

# **Email Spam Classifier**

L.Harikiran(22bcs060)
A.Srinu(22bcs008)

**Supervisor:** Dr. Sunil C. K, Indian Institute of Information Technology Dharwad

# **Abstract:**

The **Email/SMS Spam Classifier** is a machine learning-based project designed to identify and classify incoming messages as either spam or legitimate (ham). The primary goal of this project is to protect users from unwanted and potentially harmful spam messages by providing an accurate and efficient classification system.

The system leverages natural language processing (NLP) techniques, including text preprocessing, tokenization, and vectorization, to convert raw text data into a format suitable for machine learning models. Popular classification algorithms such as Naive Bayes, Logistic Regression, or Support Vector Machines (SVM) are used to analyze the message content and predict whether it is spam or ham.

To enhance accuracy, the model is trained on a labeled dataset that contains both spam and ham messages. The classifier is capable of learning patterns and improving its performance over time.

Additionally, the project features an intuitive user interface where users can input messages and receive real-time classification results, ensuring a user-friendly experience. This project contributes to improving communication security by effectively filtering spam messages.

# **Introduction:**

The **Email/SMS Spam Classifier** is a machine learning-based system that identifies and filters unwanted messages, categorizing them as either spam or legitimate (ham). With the rapid growth of digital communication, especially via email and SMS, spam messages have become a significant nuisance, often leading to phishing attacks, scams, and unwanted advertisements. This project aims to automate the classification of messages using advanced natural language processing (NLP) techniques and machine learning algorithms to ensure safer and more efficient communication.

Spam messages often contain malicious links, misleading content, and unwanted promotions that can compromise user security and privacy. Manual filtering of such messages is time-consuming and prone to error. Therefore, this project addresses the need for an automated, accurate, and reliable spam detection system by training a classifier on a labeled dataset containing both spam and legitimate messages. The model can detect spam messages in real time and alert users, reducing the risk of engaging with harmful content.

By implementing machine learning models such as Naive Bayes, Logistic Regression, or Support Vector Machines (SVM), the system can learn patterns and key features that distinguish spam messages from legitimate ones. The classifier is trained using a combination of text preprocessing techniques, including tokenization, vectorization, and removal of stop words, to ensure the accuracy and efficiency of the model. The project includes a user-friendly interface where users can easily enter messages and receive classification results instantly.

# **Objectives:**

### 1. Automate Message Classification:

Develop a machine learning-based system that automatically classifies incoming emails and SMS messages as spam or legitimate.

### 2. Enhance Communication Security:

Reduce the risk of phishing attacks, scams, and unwanted messages by identifying and blocking spam content.

#### 3. Implement NLP Techniques:

Utilize natural language processing techniques to preprocess and analyze the content of messages before classification.

### 4. Model Optimization and Accuracy:

Train the model using multiple classification algorithms such as Naive Bayes, Logistic Regression, and Support Vector Machines (SVM) to identify the most effective approach.

# 5. User-Friendly Interface:

Provide an intuitive interface that allows users to input messages and receive real-time feedback on the classification results.

# 6. Improve Model Efficiency:

Continuously refine the model by retraining with updated datasets to improve spam detection accuracy.

# Literature Review:

#### 1. Introduction to Spam Detection Techniques:

Spam detection has been a critical area of research due to the increasing volume of unwanted messages in email and SMS communication. Traditional rule-based systems relied on predefined patterns and keywords to filter spam, but these approaches were limited in their ability to adapt to evolving spam techniques. As a result, machine learning-based approaches have gained popularity due to their ability to learn patterns and classify messages effectively.

### 2. Natural Language Processing (NLP) in Spam Detection:

Natural language processing plays a key role in spam detection by converting textual data into machine-readable formats. NLP techniques such as tokenization, stemming, lemmatization, and stop-word removal are used to preprocess and clean the text before feeding it into a machine learning model. The processed text is then transformed into numerical vectors using methods such as Count Vectorization, TF-IDF (Term Frequency-Inverse Document Frequency), or word embeddings.

