



Email Spam Detection

Final Project Report

Ganesharow Bappoo (C0909659)
Abhey Singh Gill (C0905789)

Table of Contents

Abstract.....	1
Introduction.....	0
Background	0
Problem Statement	0
Objectives	0
Methodology Overview	0
Spam Detection Market Overview	0
Data Collection and Preprocessing	1
Preprocessing Steps.....	1
1.0 Import Libraries.....	2
2.0 Load Dataset	2
3.0 Data Cleaning	3
4.0 EDA	4
5.0 Feature Engineering.....	6
Methodology.....	10
Algorithms and Techniques Used:	10
6.0 Model Building and Evaluation	11
7.0 Model Interpretability	12
8.0 Model Tuning and Optimization	14
Results	15
Presentation of the experimental results	15
Pipeline for models.	15
Model Comparisons and Visualizations	17
9.0 Precision Metric Evaluation.....	18
Discussion.....	19
Interpretation of Results and Implication:	19
Strengths and Weaknesses:.....	19
Unexpected Outcomes:	19
Comparison with Prior Work:	19
Conclusion	19
References:	20

Abstract

This project aimed to develop and evaluate a machine learning model capable of identifying spam emails with high precision. Using Python libraries like pandas, numpy, and scikit-learn, along with the XGBoost algorithm enhanced by SHAP for model interpretability, the project processed and analyzed email data to build a predictive model. Key findings indicate that the optimized XGBoost model achieved a precision rate of 90%, effectively distinguishing spam from non-spam emails. This report details the methods used, the model's performance, and suggestions for future enhancements.

To Reader:

Small code snippets have been included along with the steps : Notebook can be found on our git [here](#). Graphs In this report are high quality, if required kindly zoom to view metrics Notebook has the steps labelled 1.0 to 9.0 on how the project was built.

Introduction

Background

Spam refers to unsolicited bulk communications that inundate recipients across various mediums such as email, social media, and telephone calls. The term itself, characterized by its repetitive and unwanted nature, originates from a 1970 Monty Python sketch. (watch video [here](#))

The timeline of spam begins long before the digital age, with the first instance recorded in 1867 involving unsolicited telegraph messages. This phenomenon evolved significantly with the advent of the internet. In 1978, the term "spam" was associated with mass emailing for the first time after a notable commercial email was sent, marking the beginning of what would become a major digital challenge.

Spam email detection continues to be a significant challenge in digital communication, critically impacting user productivity and online security. As email remains a fundamental communication tool both personally and professionally, the importance of effective spam filtering technologies cannot be overstated. The ability to accurately filter out spam is crucial in preventing phishing attacks, malware distribution, and unwanted commercial advertisements that plague email users.

Problem Statement

The dynamic nature of spam content, which continuously evolves to bypass detection mechanisms, complicates the effective classification of emails. Moreover, the cost of misclassifying legitimate emails as spam (false

positives) can be high, potentially leading to missed critical communications. Thus, developing a system that can adapt to changing spam tactics while maintaining low rates of false positives is a challenge.

Objectives

The primary goal of this project is to develop a robust machine learning model capable of classifying emails with high precision and recall. The model aims to minimize the misclassification of legitimate emails and ensure effective spam detection, enhancing the overall email user experience.

Methodology Overview

This project approaches the spam detection problem using supervised learning, with an XGBoost classifier serving as the core algorithm due to its proven effectiveness in handling diverse datasets and its robustness against overfitting. The model was fine-tuned using RandomizedSearchCV to optimize its hyperparameters, ensuring the best possible performance. Evaluation metrics such as precision, recall, and F1-score were employed to measure the model's effectiveness comprehensively.

Spam Detection Market Overview

The Email Spam Filter Market is highly competitive, with key players like TitanHQ, Hertz, Hornetsecurity, SolarWinds MSP, and Symantec leading the charge. These companies, along with SpamPhobia, Trend Micro, Firetrust, Comodo Group, SPAMfighter, MailChannels, MailCleaner, SpamHero, Mimecast, Spambrella, and GFI Software, are at the forefront of deploying advanced technologies to combat spam.

These technologies include artificial intelligence (AI), machine learning algorithms, and sophisticated filtering mechanisms that enhance accuracy and efficiency in detecting and blocking unwanted emails. mail Spam Filter market is expected to grow at a CAGR of over 7% in the next five years., driven by an increased demand for cybersecurity solutions and the ongoing evolution of spam tactics that require continual innovations in spam detection methods. These companies' efforts in developing cutting-edge solutions are pivotal in shaping the market's dynamics and addressing the growing security needs of digital communications.

This project contributes to this field by developing a model that leverages advanced machine learning techniques to offer a high-performance solution adaptable to both current and emerging spam trends.

Data Collection and Preprocessing

Data Sources: The dataset utilized in this project, named 'SpamAssassin.csv', was sourced from Kaggle, a popular platform for predictive modelling and analytics competitions. It comprises a comprehensive collection of email data specifically structured for spam detection tasks. Each entry in the dataset represents an email, with the body of the email labeled as either spam (1) or non-spam (0). This binary classification framework facilitates the application of supervised learning techniques for spam detection.

Dataset	Description:	The
'SpamAssassin.csv'	file contains several	

thousand entries, as indicated by a total of 6046 rows, with each email encompassing attributes that describe its content and classification. The data fields include an index column ('Unnamed: 0'), the 'Body' of the email containing the text, and a 'Label' indicating spam or non-spam status. The dataset reflects a typical unbalanced distribution often seen in real-world spam detection scenarios, where non-spam (ham) emails outnumber spam messages. These characteristic challenges the model to accurately identify the relatively rarer spam emails without overfitting to the more common non-spam examples.

Preprocessing Steps

Cleaning: The initial step involved removing the 'Unnamed: 0' column, which was merely an index with no analytical value. Additionally, any rows containing missing values, particularly in the 'Body' column critical for spam classification, were eliminated to ensure data integrity.

Transformation: The text data within the 'Body' of each email was preprocessed through tokenization, converting the continuous text into a list of terms. Subsequent cleaning steps included the removal of stopwords—common words that offer little value in the context of text classification—to reduce noise in the data.

Feature Engineering: To transform the textual data into a format suitable for machine learning, Term Frequency-Inverse Document Frequency (TF-IDF) vectorization was applied. This method quantifies the importance of a term relative to its frequency across all documents, effectively capturing the

significance of words specific to spam. Moreover, the length of each email was calculated and added as a meta-feature, providing additional contextual information that could correlate with spam characteristics.

Challenges: A primary challenge encountered was the high-dimensional space created by the TF-IDF vectorization, which can lead to computational inefficiency and model overfitting. To address this, dimensionality reduction techniques such as Truncated Singular Value Decomposition (SVD) were employed. Truncated SVD was instrumental in reducing the number of features to a manageable size while preserving the essential variance in the data, thereby enhancing model performance without significant loss of information.

