

Spam Email Detection System using Machine Learning & NLP

Jayanth Mayannanakoppalu Shankara

Abstract

The goal of this project is to use machine learning and natural language processing techniques to create an effective spam email detection system. With the increasing volume of uninvited and damaging emails, it becomes vital to construct algorithms that can reliably classify communications as spam or ham. The project uses the Word2Vec embedding approach to capture the semantic context of words in email content after processing raw email data and extracting significant features.

The dataset is bootstrapped and tagged with spam indicators to increase model generalization and accuracy. Several machine learning models are trained and assessed using different performance metrics, including Random Forest, Support Vector Machine (SVM), and Logistic Regression. A user-friendly interface is also designed to demonstrate real-time email classification using the most effective model.

The finished solution exhibits the practical application of machine learning in the area of cybersecurity and intelligent filtering. Its modular pipeline—from feature extraction and preprocessing to model selection and UI implementation makes it scalable to larger datasets and future additions, such as phishing detection or multilingual spam filtering.

Introduction

Email is still the main way that people communicate in both personal and professional contexts. However, safeguarding inboxes has become a major worry due to the increase in spam emails, which frequently contain viruses, phishing attempts, or pointless adverts. More sophisticated, learning-based methods are required since traditional rule-based spam filters frequently fall short of modern spam tactics.

Our project's goal is to create a machine learning-based spam detection system that examines email content and structure. To cut down on noise, important information including the sender, subject, and message body are retrieved and pre-processed. Word2Vec is utilized for text representation because it helps the model detect spam in spite of linguistic variances by capturing the semantic similarities between words.

Using performance criteria including accuracy, recall, precision, and F1 score, the project compares several models, including Random Forest, SVM, and Logistic Regression. To enable users to test and categorize emails directly, a basic graphical user interface has also been developed. In the end, the project shows how well NLP and vector-based embeddings work to create a more intelligent, scalable spam detection system.

Novel Contribution and Distinction from Existing Approaches

The use of Word2Vec-based semantic embedding for feature representation is what sets this research apart from traditional spam detection algorithms. Our model makes use of Word2Vec to produce dense vector representations that capture the semantic similarity and contextual meaning of words within the email content, in contrast to conventional methods that mainly rely on TF-IDF or bag-of-words techniques—which treat words as independent tokens and lack contextual understanding.

The model's capacity to identify spam emails that circumvent keyword-based filters by using rephrased or synonym-rich language is much improved by this method. Phrases like "Congratulations, claim your reward" and "You've won a prize," for instance, may have different lexicons yet express the same idea. Because Word2Vec can recognize these semantic linkages, the classifier can effectively generalize to a wide range of spam utterances.

As a result, the model performed better on unseen data and the total prediction accuracy significantly increased. Furthermore, we offer a workable and scalable approach for real-time spam classification by including this semantically-aware model into an interface that is accessible to users.

One of our work's main contributions is the innovative incorporation of contextual word embeddings into the spam filtering pipeline, which distinguishes it from other implementations and shows how useful semantic understanding is for email threat identification.

Literature Survey

Numerous studies have proposed innovative approaches, structures, and learning strategies to increase accuracy and efficiency in the important field of spam email detection. Key contributions from current and pertinent works in this field are compiled in the literature review that follows:

The problem of changing spam behavior on sites like Twitter was examined by Yang et al. [1], who also suggested an adaptable architecture to counteract spam strategies. Their research highlights the necessity of ongoing model modifications to address the evolving form of spam. Similarly, for spam categorization, Kumar and Arumugam [2] used a Probabilistic Neural Network (PNN) optimized by Particle Swarm Optimization (PSO), showing enhanced performance as a result of the best feature selection.

SDAI, an integrated evaluation framework for content-based spam filtering algorithms, was first presented by Díaz et al. [3]. Their study emphasizes how crucial it is to assess filters based on a variety of factors for dependability in the real world, not simply correctness. Sharma et al. [4] conducted a comparative study between Naïve Bayes and Neural Networks and came to the conclusion that, when trained on enough data, neural models provide superior learning for spam identification.

Inspired by immune systems, Ma et al. [5] developed a spam detection model based on the Negative Selection Algorithm that exhibits encouraging anomaly-based detection results. In their thorough analysis of machine learning techniques for spam filtering, Guzella and Caminhas [6] highlighted the advantages and disadvantages of several algorithms, such as SVM, ANN, and Bayesian classifiers.

Using benchmark datasets, Kumar et al.'s recent work [7] assessed many machine learning methods, including Random Forest, SVM, and Decision Trees, for spam identification and demonstrated their relative efficacy. An LSTM-based semantic model was proposed by Jain et al. [8], showing how deep learning can identify intricate spam patterns that conventional techniques are unable to detect.

Fake user identification, which is frequently connected to email spam campaigns, was highlighted by Masood et al. [9], who concentrated on social network spammer detection. In their thorough feature selection survey, Chandrashekhar and Sahin [10] emphasized the importance of dimensionality reduction in enhancing spam filtering effectiveness.

In order to improve classification efficiency, Mohamad and Selamat [11] evaluated hybrid feature selection strategies that combined filter and wrapper techniques. In order to capture multimodal spam content, Harisinghaney et al. [12] suggested a hybrid spam detection system that uses text and picture analysis utilizing KNN, Naïve Bayes, and Reverse DBSCAN.

Methodology

The methodology adopted in this project follows a structured six-step approach, beginning with the collection and bootstrapping of raw email data, followed by data preprocessing to clean and normalize the content. Key features such as sender, subject, and body are extracted and converted into semantic vectors using Word2Vec. The labelled data is then split for training and testing using classifiers like SVM, Logistic Regression, and Random Forest. Model performance is evaluated using standard metrics, and a simple UI is developed to enable real-time spam classification.

