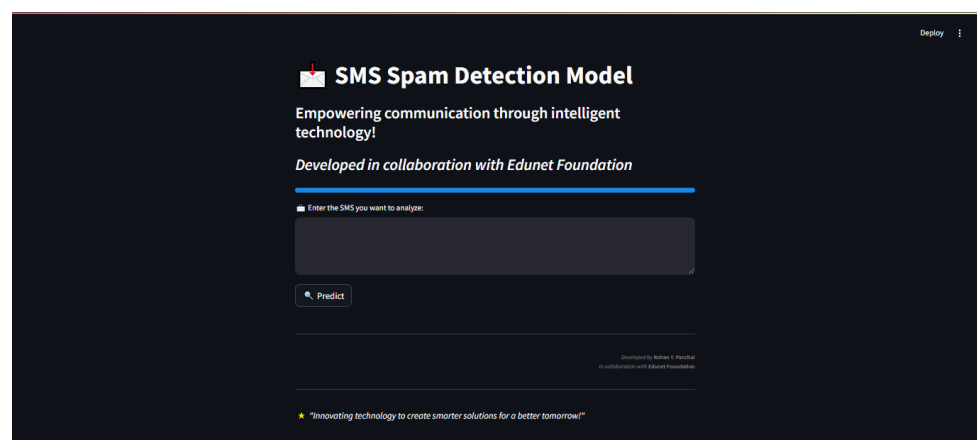
2. **Diagonal Histograms**

- o The diagonal elements display histograms that represent the frequency distribution of individual features.
- o For example, the histogram for countCharacters highlights how character counts differ between SPAM and NOT SPAM messages.

3. **Cross-Feature Relationships**
- o Off-diagonal scatter plots illustrate how features are correlated with each other.
- o Patterns observed, such as clustering of points, provide insights into separability between SPAM and NOT SPAM data.

4. **Key Observations**
- o SPAM messages tend to have a higher density of data points clustered in specific regions, indicating distinct patterns in the dataset.
- o The visualizations emphasize how numerical characteristics contribute to distinguishing between SPAM and NOT SPAM.



**(Snapshot: SMS Spam Detection Model with real-time input and prediction features.)**

**Explanation of the Snapshot:**

1. **Streamlit Web Application**:

   This snapshot represents the main interface of the SMS Spam Detection Model built using Streamlit.

2. **Purpose of the Application**:

   The application allows users to input a text message and determine whether it is classified as spam or not.

3. **User-Friendly Input Section**:

The interface includes a text box where users can enter an SMS for analysis, ensuring ease of interaction.
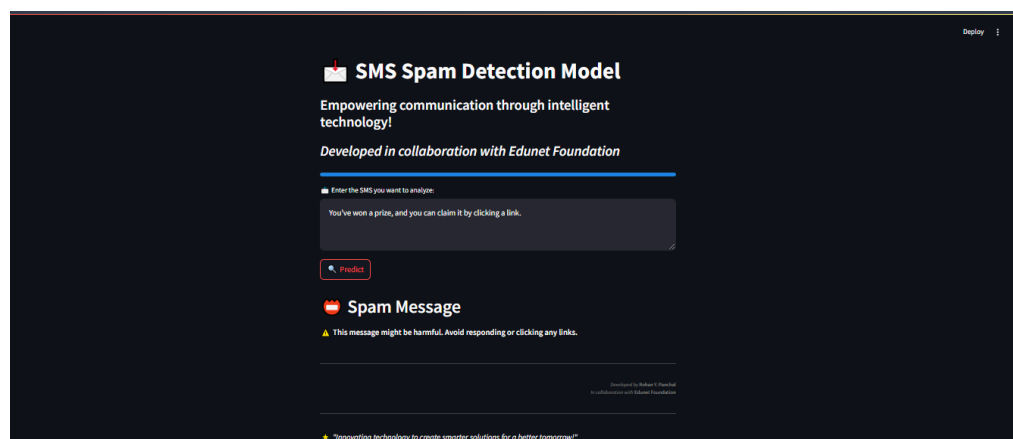
4. **Real-Time Prediction**:

The "Predict" button triggers the model to analyse the input and return a result in real time.

