

Configuration Manual

MSc Research Project
MSc Cybersecurity

Tanmay Dharmaraj Shukla
Student ID: x22112421

School of Computing
National College of Ireland

Supervisor: Eugene McLaughlin

National College of Ireland
MSc Project Submission Sheet
School of Computing



Student Name: Tanmay Dharmaraj Shukla

Student ID: x22112421

Programme: MSc Cybersecurity **Year:** 2023-2024

Module: MSc Research Project

Lecturer: Eugene McLaughlin

Submission Due Date: 25th April 2024

Project Title: Unveiling the Power of CNNs with Attention for URL Phishing Detection.

Word Count: 940 **Page Count:** 13

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: Tanmay Dharmaraj Shukla

Date: 25th April 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

Tanmay Dharmaraj Shukla
x22112421

1 INTRODUCTION

This manual is created with the stepwise implementation of how the project implement worked and how the model is created. In addition to it information and process of used libraries and tools required to carry out for project implementation is mentioned. Furthermore, information regarding the specifications of the local machine is mentioned in the manual. Each model information is mentioned in this configure manual.

.

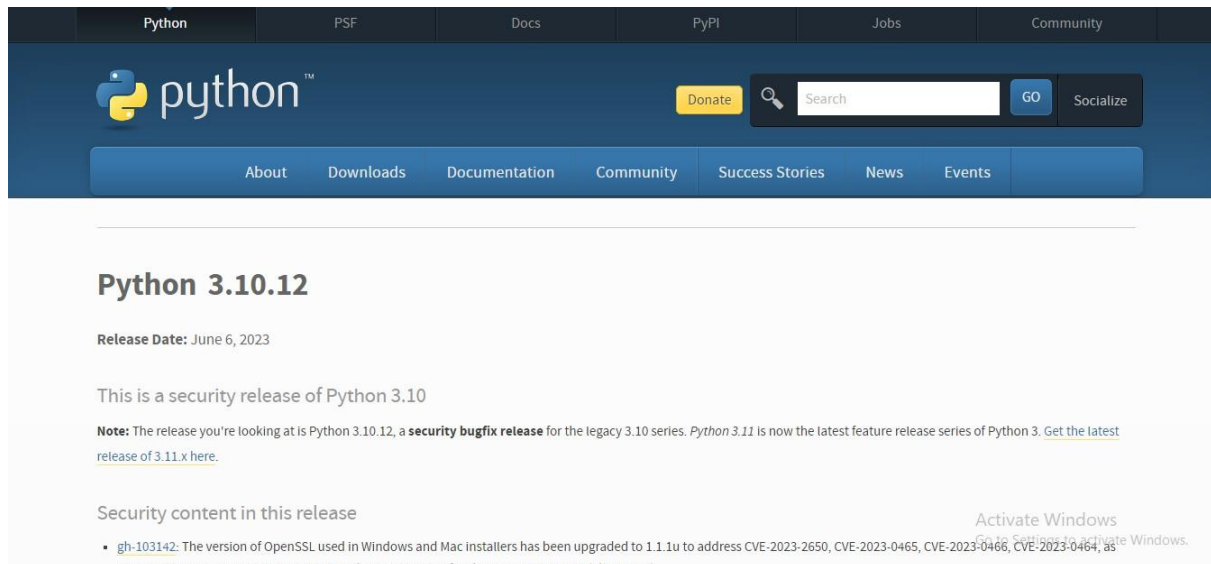
2 HARDWARE CONFIGURATION

- Operating system: Windows 11
- Processor: > 11th Gen Intel(R) Core(TM) i3-1115G4 @ 3.00GHz 3.00 GHz
- System Compatibility: 64-bit
- Hard Disk: 512GB
- RAM: 16 GB

3 SOFTWARE CONFIGURATIONS

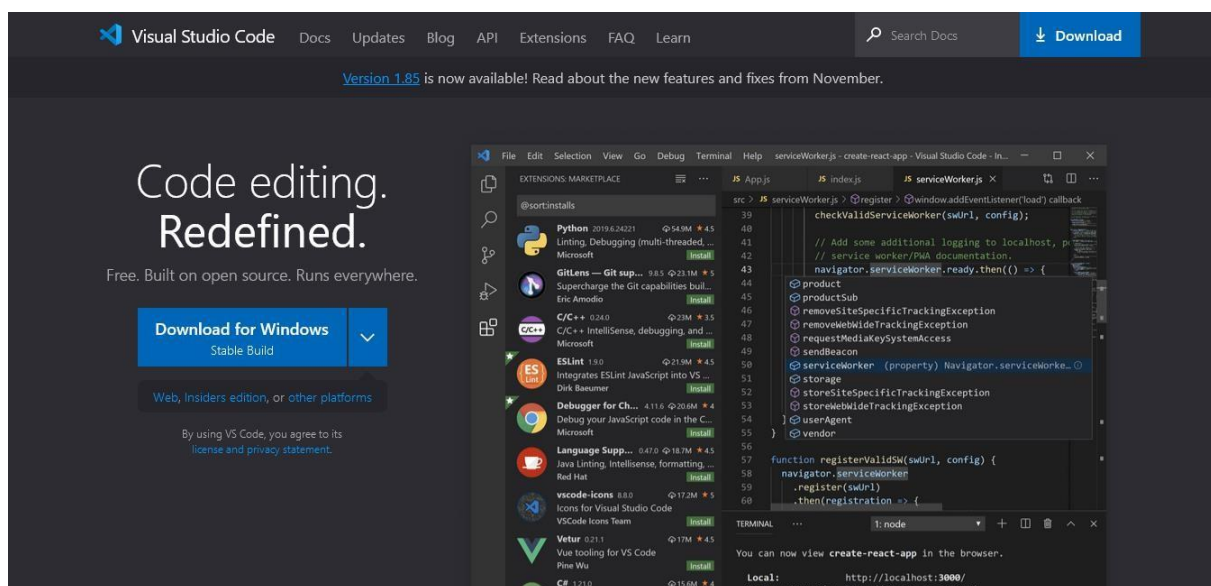
a) Python==3.10.12 :-

Python is an interpreted , high-level programming language . Its high-level built in data structures , combined with dynamic typing and dynamic binding , make it very attractive for Rapid Application Development , as well as for use as a scripting or glue language to connect existing components together. Its language concept and object-oriented approach aim to help programmers write clear, logical code for small and large-scale projects.



