

Configuration Manual

Protection Against Spear Phishing Attacks Using the Ensemble Method of Machine Learning

Jecinta Ifechukwu Fidelis

Student ID: 23148306

School of Computing
National College of Ireland

Supervisor: Imran Khan

National College of Ireland
MSc Project Submission Sheet



School of Computing

Jecinta Ifechukwu Fidelis

Student Name:
.....

Student ID:
23148306
.....

Program me: **Year**
..... :
Msc in Cyber Security 2024

Module:
.....
Msc Research Practicum

Lecturer:
.....
Imran Khan

Submission Due Date:
10/08/2024
.....

Project Title:
.....

Word Count: **Page Count:**
.....

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:
Jecinta Ifechukwu Fidelis
.....

Date:
10/08/2024
.....

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

Attach a completed copy of this sheet to each project (including multiple copies)	<input type="checkbox"/>
Attach a Moodle submission receipt of the online project submission , to each project (including multiple copies).	<input type="checkbox"/>
You must ensure that you retain a HARD COPY of the project , both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	<input type="checkbox"/>

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Jecinta Ifechukwu Fidelis
23148306

1. Introduction

This research utilizes Google Colab, a robust cloud-based platform, to facilitate the development and execution of complex machine learning models. Google Colab provides complimentary access to both CPU and GPU resources, making it an invaluable tool for handling intensive computational tasks associated with machine learning. Additionally, the platform supports efficient management and backup of datasets and codebases, ensuring data integrity and continuity throughout the research process. Google Colab also comes pre-equipped with numerous machine learning libraries, streamlining the implementation process by reducing the need for extensive setup, thereby enabling quicker and more efficient project execution.

2. System Specification

The system configuration used in this project are:

- Operating system: Linux
- Processor: Intel Xeon CPU (2)
- RAM: 13GB

This configuration provides the necessary computational power to efficiently process and analyze large datasets, essential for training and testing the machine learning models employed in this research.

3. Software Tools

The primary software tools used in the implementation of this project are centered around the Python programming language, with Google Colab serving as the development environment. The following tools were integral to the project's success:

- **Python:** The core programming language used for developing machine learning models.
- **Google Colab:** The cloud-based platform that facilitated code execution, model training, and deployment.
- **Jupyter Notebook:** Used within Google Colab for interactive code development and testing.

These tools provided a cohesive environment for the development, testing, and deployment of machine learning models.

4. Implementation

The implementation of this project involved several key Python libraries that provided the necessary functionality for data manipulation, model development, and evaluation:

- **Numpy:** Used for numerical computations and data manipulation.

- **Pandas:** Provided tools for data analysis and manipulation, particularly for handling data frames.
- **Matplotlib:** Facilitated the creation of visualizations for data analysis and model evaluation.
- **Sklearn (Scikit-learn):** A machine learning library used for model training, evaluation, and feature extraction.
- **NLTK (Natural Language Toolkit):** Employed for text preprocessing tasks, including tokenization, stop word removal, and lemmatization.
- **LightGBM:** A gradient boosting framework that was utilized for efficient model training and prediction.
- **XGBoost:** Another gradient boosting library, used alongside LightGBM for model comparison.

These libraries formed the backbone of the project's implementation, enabling sophisticated data processing and model training. Figure Descriptions:

Fig 1: Demonstrates how the Google Colab environment was set up by mounting the necessary drives.

Fig 2: Illustrates the process of importing the required Python libraries.

Fig 3: Shows the function used to calculate and display evaluation metrics for the models.

Fig 4: Depicts the function created for visualizing the confusion matrix, an important tool for evaluating model performance.

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Fig 1: Mounting the drive in google colab

```
1 ### Libraries for data pre-processing and Visualization
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import json
7 import os
8 from collections import Counter
9
10 ### Natural Language Processing Libraries
11 import re
12 from nltk.corpus import stopwords
13 from nltk.stem import WordNetLemmatizer
14 from collections import Counter
15 import nltk
16 # Download necessary NLTK data files
17 nltk.download('stopwords')
18 nltk.download('wordnet')
19
20 ## Feature Extraction
21 from sklearn.feature_extraction.text import TfidfVectorizer
22
23 ### Machine Learning Libraries
24 from sklearn.model_selection import train_test_split
25 from sklearn.metrics import accuracy_score, f1_score, matthews_corrcoef, precision_score, recall_score, confusion_matrix
26 from sklearn.linear_model import LogisticRegression
27 from sklearn.naive_bayes import MultinomialNB
28 from sklearn.tree import DecisionTreeClassifier
29 from sklearn.ensemble import RandomForestClassifier
30 import xgboost as xgb
31 import lightgbm as lgb
```

Fig 2: Importing Libraries

```

1 def evaluate_model_performance(y_true, y_pred):
2
3     metrics = {
4         'accuracy': accuracy_score(y_true, y_pred),
5         'f1_score': f1_score(y_true, y_pred, average='weighted'),
6         'mcc': matthews_corrcoef(y_true, y_pred),
7         'precision': precision_score(y_true, y_pred, average='weighted'),
8         'recall': recall_score(y_true, y_pred, average='weighted')
9     }
10
11     return metrics

```

Fig 3: Evaluation Metrics function

```

1 def plot_confusion_matrix(y_true, y_pred, class_names=None):
2     cm = confusion_matrix(y_true, y_pred)
3     if class_names is None:
4         class_names = [str(i) for i in range(cm.shape[0])]
5
6     fig, ax = plt.subplots(figsize=(8, 6))
7     im = ax.imshow(cm, interpolation='nearest', cmap=plt.cm.Blues)
8     ax.figure.colorbar(im, ax=ax)
9
10    # We want to show all ticks and label them with the respective list entries
11    ax.set(xticks=np.arange(cm.shape[1]),
12          yticks=np.arange(cm.shape[0]),
13          xticklabels=class_names, yticklabels=class_names,
14          title='Confusion Matrix',
15          ylabel='True label',
16          xlabel='Predicted label')
17
18    # Rotate the tick labels and set their alignment.
19    plt.setp(ax.get_xticklabels(), rotation=45, ha="right",
20            rotation_mode="anchor")
21
22    # Loop over data dimensions and create text annotations.
23    fmt = 'd'
24    thresh = cm.max() / 2.
25    for i in range(cm.shape[0]):
26        for j in range(cm.shape[1]):
27            ax.text(j, i, format(cm[i, j], fmt),
28                  ha="center", va="center",
29                  color="white" if cm[i, j] > thresh else "black")
30    fig.tight_layout()
31    plt.show()

```

Fig 4: Confusion Matrix visualization function

1 Data Gathering

The process of data gathering involved collecting data points that were instrumental in distinguishing between normal and phishing emails. The steps included accessing the dataset path on Google Drive, creating a function to read each JSON file, converting these files into Pandas data frames, and labeling the data accordingly. Specifically, a new column labeled "Label" was added to identify each email as normal or phishing. The dataset was further refined by removing insignificant columns, such as 'id' and 'sender name,' as shown in the figures below.