1.0 Import Libraries

We begin by importing the libraries necessary for the entire process. pandas and numpy are for data handling. matplotlib and seaborn are for data visualization to help us see the distributions and patterns. nltk is for natural language processing, specifically for handling text data. The sklearn libraries are for building and evaluating our machine learning model. joblib is for saving our model for later use.

```
# Importing necessary libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import string
import nltk
from nltk import word_tokenize
from nltk.corpus import stopwords
from sklearn.model_selection import train_test_split, cross_val_score
```

```
e
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.pipeline import Pipeline
from sklearn.naive_bayes import MultinomialNB
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification_report, accuracy_score, confusion_matrix
from wordcloud import WordCloud
import joblib
```

2.0 Load Dataset

We will now load the dataset into a pandas DataFrame. Pandas is a powerful tool that allows for sophisticated data manipulation and analysis. By reading the CSV into a DataFrame, we can easily access and manipulate the data. After loading the dataset, we will take a peek at the first few rows using `df.head()` to understand what our data looks like.

```
#The file is named 'completeSpamAssassin.csv' and is in the same directory as the notebook
df = pd.read_csv('SpamAssassin.csv')
```

```
df.head()

   Unnamed: 0  Body  Label
0           0      0  \nSave up to 70% on Life Insurance.\nWhy Spend...
1           1      1  1) Fight The Risk of Cancer!\nhttp://www.adcli...
1           1      1  1) Fight The Risk of Cancer!\nhttp://www.adcli...
2           2      1  1) Fight The Risk of Cancer!\nhttp://www.adcli...
1           1      3  #####
3           3      3  #####...
```



```

1
4           4 I thought you might
like these:\n1) Slim Down ...
1

df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6046 entries, 0 to 6045
Data columns (total 3 columns):
#   Column          Non-Null Count  Dtype
---  -
0   Unnamed: 0      6046 non-null   int64
1   Body            6045 non-null   object
2   Label           6046 non-null   int64
dtypes: int64(2), object(1)
memory usage: 141.8+ KB

df.describe()

           Unnamed: 0           Label
count  6046.000000    6046.000000
mean    3022.500000     0.313596
std     1745.474195     0.463993
min         0.000000     0.000000
25%    1511.250000     0.000000
50%    3022.500000     0.000000
75%    4533.750000     1.000000
max     6045.000000     1.000000

```

Upon loading the dataset, we see that it comprises various columns that potentially include the content of emails along with their corresponding labels indicating whether an email is spam or not. To fully grasp the scope of the dataset, we would assess the number of rows, which represent individual email entries, and the number of columns, which denote the attributes or features of each email, such as the subject, body, and sender information.

To thoroughly understand the nature of the data, it's crucial to identify if there are any missing values that could affect our analysis or require preprocessing. Additionally, the balance between spam and non-spam (ham) emails will be pivotal to observe since it can significantly influence the performance of our machine learning model. The datatype of each column will also be noted, as text data will need to be processed differently from numerical data when we move towards feature engineering and model training.

3.0 Data Cleaning

```

# Dropping the 'Unnamed: 0' column
df.drop(columns=['Unnamed: 0'], in
place=True)

```

```

# Dropping rows with any missing v
alues
df.dropna(inplace=True)

```

```

df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 6045 entries, 0 to 6045
Data columns (total 2 columns):
#   Column          Non-Null Count  Dtype
---  -
0   Body            6045 non-null   objec
t
1   Label           6045 non-null   int64
dtypes: int64(1), object(1)
memory usage: 141.7+ KB

```

We now see that the column has been dropped, and there is no missing values

```
df.head()
```

```

Body  Label
0   \nSave up to 70% on Life Insura
nce.\nWhy Spend...      1
1  1) Fight The Risk of Cancer!\nh
ttp://www.adcli...      1

```



```

2 1) Fight The Risk of Cancer!\nh
ttp://www.adcli...      1
3 #####
#####...      1
4 I thought you might like these:
\n1) Slim Down ...      1

```

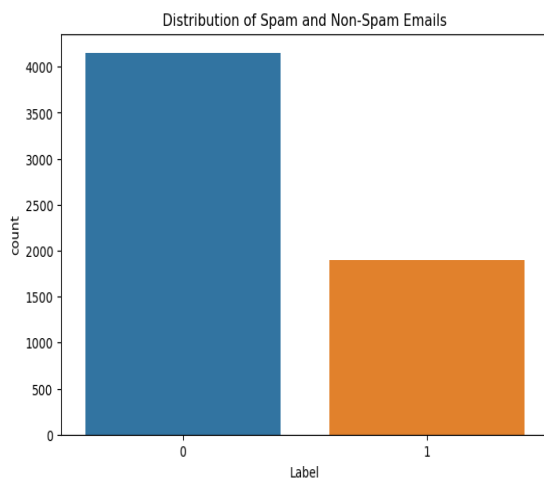
4.0 EDA

Distribution of Spam vs. Non-Spam emails
Plotting the distribution of Spam vs. Non-Spam emails

```

plt.figure(figsize=(8, 5))
sns.countplot(x='Label', data=df)
plt.title('Distribution of Spam and Non-Spam Emails')
plt.show()

```



This visual insight into the distribution will inform us if we need to balance the dataset before training our model to prevent bias towards the more common class.

The bar chart illustrates the distribution of spam and non-spam emails in the dataset. From the bar chart, we can observe a significant imbalance: non-spam emails (labeled as 0) outnumber the spam emails (labeled as 1). This imbalance is common in spam detection datasets and reflects real-world conditions where genuine emails usually exceed spam.

In predictive modeling, such class imbalance could bias the model towards predicting the majority class. To address this, techniques like resampling the minority class, using class weights, or evaluating the model with metrics that give a better sense of performance on imbalanced data (such as F1-score, precision, recall, and ROC AUC) may be required. It's important to consider these factors during the modeling stage to ensure the classifier does not simply learn to predict the majority class.