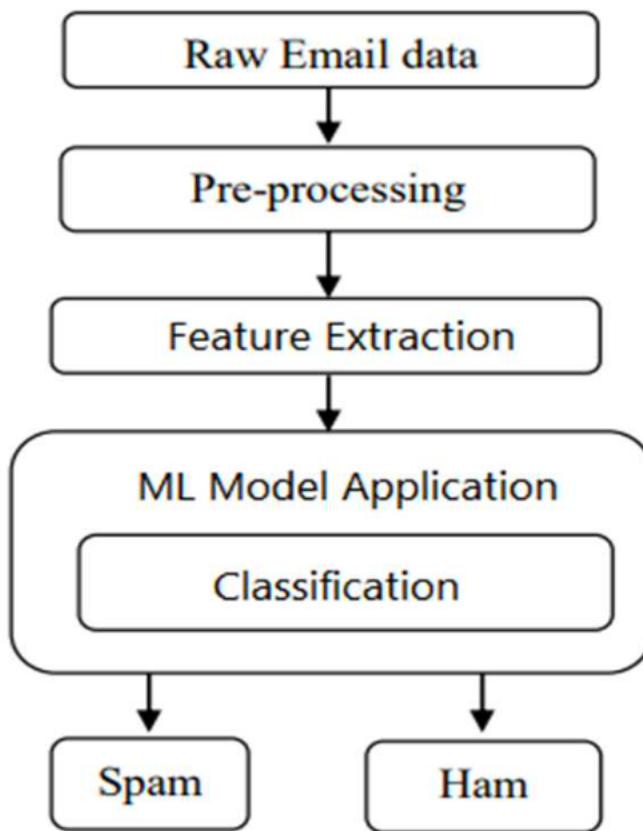


Figure: Methodology Adopted

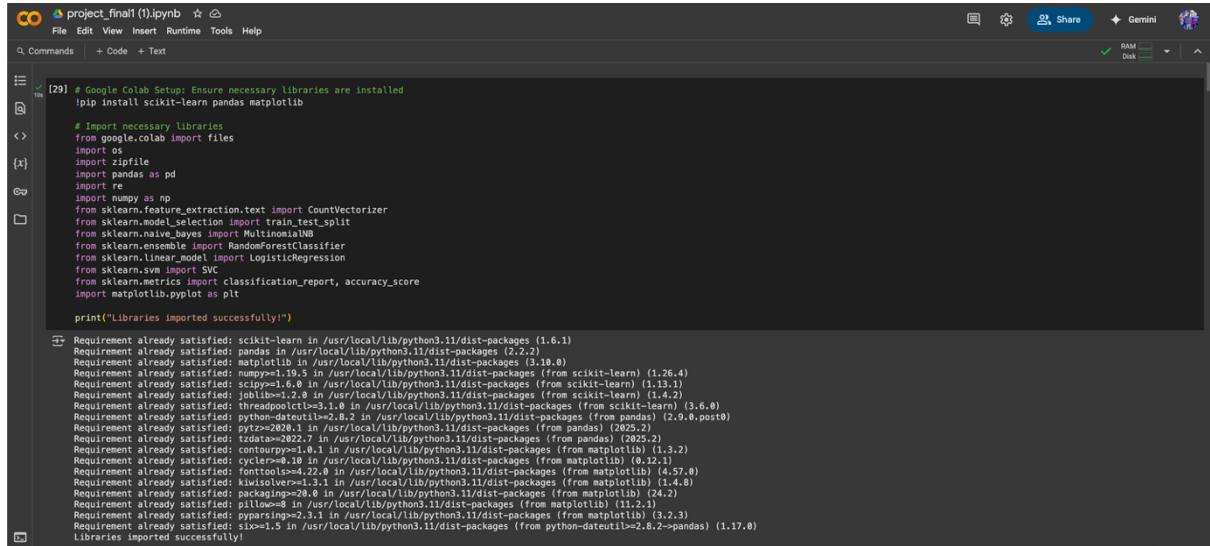
1. Platform/System Setup Procedures

To ensure smooth execution of the notebook, the development environment is set up using Google Colab follow below mentioned setup steps:

i. Environment Setup in Google Colab

Installing the necessary dependencies for the different phases of the spam email filtering project is the first step. Important Python libraries are imported, including matplotlib (used

for data visualization), pandas (used for data manipulation and analysis), and scikit-learn (used for creating machine learning models). Preprocessing, cleaning, word embedding, visualization, UI integration, and prediction tasks all depend on these libraries. The configuration guarantees that the notebook is prepared for the next steps of data loading and model development.



```
[29] # Google Colab Setup: Ensure necessary libraries are installed
!pip install scikit-learn pandas matplotlib

# Import necessary libraries
from google.colab import files
import os
import zipfile
import pandas as pd
import numpy as np
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import MultinomialNB
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
from sklearn.metrics import classification_report
import matplotlib.pyplot as plt

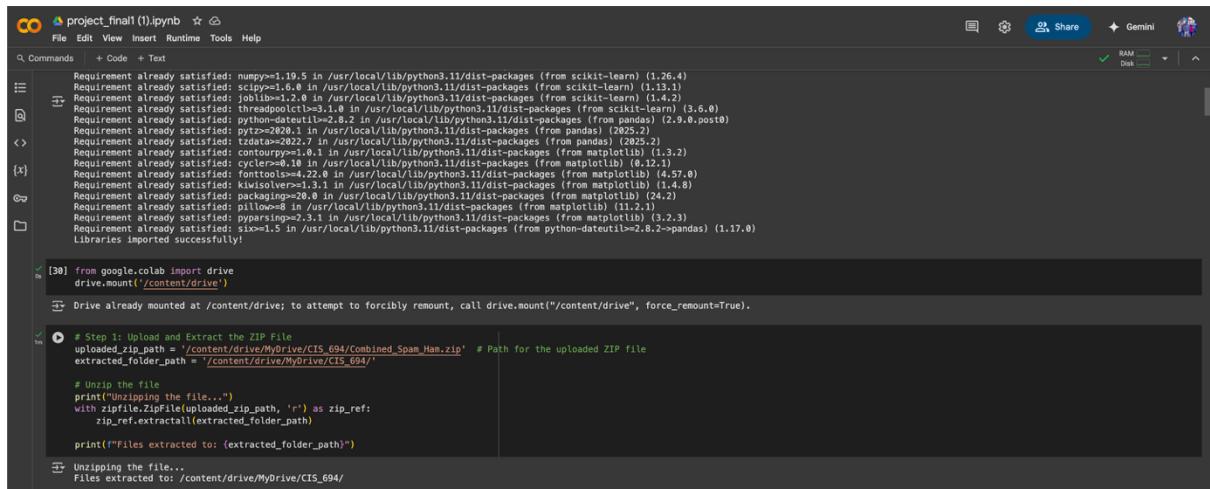
print("Libraries imported successfully!")

Requirement already satisfied: scikit-learn >=0.20.0 in /usr/local/lib/python3.11/dist-packages (from -r requirements.txt (line 1))
Requirement already satisfied: pandas >=1.0.0 in /usr/local/lib/python3.11/dist-packages (from -r requirements.txt (line 1))
Requirement already satisfied: matplotlib >=1.2.0 in /usr/local/lib/python3.11/dist-packages (from -r requirements.txt (line 1))
Requirement already satisfied: numpy>=1.19.5 in /usr/local/lib/python3.11/dist-packages (from scikit-learn (1.26.4))
Requirement already satisfied: scipy>=1.4.3 in /usr/local/lib/python3.11/dist-packages (from scikit-learn (1.26.4))
Requirement already satisfied: joblib>=0.12.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn (1.26.4))
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (3.6.0)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pandas) (2.9.0.post0)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas) (2025.2)
Requirement already satisfied: six>=1.5.2 in /usr/local/lib/python3.11/dist-packages (from python-dateutil<2.8.2->pandas) (1.17.0)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (24.2)
Requirement already satisfied: pillow>=8 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (11.2.1)
Requirement already satisfied: pyarrow>=2.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (3.2.3)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from python-dateutil>=2.8.2->pandas) (1.17.0)
Libraries imported successfully!
```