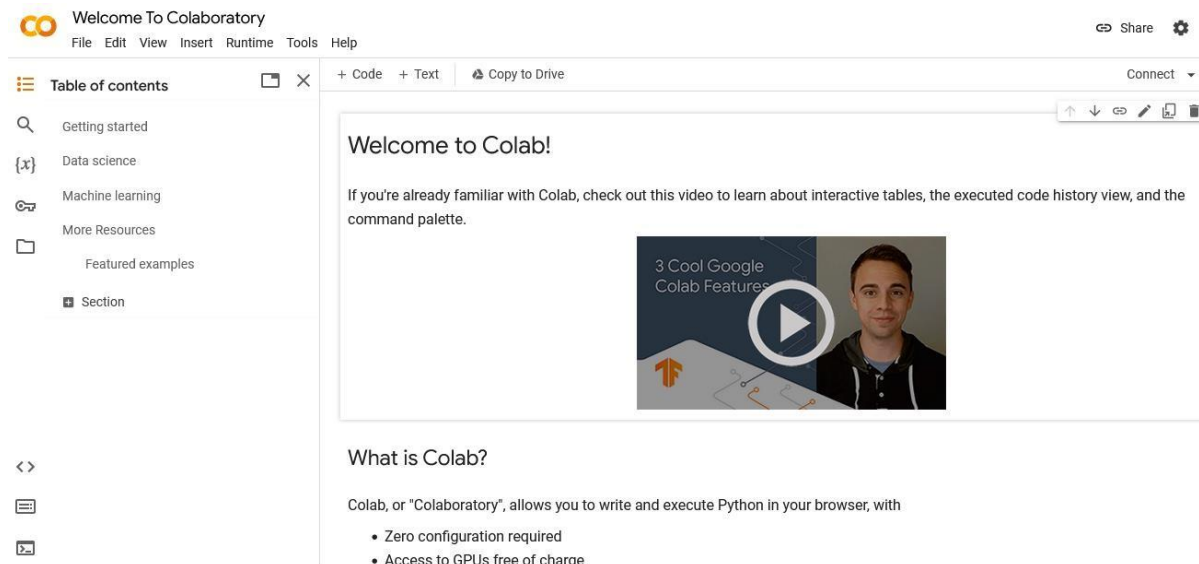
b) Visual Studio Code:-

Visual Studio Code (famously known as VS Code) is a free open source text editor by Microsoft. VS Code is available for Windows, Linux, and macOS. (Visual Studio Code - Code Editing. Redefined.)



c) Google Colab:-

Colaboratory, or “Colab” for short, is a product from Google Research. Colab allows anybody to write and execute arbitrary python code through the browser, and is especially well suited to machine learning, data analysis and education. More technically, Colab is a hosted Jupyter notebook service that requires no setup to use, while providing access free of charge to computing resources including GPUs. (colab.google,)



4 Libraries Configuration

1. Pandas
2. Numpy
3. Matplotlib
4. Seaborn
5. Plotly
6. Scikit-learn
7. Tensorflow
8. Flask
9. Keras

(*Best Python Libraries for Machine Learning - GeeksforGeeks, 2023.*)

5 Implementation

Implementation is divided into 3 modules which is mentioned below:

Module1: Visual Studio Code

First Loading Urls Dataset from each class, Extracting 19 Features and one Labelling Feature(target class) and last step is to save extracted features into csv file. i.e final_dataframe.csv.(URL 2016 | Datasets | Research | Canadian Institute for Cybersecurity | UNB, no date)

Module2: Next step is on google colab for training of model .

Saving the extracted features into csv file on drive followed by Importing & Installing all required libraries.

Dataset Loading using Pandas

Data Cleaning: null values are checked and unnecessary columns are dropped. We removed unnecessary data from the dataset, 2000 url we took from which 19+1 is feature extraction and 1 is label (target column). Then pre-processing is done where we are converting categorical data into numerical.

Data Visualization(EDA) : creating graphs and plots to understand data.

Splitting data (Features and target split in to x and y) and we also split data for training 90% and 10 % for testing, we always train model with a higher volume of data, as we have used 1000 of each benign and phishing URLs.

This data is used to train ML and DL models. (Fit trains data and predict from test data) Later Performance evaluation is done using the Confusion Matrix and Classification Report.

Module 3: Web application for detection of Phishing URL. There is web framework inside python here in this step two modules are merged. (Welcome to Flask — Flask Documentation (3.0.x),)

Step 1: Select the relevant dataset here we are using only two classes benign and phishing URLs.

Dataset Link :- <https://www.unb.ca/cic/datasets/url-2016.html>.

Step 2: After Installation of Visual Studio Code we have created new project and we have imported the necessary libraries for feature extraction and selected desired columns which are features mentioned in figure 1.

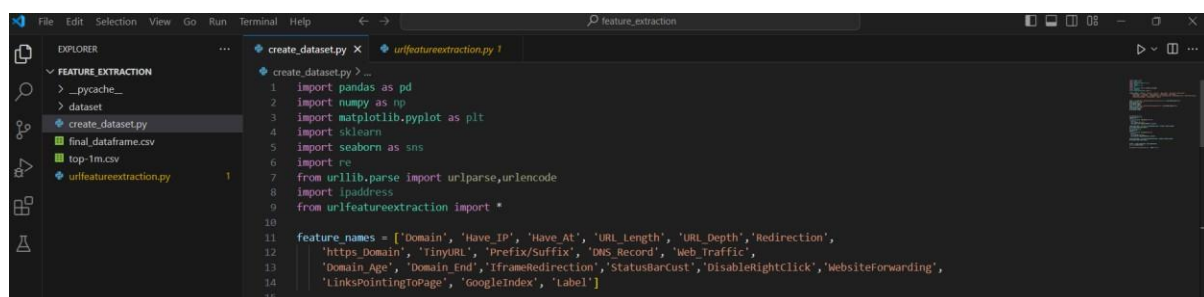


Fig. 1

Step 3: All the functions for feature extraction is in python where we are using regular expressions library which will extract features of NLP concept mentioned in figure 2 and 3.