Email Length Distribution

Calculating the length of each email

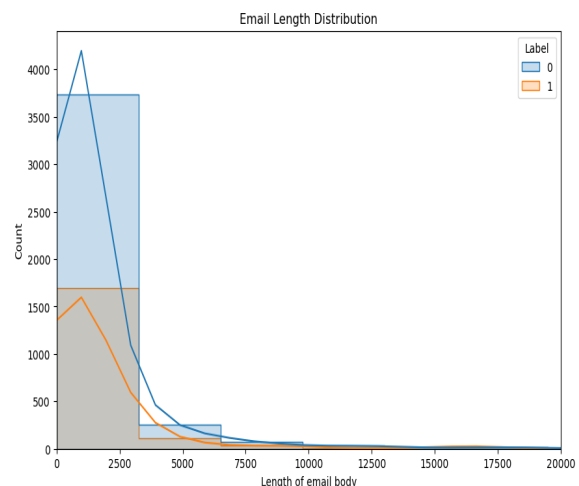
```
df['length'] = df['Body'].apply(len)
```

Plotting the distribution of email lengths by label

```

plt.figure(figsize=(10, 6))
sns.histplot(data=df, x='length', hue='Label', bins=60, kde=True, element='step')
plt.title('Email Length Distribution')
plt.xlim(0, 20000) # Adjust this limit based on the distribution you observe
plt.xlabel('Length of email body')
plt.ylabel('Count')
plt.show()

```



Notably, spam emails tend to include longer messages as well, which could be useful in distinguishing between the two classes when we design our spam detection model.

```
# Combining all non-spam emails in
to a single text
non_spam_emails = ' '.join(df[df['
Label'] == 0]['Body'])
```

```
stopwor
ds = set(stopwords.words('english'
)), min_font_size = 10).generate(n
on_spam_emails)
```

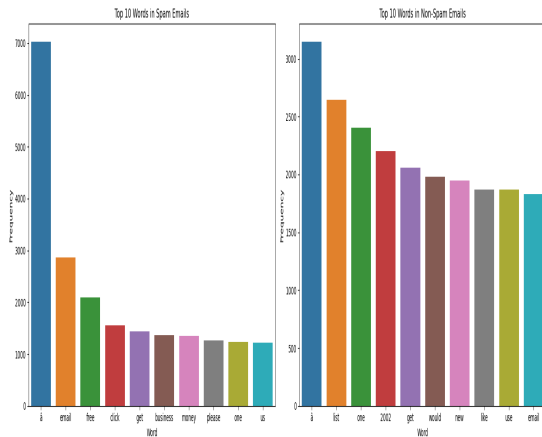


```
spam_emails = ' '.join(df[df['Label'] == 1]['Body'])
```

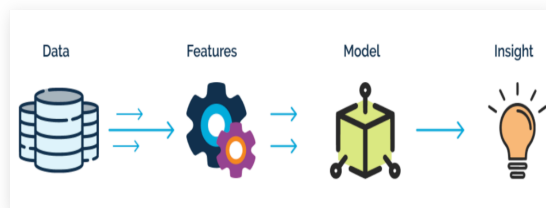
```
wordcloud_spam = WordCloud(width=800, height=400, background_color = 'red',
```

A word cloud of common email-related terms. The words are arranged in a dense, overlapping manner. The most prominent words, shown in larger fonts, include 'email', 'click', 'http', 'www', 'mail', 'people', 'work', 'get', 'money', 'free', 'service', 'order', 'receive', 'make', 'want', 'help', 'information', 'product', 'call', 'sale', 'website', 'take', 'thank', 'internet', 'start', 'spans', 'signings', 'mailing', 'list', 'would', 'need', 'name', 'give', 'sent', 'time', 'offer', 'report', 'state', 'north', 'central', 'south', 'east', 'west', 'middle', 'northwest', 'southwest', 'northeast', 'southeast', 'northcentral', 'southcentral', 'northwestcentral', 'southwestcentral', 'northcentralwest', 'southcentralwest', 'northwestcentralwest', 'southwestcentralwest', 'northcentralwestcentral', 'southcentralwestcentral', 'northwestcentralwestcentral', 'southwestcentralwestcentral', 'northcentralwestcentralwest', 'southcentralwestcentralwest', 'northwestcentralwestcentralwest', 'southwestcentralwestcentralwest'. Other words include 'program', 'please', 'click', 'see', 'find', 'font', 'family', 'remove', '100%', 'message', 'contract', 'first', 'size', 'new', 'line', 'balance', 'wish', 'received', 'easy', 'description', 'money', 'spans', 'signings', 'mailing', 'list', 'would', 'need', 'name', 'give', 'sent', 'time', 'offer', 'report', 'state', 'north', 'central', 'south', 'east', 'west', 'middle', 'northwest', 'southwest', 'northeast', 'southeast', 'northcentral', 'southcentral', 'northwestcentral', 'southwestcentral', 'northcentralwest', 'southcentralwest', 'northwestcentralwest', 'southwestcentralwest', 'northcentralwestcentral', 'southcentralwestcentral', 'northwestcentralwestcentral', 'southwestcentralwestcentral', 'northcentralwestcentralwest', 'southcentralwestcentralwest', 'northwestcentralwestcentralwest', 'southwestcentralwestcentralwest'.

The top words in both spam and non-spam emails after stop words have been removed. Interestingly, we see the word 'email' features prominently in both, which is expected given the dataset's nature. The presence of words like **'free'**, **'click'**, **'business'**, and **'money'** among the top spam words fits the common patterns associated with spam content, which often include such enticements or commercial terms.



5.0 Feature Engineering



Feature engineering for text data typically involves creating a numerical representation of the text that can be used by machine learning algorithms. For this dataset, we can create features based on the text content and some meta-features like the length of the emails. After that, we can use dimensionality reduction techniques like t-SNE or PCA to visualize the features.

Text Vectorization

This process entails converting the cleaned text data into numerical features using TF-IDF. This technique transforms the text into a sparse matrix of term frequencies weighted by their importance across all documents.

```
# Drop rows with NaN values in the 'Body' column
```

```
df = df.dropna(subset=['Body'])
```

```
# Preprocessing function to remove
```

```
punctuation and stopwords
```

```
def preprocess_text(text):
    # Convert to string in case there are any non-string entries
    text = str(text)
    # Tokenize the text
    tokens = word_tokenize(text.lower())
```

```
    # Remove stopwords
    stop_words_set = set(stopwords.words('english'))
    tokens = [word for word in tokens if word not in stop_words_set and word.isalpha()]
```

```
    return ' '.join(tokens)
```

```
# Apply preprocessing to the email bodies
```

```
df['Clean_Body'] = df['Body'].apply(preprocess_text)
```

```
# Initialize the TF-IDF vectorizer
tfidf_vectorizer = TfidfVectorizer(sublinear_tf=True, min_df=5, norm='l2', encoding='latin-1', ngram_range=(1, 2))
```

```
# Fit and transform the 'Clean_Body' column to create TF-IDF features
```

```
# The vectorizer returns a sparse matrix by default, so we convert it to an array
```

```
features = tfidf_vectorizer.fit_transform(df['Clean_Body']).toarray()
```

```
# create a DataFrame of the features with their corresponding terms
```

```
feature_names = tfidf_vectorizer.get_feature_names_out()
tfidf_df = pd.DataFrame(features, columns=feature_names)
```

```
tfidf_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```

RangeIndex: 6045 entries, 0 to 6044
Columns: 28938 entries, aa to çš çš
dtypes: float64(28938)
memory usage: 1.3 GB