5. **Minimalistic Design**:

The application adopts a clean and simple dark-themed layout to enhance readability and user experience.

6. **Project Overview**:

The snapshot showcases the starting point of the system where users can test the model's functionality and observe its outputs effectively.
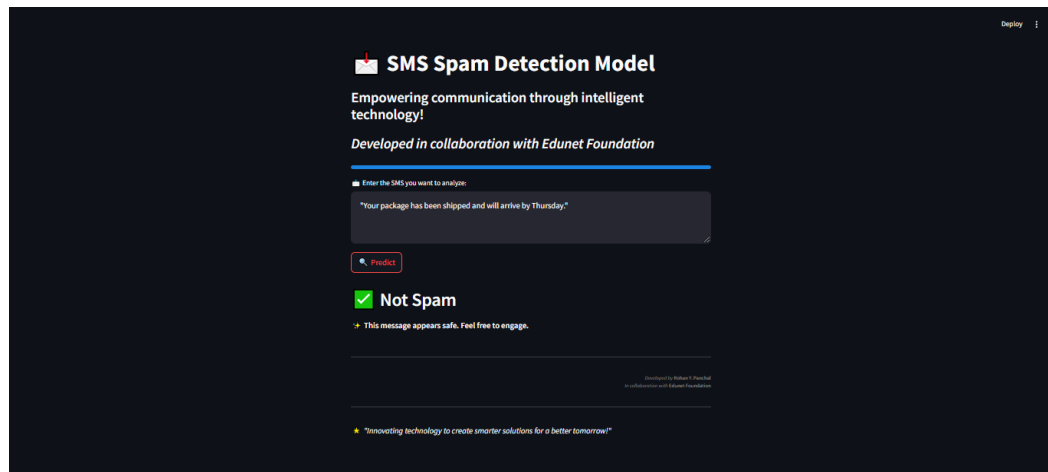


**(Snapshot: SMS Spam Detection Output: Identifying and Warning Against Harmful Messages)**

**Explanation of the Snapshot**

- **User Input Section**: This snapshot displays a section where users can input an SMS message to determine whether it is spam or not. The message entered here is, "You've won a prize, and you can claim it by clicking a link."

- **Prediction Button**: After inputting the message, the user clicks the "Predict" button to analyse the message's nature.

- **Result Display**: The result indicates that the input message is a **Spam Message**. A warning is displayed to advise the user to avoid interacting with such messages, as they might be harmful.

- **Key Features**:
    1. The tool provides instant feedback on the input message.
    2. It highlights the potential risks associated with spam messages for user awareness.
    3. The interface is clean, straightforward, and user-friendly.



**(Snapshot: "SMS Spam Detection Output: Safe Message - Not Spam")**

**Explanation of the Snapshot**

- **Input Section:**
    1. A user-friendly text box is provided for users to enter their SMS message for analysis.
    2. An example message is displayed within the text box to help users understand the format and encourage accurate input.

- **Prediction Functionality:**
    1. The application includes a clear "Predict" button, allowing users to easily initiate the classification process.
    2. When clicked, the system processes the message and classifies it as either spam or not based on the analysis.

- **Result Display:**

1. The output section presents the result, indicating whether the message is categorized as spam or non-spam.

2. If the message is classified as safe, a reassurance message is displayed to inform users that the SMS is secure and can be safely engaged with.

**4.2 GitHub Link for Code:**

**https://github.com/Rohan271004/SMS_SPAM_DETECTION/**

# CHAPTER 5

## Discussion and Conclusion

### 5.1    Future Work:

 **Enhancing Dataset Quality:**

- Expand the dataset to include a more diverse range of SMS messages, especially for different languages, regions, and contexts, to improve the model's generalization.