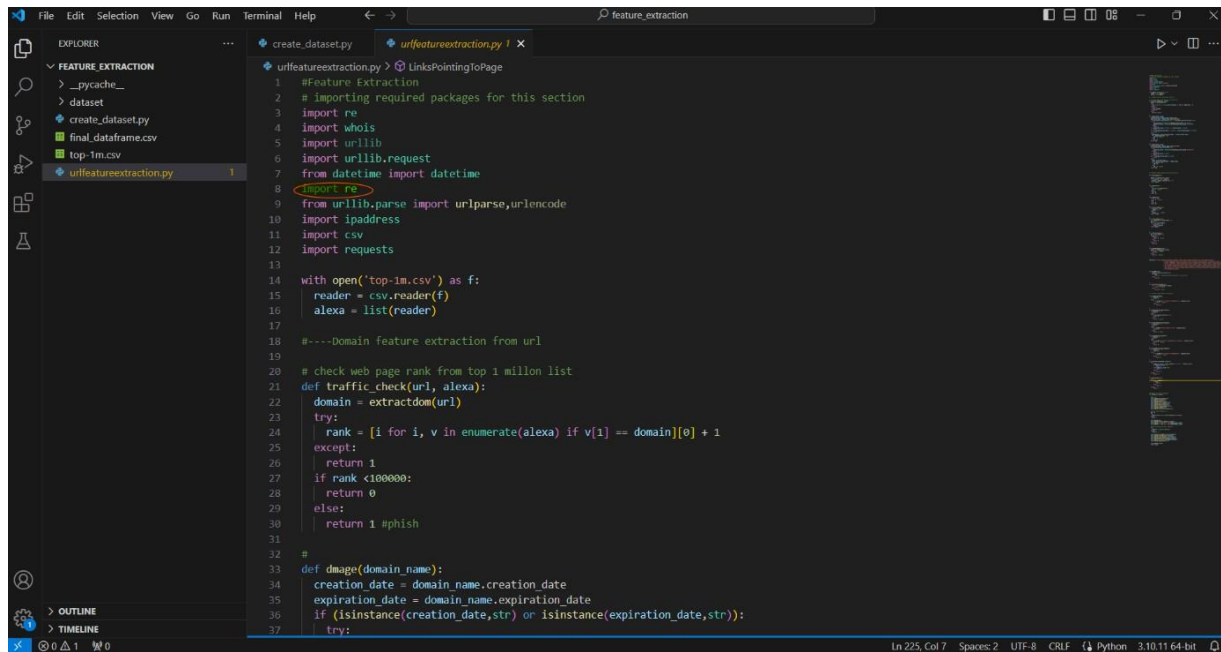


Fig. 2

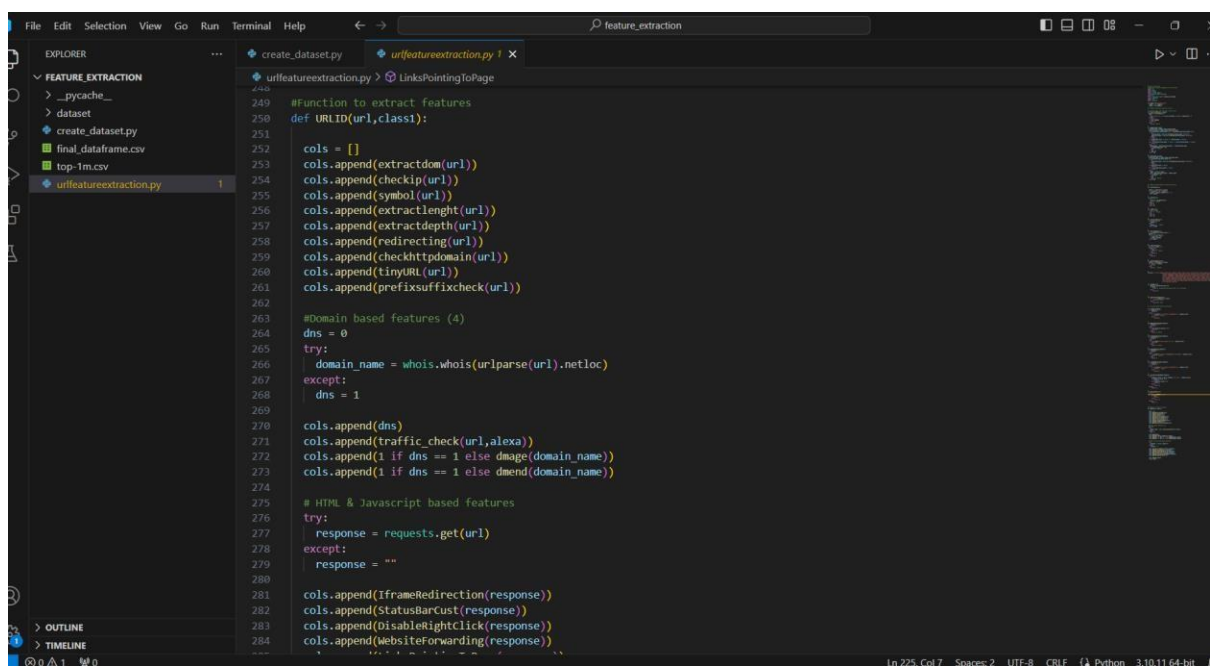


Fig.3

Step 4.After running the create_dataset ,it will start checking each URL and later it will save the final data as mentioned in below image.