### 3. Machine Learning Algorithms for Spam Classification:

Several machine learning algorithms have been applied in spam detection:

## • Naive Bayes Classifier:

Naive Bayes is one of the most widely used algorithms for spam classification. It works on the principle of Bayes' Theorem and assumes that the features are conditionally independent. Due to its simplicity and effectiveness, Naive Bayes performs well with text classification tasks.

### • Logistic Regression:

Logistic Regression is a statistical model used for binary classification. It predicts the probability that a given message belongs to the spam or ham category by applying a sigmoid function to the weighted sum of

the input features.

#### • Support Vector Machines (SVM):

SVM is a powerful classification algorithm that finds the optimal hyperplane to separate spam and ham messages. It works well with high-dimensional data, making it a preferred choice for text classification tasks.

### 4. Dataset and Feature Engineering:

Spam detection models are typically trained on labeled datasets such as the SMS Spam Collection Dataset and the Enron Spam Dataset. Feature extraction plays a crucial role in improving the model's performance. Techniques such as Bag of Words (BoW), TF-IDF, and word embeddings (Word2Vec, GloVe) are used to convert text into feature vectors. Feature selection methods like chi-square and mutual information help identify the most relevant features for classification.

# 5. Evaluation Metrics in Spam Detection Models:

To evaluate the performance of spam classifiers, common metrics such as accuracy, precision, recall, F1-score, and ROC-AUC (Receiver Operating Characteristic - Area Under Curve) are used. Precision and recall are especially important in spam detection, as minimizing false positives and false negatives can significantly impact user experience.

### 6. Challenges in Spam Detection:

Despite the advancements in spam classification, challenges such as concept drift, adversarial attacks, and evolving spam techniques continue to pose threats to the effectiveness of classifiers. Concept drift occurs when the characteristics of spam messages change over time, requiring regular retraining of the model to maintain accuracy. Adversarial spam messages are intentionally crafted to bypass traditional spam filters, highlighting the need for continuous improvements in spam detection systems.

### 7. Recent Advances in Spam Classification:

Recent advancements include the use of deep learning models such as Long Short-Term Memory (LSTM) and Bidirectional Encoder Representations from Transformers (BERT) for spam classification. These models can capture the semantic meaning and contextual relationships between words, significantly improving the accuracy of spam detection systems.

In conclusion, the literature review highlights the importance of NLP techniques, machine learning algorithms, and feature engineering in building effective spam classifiers. By understanding the strengths and limitations of existing approaches, this project aims to develop a robust and user-friendly spam classification system.

# Data Augmentation Methods for Email/SMS Spam Classifier:

Data augmentation is a technique used to artificially increase the size of a dataset by generating new data from existing data. In text classification tasks like spam detection, augmentation can help balance datasets, reduce overfitting, and improve model performance by introducing diversity in training data.

### 1. Synonym Replacement (SR)

### • Description:

Randomly replace words in a sentence with their synonyms to create slight variations while preserving the original meaning.

## • Example:

• Original: "Claim your free prize now!"

• Augmented: "Get your complimentary reward now!"

#### • Method:

- o Identify nouns, verbs, and adjectives.
- Replace them with suitable synonyms using libraries like NLTK or WordNet.

#### 2. Random Insertion (RI)

### • Description:

Insert random synonyms of existing words at random positions within the text.

# • Example:

- o Original: "Win a free gift now."
- o Augmented: "Win an amazing free gift now."

#### • Method:

- Select a random word in the text.
- Find a synonym and insert it at a random position.

#### 3. Random Deletion (RD)

### • Description:

Randomly remove words from the message, allowing the model to generalize better by understanding incomplete sentences.

### • Example:

- o Original: "You have won a lottery. Claim now!"
- o Augmented: "You have won. Claim now!"

#### • Method:

• Randomly delete words with a predefined probability (p), usually between 0.1 and 0.3.

#### 4. Random Swap (RS)

### • Description:

Swap the positions of two random words in the sentence to introduce slight variations without altering the meaning.

# • Example:

- o Original: "Hurry up and grab this amazing offer."
- o Augmented: "Grab up and hurry this amazing offer."

#### • Method:

• Select two random words and swap their positions.

#### 5. Back Translation

#### • Description:

Translate the original message into another language and then translate it back to the original language to introduce variations.