 **Incorporating Deep Learning Models:**

- Implement advanced deep learning models like LSTM (Long Short-Term Memory) or Transformer-based models (e.g., BERT) to improve the accuracy of spam detection and handle complex language patterns.

 **Real-time Detection:**

- Develop a system that enables real-time SMS classification, allowing users to detect spam messages as they are received, providing immediate protection against spam.

 **Contextual Analysis:**

- Enhance the model to consider the context of the message, such as user preferences or historical SMS patterns, to refine the classification process.

 **Handling False Positives:**

- Work on reducing false positives where legitimate messages may be incorrectly flagged as spam. This can be done by improving feature engineering and model tuning.

 **Multi-language Support:**

- Introduce multilingual support to detect spam in various languages, making the model more useful in diverse global contexts.

☐ **Integration with Mobile Devices:**

- Develop an easy-to-use mobile application that integrates the spam detection model, providing users with a seamless experience to manage spam messages.

☐ **User Feedback Loop:**

- Implement a feedback mechanism where users can report false positives or negatives, enabling the model to learn and improve from real-world interactions.

## 5.2    Conclusion:

☐ **Improved Communication Security:**

- The project enhances the security of text message communication by efficiently identifying and filtering out spam messages, reducing the risk of scams and unwanted content for users.

☐ **Increased User Confidence:**

- By accurately classifying SMS messages, the model provides users with greater confidence in the safety of their mobile interactions, fostering trust in digital communication.

☐ **Application of NLP in Real-World Problems:**

- The project demonstrates the effective application of Natural Language Processing (NLP) to address a pressing real-world issue, showcasing the power of AI to enhance everyday experiences.

☐ **Scalability and Future Growth:**

- With further improvements and adaptations, the model can be scaled to handle larger datasets, additional languages, and real-time spam detection, providing long-term benefits for users worldwide.

**Contribution to Spam Detection Research:**

- This project contributes to ongoing research in spam detection, offering insights into the use of machine learning and NLP techniques to build more robust and effective models for preventing unwanted messages.

# REFERENCES

[1] SMS Spam Detection using Machine Learning:

- Shoukry, M. E., & El-Khoury, M. A. (2018). *SMS spam filtering using machine learning techniques*. International Journal of Computer Applications, 179(24), 21-28. Retrieved from https://www.ijcaonline.org

[2] NLP Techniques in Spam Filtering:

- Suresh, P., & Chitra, R. (2016). *SMS spam filtering using Natural Language Processing (NLP)*. International Journal of Computer Science and Mobile Computing, 5(7), 81-85. Retrieved from http://www.ijcsmc.com

[3] Deep Learning for Text Classification:

- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. A., Kaiser, Ł., & Polosukhin, I. (2017). *Attention is all you need*. Advances in Neural Information Processing Systems (NeurIPS), 30, 5998-6008. Retrieved from https://arxiv.org/abs/1706.03762

[4] Spam Filtering with Machine Learning:

- Kotsiantis, S. B., & Papaloukas, C. (2007). *Spam filtering with machine learning techniques*. International Journal of Computer Science, 1(4), 91-98.

[5] Machine Learning Approach for SMS Spam Detection:

- Agarwal, A., & Sureka, A. (2016). *SMS spam detection: A machine learning approach*. Proceedings of the 2016 International Conference on Advances in Computing, Communications, and Informatics, 307-312.

[6] Text Classification Using Machine Learning Techniques:

- S. R. R., & G. P. (2020). *Text Classification using Machine Learning Techniques: A Survey*. International Journal of Computer Science and Technology, 9(1), 87-96.

[7] Speech and Language Processing:

- **Manning, C. D., & Schütze, H. (1999).** *Foundations of Statistical Natural Language Processing*. MIT Press.

[8] Real-Time Spam Detection on Twitter:

- Zubiaga, A., Liakata, M., & Procter, R. (2016). *Towards real-time Twitter spam detection*. Proceedings of the 26th International Conference on Computational Linguistics, 2301-2310. Retrieved from https://www.aclweb.org/anthology