Fig 5: Illustrates the process of reading and processing the datasets.

Fig 6: Shows the data gathering process, including the labeling and cleaning of the dataset.

```

1 path = '/content/drive/MyDrive/Email-project'

1 def json_files_to_dataframe(directory):
2     # Initialize an empty list to hold the data from each file
3     data_list = []
4
5     # Loop through each file in the directory
6     for filename in os.listdir(directory):
7         if filename.endswith('.json'):
8             file_path = os.path.join(directory, filename)
9
10            # Open and read the JSON file
11            with open(file_path, 'r') as file:
12                data = json.load(file)
13
14            # Append the data to the list
15            data_list.append(data)
16
17            # Create a DataFrame from the list of data
18            df = pd.DataFrame(data_list)
19            return df

1 ## The path of the dataset
2 normal_path = path + '/enron_ham'
3 spam_path = path + '/spear_phishing'

1 ### Normal (Non-Spam) Email
2 normal_df = json_files_to_dataframe(normal_path)
3 ### SPam Email
4 spam_df = json_files_to_dataframe(spam_path)

```

Fig 5: Reading the datasets

```

1 ### Normal Email Target
2 normal_df['Label'] = 'normal'

1 ### Phising Email Target
2 spam_df['Label'] = 'phishing'

1 data = pd.concat([normal_df, spam_df]).reset_index(drop=True)

1 ## Drop insignificant features
2 data = data.drop(['id', 'sender_name'], axis=1)

```

Fig 6: Data Gathering and Labelling

2 Data Understanding

Understanding the dataset involved gaining insights into its structure and composition. The dataset's basic characteristics, such as its head, shape, and label frequency, were examined. The dataset comprises 3,334 rows and three features, including the target feature (label). It was noted that approximately 90% of the emails were normal, while 10% were phishing. These insights were critical for guiding the subsequent data processing steps.

- **Fig 7:** Provides information about the dataset, including the data types of each feature.
- **Fig 8:** Displays the distribution of the target feature, highlighting the imbalance between normal and phishing emails.

```
1 data.shape
(3334, 3)

1 data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3334 entries, 0 to 3333
Data columns (total 3 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   email_subject   3334 non-null   object
1   email_body      3334 non-null   object
2   Label           3334 non-null   object
dtypes: object(3)
memory usage: 78.3+ KB

1 data['Label'].value_counts(normalize=True)*100
Label
normal      89.982004
phishing     10.017996
Name: proportion, dtype: float64
```

Fig 7: Data Info

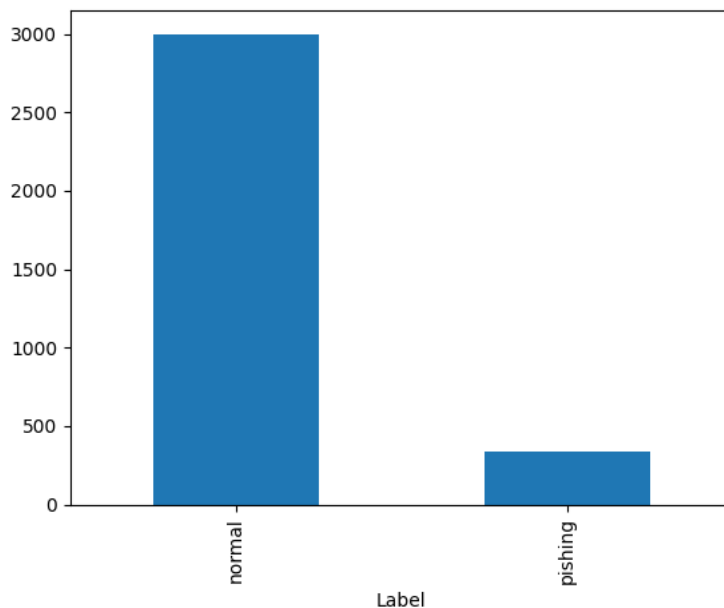


Fig 8: Target feature distribution

3 Data Pre-processing

Data pre-processing is a critical step where the text data from email bodies, which serves as the feature (X) in the dataset, is cleaned and prepared for analysis. This involved applying functions for stop-word removal and lemmatization, converting text to lowercase, and removing unnecessary elements such as URLs, punctuation, special characters, digits, and repetitive words. The pre-processing steps were designed to enhance the quality of the text data, making it more suitable for machine learning model training.

- **Fig 9:** Shows the pre-processing steps applied to the text data, illustrating the improvement in data quality.

```

1 # Initialize stopwords and lemmatizer
2 stop_words = set(stopwords.words('english'))
3 lemmatizer = WordNetLemmatizer()

1 top_10_words = ["subject", "cc", "please", "phillip", "pm", "email", "k", "forwarded", "would", "john"]
2 underscore_patterns = ['_', '____', '_____']
3
4 def preprocess_text(text):
5     # Convert to lowercase
6     text = text.lower()
7     # Remove URLs
8     text = re.sub(r'http\S+|www\S+|https\S+', '', text, flags=re.MULTILINE)
9     # Remove punctuation, special characters, and numbers
10    text = re.sub(r'\n', ' ', text) # Replace newline characters with space
11    text = re.sub(r'[\W\s]', '', text)
12    text = re.sub(r'\d+', '', text) # Remove numbers
13    # Remove underscore patterns
14    for pattern in underscore_patterns:
15        text = text.replace(pattern, '')
16    # Remove stopwords and top 10 words
17    text = ' '.join([word for word in text.split() if word not in stop_words and word not in top_10_words])
18    # Remove repeated words
19    text = ' '.join(sorted(set(text.split()), key=text.split().index))
20    # Lemmatize words
21    text = ' '.join([lemmatizer.lemmatize(word) for word in text.split()])
22
23    return text

1 data['clean_email_body'] = data['email_body'].apply(preprocess_text)

```

Fig 9: Text Data Pre-processing

4 Data Visualization and Data Transformation

Data visualization was used to identify the most frequent words in the cleaned email bodies, while data transformation involved converting the labels into binary values (0 for normal emails and 1 for phishing emails). The transformed dataset, consisting of the "clean email body" and "Label" features, was then prepared for building the machine learning models.