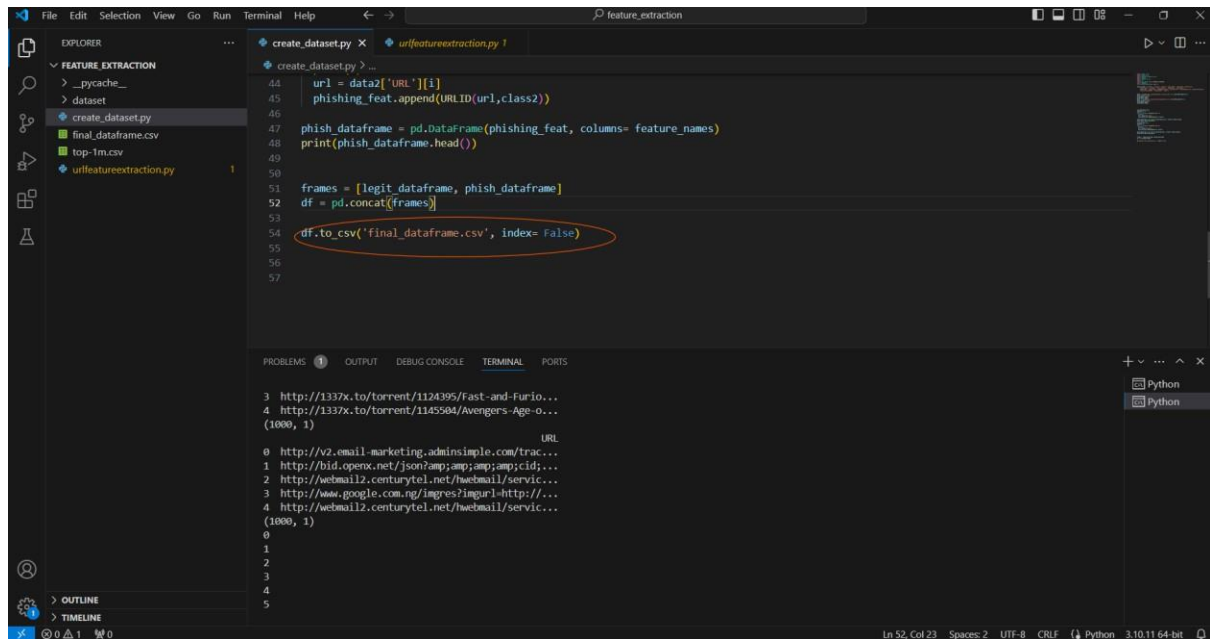
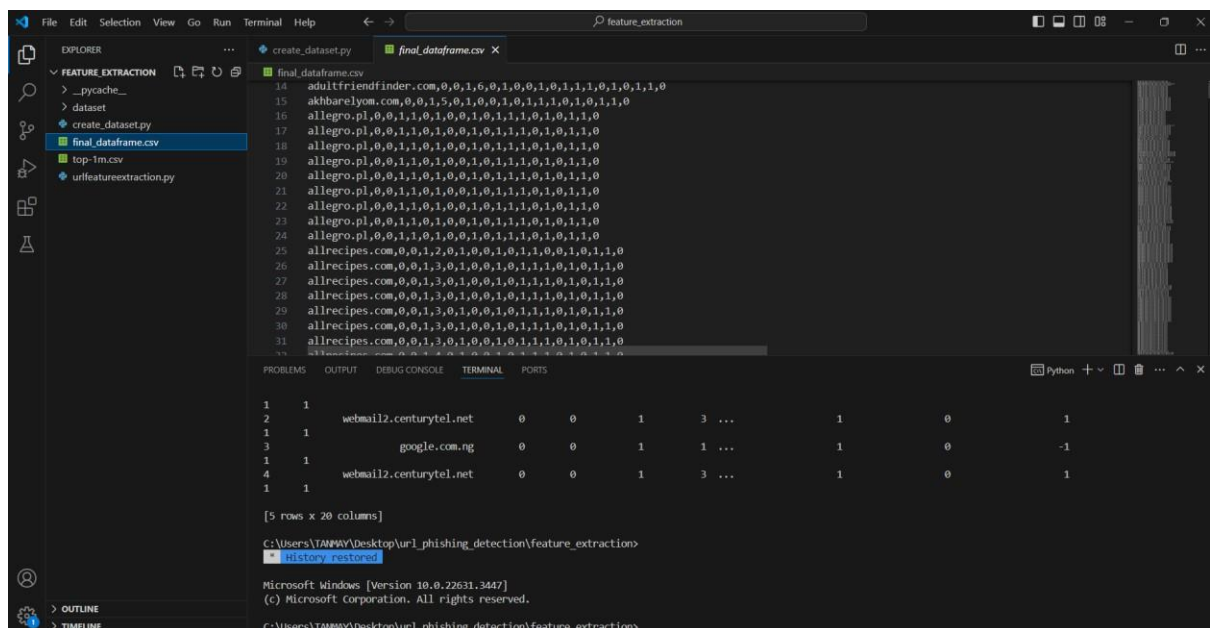


Fig. 4

In Below figure **Final feature extraction.csv** file of visual studio code is mentioned below. Which is final .csv file after feature extraction which will be used on google colab.



Step 5: All the below steps and image illustrates the work done in google colab which is explained in module 2 mentioned above. Like figure 5 states that we are Importing Libraries for visualization like plotly, sklearn, and seaborn. TensorFlow is for deep learning to speed up the process.

Code.ipynb

File Edit View Insert Runtime Tools Help Last saved at April 24

+ Code + Text

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

Mounted at /content/drive
```

Importing Libraries

```
import pickle
import numpy as np
import pandas as pd
import seaborn as sns
import keras.backend as K
import keras.backend as K
import plotly.express as px
from sklearn import ensemble
import matplotlib.pyplot as plt
import plotly.graph_objects as go
import plotly.figure_factory as ff
from tensorflow.keras.utils import to_categorical
from sklearn.model_selection import train_test_split
from sklearn.metrics import precision_recall_fscore_support
from tensorflow.keras.models import Sequential, ModelFromJSON
from sklearn.ensemble import RandomForestClassifier, ExtraTreesClassifier
from sklearn.preprocessing import LabelBinarizer, LabelEncoder, MinMaxScaler
from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
from tensorflow.keras.layers import Conv1D, Dropout, Flatten, Dense, Bidirectional, LSTM, Activation, Layer, MaxPooling1D
```

Fig. 5

Data Cleaning

```
[ ] #dropping unnecessary column
dataframe.drop(['Domain'], axis='columns', inplace=True)