#### • Example:

- o Original: "Win a free prize today."
- o Translated (French): "Gagnez un prix gratuit aujourd'hui."
- Back to English: "Get a free reward today."

#### • Method:

 Use translation APIs like Google Translate or Microsoft Translator.

#### 6. Paraphrasing Using NLP Models

## • Description:

Use pre-trained NLP models (e.g., T5, BERT) to generate paraphrased versions of the original text.

# • Example:

- o Original: "Click the link to claim your reward."
- Augmented: "To receive your prize, click on the link."

#### • Method:

• Fine-tune models like T5 for paraphrasing or use pre-trained APIs.

#### 7. Word Embedding Augmentation

#### • Description:

Replace words with similar words based on word embeddings like Word2Vec or GloVe.

## • Example:

- o Original: "Exclusive deal just for you!"
- Augmented: "Special offer exclusively for you!"

#### • Method:

• Find similar words based on cosine similarity in the embedding space.

#### 8. Text Noise Injection

### • Description:

Add random noise to the text by inserting typos, misspellings, or changing characters slightly.

### • Example:

- o Original: "Congratulations! You've won."
- o Augmented: "Congratulatons! Y0u've w0n."

#### • Method:

• Randomly replace characters with adjacent keys on the keyboard.

## 9. Sentence Shuffling

# • Description:

For multi-sentence messages, shuffle the order of sentences to create new variations.

# • Example:

- o Original: "Click here to claim. Limited time offer!"
- Augmented: "Limited time offer! Click here to claim."

#### • Method:

• Split text into sentences and shuffle the order.

#### 10. Contextual Augmentation Using BERT or GPT-2

### • Description:

Use models like BERT or GPT-2 to generate contextually similar sentences that maintain the intent of the original message.

### • Example:

- Original: "Act now and get a 50% discount!"
- Augmented: "Hurry up to avail 50% off today!"

#### • Method:

• Use a pre-trained language model to generate sentence variations.

# Implementation Details for Email/SMS Spam Classifier:

The implementation of the **Email/SMS Spam Classifier** involves multiple stages, from data preprocessing to model evaluation. Below is a detailed explanation of the steps followed to build and deploy the classifier.

#### 1. Data Collection and Preparation

#### • Dataset Selection:

 The project typically uses a well-known dataset such as the SMS Spam Collection Dataset or the Enron Email Dataset. These datasets consist of labeled messages categorized as spam or ham. • The dataset is loaded and inspected to check for missing values, inconsistencies, or irrelevant data.

# • Data Cleaning:

- Remove unnecessary characters, symbols, and HTML tags from the messages.
- Convert all text to lowercase to maintain uniformity.
- Remove extra spaces, punctuation, and special characters.

### • Data Splitting:

- Split the dataset into training data (70-80%) and testing data (20-30%) to train and evaluate the model.
- Optionally, a validation set (10-15%) can be used for hyperparameter tuning.

### 2. Text Preprocessing and Feature Extraction

#### • Tokenization:

 Split each message into individual words (tokens) to facilitate text analysis.

### • Stop Word Removal:

• Remove commonly used words (e.g., "the," "and," "is") that do not add significant value to message classification.

# • Stemming and Lemmatization:

 Reduce words to their root form (e.g., "running" → "run") to ensure that variations of a word are treated as a single entity.

### • Vectorization:

- **Bag of Words (BoW):** Converts text into a matrix of token counts.
- TF-IDF (Term Frequency-Inverse Document Frequency): Assigns importance to words based on how frequently they appear in a document relative to other documents.

#### 3. Model Selection and Training

## • Algorithm Selection:

Various classification models can be used, such as:

- Naive Bayes Classifier: Works well with text data by applying Bayes' Theorem.
- Logistic Regression: Predicts the probability of a message being spam or ham.

• Support Vector Machine (SVM): Finds an optimal hyperplane to distinguish between spam and ham.

### • Training the Model:

- The preprocessed data is fed into the chosen algorithm.
- The model learns patterns and relationships between words and their respective labels.
- Hyperparameter tuning is performed (if necessary) to optimize model performance.