- **Fig 10:** Visualizes the top words in the cleaned dataset.
- **Fig 11:** Lists the top 10 words, showing their frequency.
- **Fig 12:** Demonstrates the data transformation process, preparing the dataset for model training.

```

1 # Combining all email bodies into a single string
2 all_text = ' '.join(data['clean_email_body'])
3 # Splitting the text into words
4 words = all_text.split()
5
6 # Counting the frequency of each word
7 word_counts = Counter(words)
8 # Getting the top 10 most common words
9 top_10_words = word_counts.most_common(10)
10 words, counts = zip(*top_10_words)
11
12 # Creating the bar chart
13 plt.figure(figsize=(10, 5))
14 plt.bar(words, counts, color='skyblue')
15 plt.xlabel('Words')
16 plt.ylabel('Frequency')
17 plt.title('Top 10 Words in Email Bodies')
18 plt.xticks(rotation=45)
19 plt.show()

```

Fig 10: Data Visualization

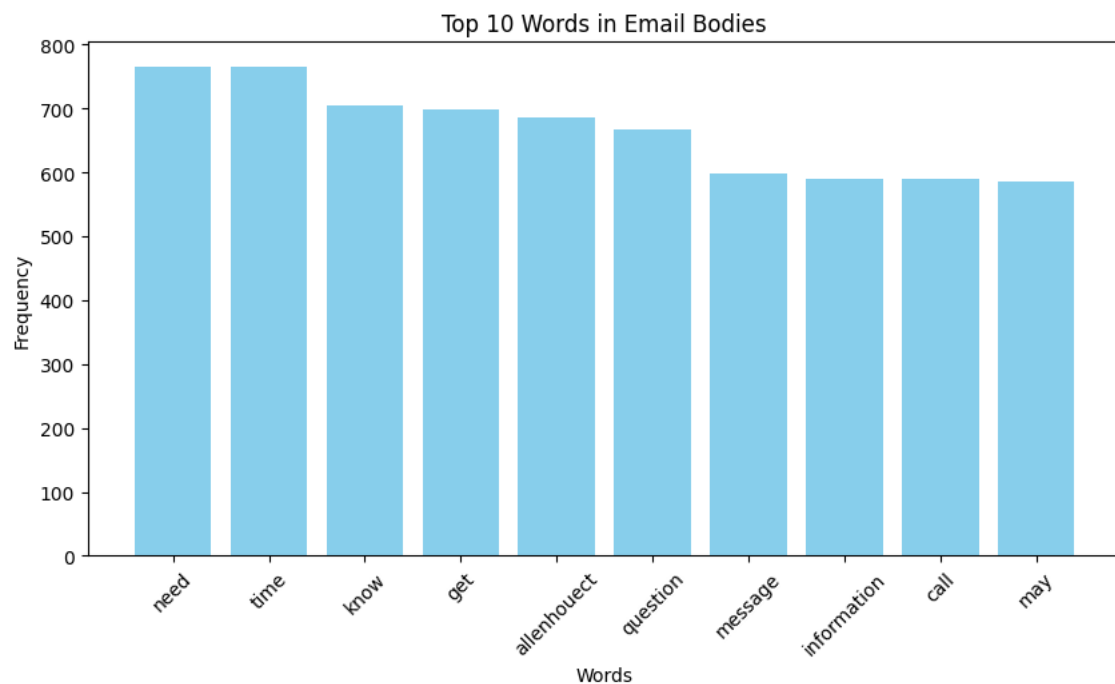


Fig 11: Top 10 words

```

1 data['Label'] = data['Label'].replace(to_replace=['normal', 'pishing'], value=[0, 1])

1 data_new = data[['clean_email_body', 'Label']]

```

Fig 12: Data Transformation

5 Data Splitting

The dataset was split into training and testing sets using Scikit-learn's train-test split function. This stratified approach ensured that the label distribution in both the training and testing sets mirrored the original dataset. The training set comprised 70% of the data, while the testing set included the remaining 30%, with a random seed of 42 used for reproducibility.

- **Fig 13:** Illustrates the data splitting process.

```

1 X = data_new['clean_email_body']
2 y = data_new['Label']

1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42, shuffle=True, stratify=y)

1 print(X_train.shape)
2 print(X_test.shape)
3 print(y_train.shape)
4 print(y_test.shape)

(2333,)
(1001,)
(2333,)
(1001,)

```

Fig 13: Data Splitting

6 Feature Extraction

Feature extraction involved converting the text data into numerical vectors using TF-IDF (Term Frequency-Inverse Document Frequency). This technique was implemented using Scikit-learn libraries and involved setting parameters to filter text, remove rare terms, and handle common words effectively. The TF-IDF was fitted to the training data and subsequently applied to both the training and testing datasets.

- **Fig 14:** Shows the feature extraction process using TF-IDF.

```

1 tfv = TfidfVectorizer(min_df=5, max_features=None,
2 strip_accents='unicode', analyzer='word', token_pattern=r'\w{1,}',
3 ngram_range=(1, 4), use_idf=1, smooth_idf=1, sublinear_tf=1,
4 stop_words = 'english')

1 # Fitting TF-IDF to both training and test sets
2 tfv.fit(X_train)
3 X_train_tfv = tfv.transform(X_train)
4 X_test_tfv = tfv.transform(X_test)

```

Fig 14. Feature Extraction using Tf-idf

7 Machine Learning Model

Various Machine learning models are used to predict the test set.

7.1 Logistics Regression

The logistics regression confusion matrix shows that we have two false negatives.

```

1 lr_model = LogisticRegression() #initialize the logistic regression model
2 lr_model.fit(X_train_tfv, y_train) # fit the logistic regression model to the training data
3 predictions = lr_model.predict(X_test_tfv) # make prediction on the test data
4 lr_eval = evaluate_model_performance(y_test, predictions) # evaluate the performance of the model
5 plot_confusion_matrix(y_test, predictions, class_names=['normal', 'pishing']) ## Confusion Matrix (TP, FP, FN, TN)