[ ] dataframe.columns

Index(['Have_IP', 'Have_At', 'URL_Length', 'URL_Depth', 'Redirection',
      'https_Domain', 'TinyURL', 'Prefix/Suffix', 'DNS_Record', 'Web_Traffic',
      'Domain_Age', 'Domain_End', 'IframeRedirection', 'StatusBarCust',
      'DisableRightClick', 'WebsiteForwarding', 'LinksPointingToPage',
      'GoogleIndex', 'Label'],
      dtype='object')
```

Fig.6

```
[ ] #statistic of dataframe
dataframe.describe()
```

	Have_IP	Have_At	URL_Length	URL_Depth	Redirection	https_Domain	TinyURL	Prefix/Suffix	DNS_Record	Web_Traffic	Domain_Age	Domain_End	IframeRedirection	StatusBarCust
count	2000.0	2000.000000	2000.0	2000.000000	2000.000000	2000.0	2000.000000	2000.000000	2000.0	2000.000000	2000.0	2000.0	2000.000000	2000.000000
mean	0.0	0.005500	1.0	2.782000	0.01250	1.0	0.042500	0.337000	1.0	0.596500	1.0	1.0	0.535000	0.459000
std	0.0	0.073976	0.0	2.101591	0.11113	0.0	0.201777	0.472803	0.0	0.490722	0.0	0.0	0.498898	0.498898
min	0.0	0.000000	1.0	0.000000	0.00000	1.0	0.000000	0.000000	1.0	0.000000	1.0	1.0	0.000000	0.000000
25%	0.0	0.000000	1.0	1.000000	0.00000	1.0	0.000000	0.000000	1.0	0.000000	1.0	1.0	0.000000	0.000000
50%	0.0	0.000000	1.0	2.000000	0.00000	1.0	0.000000	0.000000	1.0	1.000000	1.0	1.0	1.000000	0.000000
75%	0.0	0.000000	1.0	4.000000	0.00000	1.0	0.000000	1.000000	1.0	1.000000	1.0	1.0	1.000000	1.000000
max	0.0	1.000000	1.0	16.000000	1.00000	1.0	1.000000	1.000000	1.0	1.000000	1.0	1.0	1.000000	1.000000

Fig. 7

CountPlot of Url_Depth Column

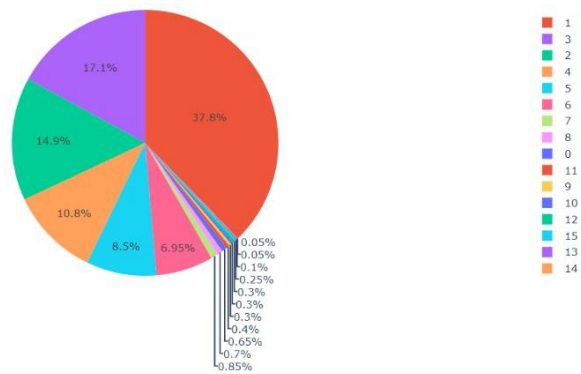


Fig . 8

```
[ ] #model feature importance
feature_name = X_train.columns.values
model = ensemble.ExtraTreesClassifier()
model.fit(X_train,y_train)
#plot imp
importance = model.feature_importances_
std = np.std([tree.feature_importances_ for tree in model.estimators_],axis=0)
indices = np.argsort(importance)[::-1][:5]
plt.figure(figsize=(7,7))
plt.title("Top 5 Important Features")
plt.bar(range(len(indices)), importance[indices], color="lightgreen")
plt.xticks(range(len(indices)), feature_name[indices], rotation='vertical')
plt.xlim([-1, len(indices)])
plt.show()
```

Top 5 Important Features

Fig . 9

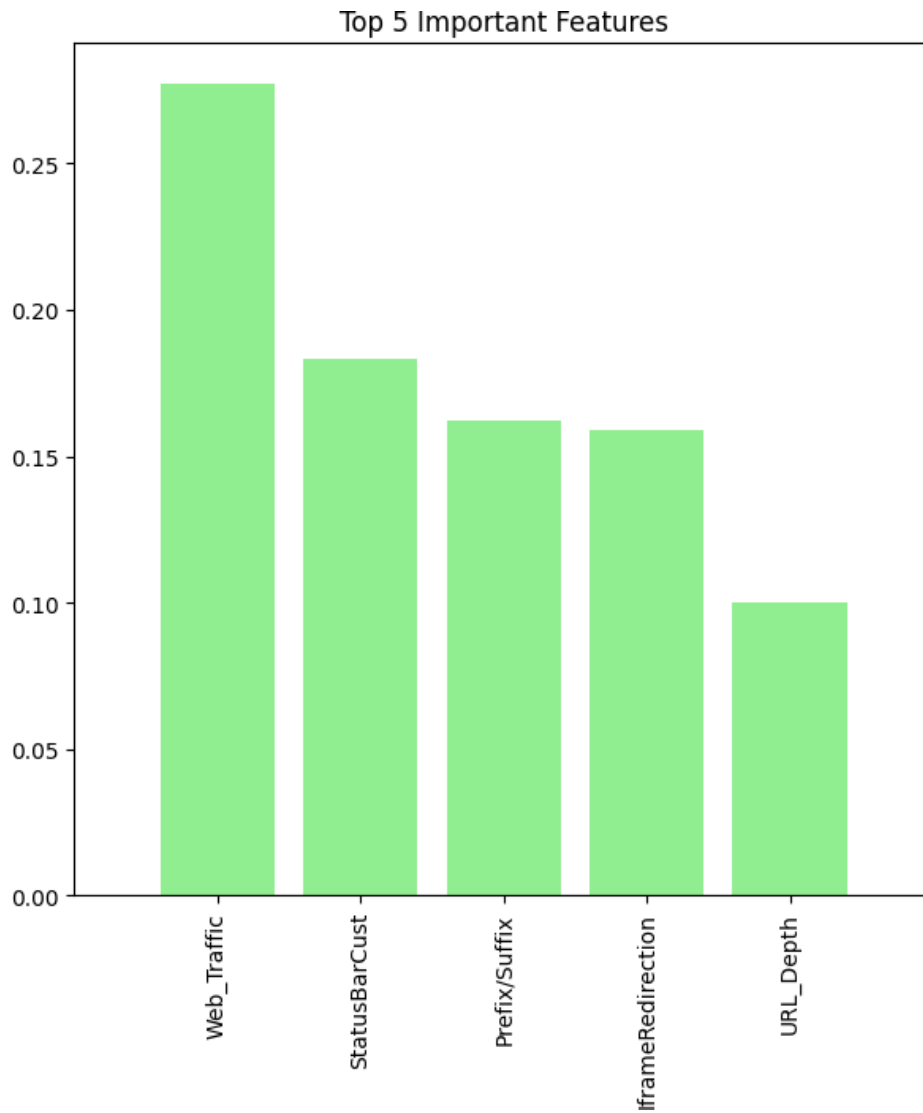
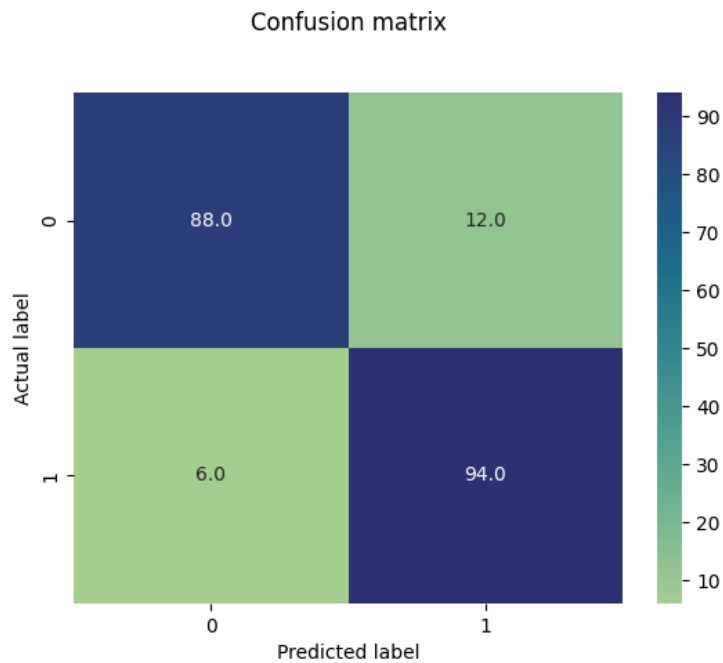


Fig. 10

Step 6:figure 9 and 10 tells us that we what are the impactful fature for detecting phishing URL which are correlated from 19 features.

Step 7: We get 91% accuracy using CNN model with attention mechanism which is evaluated using confusion matrix and classification report mentioned in below image.