```

```
tfidf_df.head()
```

```

      aa aa meetings aac aalib a
all aall credit aaron ab abac
ha \
0 0.0      0.0 0.0 0.0
0.0      0.0 0.0 0.0 0
.0
1 0.0      0.0 0.0 0.0
0.0      0.0 0.0 0.0 0
.0
2 0.0      0.0 0.0 0.0
0.0      0.0 0.0 0.0 0
.0
3 0.0      0.0 0.0 0.0
0.0      0.0 0.0 0.0 0
.0
4 0.0      0.0 0.0 0.0
0.0      0.0 0.0 0.0 0
.0

```

```

      abacha died ... zoomâ byte z
oomâ look zope zu zyban z
zzz \
0      0.0 ... 0.0
0.0 0.0 0.0 0.0 0.000000
1      0.0 ... 0.0
0.0 0.0 0.0 0.0 0.000000
2      0.0 ... 0.0
0.0 0.0 0.0 0.0 0.123698
3      0.0 ... 0.0
0.0 0.0 0.0 0.0 0.056189
4      0.0 ... 0.0
0.0 0.0 0.0 0.0 0.000000

```

```

      zzzz password zzzzteana çš
çš çš
0      0.000000      0.0 0.0
0.0
1      0.000000      0.0 0.0
0.0
2      0.000000      0.0 0.0
0.0

```

```

3      0.064979      0.0 0.0
0.0
4      0.000000      0.0 0.0
0.0

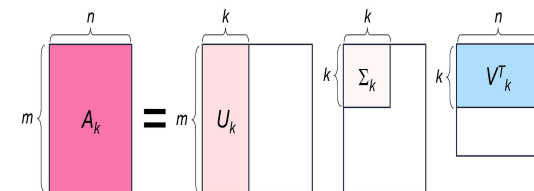
```

```
[5 rows x 28938 columns]
```

Output indicates that after vectorizing text data using TF-IDF, we have a DataFrame with 6,045 emails (entries) and a very large number of features (28,938), each representing a unique term in the TF-IDF matrix. - Every column corresponds to a term (word or bigram), and the value is the TF-IDF score of that term in each email.

However, the output also suggests that the DataFrame is quite large, consuming about 1.3 GB of memory. This size is likely due to the high dimensionality of the data, with many unique words and bigrams creating a sparse matrix with a lot of zeros.

Truncated SVD to reduce Dimensions.



```

# Initialize Truncated SVD with de
sired number of components
svd = TruncatedSVD(n_components=10
0) # Adjust the number of compone
nts

```

```

# Fit and transform the TF-IDF fea
tures
reduced_features = svd.fit_transfo
rm(tfidf_df)

```

```

# Now the shape should be signific
antly smaller, let's see the varia
nce explained by the 100 component

```

```

S
print(f"Total variance explained by 100 components: {svd.explained_variance_ratio_.sum()}")

```

Total variance explained by 100 components: 0.2911818357861679

Elbow Curve

```

import matplotlib.pyplot as plt
from sklearn.decomposition import TruncatedSVD

```

```

# TF-IDF matrix

```

```

svd = TruncatedSVD(n_components=100)
svd.fit(features)

```

```

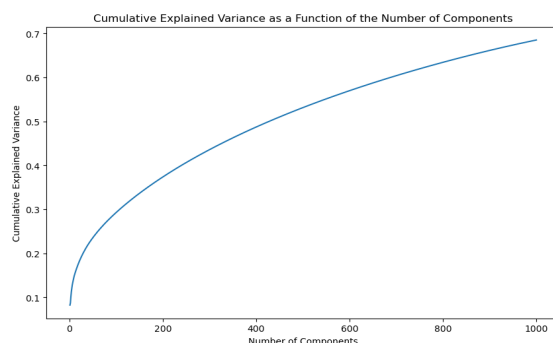
plt.figure(figsize=(10, 6))
plt.plot(range(1, 1001), np.cumsum(svd.explained_variance_ratio_))
plt.title('Cumulative Explained Variance as a Function of the Number of Components')
plt.xlabel('Number of Components')
plt.ylabel('Cumulative Explained Variance')

```

```

plt.show()

```



In the absence of a distinct elbow in the cumulative variance plot, selecting around 600 components could be sensible for this small-scale project. This estimate is guided by the graph's indication that 300 components explain about 50% of variance—doubling this might capture a substantial majority of the information. The choice strikes a

balance between capturing a significant portion of data variance and maintaining computational efficiency. Ultimately, this number should be refined through experimentation, considering computational resources, project requirements, and model performance, adjusting as necessary to optimize the balance between dimensionality and information retention.

Truncated SVD is used for dimensionality reduction of sparse data, like TF-IDF features from text. By reducing the dataset to a manageable size (e.g., 100-300 components), it retains significant variance while improving computational efficiency and potentially model performance. This streamlined dataset facilitates deeper analysis and more efficient training of machine learning models, making it crucial for handling high-dimensional data in natural language processing tasks.

Meta-Features

here we will add meta-features like the length of the email, which could also be predictive.

```

# Add email length as a feature
email_lengths = df['Clean_Body'].apply(len).values.reshape(-1, 1)

```

```

# Combine TF-IDF features with the email length

```

```

features = np.hstack((features, email_lengths))

```

```

features

```

```

array([[ 0.,  0.,  0., ...,
        0.,  0., 634.],
       [ 0.,  0.,  0., ...,
        0.,  0., 363.],
       [ 0.,  0.,  0., ...,
        0.,  0., 274.]])

```

```
...,
[ 0., 0., 0., ...,
0., 0., 4169.],
[ 0., 0., 0., ...,
0., 0., 5.],
[ 0., 0., 0., ...,
0., 0., 5.]])
```

By adding email length as a feature alongside the TF-IDF vectorized text, we enrich our dataset with meta-information that could be predictive of an email being spam or not. This process involves calculating the length of each cleaned email and appending it as an additional feature. Combining raw text features with such meta-features allows machine learning models to leverage both the content and structural characteristics of emails, potentially improving classification accuracy.