```

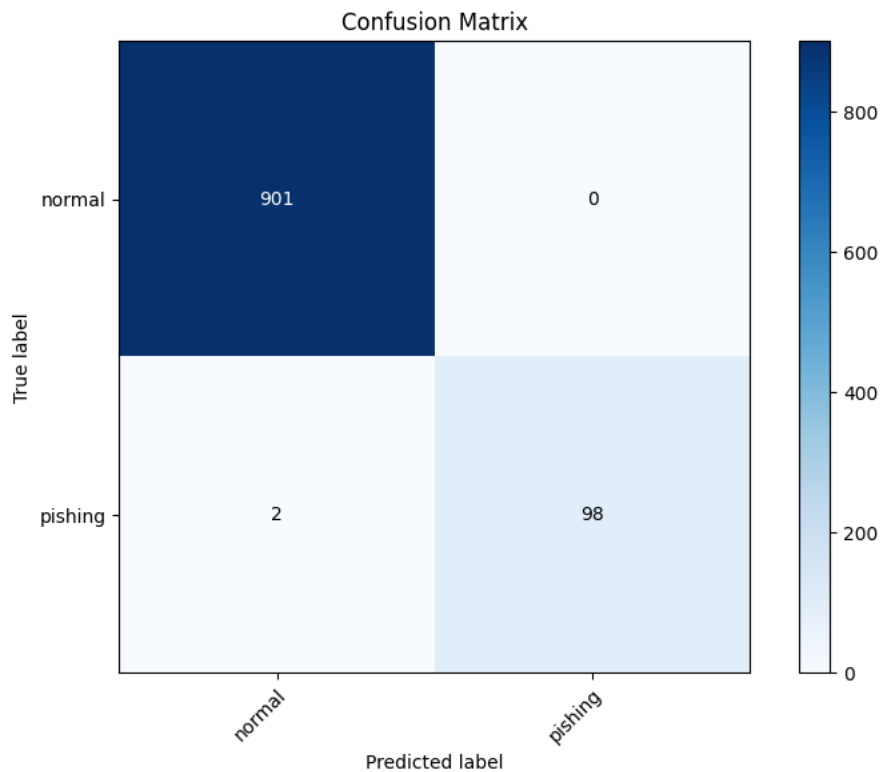


Fig. 15: Logistics Regression Confusion Matrix

7.2 Naïve Bayes

The naïve bayes model confusion matrix shows that we have eight false positives.

```
1 nb_model = MultinomialNB()
2 nb_model.fit(X_train_tfv, y_train)
3 predictions = nb_model.predict(X_test_tfv)
4 nb_eval = evaluate_model_performance(y_test, predictions)
5 plot_confusion_matrix(y_test, predictions, class_names=['normal', 'pishing']) ## The TP, FP, FN, TN
```

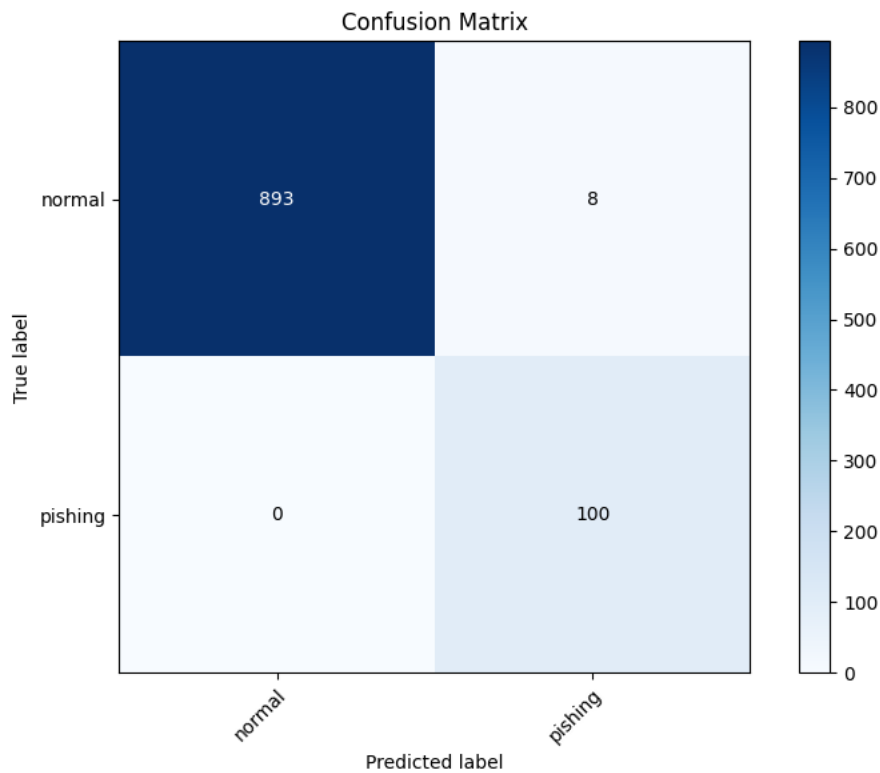


Fig. 16: Naïve Bayes Confusion Matrix

7.3 Decision Tree

The decision tree model confusion matrix shows that we have one false negative.

```
1 dtr_model = DecisionTreeClassifier(random_state=101)
2 dtr_model.fit(X_train_tfv, y_train)
3 predictions = dtr_model.predict(X_test_tfv)
4 dtr_eval = evaluate_model_performance(y_test, predictions)
5 plot_confusion_matrix(y_test, predictions, class_names=['normal', 'pishing']) ## The TP, FP, FN, TN
```

