

Methodologies

3.1 Introduction .

In today's digital age, email communication has become an integral part of our daily lives, serving as a primary means of communication for personal and professional purposes. However, alongside legitimate communication, the proliferation of spam emails poses significant challenges, including security risks, privacy concerns, and productivity losses. To address these challenges, the development of effective spam email classification systems is essential.

The objective of this mini-project is to design, implement, and deploy a spam email classifier using state-of-the-art machine learning techniques. Specifically, we employ the Multinomial Naive Bayes algorithm, a probabilistic classifier known for its simplicity, efficiency, and effectiveness in text classification tasks.

Multinomial Naive Bayes

The Multinomial Naive Bayes (MNB) algorithm is a variant of the Naive Bayes classifier, which is based on Bayes' theorem and the assumption of feature independence. In text classification tasks, such as spam email detection, MNB is particularly well-suited due to its ability to handle high-dimensional and sparse feature spaces typical of textual data.

In MNB, each feature represents the frequency of a term (word) in a document, and the classifier calculates the conditional probability of a class (spam or ham) given the document's feature vector. Despite its simplicity, MNB often achieves competitive performance and is widely used in various natural language processing (NLP) applications.

Streamlit Deployment

To make our spam email classifier accessible and user-friendly, we leverage the Streamlit framework for deployment. Streamlit is an open-source Python library that allows for the creation of interactive web

applications with minimal effort. It provides a straightforward and intuitive interface for building data-driven applications, making it ideal for showcasing machine learning models to a broader audience.

With Streamlit, we create a user-friendly interface where users can input an email text, and the classifier provides real-time predictions on whether the email is spam or ham. The Streamlit application streamlines the deployment process, enabling seamless integration of the trained model into a web-based environment without extensive web development experience.

3.2 Research Design .

The research design for this project follows a structured approach to develop and deploy a spam email classifier using the Multinomial Naive Bayes algorithm. The CRISP-DM methodology provides a framework for conducting data mining projects, guiding the process from initial data exploration to model deployment.

Naïve Bayes Algorithm (NBA) is also known NB classifier is kind of machine learning classifier which is Bayes theorem related to probabilistic classifier that used for text classification. Bayes theorem has some attributes such as ability to arrange scale-large datasets, strong independence and probability distribution. This is one types of popular machine learning algorithm that called supervised. This has been improved during the Bayes' rule which attempt to derive the probability of an event occurrence based on even associated prior knowledge and conditions [10]. Thus, Naïve Bayes algorithm utilized Bayes theorem to determine that probabilities spam e-mail. some words contain a chance of occurring in either spam or non-spam. for instance, there are some words related to know spam such as Free, which never occur in a non-spam email. However, when we find a message containing this word [11]. We already say that is spam. In addition, there are also some words that are known as non-spam such as the name of family members and friends, So the the naive Bayes technique was trained to recognize and categorize words like "free" and "earn money," but it had a very low spam probability for words found in non-spam e-mail, like friend and family member names. Hence, to determine whether an email is likely spam or ham (non-spam). As shown in the formula below, the naive Bayes technique applied the Bayes theorem. [12]:

Multinomial Naive Bayes Algorithm

The Multinomial Naive Bayes (MNB) algorithm is a probabilistic classifier that is widely used for text classification tasks, such as spam email detection. It is based on Bayes' theorem with the "naive" assumption of independence between features. Here's a detailed explanation of how the Multinomial Naive Bayes algorithm works, along with its mathematical formulation:

Given a set of features $X = \{x_1, x_2, \dots, x_n\}$ where each x_i represent a features. We have for the multinomial navye bayes the formula as follows

$$P(C_k|X) = P(X|C_k) \cdot P(C_k) / P(X);$$

Data Understanding:

Dataset Selection: The dataset used for this project is the "Spam-Ham Dataset" obtained from Kaggle. It consists of a collection of labeled emails, with labels indicating whether an email is spam (1) or ham (0).

Data Exploration: Initial exploration of the dataset involves understanding the distribution of spam vs. ham emails, examining the length and content of emails, and identifying any patterns or anomalies in the data.

Data Preparation:

Cleaning and Preprocessing: The raw email data undergoes preprocessing steps to clean and standardize the text. This includes removing non-ASCII characters, eliminating short words, and lemmatizing the text to reduce noise and improve model performance.

Feature Engineering: The email text is transformed into numerical features using the CountVectorizer from the Scikit-Learn library. This step converts the text data into a format suitable for training the Multinomial Naive Bayes model.

Modeling:

Algorithm Selection: The Multinomial Naive Bayes algorithm is chosen for its simplicity, efficiency, and effectiveness in text classification tasks. It is well-suited for handling the high-dimensional and sparse feature space of text data.

Training and Validation: The model is trained on a subset of the data and validated using cross-validation techniques to assess its generalization performance.

Evaluation:

Performance Metrics: The performance of the trained model is evaluated using standard classification metrics such as accuracy, precision, recall, F1-score, and confusion matrix. These metrics provide insights into the model's predictive capabilities and help identify areas for improvement.