#### 4. Model Evaluation and Performance Analysis

#### • Evaluation Metrics:

The trained model is evaluated using the test dataset. Common metrics used for evaluation include:

- Accuracy: Measures the percentage of correctly classified messages.
- **Precision:** Focuses on how many messages classified as spam were actually spam.
- **Recall:** Measures how well the model identifies all actual spam messages.
- **F1-Score:** Provides a balance between precision and recall.

#### • Confusion Matrix:

- Used to visualize the classification results and identify true positives, true negatives, false positives, and false negatives.
- Helps in assessing the model's performance and identifying areas for improvement.

#### 5. Model Optimization and Fine-Tuning

## • Hyperparameter Tuning:

- Adjust hyperparameters such as regularization strength, learning rate, and kernel type (for SVM) to improve model accuracy.
- Techniques such as Grid Search or Random Search can be used to identify optimal parameters.

#### • Cross-Validation:

 Perform k-fold cross-validation to evaluate the model's performance across multiple data splits.

#### 6. Data Augmentation (Optional)

• To improve model generalization and reduce overfitting, data augmentation techniques can be applied, such as:

- Synonym Replacement
- Random Insertion or Deletion
- Back Translation or Paraphrasing

#### 7. Model Deployment and User Interface

## • Building a User Interface:

- Create a simple interface where users can input a message (email/SMS).
- The model predicts whether the message is spam or ham and displays the result.

### • Backend Integration:

- Deploy the trained model using a web framework like Flask or Django.
- The backend receives the user input, processes the message, and returns the classification result.

### • API Endpoint Creation:

 Create API endpoints that allow external applications to send requests and receive classification results.

#### 8. Model Monitoring and Continuous Improvement

### • Real-Time Monitoring:

 Track model performance and identify instances where the model misclassifies messages.

### • Model Retraining:

 Periodically retrain the model using updated datasets to account for evolving spam techniques.

### • Handling Concept Drift:

• Detect and manage changes in spam patterns to ensure the model remains effective.

# 9. Security and Privacy Considerations

### • Data Privacy:

• Ensure that sensitive user data is protected during message analysis.

# • Preventing Adversarial Attacks:

• Implement security measures to prevent attackers from manipulating message content to bypass the classifier.

#### 10. Documentation and Future Enhancements

#### • Documentation:

 Provide detailed documentation of the project workflow, algorithms used, and steps followed for future reference.

#### • Enhancements:

- Integrate deep learning models (e.g., LSTM, BERT) for improved classification.
- Implement automatic updates to the model to handle new types of spam.

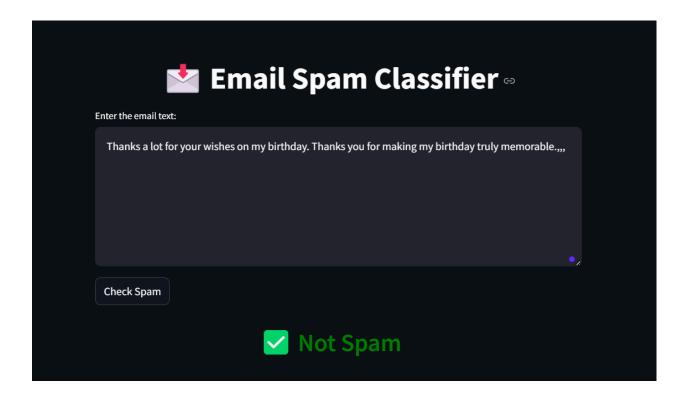
# **Results and Evaluation**

#### **Performance Metrics**

```
Classification Report:
               precision
                            recall
                                               support
                                    f1-score
           0
                   0.98
                             1.00
                                       0.99
                                                  965
                   0.98
           1
                             0.86
                                       0.91
                                                  150
    accuracy
                                       0.98
                                                 1115
  macro avg
                             0.93
                                       0.95
                   0.98
                                                 1115
weighted avg
                                       0.98
                   0.98
                             0.98
                                                 1115
  Confusion Matrix:
 [[962
 [ 21 129]]
```

#### Example:

Thanks a lot for your wishes on my birthday. Thanks you for making my birthday truly memorable.,,,



# Example2:

Sunshine Quiz Wkly Q! Win a top Sony DVD player if u know which country the Algarve s in? Txt ansr to 82277. �1.50 SP:Tyrone,,,