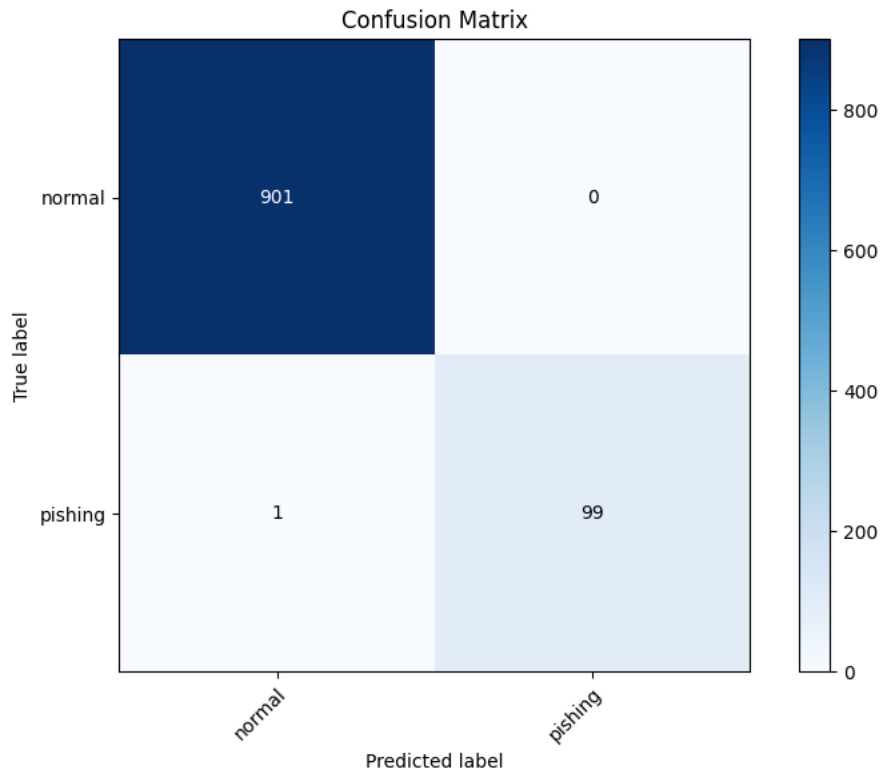


Fig. 17: Decision Tree Confusion Matrix

7.4 Random Forest

The Random Forest model confusion matrix shows that we have three false negative.

```
1 rf_model = RandomForestClassifier(random_state=0, max_depth=7) ## Restricting the Max depth to handle the Model overfitting
2 rf_model.fit(X_train_tfv, y_train)
3 predictions = rf_model.predict(X_test_tfv)
4 rf_eval = evaluate_model_performance(y_test, predictions)
5 plot_confusion_matrix(y_test, predictions, class_names=['normal', 'pishing']) ## The TP, FP, FN, TN
```

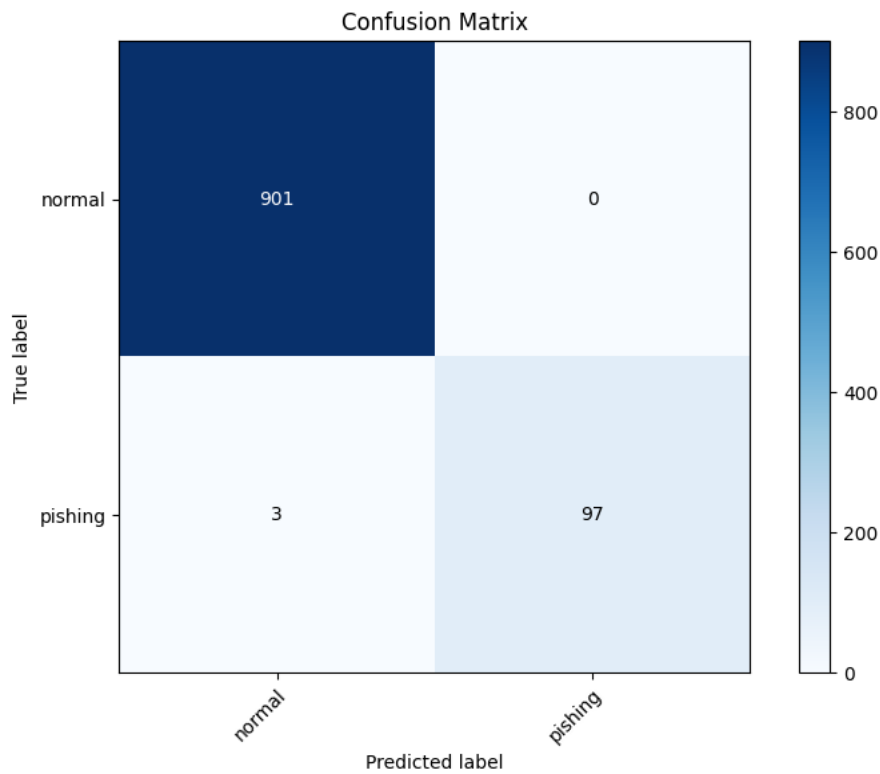


Fig. 18: Random Forest Confusion Matrix

7.5 Xgboost

The Xgboost model confusion matrix shows that we have one false negative.

```
1 xgb_model = xgb.XGBClassifier()
2 xgb_model.fit(X_train_tfv, y_train)
3 predictions = xgb_model.predict(X_test_tfv)
4 xgb_eval = evaluate_model_performance(y_test, predictions)
5 plot_confusion_matrix(y_test, predictions, class_names=['normal', 'pishing']) ## The TP, FP, FN, TN
```

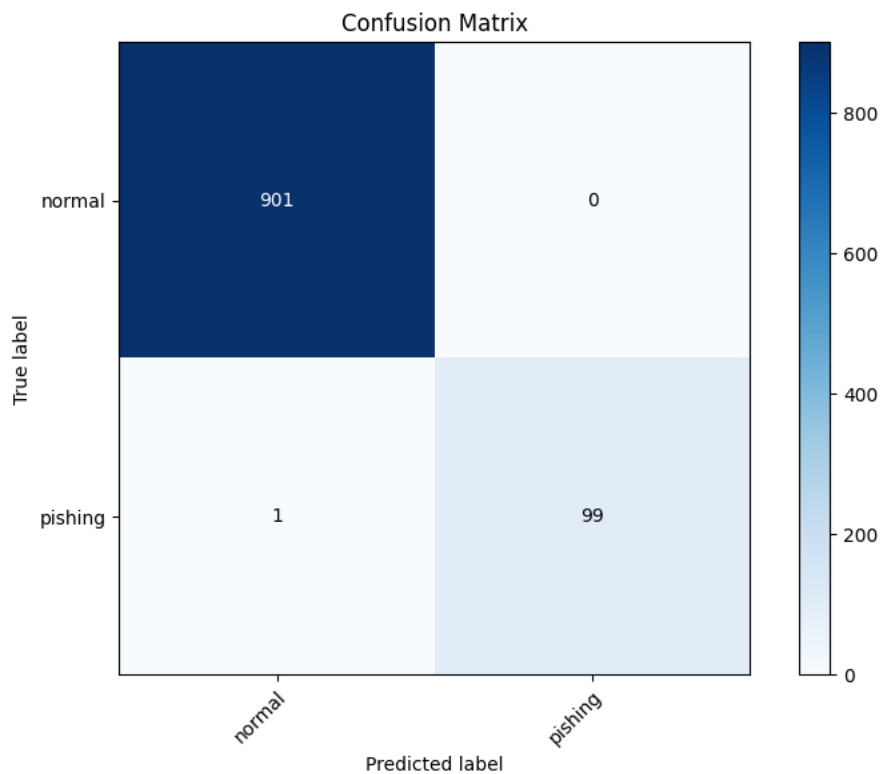


Fig. 19: Xgboost Confusion Matrix

7.6 LightGbm

The Lightgbm model confusion matrix shows that we have one false negative.

```
1 lgb_model = lgb.LGBMClassifier(silent = True)
2 lgb_model.fit(X_train_tfv, y_train)
3 predictions = lgb_model.predict(X_test_tfv)
4 lgb_eval = evaluate_model_performance(y_test, predictions)
5 plot_confusion_matrix(y_test, predictions, class_names=['normal', 'pishing']) ## The TP, FP, FN, TN
```