Deployment:

Deployment Platform: Utilize Streamlit, a Python library for creating interactive web applications, to deploy the trained model.

User Interface Design: Develop a user-friendly interface where users can input email text and receive instant predictions on whether the email is classified as spam or ham.

Integration: Embed the model prediction logic within the Streamlit application to ensure seamless interaction and real-time feedback.

3.3 Data Collection Methods .

The data collection involved acquiring the Spam-Ham Dataset from Kaggle, ensuring it met the project's requirements for labeled email data. This dataset provided a diverse range of emails, enabling comprehensive training and testing of the spam email classifier.

3.4 Data Analysis Techniques .

3.4.1 Text Preprocessing

Text preprocessing is crucial to ensure that the data is clean, standardized, and ready for analysis. The steps involved in text preprocessing for this project include:

Cleaning Text:

Non-ASCII Character Removal: Many email datasets, including the Spam-Ham Dataset used here, may contain non-ASCII characters that can interfere with text analysis and modeling. Cleaning involves removing these characters to ensure uniformity and readability of the text.

Function Example

```
import string
import re

def clean_text(s):
    # Remove non-ASCII characters
    s = ''.join(filter(lambda x: x in string.printable, s))
    return s.rstrip('\r\n')
```

Removing Little Words:

Short Word Removal: Words with fewer than a specified number of characters (e.g., less than 3) are often removed during preprocessing to reduce noise and improve the quality of feature extraction.

Function Example:

```
def remove_little(s):
    wordsList = s.split()
    k_length = 2 # Example threshold for minimum word length
    resultList = [element for element in wordsList if len(element) > k_length]
    resultString = ''.join(resultList)
    return resultString
```

Lemmatization:

Normalization: Lemmatization reduces words to their base or root form to handle variations like pluralization or verb conjugation. This process helps in standardizing the vocabulary and improving the model's understanding of the text.

Function Example:

```
import spacy

nlp = spacy.load("en_core_web_sm")

def lemmatize_text(text):
    doc = nlp(text)
    lemmatized_text = " ".join([token.lemma_ for token in doc])
    return lemmatized_text
```

3.4.2 Exploratory Data Analysis (EDA)

Exploratory Data Analysis (EDA) involves visual and statistical methods to understand the dataset's characteristics and distributions. For this project, EDA focused on:

Word Cloud Generation:

Visualization: Word clouds provide a visual representation of the most frequent words in the dataset, helping to identify common terms and themes in both spam and ham emails.

Library Used: WordCloud from the wordcloud library was utilized for generating word clouds.

```
from wordcloud import WordCloud

# Generating a word cloud for spam emails

word_cloud = WordCloud(stopwords=STOPWORDS, background_color='white').generate("
".join(df[df['label']==1]['text']))
```

Frequency Analysis:

Counting Words: Analyzing the frequency of words in the dataset to determine prevalent terms. This analysis aids in feature selection and understanding the vocabulary used in spam and ham emails.

Example:

```
def count_words(s):
    global dic_all
    wordsList = s.split()
    for w in wordsList:
        if not w in dic_all:
            dic_all[w] = 1
        else:
            dic_all[w] += 1
```

3.4.3 Feature Extraction

Feature extraction transforms raw text data into numerical features that machine learning models can process. In this project, the primary technique used for feature extraction is:

CountVectorizer:

Purpose: Convert text data into a matrix of token counts, where each row represents an email and each column represents a word in the vocabulary.

Functionality: CountVectorizer tokenizes the text, builds a vocabulary of known words, and counts the occurrences of each word in the text documents.

Example:

```
from sklearn.feature_extraction.text import CountVectorizer

vectorizer = CountVectorizer()

X = vectorizer.fit_transform(df['text'])
```

Stopword Removal:

Enhancement: Incorporating additional stopwords (common words like 'the', 'and', etc.) into the analysis to filter out irrelevant words that do not contribute to the classification of emails.

Example:

```
from wordcloud import STOPWORDS

more_stopwords = {'re', 's', 'subject', 'hpl', 'hou', 'enron'}

STOPWORDS = STOPWORDS.union(more_stopwords)
```

3.5 Ethical Considerations .

Ethical considerations are paramount in developing and deploying a spam email classifier:

Privacy: Ensured that the dataset used does not contain sensitive or personally identifiable information, prioritizing user confidentiality.

Bias Mitigation: Strived for a balanced dataset representation to avoid biased model predictions that may disproportionately impact certain email types.

Transparency: Documented the model's limitations and performance metrics clearly, providing users with insights into its reliability and potential errors.

3.6 Limitations .

Despite rigorous methodology, several limitations were identified:

Data Quality: Potential noise or mislabeling in the dataset may affect model accuracy.

Model Generalization: The model's performance may vary on unseen data not represented in the training dataset.

Feature Representation: CountVectorizer may not capture semantic nuances in email text, potentially limiting the model's understanding of context and content variations.