```
[ ] print(classification_report(yorg, ypred))
```

	precision	recall	f1-score	support
0	0.94	0.88	0.91	100
1	0.89	0.94	0.91	100
accuracy			0.91	200
macro avg	0.91	0.91	0.91	200
weighted avg	0.91	0.91	0.91	200

Step 8. We have analysed from the evaluation that CNN with attention mechanism is best model and we save it for further integration with flask web application.

```

history = model.fit(X_train, y_train, batch_size=32, epochs=5, validation_data=(X_test, y_test), verbose=1)

Epoch 1/5
57/57 [=====] - 3s 14ms/step - loss: 0.6467 - accuracy: 0.6689 - val_loss: 0.4941 - val_accuracy: 0.8480
Epoch 2/5
57/57 [=====] - 0s 8ms/step - loss: 0.3926 - accuracy: 0.8350 - val_loss: 0.2541 - val_accuracy: 0.8800
Epoch 3/5
57/57 [=====] - 1s 11ms/step - loss: 0.2793 - accuracy: 0.8744 - val_loss: 0.2563 - val_accuracy: 0.8950
Epoch 4/5
57/57 [=====] - 1s 13ms/step - loss: 0.2641 - accuracy: 0.8894 - val_loss: 0.2400 - val_accuracy: 0.8950
Epoch 5/5
57/57 [=====] - 0s 8ms/step - loss: 0.2565 - accuracy: 0.8922 - val_loss: 0.1936 - val_accuracy: 0.9100

[ ] ypred = model.predict(X_test)
ypred = np.argmax(ypred,axis=1)
yorg = np.argmax(y_test,axis=1)

7/7 [=====] - 0s 3ms/step

[ ] #Train and Validation Accuracy
plt.plot(history.history['accuracy'], 'b', marker='*', label='Train Accuracy')
plt.plot(history.history['val_accuracy'], 'g', marker='*', label='Val Accuracy')
plt.title('Train and Val Accuracy')
plt.legend()
plt.figure()

#Train and Validation Loss
plt.plot(history.history['loss'], 'b', marker='*', label='Train Loss')
plt.plot(history.history['val_loss'], 'g', marker='*', label='Val Loss')
plt.title('Train and Val Loss')
plt.legend()
plt.show()

```

Step 9: Web Application for testing the phishing URLs is run on visual studio code where we load the best model and integrate all 3 modules. Here in assets we have screenshot of web application, models have best model selected in google colab i.e. CNN with attention mechanism. Which focuses constantly on impactful layers features. Static is having design.

All this is shared in artefacts upload.

```

app.py 2 X
app.py > ...
16 class Attention(Layer):
31
32 def compute_output_shape(self, input_shape):
33     return (input_shape[0], input_shape[-1])
34
35 def get_config(self):
36     return super(Attention, self).get_config()
37
38
39 # load json and create model
40 json_file = open('./Models/bestmodel.json', 'r')
41 loaded_model_json = json_file.read()
42 json_file.close()
43 loaded_model = model_from_json(loaded_model_json, custom_objects={"Attention": Attention})
44 # load weights into new model
45 loaded_model.load_weights("./Models/bestmodel.h5")
46 print("Loaded model from disk")
47
48 #
49 col = ['Domain', 'Have_IP', 'Have_AT', 'URL_Length', 'URL_Depth', 'Redirection', 'https_Domain', 'TinyURL', 'Prefix/Suffix', 'DNS_Record',
50
51
52
53
54 @app.route('/')
55 def home():
56     return render_template('index.html')
57
58 @app.route('/index')
59 def index():
60     return render_template('index.html')
61
62
63 @app.route('/service')
64 def service():
65     return render_template('detection.html')
66

```

Step 10: Output of the web page once we receive the IP after running app.py.