Figure: Screenshot of importing libraries

ii. File Preparation and Verification

First, the dataset file Combined_Spam_Ham.zip is uploaded to Google Colab's working directory. The contents are extracted into a folder called Combined_Spam_Ham when it has been uploaded and unzipped. Because the emails are encoded and not saved in the common.txt format, it is crucial to confirm that the files in the folder have been extracted correctly. In order to verify successful extraction, a sample file is loaded into the environment as part of the verification procedure. The extracted folder is shown in Colab's left panel, and the output window verifies that a raw email file was read, signifying that the setup and file preparation processes are finished.



```
[30] from google.colab import drive
drive.mount('/content/drive')

# Step 1: Upload and Extract the ZIP File
uploaded_zip_path = '/content/drive/MyDrive/CIS_694/Combined_Spam_Ham.zip' # Path for the uploaded ZIP file
extracted_folder_path = '/content/drive/MyDrive/CIS_694/'

# Unzip the file
print("Unzipping the file...")
with zipfile.ZipFile(uploaded_zip_path, 'r') as zip_ref:
    zip_ref.extractall(extracted_folder_path)

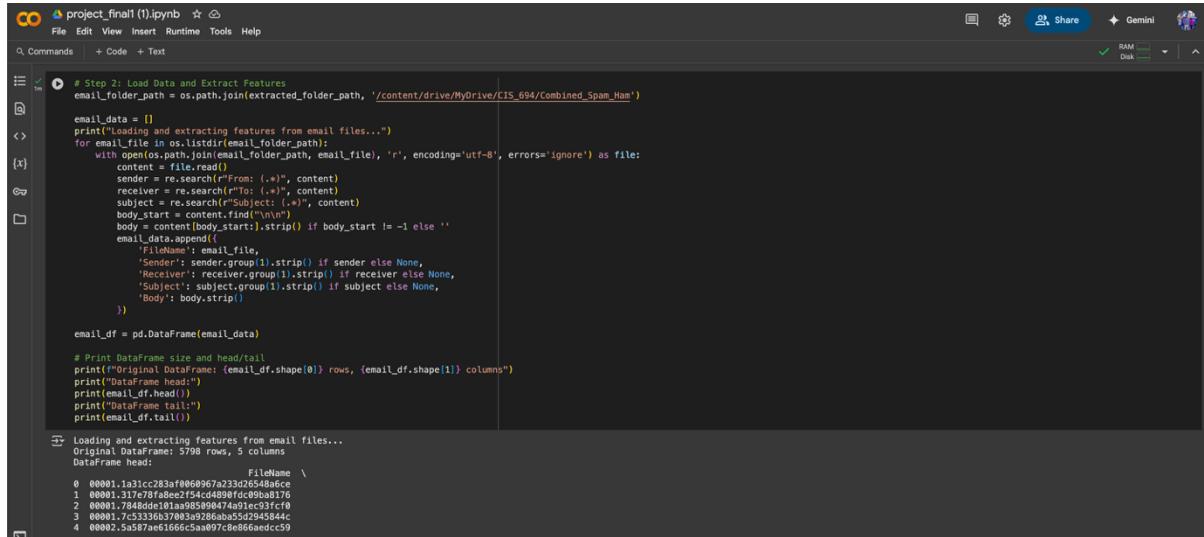
print(f"Files extracted to: {extracted_folder_path}")

# Unzipping the file...
Files extracted to: /content/drive/MyDrive/CIS_694/
```

Figure: Screenshot of extraction of Raw files

2. Load Data and extract features

The raw email files in the Combined_Spam_Ham folder are put into a single DataFrame for analysis following the successful extraction of the dataset. In order to streamline the dataset for spam identification, a number of superfluous elements, including HTML tags, CSS content, and message IDs, were eliminated, as was illustrated in the class presentation. Extracting the most pertinent information from each email—the file name, sender, recipient, email subject, and body—was the main goal. The following phases of preprocessing and model creation are built upon these features that have been cleaned and organized.



```
# Step 2: Load Data and Extract Features
email_folder_path = os.path.join(extracted_folder_path, '/content/drive/MyDrive/CIS_694/Combined_Spam_Ham')

email_data = []
print("Loading and extracting features from email files...")
for email_file in os.listdir(email_folder_path):
    with open(os.path.join(email_folder_path, email_file), 'r', encoding='utf-8', errors='ignore') as file:
        content = file.read()
        sender = re.search("From: (.*)", content)
        receiver = re.search("To: (.*)", content)
        subject = re.search("Subject: (.*)", content)
        body_start = content.find("\n\n")
        body = content[body_start:1].strip() if body_start != -1 else ''
        email_data.append({
            'FileName': email_file,
            'Sender': sender.group(1).strip() if sender else None,
            'Receiver': receiver.group(1).strip() if receiver else None,
            'Subject': subject.group(1).strip() if subject else None,
            'Body': body.strip()
        })
email_df = pd.DataFrame(email_data)

# Print DataFrame size and head/tail
print("Original DataFrame: {} rows, {} columns".format(email_df.shape[0], email_df.shape[1]))
print("DataFrame head:")
print(email_df.head())
print("DataFrame tail:")
print(email_df.tail())

# Loading and extracting features from email files...
Original DataFrame: 5798 rows, 5 columns
DataFrame head:
   File Name \
0  00001.la1lc263r70659567a233d055486cc
1  00001.317c7df8eex2754d4899fd409a8176
2  00001.7848dfe101a99589847491ec03fcf76
3  00001.7c53336b370839286aba55d2945844c
4  00002.5a587ae8168c5aa097cbe806aedcc59
```