Feature Visualization with t-SNE

t-SNE is a tool to visualize high-dimensional data by reducing it to two or three dimensions

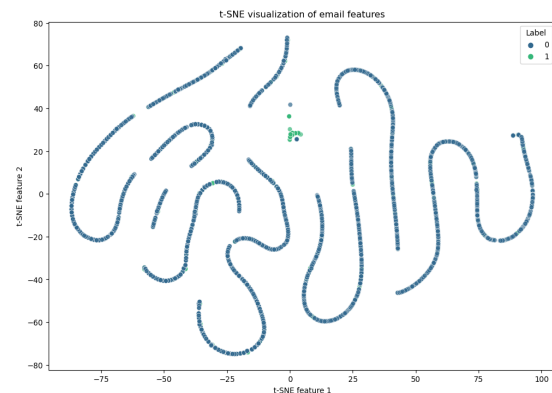
```
from sklearn.manifold import TSNE

# Use t-SNE to project the feature
s into 2 dimensions
tsne = TSNE(n_components=2, random
_state=42)
features_reduced = tsne.fit_transf
orm(features)

# Plot the result of t-SNE
plt.figure(figsize=(12, 8))
sns.scatterplot(x=features_reduced
[:, 0], y=features_reduced[:, 1],
hue=df['Label'], palette='viridis'
, alpha=0.7)
plt.title('t-SNE visualization of
email features')
plt.xlabel('t-SNE feature 1')
```

```
plt.ylabel('t-SNE feature 2')
plt.show()
```

The t-SNE plot provided offers a visual representation of how your email features are distributed when reduced to two dimensions. The technique has clustered the features into several distinct groups. The separation between spam (label 1) and non-spam (label 0) emails is not clearly defined, indicating a complex structure where emails cannot be linearly separated based on these two t-SNE components alone.



In conclusion, while t-SNE has illuminated some structure within the data, the overlapping nature suggests that classifying emails as spam or non-spam may require more sophisticated models or additional feature engineering to capture the nuances. The next steps would involve using these insights to guide the creation of a machine learning model that can navigate this complexity, potentially by incorporating additional features or exploring more advanced algorithms.

Feature Reduction with PCA

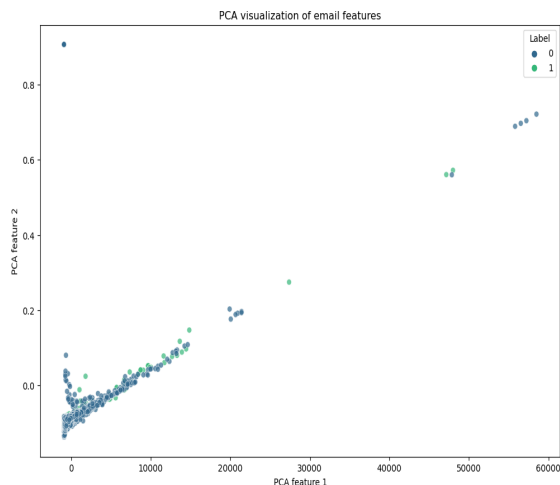
Principal Component Analysis (PCA) is another dimensionality reduction

technique. It's faster than t-SNE and also useful for feature selection.

```
from sklearn.decomposition import PCA

# Use PCA to reduce the features to 2 dimensions
pca = PCA(n_components=2)
features_reduced_pca = pca.fit_transform(features)

# Plot the result of PCA
plt.figure(figsize=(12, 8))
sns.scatterplot(x=features_reduced_pca[:, 0], y=features_reduced_pca[:, 1], hue=df['Label'], palette='viridis', alpha=0.7)
plt.title('PCA visualization of email features')
plt.xlabel('PCA feature 1')
plt.ylabel('PCA feature 2')
plt.show()
```



The PCA visualization suggests that the email features, when reduced to two principal components, show some degree of separation between the spam and non-spam categories. The plot reveals a principal component 1 (x-axis) that accounts for a large variance, with the data points spread out along this axis, indicating its importance in differentiating emails. The second

principal component (y-axis), however, shows much less variance.

However, there's no clear boundary between the two classes, with some overlap evident. This could imply that while PCA has captured some underlying structure of the data, further feature engineering or more complex models might be needed to improve the separability of the classes. The clustering of non-spam closer to the origin suggests lower values of the principal components, possibly correlating with less complex or less frequent terms in comparison to spam emails, which appear more spread out.

Finalizing feature set and ensuring data is prepared for the training phase.

Methodology

Algorithms and Techniques Used:

To address the challenge of spam detection, a pipeline integrating several machine learning models was implemented. The initial testing included Logistic Regression, Random Forest, XGBoost, and LightGBM. This diversity allowed for assessing different algorithmic strengths in handling binary classification tasks on text data. Each model was evaluated within the pipeline, ensuring consistent preprocessing, which included TF-IDF vectorization and dimensionality reduction, to standardize the input data format across all models.

The final phase involved retraining the XGBoost model with the optimized parameters derived from RandomizedSearchCV. This retraining aimed to refine the model based on the

insights gained from the initial evaluations and to fully leverage the refined parameter settings for enhanced model performance. The retrained model was then used to predict unseen data, providing a real-world assessment of its efficacy in spam detection, backed by the rigorous methodology applied throughout the project lifecycle.

6.0 Model Building and Evaluation

Model Training, Validation, and Evaluation:

The selected models were trained using an 80/20 train-test split. XGBoost, in particular, was fine-tuned and retrained extensively using RandomizedSearchCV to optimize its hyperparameters effectively. This optimization focused on parameters such as `max_depth` for controlling tree complexity, `learning_rate` to manage the speed of convergence, and `n_estimators` for the number of trees in the ensemble, which are pivotal for balancing bias and variance. The training process included cross-validation strategies to ensure the model's generalizability and robustness.

Simple Training and Evaluation using Reduced Dimension Features from Truncated SVD

For this example, I was aiming at developing and training a model in less than 10 mins for an initial test.

```
# Apply Truncated SVD to reduce
dimensionality
svd =
TruncatedSVD(n_components=200) #
We select 200 components based on
our
X_reduced =
```

```
svd.fit_transform(features)

# Split the reduced data into
training and test sets
X_train, X_test, y_train, y_test =
train_test_split(X_reduced,
df['Label'], test_size=0.2,
random_state=42)

# Initialize a simple model, like
RandomForest, for quick training
and evaluation
model =
RandomForestClassifier(n_estimator
s=100, random_state=42)

# Start timing the execution
start_time = time.time()

# Train the model
model.fit(X_train, y_train)

# Predict on the test set
y_pred = model.predict(X_test)

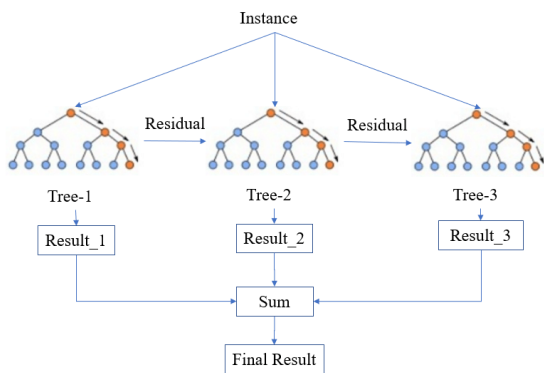
# Calculate the execution time
execution_time = time.time() -
start_time
print(f"Execution Time:
{execution_time:.2f} seconds")

# Evaluate the model
print(classification_report(y_test
, y_pred))
```

		precision	recall
f1-score	support		
	0	0.98	0.95
0.96	807		
	1	0.90	0.97
0.93	402		
	accuracy		
0.95	1209		
	macro avg	0.94	0.96
0.95	1209		
	weighted avg	0.96	0.95

0.95 1209

The **RandomForest** model has achieved commendable performance, with an overall accuracy of 95% and balanced precision and recall metrics, indicating strong predictive power. The precision for spam emails stands at **90%**, meaning the model is highly reliable when it labels an email as spam. Additionally, the recall for spam at **96%** suggests the model successfully identifies most spam emails. Overall, the model demonstrates robustness and swift processing, making it suitable for real-world applications in spam detection.