```

class Attention(Layer):
    def compute_output_shape(self, input_shape):
        return (input_shape[0], input_shape[-1])

    def get_config(self):
        return super(Attention, self).get_config()

# load json and create model
json_file = open('./Models/bestmodel.json', 'r')
loaded_model_json = json_file.read()
json_file.close()
loaded_model = model_from_json(loaded_model_json, custom_objects={"Attention": Attention})
# load weights into new model
loaded_model.load_weights("./Models/bestmodel.h5")
print("loaded model from disk")

```

```

Microsoft Windows [Version 10.0.22631.3447]
(c) Microsoft Corporation. All rights reserved.

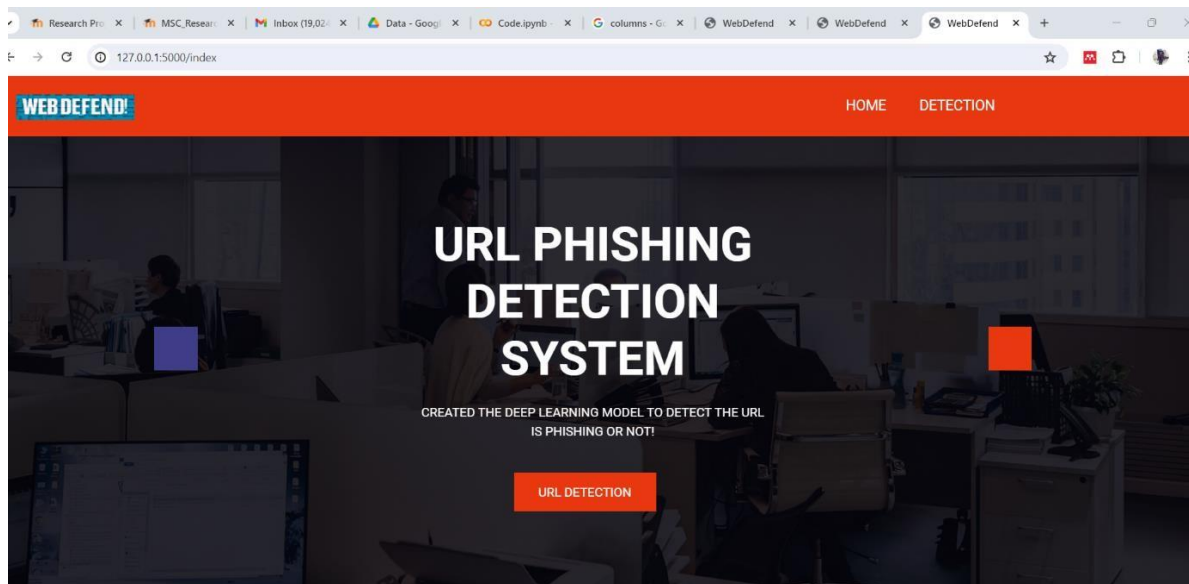
C:\Users\TANMAY\Desktop\url_phishing_detection\flaskwebapp>"C:/Program Files/Python310/python.exe" c:/Users/TANMAY/Desktop/url_phishing_detection/flaskweb
app/app.py
2024-04-25 12:31:11.370698: I tensorflow/core/platform/cpu_feature_guard.cc:182] This TensorFlow binary is optimized to use available CPU instructions in
performance-critical operations.
To enable the following instructions: SSE SSE2 SSE3 SSE4.1 SSE4.2 AVX AVX2 AVX512F AVX512_VNNI FMA, in other operations, rebuild TensorFlow with the appro
priate compiler flags.
loaded model from disk
* Serving Flask app 'app'
* Debug mode: off
INFO:werkzeug:WARNING: This is a development server. Do not use it in a production deployment. Use a production WSGI server instead.
* Running on http://127.0.0.1:5000
INFO:werkzeug:Press CTRL+C to quit