Figure: Load and extract features

3. Putting all the data into csv format for next process of data cleaning

The data was saved to a CSV file called email_features.csv for later usage once the important features were extracted and organized into a DataFrame. To ensure portability and accessibility, this CSV file is part of the project's zip folder. In order to verify that the structure is consistent and prepared for additional preprocessing and analysis, a screenshot was also taken to show the data types of each feature in the DataFrame.

```
project_final1(1).ipynb  ⋆ ⓘ
File Edit View Insert Runtime Tools Help
Commands + Code + Text
5797 None
5793 Subject \
    Re: [ILUG] Linux Install \
5794 [Spambayes] Re: New Application of Spambayes...
5795 Re: [ILUG] Linux Install \
5796 None \
5797 None
{x}
5793 Gianni Ponzi wrote: I have a prob when tryi...
5794 Heile Pickett <heale@ngole.org> writes:\n> ...
5795 Hi,\n\nI think you need to give us a little mo...
5796 ...
5797

[35] # Save the DataFrame to a CSV file
output_csv_path = os.path.join(extracted_folder_path, '/content/drive/MyDrive/CIS_694/email_features.csv')
email_df.to_csv(output_csv_path, index=False, encoding='utf-8')

# Confirm the file has been saved
print(f"CSV file saved at: {output_csv_path}")

CSV file saved at: /content/drive/MyDrive/CIS_694/email_features.csv

[36] # Download the CSV file
files.download(output_csv_path)

[37] print("Data types of features:")
print(email_df.dtypes)

Data types of features:
FileName    object
Sender     object
Receiver   object
Subject    object
Body       object
dtype: object
```

Figure: Converting data frame to CSV file

4. Removing all the unnecessary tags and creating cleaned body

A number of cleaning procedures were carried out in order to get the email body content ready for analysis. These included eliminating superfluous spaces, digits, URLs, HTML tags, and non-alphanumeric characters from the text. Additionally, the email body was cleared of superfluous or extraneous terms by employing the STOPWORDS function to remove common English stop words. Following cleaning, the refined text data was stored as preprocessed_email.csv, a new CSV file that is also included in the project's zip folder for future processing and reference.

```
project_final1(1).ipynb  ⋆ ⓘ
File Edit View Insert Runtime Tools Help
Commands + Code + Text
import re
from nltk.corpus import stopwords
import nltk

# Download stopwords from NLTK (if not already done)
nltk.download('stopwords')

def clean_text(text):
    # Step 1: Remove HTML tags
    text = re.sub(r'<[^>]+>', '', text) # Removes HTML tags like <html>, <body>, etc.
    # Step 2: Remove non-alphanumeric characters (except spaces)
    text = re.sub(r'[^a-zA-Z\s]', '', text)
    # Step 3: Remove extra spaces
    text = re.sub(r'\s+', ' ', text)
    # Step 4: Remove URLs
    text = re.sub(r'^https?:\/\/.*[\r\n]*', '', text)
    # Step 5: Remove numbers
    text = re.sub(r'\d+', '', text)
    # Step 6: Remove programming keywords or system-specific tokens
    text = re.sub(r'\bpython\b|csharp|java|c|c++|c\+\+|c\#\#\b', '', text, flags=re.IGNORECASE)
    # Step 7: Convert text to lowercase
    text = text.lower()
    # Step 8: Remove stopwords
    stop_words = set(stopwords.words('english'))
    words_set = set(text.split())
    filtered_words = [word for word in words_set if word not in stop_words]
    # Step 10: Join the cleaned words back into a single string
    cleaned_text = ' '.join(filtered_words)
    return cleaned_text.strip()

# Apply the updated clean_text function to the 'Body' column
email_df['Cleaned_Body'] = email_df['Body'].apply(clean_text)

# Create a new DataFrame with preprocessed columns
preprocessed_email_df = email_df[['FileName', 'Sender', 'Receiver', 'Subject', 'Cleaned_Body']].copy()

# Print preprocessed DataFrame size and head/tail for verification
print(f"Preprocessed DataFrame: ({preprocessed_email_df.shape[0]} rows, {preprocessed_email_df.shape[1]} columns")
print(f"Preprocessed DataFrame head:")
print(preprocessed_email_df.head())
print(f"Preprocessed DataFrame tail:")
print(preprocessed_email_df.tail())
```

Figure: Cleaned body with features

5. Bootstrapping to create a larger dataset

The dataset's initial size of 6,000 emails was insufficient to create a reliable spam detection algorithm. The bootstrapping technique, which uses replacement sampling and random sampling from the available data to artificially expand the dataset size, was used to overcome this constraint. Consequently, the dataset was increased to 15,000 rows, which gave machine learning algorithms a bigger and more balanced input, improving model training efficacy and generalization.

The screenshot shows a Jupyter Notebook interface with the following code:

```
[39] preprocess_email_df.to_csv('/content/drive/MyDrive/CIS_694/Preprocessed_emails.csv', index=False)
print("Preprocessed DataFrame saved as 'Preprocessed_emails.csv'")
files.download('/content/drive/MyDrive/CIS_694/Preprocessed_emails.csv')

[40] # Step 4: Bootstrapping to Create a Larger Dataset
desired_rows = 15000
bootstrapped_df = preprocess_email_df.sample(n=desired_rows, replace=True, random_state=42).reset_index(drop=True)

print(f"Bootstrapped DataFrame: {bootstrapped_df.shape[0]} rows, {bootstrapped_df.shape[1]} columns")
print(bootstrapped_df.head())
print(bootstrapped_df.tail())
print(bootstrapped_df.dtypes)
```

Output:

```
Preprocessed DataFrame saved as 'Preprocessed_emails.csv'

Bootstrapped DataFrame: 15000 rows, 5 columns
Bootstrapped DataFrame:
   FileName \
0  00216.53a4d271ae7b0752fe5f21c2d5789ff \
1  02095.3d93919db01a74951fa2302b39c23a \
2  01931.95b9c125f167b05a7fe51796fc69a784c \
3  01896.516f25e3f181855da0e02a2d896ddff \
4  01093.5089c93d291361f2d0ee9c30a670

   Sender \
0  Wesley Burlington <esleywyle@gmail.com>
1  fark <rssfeeds@spamassassin.taint.org>
2  "hyatt@mozilla" <rssfeeds@spamassassin.taint.org>
3  boingboing <rssfeeds@spamassassin.taint.org>
4  bbridge@hotmail.com

   Receiver \
0  yyyy@localhost.netnoteinc.com \
1  yyyy@localhost.spamassassin.taint.org \
2  yyyy@localhost.spamassassin.taint.org \
3  yyyy@localhost.spamassassin.taint.org \
4  yyyy@localhost.netnoteinc.com

   Subject \
0  Re: [ILUG] Strange ssh problem
```