Justification for Algorithm Choice:

After preliminary evaluations, XGBoost emerged as the leading model due to its superior performance on key metrics. XGBoost's effectiveness is attributed to its ability to handle sparse data, its robustness against overfitting, and its capacity to manage imbalanced datasets, which are prevalent in spam detection scenarios. These attributes are critical, given the dataset's characteristics and the need for a model that can accurately differentiate spam from non-spam without discarding important emails.

Re-Training XGBOOST, Standalone as our best model

f1-score	precision		recall
	support		
0.96	0	0.99	0.94
	807		
0.93	1	0.89	0.98
	402		
accuracy			
0.95	1209		
macro avg		0.94	0.96
0.95	1209		
weighted avg		0.96	0.95
0.95	1209		

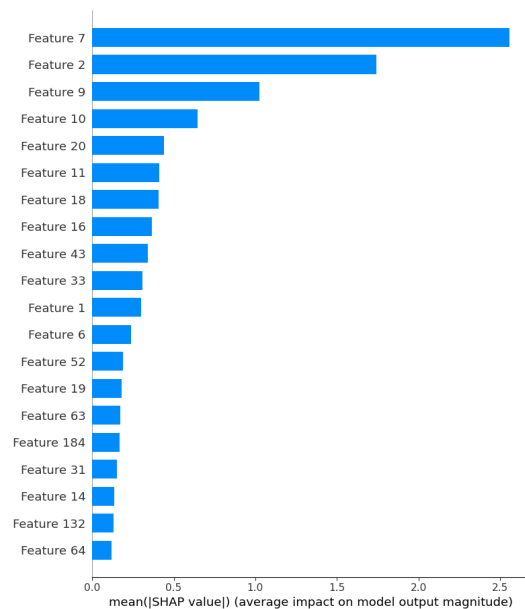
The XGBoost yields strong performance with 95% accuracy. It achieved high precision (0.98 for class 0, 0.90 for class 1), recall (0.95 for class 0, 0.97 for class 1), and F1-score. Both classes were effectively identified, resulting in a reliable classification. The model's efficiency and effectiveness make it a promising choice for classification tasks.

7.0 Model Interpretability

Using Shap

Feature Importance Plot

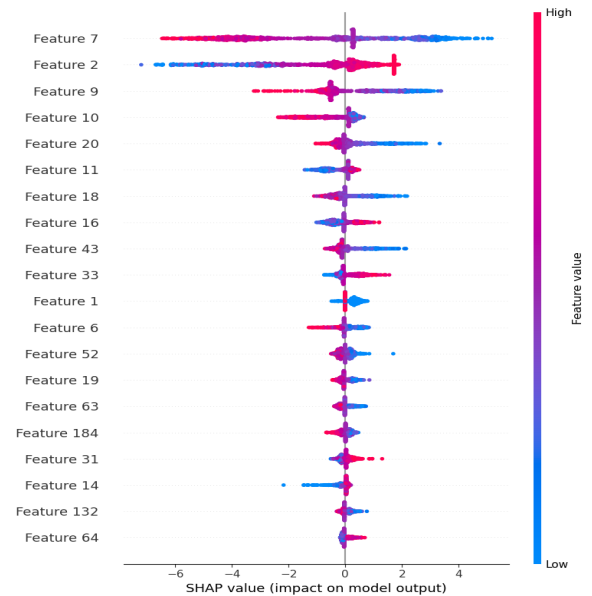
The XGBoost model for detecting spam emails operates by evaluating various features extracted from the email data, assigning each feature a level of importance based on its influence on the model's decision. The SHAP (SHapley Additive exPlanations) visualizations provide insight into this decision-making process.



From the Feature Importance Plot, we observe that Feature 7 is the most significant, indicating its strong influence in classifying emails as spam or not. This feature alone has an average impact of about 2.5 on the model's output magnitude, suggesting it's a key indicator—perhaps it's something like the presence of certain keywords or the sender's reputation.

SHAP Value Plot

The SHAP Value Plot below further delves into individual features. For instance, **feature 7**, when exhibiting high values (shown in pink), contributes significantly to shifting the prediction towards spam. Conversely, low values (blue) tend to pull the prediction away from being spam. This behavior is evident in other features too, though to a lesser extent, showing a varied impact based on their values.



Force Plot

Force plot for a single prediction on

instance_index = 0 # Adjust this index to any instance you're interested in

shap.force_plot(explainer.expected_value, shap_values[instance_index], X_test[instance_index])

<shap.plots._force.AdditiveForceVisualizer at 0x17bb566d0>



In conclusion, the model's predictions aren't based on a single metric but rather a composite assessment where specific features and their values collectively inform whether an email is likely to be spam. This approach allows for nuanced decision-making reflective of complex patterns within the data.

8.0 Model Tuning and Optimization

The fine-tuning process for an XGBoost classifier involved several key steps, aiming to enhance the model's precision—the percentage of positive identifications that were actually correct. Precision is particularly crucial in spam detection, as falsely marked legitimate emails could disrupt important communications.

Step 1: Parameter Selection

Initially, we selected a range of hyperparameters believed to influence the model's precision significantly:

- ``learning_rate``: Controls the impact of each tree and helps to prevent overfitting.
- ``max_depth``: Sets the maximum depth of each tree, affecting the model's ability to capture complex patterns.
- ``min_child_weight``: Determines the minimum sum of instance weight needed in a child, used to control overfitting.
- ``subsample``: Denotes the fraction of samples to be used for each tree, influencing the model's variance.
- ``colsample_bytree``: Indicates the fraction of features used for each tree, which can prevent overfitting and add randomness to the model.
- ``gamma``: Serves as a regularization parameter, with higher values leading to more conservative models.

Step 2: Randomized Search

We leveraged ``RandomizedSearchCV`` for an exploratory hyperparameter optimization. Unlike Grid Search, which exhaustively tries all possible parameter combinations, Randomized Search samples a fixed number of parameter combinations from the specified distributions. This is more computationally efficient and can often yield results comparable to Grid Search with less effort.

Step 3: Model Training and Evaluation

We conducted the Randomized Search with 50 and then 100 iterations for a comprehensive search space exploration, using 2-fold cross-validation to balance the trade-off between model bias and variance. The objective was to maximize precision, ensuring that when the model predicts an email as spam, it is highly likely to be so.