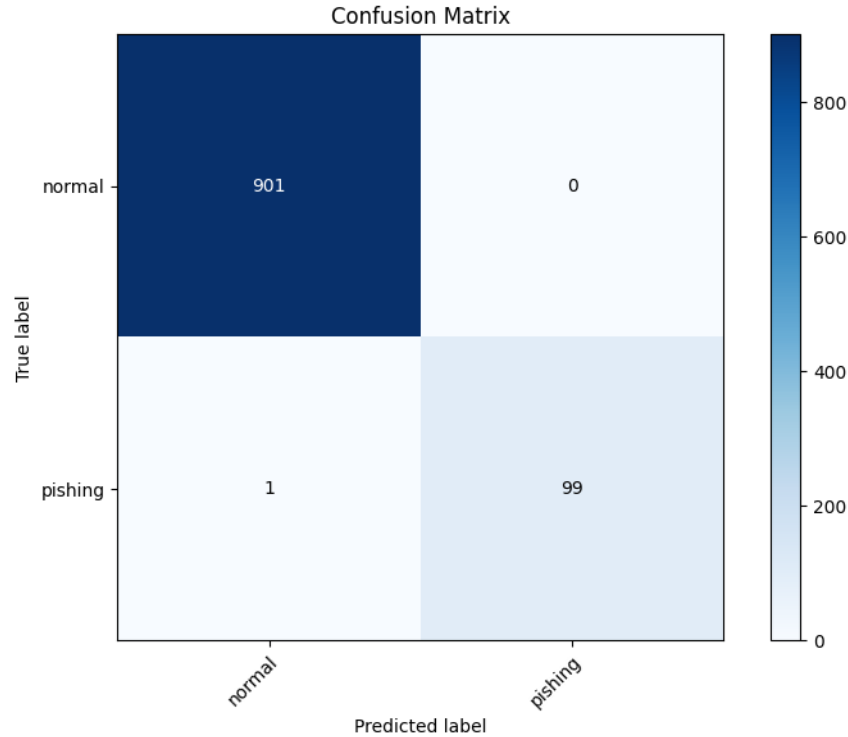


Fig. 20: LightGbm Confusion Matrix

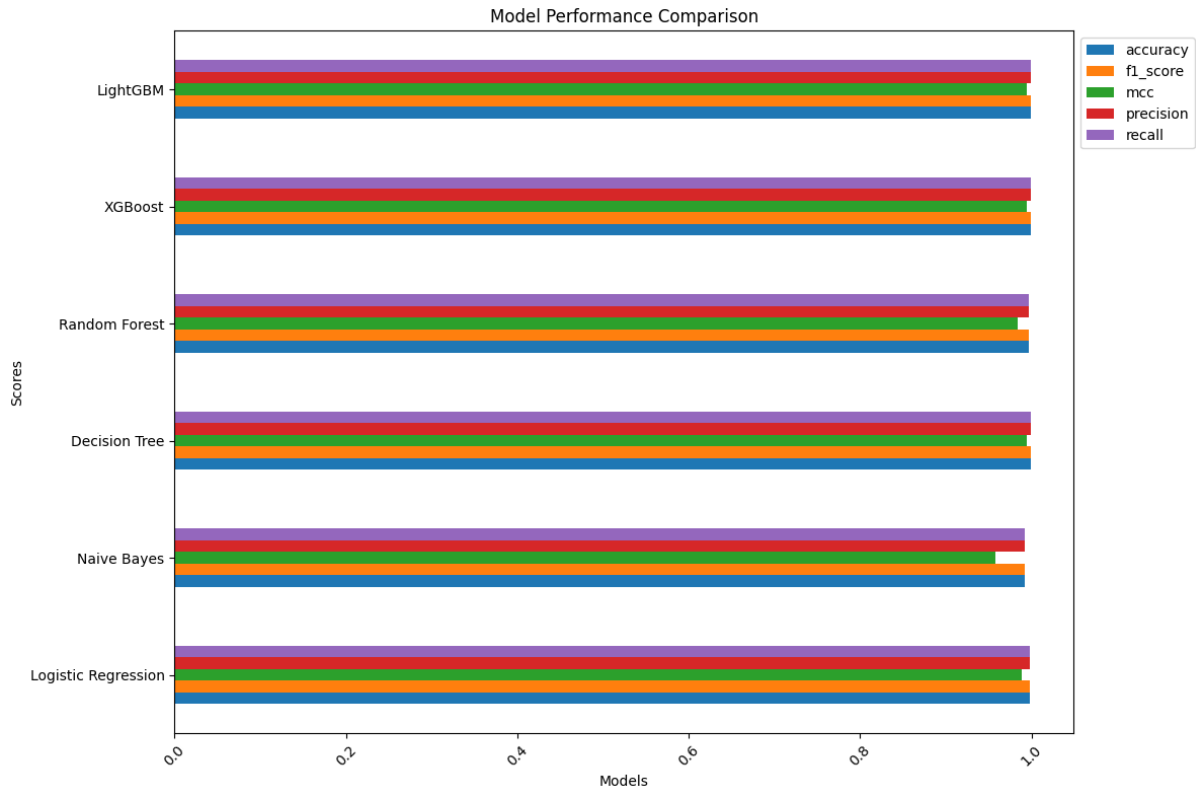
Results and Evaluation

The XGBoost and LightGBM models emerged as the top performers, each achieving an accuracy rate of 99.91%. These models also scored the highest across various evaluation metrics, indicating their effectiveness in classifying phishing emails accurately.

Fig 21: Comparison of model performance results.

	Model	accuracy	f1_score	mcc	precision	recall
0	Logistic Regression	0.998002	0.997993	0.988853	0.998006	0.998002
1	Naive Bayes	0.992008	0.992144	0.957969	0.992600	0.992008
2	Decision Tree	0.999001	0.998999	0.994436	0.999002	0.999001
3	Random Forest	0.997003	0.996983	0.983250	0.997013	0.997003
4	XGBoost	0.999001	0.998999	0.994436	0.999002	0.999001
5	LightGBM	0.999001	0.998999	0.994436	0.999002	0.999001

Fig. 21: Models Comparison results



Email Prediction Interface

The final model was saved using a joblib file, allowing it to be utilized in real-time to predict whether an email is normal or phishing. The interface, built with Gradio, enables users to classify emails interactively.

- **Fig 22:** Shows the real-time email classification interface.
- **Fig 23:** Demonstrates the saved machine learning model and feature extraction function.

The trained model was tested with new email inputs to validate its performance. The results indicated whether the emails were normal or phishing.