Figure: Bootstrapping to create larger dataset

The screenshot shows a Jupyter Notebook interface with the following code:

```
print("Bootstrapped DataFrame tail:")
bootstrapped_df.tail()

[41] print(f"Bootstrapped DataFrame: {bootstrapped_df.shape[0]} rows, {bootstrapped_df.shape[1]} columns")
```

Output:

```
Bootstrapped DataFrame tail:
Bootstrapped DataFrame:
   FileName \
14995 01153.e0954ec16f914973d54090f5958ed9b02 \
14996 000193.30ff42065811bd9ad262032953001415 \
14997 01385.508a461a95c42be52a29c92c2ca912 \
14998 00746.acf988d97eccf01a3d5f2d2d48882e \
14999 02098.96a584164de1919fe146c1ef4012182e

   Sender \
14995 <latest7687c2@yahoo.com> \
14996 Gary Lawrence Murphy <gary@canada.com> \
14997 Frank Burkhardt <burk@dns.mpd.de> \
14998 Tom <tomhore@silck.net> \
14999 newscientist <rssfeeds@spamassassin.taint.org>

   Receiver \
14995 <latest7687c2@yahoo.com> \
14996 yyyy@localhost.netnoteinc.com \
14997 yyyy@localhost.netnoteinc.com \
14998 yyyy@localhost.spamassassin.taint.org \
14999 yyyy@localhost.spamassassin.taint.org

   Subject \
14995 Hello ! \
14996 Re: My brain hurts \
14997 spamsassassin mailbox delivery problems \
14998 Here Come Da Feds \
14999 Space shuttle ready for return

   Cleaned Body \
14995 .nextpart_.b_dcc.ca content type text/plain charset \
14996 luck mon sends time never anything else old na... \
14997 hi testing spamsassassin found little problem k... \
14998 damn federalists another bumrush oregons stat... \
14999 url http://www.newsiscfree.com click date supplie...

[41] print(f"Bootstrapped DataFrame: {bootstrapped_df.shape[0]} rows, {bootstrapped_df.shape[1]} columns")

Bootstrapped DataFrame: 15000 rows, 5 columns
```

Figure: 15000 rows of data

6. Adding class label (spam/ham) based on spam keywords

Each email was given a binary class label, with 0 denoting ham (non-spam) and 1 denoting spam, in order to prepare the dataset for training a spam detection algorithm. In addition to a comprehensive manual assessment of spam and ham emails, search engine research was done to find frequently used keywords and indicators linked to spam content. These keywords formed the basis for creating a more precise and trustworthy model and assisted in informing the categorization process. The recognized spam-related phrases that help differentiate between spam and ham emails are displayed in a screenshot.

Figure: Adding class label on spam indicators

Figure: CSV file with class label on spam indicators

7. Word2Vec Model for Spam Email Detection

By capturing the semantic and contextual meaning of words, the Word2Vec model was used in this experiment to improve spam email identification. Three main benefits set Word2Vec apart from conventional keyword-based or bag-of-words strategies.

- **Semantic Understanding:** By rewording typical spam terms, spam emails frequently evade detection. The semantic similarity between these variations is captured by Word2Vec, which recognizes, for instance, that "Congratulations! You're the winner!" and "Well done! You've won!" expresses the same idea.
- **Contextual Features:** Word2Vec creates word embeddings that comprehend the context in which concepts arise by taking into account the surrounding words. In sentences like "Limited time offer," for instance, the model recognizes that the word "offer" is spam-indicative when combined with specific nouns.
- **Dimensionality Reduction:** By transforming words into low-dimensional dense vectors, Word2Vec preserves substantial linguistic patterns while drastically lowering memory consumption and processing costs as compared to sparse representations like TF-IDF.

```
File Edit View Insert Runtime Tools Help
Q Commands + Code + Text
14979 url http://www.newsisfree.com click date supplie...
[45] bootstrapped_df.to_csv('/content/drive/MyDrive/CIS_694/Class_label_emails.csv', index=False)
print("Bootstrapped DataFrame with class labels saved as 'Class_label_emails.csv'")
files.download('/content/drive/MyDrive/CIS_694/Class_label_emails.csv')
[46] #!pip install --upgrade --force-reinstall numpy==1.26.4 scipy==1.13.1 gensim
[47] !pip install gensim
from gensim.models import Word2Vec
from gensim.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.naive_bayes import MultinomialNB
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score, classification_report
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

# Step 6: Prepare Data for Model Training using Word2Vec
tokenized_texts = bootstrapped_df['Combined_Text'].apply(lambda x: x.split())

# Train a Word2Vec model on the tokenized texts
word2vec_model = Word2Vec(sentences=tokenized_texts, vector_size=100, window=5, min_count=1, workers=4)

# Function to vectorize text using Word2Vec
def vectorize_texts(texts, model):
    vectors = []
    for text in texts:
        tokens = text.split()
        token_vectors = [model.wv[word] for word in tokens if word in model.wv]
        if token_vectors:
            vectors.append(np.mean(token_vectors, axis=0))
        else:
            vectors.append(np.zeros(model.vector_size))
    return np.array(vectors)

# Vectorize the data
X_bootstrapped = vectorize_texts(bootstrapped_df['Combined_Text'], word2vec_model)
y_bootstrapped = bootstrapped_df['Spam']

Requirement already satisfied: gensim in /usr/local/lib/python3.11/dist-packages (4.5.3)
Requirement already satisfied: numpy<2.0,!=1.19.* in /usr/local/lib/python3.11/dist-packages (from gensim) (1.26.4)
Requirement already satisfied: scipy<1.14.0,!=1.7.0 in /usr/local/lib/python3.11/dist-packages (from gensim) (1.13.1)
Requirement already satisfied: smart-open>=1.8.1 in /usr/local/lib/python3.11/dist-packages (from gensim) (7.1.0)
Requirement already satisfied: wrapt in /usr/local/lib/python3.11/dist-packages (from smart-open>=1.8.1->gensim) (1.17.2)
```

Figure: Word2Vec model for word embedding to predict spam

8. Training and Testing

To build and evaluate the spam detection model effectively, the dataset was split into two subsets: 70% for training and 30% for testing. The training set is used to teach the model underlying patterns and associations in the labelled data, while the testing set is used to measure how well the model performs on previously unseen emails. This split maintains the overall distribution of the original dataset to ensure fairness and accuracy in evaluation. Additionally, a random state parameter was applied during splitting to ensure reproducibility—without it, each run would produce a different data split, potentially affecting the consistency of results.