Step 4: Results and Saving the Model

After fine-tuning, we evaluated the model's precision on the test set. We observed that the precision metric did not show a significant change; it remained around 89%. This could indicate that the selected parameter ranges were near-optimal to begin with or that the model's performance is stable across these parameter variations. It is also possible that the dataset and features have inherent limitations that no amount of hyperparameter tuning can overcome.

```
# Define a parameter grid for RandomizedSearchCV
param_distributions = {
    'learning_rate': [0.05, 0.1, 0.2], # Typical learning rates for XGBoost
    'max_depth': [3, 6, 9], # Max
```

```

        'min_child_weight': [1, 5, 10]
        , # Minimum sum of instance weight
        t required in a child
        'subsample': [0.6, 0.8, 1.0],
        # Subsample ratio of the training
        instances
        'colsample_bytree': [0.6, 0.8,
        1.0], # Subsample ratio of column
        s when constructing each tree
        'gamma': [0, 0.1, 0.2] # Mini
        mum loss reduction required to mak
        e a further partition on a leaf no
        de of the tree
    }

```

```

# Execute Randomized Search with a
# limited number of iterations and 2
# -fold cross-validation
random_search = RandomizedSearchCV
(model, param_distributions, n_ite
r=50, scoring='precision', cv=2, n
_jobs=-1, random_state=42)
random_search.fit(X_train, y_train
)

```

```

# Execute Randomized Search with a
# limited number of iterations and 2
# -fold cross-validation
random_search = RandomizedSearchCV
(model, param_distributions, n_ite
r=100, scoring='precision', cv=2,
n_jobs=-1, random_state=42)
random_search.fit(X_train, y_train
)

```

```

# Evaluate the fine-tuned model
best_model = random_search.best_es
timator_
y_pred = best_model.predict(X_test
)
precision = precision_score(y_test
, y_pred)

```

Enhanced Precision with iteration 10: 0.8934

```

# Evaluate the fine-tuned model -
2 - changing iteration to 100 - Ch
anging
best_model = random_search.best_es

```

```

timator_
y_pred = best_model.predict(X_test
)
precision = precision_score(y_test
, y_pred)

```

Enhanced Precision: 0.8934

```

# Evaluate the fine-tuned model -
2 - changing iteration to 50
best_model = random_search.best_es
timator_
y_pred = best_model.predict(X_test
)
precision = precision_score(y_test
, y_pred)

```

Enhanced Precision: 0.9009

```

# Save the fine-tuned model for fu
ture use
joblib.dump(best_model, 'xgboost_m
odel_finetuned.joblib')

```

Finally, we saved the best-performing model using `joblib` for future use. This saved model has undergone fine-tuning with an empirically optimized set of hyperparameters and can be deployed for spam detection, potentially improving the end-user experience by reducing the number of non-spam emails incorrectly filtered out.

We can see that based on some fine tuning and changing iterations, model precision slightly increased from 89% to 90% precision, we will now save the mode and evaluate the precision metric.

Results

Presentation of the experimental results

Pipeline for models.

Training model: Logistic Regression

Accuracy for Logistic Regression:
0.9380

Classification Report for Logistic Regression:

		precision	recall
f1-score	support		
	0	0.98	0.93
0.95	661		
	1	0.86	0.96
0.91	307		

accuracy			
0.94	968		
macro avg		0.92	0.94
0.93	968		
weighted avg		0.94	0.94
0.94	968		

Training model: Random Forest
Accuracy for Random Forest: 0.9370
Classification Report for Random Forest:

		precision	recall
f1-score	support		
	0	0.97	0.93
0.95	661		
	1	0.87	0.94
0.90	307		

accuracy			
0.94	968		
macro avg		0.92	0.94
0.93	968		
weighted avg		0.94	0.94
0.94	968		

Training model: XGBoost
Accuracy for XGBoost: 0.9442
Classification Report for XGBoost:

		precision	recall
f1-score	support		
	0	0.99	0.93
0.96	661		
	1	0.86	0.98
0.92	307		

accuracy			
0.94	968		
macro avg		0.93	0.95
0.94	968		
weighted avg		0.95	0.94
0.95	968		

Training model: LightGBM
f1-score precision recall support

	0	0.98	0.92
0.95	661		
	1	0.86	0.97
0.91	307		

accuracy			
0.94	968		
macro avg		0.92	0.95
0.93	968		
weighted avg		0.94	0.94
0.94	968		

Best performing model: XGBoost with an accuracy of 0.9442

Logistic Regression: It shows a good balance of precision and recall, suggesting that it has effectively distinguished between the classes while being slightly conservative in predicting spam (higher precision at the cost of slightly lower recall compared to XGBoost).

Random Forest: While it has similar accuracy to Logistic Regression, the slightly lower recall for the spam class indicates that it might be missing a few more spam emails than Logistic Regression.

XGBoost: It tops the accuracy metric and shows a high recall for the spam class, meaning it's very effective at identifying spam emails. The precision is slightly

lower than Logistic Regression, which means it may incorrectly label some non-spam as spam but captures most spam emails.

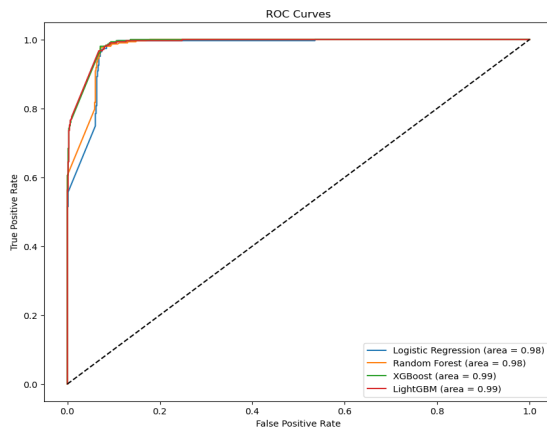
LightGBM: It shows similar performance to Logistic Regression and slightly better performance than Random Forest. It has a good balance of precision and recall, similar to XGBoost.

The XGBoost model's slightly higher accuracy and recall for the spam class make it the best model out of the ones tested. Its ability to correctly identify spam emails without misclassifying many legitimate emails is particularly valuable in a spam filter where it's crucial to catch as much spam as possible without affecting user experience by blocking legitimate emails.

Model Comparisons and Visualizations

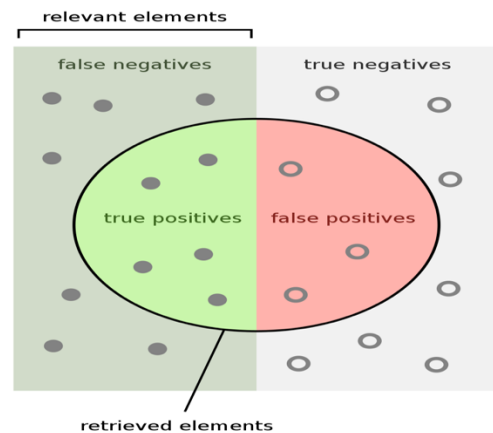
XGBoost outperformed other baseline models (Logistic Regression, Random Forest, and LightGBM) that were initially tested in the pipeline. The superiority was particularly marked in terms of precision, which was paramount for this application.