- **Fig 24:** Interface displaying a phishing email prediction.
- **Fig 25:** Interface displaying a normal email prediction.


```

1 def classify_email(model_path, vectorizer_path, email_text):
2     """
3     Classifies an input email text to determine if it is a phishing email or not.
4
5     Parameters:
6     model_path (str): Path to the saved model.
7     vectorizer_path (str): Path to the saved vectorizer.
8     email_text (str): The email text to be classified.
9
10    Returns:
11    str: 'Phishing' if the email is classified as phishing, 'Not Phishing' otherwise.
12    """
13    # Load the saved model and vectorizer
14    model_saved = joblib.load(model_path)
15    vector_saved = joblib.load(vectorizer_path)
16
17    # Preprocess the text
18    clean_sample_text = preprocess_text(email_text)
19
20    # Transform the text using the vectorizer
21    sample_vec = vector_saved.transform([clean_sample_text])
22
23    # Predict using the model
24    prediction = model_saved.predict(sample_vec)
25
26    # Map the prediction to a human-readable label
27    if prediction[0] == 1:
28        return 'Phishing'
29    else:
30        return 'Normal'

```

Fig. 22: Email classify near real time prediction

```

1 model_saved = joblib.load("lgb_model.joblib")
2 vector_saved = joblib.load("tfv.joblib")

```

Fig. 23: The saved ML model and feature extraction function

The feature extraction model (TF-IDF) and the saved machine learning model. Whether an email is normal or phishing, the model is utilized to ascertain its input. The email message is phishing, as shown in Fig. 24. When we run the cell and enter a new email message, the text input message will refresh.

Gradio was used to develop the user interface for the email classification system. Its ease of integration with machine learning models and intuitive design features made it an excellent choice for creating a user-friendly interface.

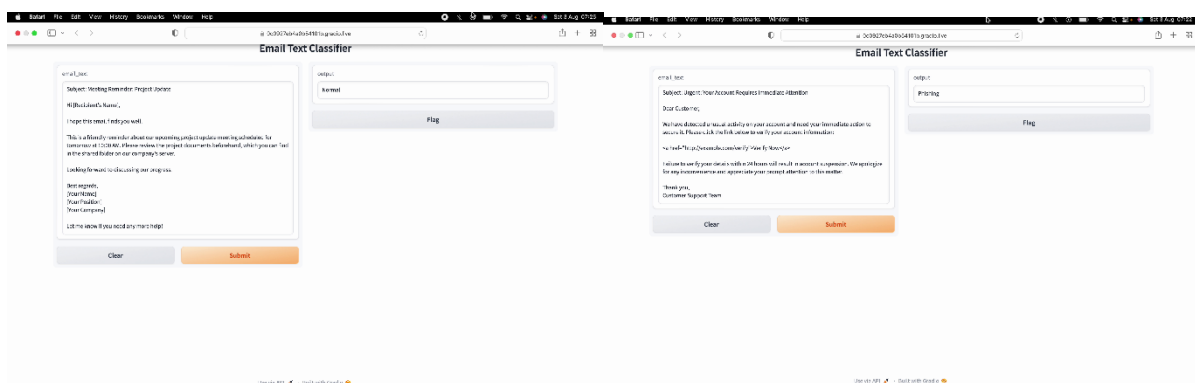


Fig. 24: Phishing Email interface

The architecture comprises several components: user interface, data pre-processing, feature extraction, model training and evaluation, deployment environment, and real-time classification. The system architecture diagram is presented in Figure 25 below.