The screenshot shows a Jupyter Notebook interface with the following code cells:

```
File Edit View Insert Runtime Tools Help
Q Commands + Code + Text
14979 url http://www.newsisfree.com click date supplie...
[45] <latestd7607c02@yahoo.com> Hello ! ... 1
14979 [REDACTED] Murphy [REDACTED]@chicago.com> Re: My... 1
14978 Frank Burkhardt [REDACTED]@blacknet.de> spammassassin... 0
14978 Tom <tomwhore@black.net> Here Come Da Feds dam... 0
14979 newscientist <rssfeeds@spammassassin.taint.org>... 1
[46] bootstrapped_df.to_csv('/content/drive/MyDrive/CIS_694/Class_label_emails.csv', index=False)
print("Bootstrapped DataFrame with class labels saved as 'Class_label_emails.csv'")
files.download('/content/drive/MyDrive/CIS_694/Class_label_emails.csv')
[47] #!pip install --upgrade --force-reinstall numpy==1.26.4 scipy==1.13.1 gensim
[48] # Step 6: Prepare Data for Model Training using Word2Vec
from gensim.models import Word2Vec
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.naive_bayes import MultinomialNB
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score, classification_report
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
# Step 6: Prepare Data for Model Training using Word2Vec
tokenized_texts = bootstrapped_df['Combined_Text'].apply(lambda x: x.split())
# Train a Word2Vec model on the tokenized texts
word2vec_model = Word2Vec(tokenized_texts, vector_size=100, window=5, min_count=1, workers=4)
# Function to vectorize text using Word2Vec
def vectorize_texts(texts, model):
    vectors = []
    for text in texts:
        tokens = text.split()
        token_vectors = [model.wv[word] for word in tokens if word in model.wv]
        if token_vectors:
            vectors.append(np.mean(token_vectors, axis=0))
        else:
            vectors.append(np.zeros(model.vector_size))
    return np.array(vectors)
# Vectorize the data
X_bootstrapped = vectorize_texts(bootstrapped_df['Combined_Text'], word2vec_model)
y_bootstrapped = bootstrapped_df['Spam']
Requirement already satisfied: gensim in /usr/local/lib/python3.11/dist-packages (4.3.3)
Requirement already satisfied: numpy<2.0,>=1.18.5 in /usr/local/lib/python3.11/dist-packages (from gensim) (1.26.4)
Requirement already satisfied: scipy<1.14.0,>=1.7.0 in /usr/local/lib/python3.11/dist-packages (from gensim) (1.13.1)
Requirement already satisfied: smart-open<1.8.1 in /usr/local/lib/python3.11/dist-packages (from gensim) (7.1.0)
Requirement already satisfied: wrapt in /usr/local/lib/python3.11/dist-packages (from smart-open<1.8.1>gensim) (1.17.2)
[48] # Step 7: Split Data for Training and Testing
X_train, X_test, y_train, y_test = train_test_split(
    X_bootstrapped, y_bootstrapped, test_size=0.3, random_state=42)
```

Figure: Split data for training and testing

Result and Discussion

1. Model Training and Evaluation

Three machine learning models—Support Vector Machine (SVM), Logistic Regression, and Random Forest—were used to classify spam emails. These models were chosen because they work well with text data and are compatible with Word2Vec and other word embeddings. To identify the patterns that differentiate spam from non-spam emails, each model was trained using the processed dataset. Then, criteria including recall, accuracy, and precision were used to assess their performance. A screenshot demonstrating the use of all three models and demonstrating how they were used to successfully forecast email categories is given.

```
# Step 8: Train and Evaluate Multiple Models
models = {
    "Random Forest": RandomForestClassifier(random_state=42, n_estimators=100),
    "Logistic Regression": LogisticRegression(random_state=42, max_iter=1000),
    "SVM": SVC(random_state=42)
}

results = []

for model_name, model in models.items():
    print(f"\nTraining {model_name}...")
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    acc = accuracy_score(y_test, y_pred)
    report = classification_report(y_test, y_pred, output_dict=True)
    results.append({
        'Model': model_name,
        'Accuracy': acc,
        'Precision': report['weighted avg']['precision'],
        'Recall': report['weighted avg']['recall'],
        'F1-Score': report['weighted avg']['f1-score']
    })

# Convert results to DataFrame
results_df = pd.DataFrame(results)

# Print and visualize the results
print("\nModel Performance Comparison:")
print(results_df)

plt.figure(figsize=(10, 6))
for metric in ['Accuracy', 'Precision', 'Recall', 'F1-Score']:
    plt.plot(results_df['Model'], results_df[metric], marker='o', label=metric)
plt.title("Model Performance Metrics")
plt.xlabel("Model")
plt.ylabel("Score")
plt.legend()
plt.grid()
plt.show()
```

Training Random Forest...
Training Logistic Regression...
Training SVM...
Model Performance Comparison:

Model	Accuracy	Precision	Recall	F1-Score
Random Forest	0.969968	0.969826	0.969968	0.969857

Figure: Evaluation major for spam email detection

2. Model Performance Visualization

Visualizations of the accuracy, precision, recall, and F1 score of the three trained models—SVM, Logistic Regression, and Random Forest—were made in order to compare their efficacy. The capacity of each model to accurately distinguish between spam and non-spam emails was thoroughly assessed by these performance criteria. The best-performing model was chosen for deployment with the aid of the visual comparison. The final spam detection model was selected based on its balanced metric scores and overall accuracy. A clean and succinct screenshot shows how well each of the three models performs in comparison.

Based on the outcomes displayed in the screenshot:

- With roughly 96.96% accuracy, precision, recall, and F1-score, Random Forest performed best across all measures, demonstrating its consistent and accurate email classification.
- With an accuracy of about 82.47%, logistic regression performed moderately well and was less dependable than random forest.
- With an accuracy of 84.69%, SVM outperformed Logistic Regression by a small margin, but Random Forest was still far superior.

Each model's performance across the four metrics is shown in detail in the line graph. The models are shown on the x-axis, and the metric score (which ranges from 0.80 to 0.97) is shown on the y-axis. Every coloured line represents a distinct metric:

- Blue for Precision
- Orange for Accuracy
- Green to Help You Remember
- F1-Score in red

The peak at the leftmost position of the graph indicates that Random Forest routinely receives the highest scores across all metrics. SVM performs marginally better than Random Forest but is still not competitive, while Logistic Regression has the lowest performance. The best model for this spam detection task is Random Forest, as can be easily seen from this visualization.