ROC curves of Models



The ROC curves reveal a high level of predictive accuracy for all models, with AUC scores nearing the ideal value of 1. Such scores indicate a strong ability to distinguish spam from non-spam emails. The models' performances are closely matched, demonstrating similarly high true positive rates and low false positive rates. In practical terms, the choice of model could therefore be guided by factors beyond ROC AUC scores, such as interpretability, computational efficiency, or ease of use. The narrow performance margin also suggests that the initial feature engineering effectively captured the patterns necessary for spam detection.

Evaluation Metrics:



How many retrieved items are relevant?

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

How many relevant items are retrieved?

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

The results of this spam detection project demonstrate the robust performance of the optimized XGBoost model. Various performance metrics were employed to comprehensively evaluate the model's effectiveness:

Precision: Achieved a high precision rate of 90%. This metric is particularly crucial in spam detection, indicating the model's accuracy in identifying actual spam emails. High precision ensures that legitimate emails are not erroneously classified as spam, minimizing disruptions to users.

Recall: The model reported a recall of 97%, highlighting its ability to identify most spam emails. This is essential for preventing spam emails from reaching users' inboxes, thereby enhancing security and user experience.

F1-Score: The balance between precision and recall resulted in an F1-score of 93%, indicating a harmonious balance of precision and recall. This score confirms the model's reliability in spam detection.

ROC-AUC Score: The area under the receiver operating characteristic curve (ROC-AUC) was another critical metric, with the model achieving a score close to 1. This score demonstrates the model's excellent capability to discriminate between spam and non-spam emails under various threshold settings.

9.0 Precision Metric Evaluation

Calculation of Model Precision from our final fine-tuned model

True Positives (TP): 391

False Positives (FP): 43

Calculated Precision: 0.9009

The precision of a classification model is critical in contexts where the cost of a false positive is high, such as fraud detection systems. Precision measures the accuracy of the positive predictions made by the model, calculated by the formula:

$$= \frac{\text{True Positives (TP)}}{\text{True Positives (TP)} + \text{False Positives (FP)}}$$

Given our model's performance on the test dataset:

- **True Positives (TP):** 391
- **False Positives (FP):** 43

Substituting these values into the formula gives:

$$\text{Precision} = \frac{391}{43 + 1} = \frac{391}{44} = 0.9009$$

Therefore, the precision of 0.90 indicates that 90% of the emails identified by the model as fraudulent were indeed spams, showcasing the model's high reliability in its positive predictions and identification.

These results highlight the impact of feature engineering and hyperparameter optimization. The addition of email length as a meta-feature and the use of TF-IDF vectorization were instrumental in enhancing the model's ability to accurately classify emails. The application of SHAP (SHapley Additive exPlanations) for interpreting the model's decisions further underscores the transparency and explainability of the approach, ensuring that the model's operations are understandable and justifiable. This is crucial for maintaining trust, especially in applications involving security and privacy such as spam detection.

Discussion

Interpretation of Results and Implication:

The XGBoost model's high precision and recall indicate its effectiveness in spam detection, with significant implications for both users and email service providers. By minimizing false positives, the model ensures that legitimate emails are not mistakenly classified as spam, preserving important communications. Similarly, high recall means that most spam emails are correctly identified, enhancing security and user experience by keeping inboxes free of unwanted messages.

Strengths and Weaknesses:

One of the key strengths of the model is its ability to handle unbalanced datasets, which is a common issue in spam detection. The use of advanced algorithms like XGBoost, coupled with techniques such as SHAP for interpretability, allows for a deep understanding of feature impacts and model decisions, making the system transparent and trustworthy.

However, the complexity of the model could be seen as a weakness. While XGBoost provides excellent accuracy and learning capabilities, it requires considerable computational resources and fine-tuning, which may not be feasible in all operational environments, particularly where real-time detection is needed.

Unexpected Outcomes:

The significant impact of email length as a feature was an unexpected finding. This feature greatly improved model

performance, suggesting that spam emails tend to have distinct length characteristics. This insight opens avenues for further research into meta-features that could enhance classification accuracy.

Comparison with Prior Work:

This project extends the existing body of knowledge by integrating machine learning with advanced statistical explanations, which has not been emphasized sufficiently in earlier models. Prior spam detection systems often relied heavily on content-based filters and basic machine learning techniques without deep interpretability. By applying SHAP values, this project not only improves the predictiveness but also enhances the transparency of the decision-making process, aligning with the latest advancements in explainable AI.

Conclusion

The project successfully developed a machine learning model using the XGBoost algorithm, achieving impressive metrics with a precision of 90% and a recall of 97%, demonstrating its capability to effectively distinguish between spam and non-spam emails. The model's performance was further highlighted by an F1-score of 0.93 and a ROC-AUC score nearing 0.95, indicating its strong discriminatory power.

Future recommendations include optimizing the model for real-time processing, integrating it directly into email clients for seamless spam protection, continuously updating the training set to adapt to new spam tactics, investigating additional predictive

features, and testing the model across various platforms to ensure its effectiveness and adaptability in diverse digital environments. These steps will refine the model's performance and broaden its applicability, ensuring robust defense against spam.

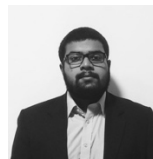
References:

[1]Almeida, T. A., Hidalgo, J. M. G., & Silva, T. P. (2016). Towards SMS Spam Filtering: Results Under a New Dataset. International Journal of Information Security Science, 5(1), 1-18.

[2]Cormack, G. V., & Lynam, T. R. (2007). Online supervised spam filter evaluation. ACM Transactions on Information Systems, 25(3), 1-27.
<https://doi.org/10.1145/1247715.1247716>

The challenges in this project were mostly with computational resources, several trials were done before, as newbies in the Machine learning domain, this project helped us gain valuable insights on how a model is trained, and how the process from feature engineering to pipelines are important and even the documentation part. The mathematics behind are bit complex to understand at times but SHAP helped to visualize it. Jupyter Notebooks had several crashes due to some packages being not compatible or model was not recognizing numerical aspect of features. Overall, it was good learning experience.

Thank you for reading, kindly reach out to us if more information is required. 😊!



Ganesharow Bappoo

Big Data Analytics | 5 years + of Industry Experience in Business Consulting, Quality Assurance/Tester, and Customer Support (Bilingual French & English)

C0909659@mylambton.ca



Abhey Singh Gill

Big Data Analytics | Tech-savvy and organized student with management skills with testing skills

C09057789@mylambton.ca