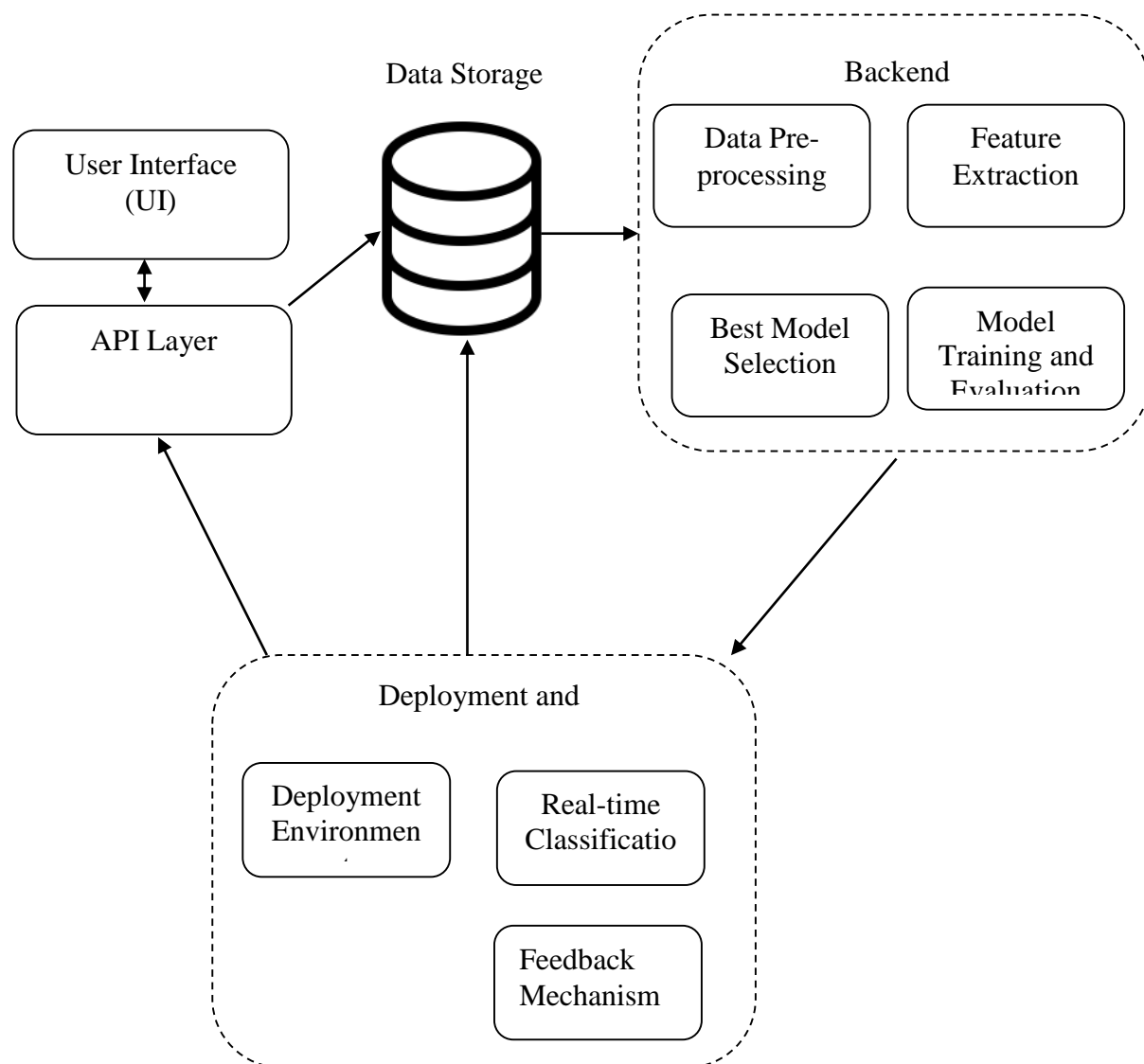


Figure 25: Phishing email classification system architecture

User Interface

The system leverages Gradio to create an intuitive web-based interface, allowing users to input email text, submit it for analysis, and receive classification results immediately. Gradio is an open-source library that simplifies the process of building and deploying user interfaces for machine learning models. It enables developers to create interactive and accessible interfaces without requiring extensive web development experience, making it an ideal choice for the email classification system. A key advantages of using Gradio is its ease of integration with machine learning models. With Gradio, developers can rapidly prototype and deploy

applications by linking pre-trained models directly to user interfaces. This capability is particularly valuable in the email classification system, where the model's ability to process and classify data in real-time is crucial for user satisfaction.

The UI includes features such as:

- **Text Input Field:** Users can enter email content directly into the system.
- **Submit Button:** Initiates the classification process, sending the text to the backend model for analysis.
- **Clear Button:** Resets the input field for new entries.
- **Output Display:** Shows the classification result (e.g., "Phishing" or "Non-Phishing") along with processing time.
- **Feedback Mechanism:** Users can flag emails for further review or feedback, enabling continuous improvement of the model. This integration of user feedback allows the system to learn from real-world usage and refine its detection capabilities.

The screenshot shows a web browser window with the title "Email Text Classifier". The interface is divided into two main sections. On the left, under the label "email_text", there is a text input area containing a sample email body. The email text is: "Subject: Meeting Reminder: Project Update", "Hi [Recipient's Name],", "I hope this email finds you well.", "This is a friendly reminder about our upcoming project update meeting scheduled for tomorrow at 10:00 AM. Please review the project documents beforehand, which you can find in the shared folder on our company's server.", "Looking forward to discussing our progress.", "Best regards,", "[Your Name]", "[Your Position]", "[Your Company]", and "Let me know if you need any more help!". Below the input area are two buttons: "Clear" and "Submit". On the right, under the label "output", there is a text display area showing the classification result "Normal". Below the output area is a "Flag" button. At the bottom of the page, there is a footer that reads "Use via API" and "Built with Gradio".

email_text

Subject: Meeting Reminder: Project Update

Hi [Recipient's Name],

I hope this email finds you well.

This is a friendly reminder about our upcoming project update meeting scheduled for tomorrow at 10:00 AM. Please review the project documents beforehand, which you can find in the shared folder on our company's server.

Looking forward to discussing our progress.

Best regards,
[Your Name]
[Your Position]
[Your Company]

Let me know if you need any more help!

output

Normal

Flag

Clear Submit

Use via API · Built with Gradio

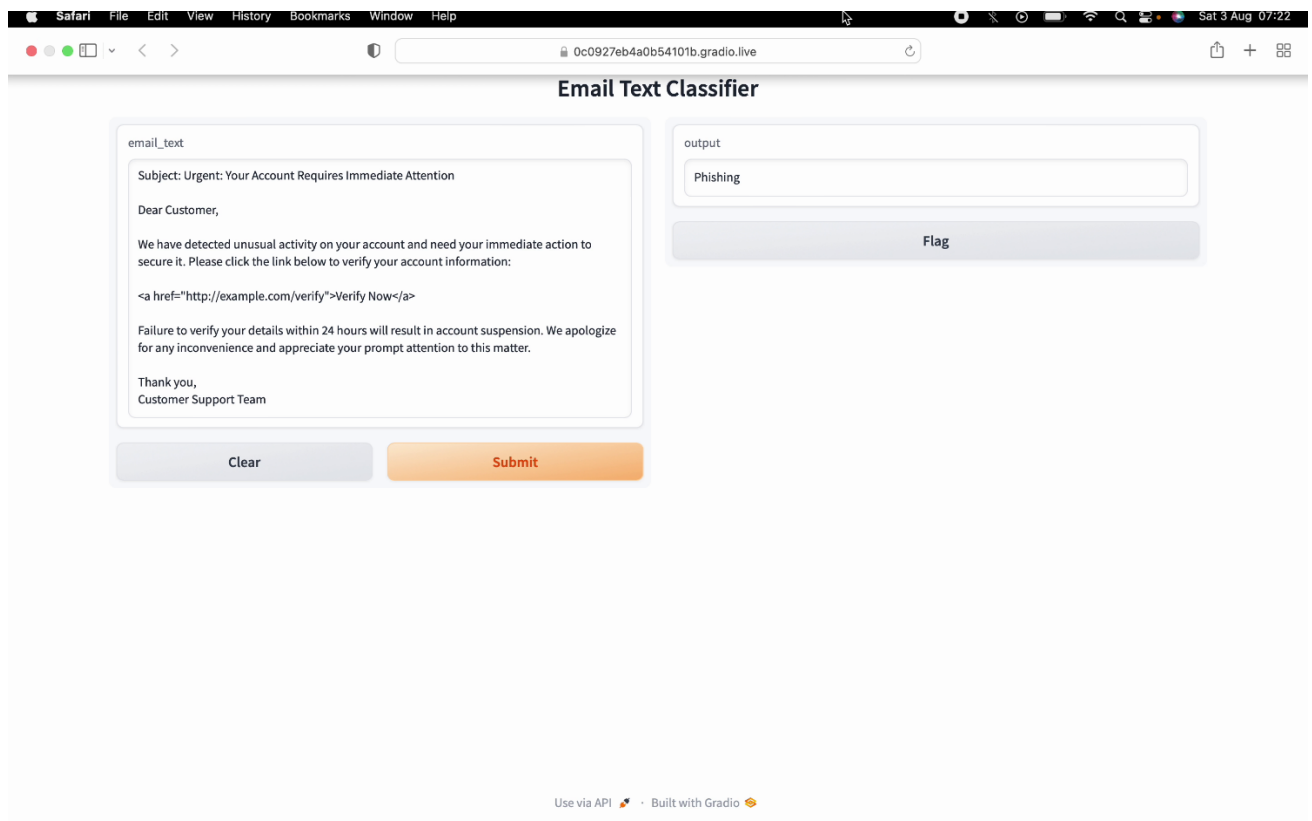


Figure 26: User evaluation with normal and phishing emails