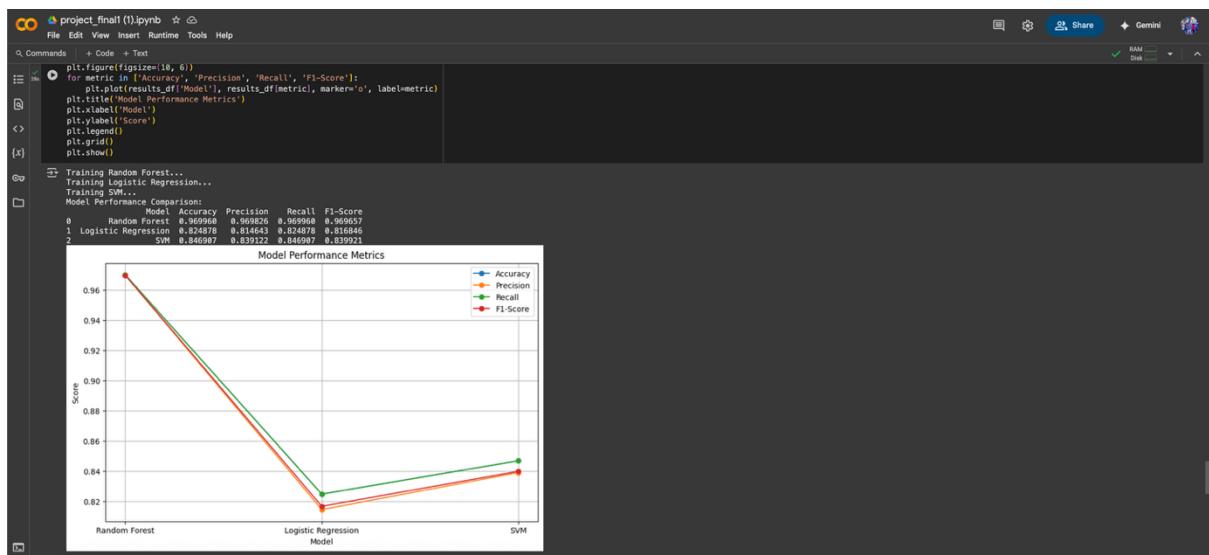


Figure: Visualization of all three models with accuracy, precision, F1 and recall

3. Future Predictions and Basic User Interface

A simple User Interface (UI) for determining if an email is spam or not was created to show practical application. Users can enter new email text into this interface and get immediate classification results using the trained model. To generate future predictions on unknown data, it incorporates the model's prediction logic. Designed for testing and display, the design is simple and practical. The underlying user interface programming that drives this interactive spam classification capability is displayed in two screenshots.

```

[50] # Step 9: Use the Best Model for Future Predictions
best_model_name = results_df.sort_values(by='Accuracy', ascending=False).iloc[0]['Model']
print("The best performing model is: " + best_model_name)
best_model = models[best_model_name]

# Example Prediction
#example_email = "Ok it's Thanksgiving - so you deserve a FEAST today. And like at Thanksgiving dinner, instead of one main course (like our usual weekly essays), today you're getting a full spread: 10 quick-hitter ideas that will make your class more interesting and fun."
#example_email = "Some Grading Errors Have Been Found in Short Exam 2. They all have been corrected now. Some of scores of the Short Exam2 have been updated (mostly lowered) now"
#example_vectorized = vectorizer.transform([example_email])
#prediction = best_model.predict(example_vectorized)
#print("The example email is classified as: 'Ham' if prediction[0] == 1 else 'Spam'")

# The best performing model is: Random Forest

from pywinauto import Text, Button, Output, Vbox
from PySide.QtCore import *
import numpy as np

# Function to vectorize a single input text using Word2Vec
def vectorize_single_text(text, model):
    tokens = text.split()
    token_vectors = [model.wv[word] for word in tokens if word in model.wv]
    if token_vectors:
        return np.mean(token_vectors, axis=0)
    else:
        return np.zeros(model.vector_size)

# Create widgets for user interaction
text_input = Text(description="Email:", placeholder="Enter email content here...")
predict_button = Button(description="Classify")
output = Output()

# Define prediction function
def predict_email_type(change):
    email_content = text_input.value
    if not email_content.strip():
        with output:
            output.clear_output()
            print("Please enter a valid email.")
    return

# Vectorize the input email
email_vectorized = vectorize_single_text(email_content, word2vec_model).reshape(1, -1)

```

Figure: UI code for spam classification

As part of the project, a basic user interface was created to show off the spam detection system's usefulness. Users can input a sample email into this interface's input box, and the trained model will use the content to determine whether the email is categorized as "spam" or "ham" (non-spam). Users can evaluate the efficacy of the spam filter on custom email text using this interactive feature, which demonstrates how the model can be used in real-time applications.

The Input Email: Your refund is on its way Hi Jayanth, Sorry things didn't work out. We've initiated a refund and wanted to provide a quick overview for you. If you have any follow-up questions, please visit our Help Center. Refund overview Refund Initiated card image AMEX ending in 1006\$19.91 Credited within 10 business days Thanks for shopping with us, Team Walmart Order number 2000130-07764860 Refund initiated Sold and shipped by Walmart item image Jif Creamy Peanut Butter, 16-Ounce Jar Qty: 1 \$3.12 item image Great Value 100% Whole Wheat Round Top Bread, 20 oz Qty: 1 \$1.97 item image Fresh Roma Tomato, Each \$1.65/LB 1.693LB \$1.65 item image Fresh Whole Red Onion, Each \$3.69/LB 2.078LB \$3.69 item image Cetaphil Daily Facial Cleanser for Sensitive, Combination to Oily Skin, 8 oz Qty: 1 \$8.78 Refund summary Subtotal \$19.21 Taxes \$0.70 Refund total \$19.91

Output: Ham

```

[50] # Step 9: Use the Best Model for Future Predictions
best_model_name = results_df.sort_values(by='Accuracy', ascending=False).iloc[0]['Model']
print("The best performing model is: " + best_model_name)
best_model = models[best_model_name]

# Example Prediction
#example_email = "Ok it's Thanksgiving - so you deserve a FEAST today. And like at Thanksgiving dinner, instead of one main course (like our usual weekly essays), today you're getting a full spread: 10 quick-hitter ideas that will make your class more interesting and fun."
#example_email = "Some Grading Errors Have Been Found in Short Exam 2. They all have been corrected now. Some of scores of the Short Exam2 have been updated (mostly lowered) now"
#example_vectorized = vectorizer.transform([example_email])
#prediction = best_model.predict(example_vectorized)
#print("The example email is classified as: 'Ham' if prediction[0] == 1 else 'Spam'")

# The best performing model is: Random Forest

from pywinauto import Text, Button, Output, Vbox
from PySide.QtCore import *
import numpy as np

# Function to vectorize a single input text using Word2Vec
def vectorize_single_text(text, model):
    tokens = text.split()
    token_vectors = [model.wv[word] for word in tokens if word in model.wv]
    if token_vectors:
        return np.mean(token_vectors, axis=0)
    else:
        return np.zeros(model.vector_size)

# Create widgets for user interaction
text_input = Text(description="Email:", placeholder="Enter email content here...")
predict_button = Button(description="Classify")
output = Output()

# Define prediction function
def predict_email_type(change):
    email_content = text_input.value
    if not email_content.strip():
        with output:
            output.clear_output()
            print("Please enter a valid email.")
    return

# Vectorize the input email
email_vectorized = vectorize_single_text(email_content, word2vec_model).reshape(1, -1)
# Predict using the trained model
prediction = best_model.predict(email_vectorized)
with output:
    output.clear_output()
    print("The email is classified as: 'Ham' if prediction[0] == 1 else 'Spam'")

# Attach the handler to the button
predict_button.on_click(predict_email_type)

# Display the UI
UI = Vbox([text_input, predict_button, output])
display(UI)

Email: Your refund is on its way Hi Jayanth
Classify
The email is classified as: Ham

```

Figure: Email Classification using best model

4. Problems Encountered and Resolutions

A number of difficulties were encountered and successfully overcome throughout the spam detection system's development:

Error 1: Problems with File Encoding

Issue: When reading some email files, encoding issues occurred.

Resolution: To get around troublesome characters, use errors='ignore' in the open() method.

Error 2: Difficulties with Data Parsing

Issue: Structured headers like From, To, and Subject were missing from a lot of email files.

Resolution: To extract data even in cases when headers were missing, regular expressions with fallback techniques were used.

Error 3: Problems with CSV Files Following Bootstrapping: Several records had null values added by bootstrapped data.

Resolution: Prior to additional processing, null values were managed programmatically using data cleaning techniques.

```
✓ # Handle missing values by removing rows with null values
✓ print("Handling missing values by removing rows with null values...")  
✓  
✓ # Check for any null values in required columns
✓ print("Before dropping null values:")
✓ print(bootstrapped_df.isnull().sum())
✓
```

Error 4: TF-IDF and Other Vectorizers' Limitations

Issue: Inconsistencies were produced by TF-IDF and related techniques because of improper word weighting, which resulted in incorrect classifications.

Resolution: To improve performance and interpretability, manually recognized spam indicators were used in place of TF-IDF.

Error 5: Class Imbalance in Dataset

Issue: There were disproportionately high and low amounts of spam and ham emails in the dataset.

Resolution: To preserve class distribution in both training and testing datasets, stratified sampling was used during the train_test_split.

Conclusion

The goal of this research was to create a reliable spam email detection system by combining machine learning models, semantic word embedding, and data preparation. After extracting significant features from unprocessed, encoded email files, the dataset was cleaned to eliminate extraneous components such as HTML tags, special characters, and stop words. As a result, just the most crucial elements—such as the sender, subject, and body—were kept for examination.

Stratified sampling was utilized to guarantee a balanced class distribution during model training and testing, and bootstrapping techniques were employed to increase the data from 6,000 to 15,000 samples in order to overcome dataset restrictions. In order to capture the semantic relationships between words and their context within the email content, Word2Vec was selected above more conventional techniques like TF-IDF.

Key performance parameters like accuracy, precision, recall, and F1 score were used to train and assess three machine learning models: Random Forest, SVM, and Logistic Regression. Finding the best model for deployment was made easier by the visualization of these findings. To enable users to test new email content for spam prediction, a basic user interface was also put into place.

Multiple technical problems, such as inconsistent email formats, file encoding mistakes, and data sparsity issues, were encountered and fixed during the project. By overcoming these obstacles, the study not only succeeded in detecting spam but also showed how well machine learning and natural language processing can be combined for practical uses.

References

- [1] C. Yang, R. Harkreader, and G. Gu, "Die free or live hard? empirical evaluation and new design for fighting evolving twitter spammers", In: Recent Advances in Intrusion Detection, Springer Berlin/Heidelberg, pp.318-337, 2011.
- [2] S. Kumar, and S. Arumugam, "A Probabilistic Neural Network Based Classification of Spam Mails Using Particle Swarm Optimization Feature Selection", Middle-East Journal of Scientific Research, Vol.23, No.5, pp.874-879, 2015.
- [3] N. P. DíAz, D. R. OrdáS, F. F. Riverola, and J. R. MéNdez, "SDAI: An integral evaluation methodology for content-based spam filtering models", Expert Systems with Applications, Vol.39, No.16, pp.12487-12500, 2012.
- [4] A. K. Sharma, S. K. Prajapat, and M. Aslam, "A Comparative Study Between Naive Bayes and Neural Network (MLP) Classifier for Spam Email Detection", In IJCA Proceedings on National Seminar on Recent Advances in Wireless Networks and Communications. Foundation of Computer Science (FCS), pp.12- 16, 2014.
- [5] W. Ma, D. Tran, and D. Sharma, "A novel spam email detection system based on negative selection", In: Proc. of Fourth International Conference on Computer Sciences and Convergence Information Technology, ICCIT'09, Seoul, Korea, pp.987-992, 2009.
- [6] T. S. Guzella, and W. M. Caminhas, "A review of machine learning approaches to spam filtering", Expert Systems with Applications, Vol.36, No.7, pp.10206-10222, 2009.
- [7] N. Kumar, S. Sonowal, and Nishant, "Email spam detection using machine learning algorithms," in Proceedings of the 2020 Second International Conference on Inventive Research in Computing Applications (ICIRCA), pp. 108–113, IEEE, Coimbatore, India, July 2020.
- [8] G. Jain, M. Sharma, and B. Agarwal, "Optimizing semantic lstm for spam detection," International Journal of Information Technology, vol. 11, no. 2, pp. 239–250, 2019.
- [9] F Masood, G. Ammad, A. Almogren et al., "Spammer detection and fake user identification on social networks," IEEE Access, vol. 7, pp. 68140–68152, 2019.
- [10] G. Chandrashekhar, and F. Sahin, "A survey on feature selection methods", Computers & Electrical Engineering, Vol.40, No.1, pp.16-28, 2014.
- [11] M. Mohamad, and A. Selamat, "An evaluation on the efficiency of hybrid feature selection in spam email classification", In: Proc. of 2015 International Conference on Computer, Communications, and Control Technology (I4CT), Kuching, Sarawak, Malaysia, pp.227- 231, 2015.
- [12] M. Sahami, S. Dumais, D. Heckerman, E. Horvitz: A Bayesian approach to filtering junk e-mail. In: AAAI-98 Workshop on Learning for Text Categorization. 1998. Volume 460.