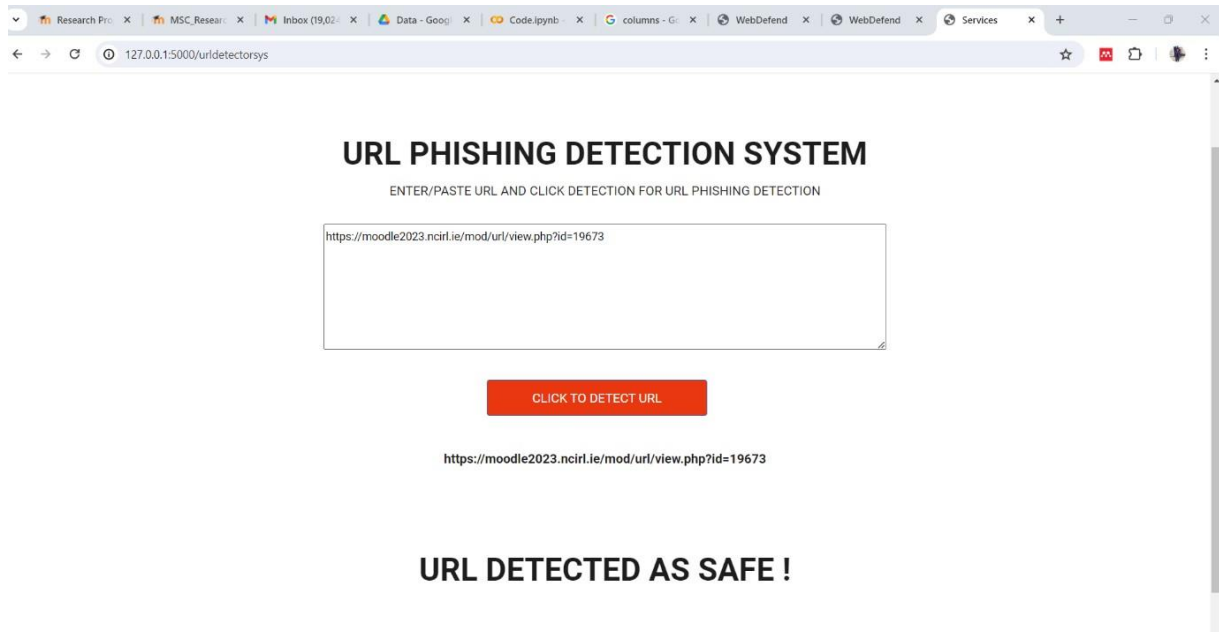
```

Step 11:

Copy URL from data set to identify whether its safe or phishing.

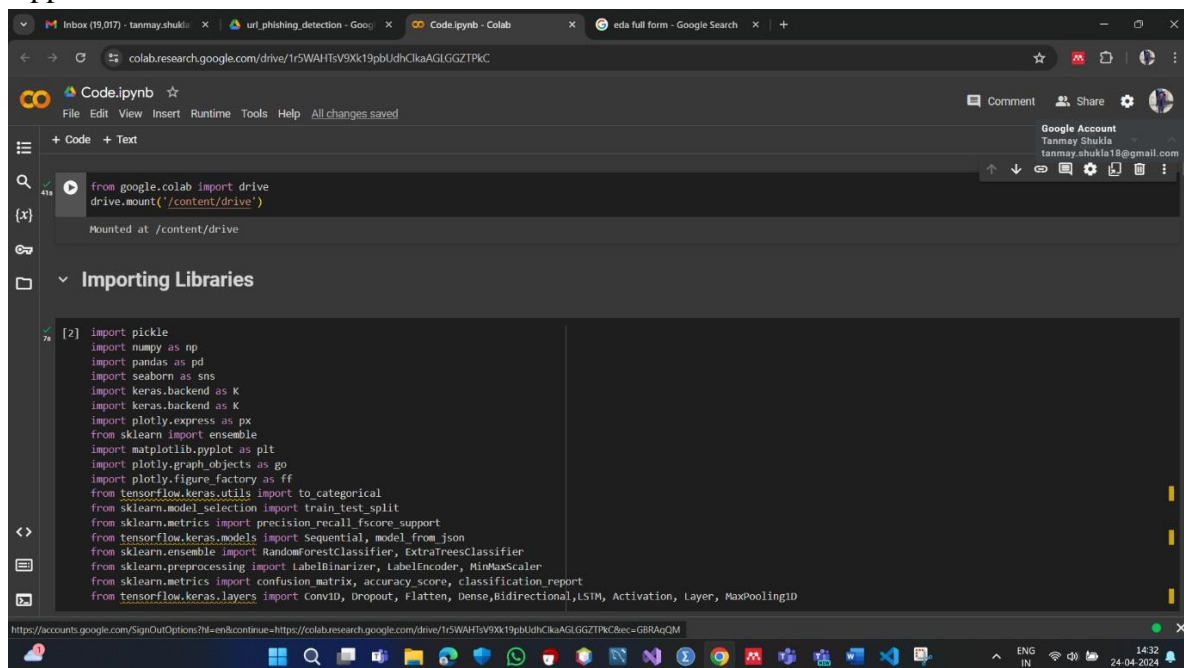
Below figure represents GUI of this App.





The above figure represents that the URL is safe.

Appendix:



Code.ipynb

Data Cleaning

```
[7] #dropping unnecessary column
dataframe.drop(['Domain'], axis='columns', inplace=True)
```

```
[8] dataframe.columns
```

```
Index(['Have_IP', 'Have_At', 'URL_Length', 'URL_Depth', 'Redirection',
       'https_domain', 'tinyURL', 'Prefix/Suffix', 'DNS_Record', 'Web_Traffic',
       'Domain_Age', 'Domain_End', 'IframeRedirection', 'StatusBarCust',
       'DisableRightClick', 'WebsiteForwarding', 'LinksPointingToPage',
       'GoogleIndex', 'Label'],
      dtype='object')
```

```
[9] #checking for null values
dataframe.isna().sum()
```

```
Have_IP      0
Have_At      0
URL_Length   0
URL_Depth    0
Redirection   0
https_Domain  0
TinyURL      0
Prefix/Suffix 0
DNS_Record   0
Web_Traffic  0
Domain_Age   0
```

Connected to Python 3 Google Compute Engine backend

Code.ipynb

```
Domain_End      0
IframeRedirection 0
StatusBarCust    0
DisableRightClick 0
WebsiteForwarding 0
LinksPointingToPage 0
GoogleIndex     0
Label           0
dtype: int64
```

```
#statistic of dataframe
dataframe.describe()
```

	Have_IP	Have_At	URL_Length	URL_Depth	Redirection	https_Domain	TinyURL	Prefix/Suffix	DNS_Record	Web_Traffic	Domain_Age	Domain_End	IframeRedirection	StatusBar
count	2000.0	2000.000000	2000.0	2000.000000	2000.000000	2000.0	2000.000000	2000.000000	2000.0	2000.000000	2000.0	2000.0	2000.000000	2000.00
mean	0.0	0.005500	1.0	2.782000	0.01250	1.0	0.042500	0.337000	1.0	0.596500	1.0	1.0	0.535000	0.48
std	0.0	0.073976	0.0	2.101591	0.11113	0.0	0.201777	0.472803	0.0	0.480722	0.0	0.0	0.498898	0.49
min	0.0	0.000000	1.0	0.000000	0.00000	1.0	0.000000	0.000000	1.0	0.000000	1.0	1.0	0.000000	0.00
25%	0.0	0.000000	1.0	1.000000	0.00000	1.0	0.000000	0.000000	1.0	0.000000	1.0	1.0	0.000000	0.00
50%	0.0	0.000000	1.0	2.000000	0.00000	1.0	0.000000	0.000000	1.0	1.000000	1.0	1.0	1.000000	0.00
75%	0.0	0.000000	1.0	4.000000	0.00000	1.0	0.000000	1.000000	1.0	1.000000	1.0	1.0	1.000000	1.00
max	0.0	1.000000	1.0	16.000000	1.00000	1.0	1.000000	1.000000	1.0	1.000000	1.0	1.0	1.000000	1.00

Connected to Python 3 Google Compute Engine backend

Code.ipynb

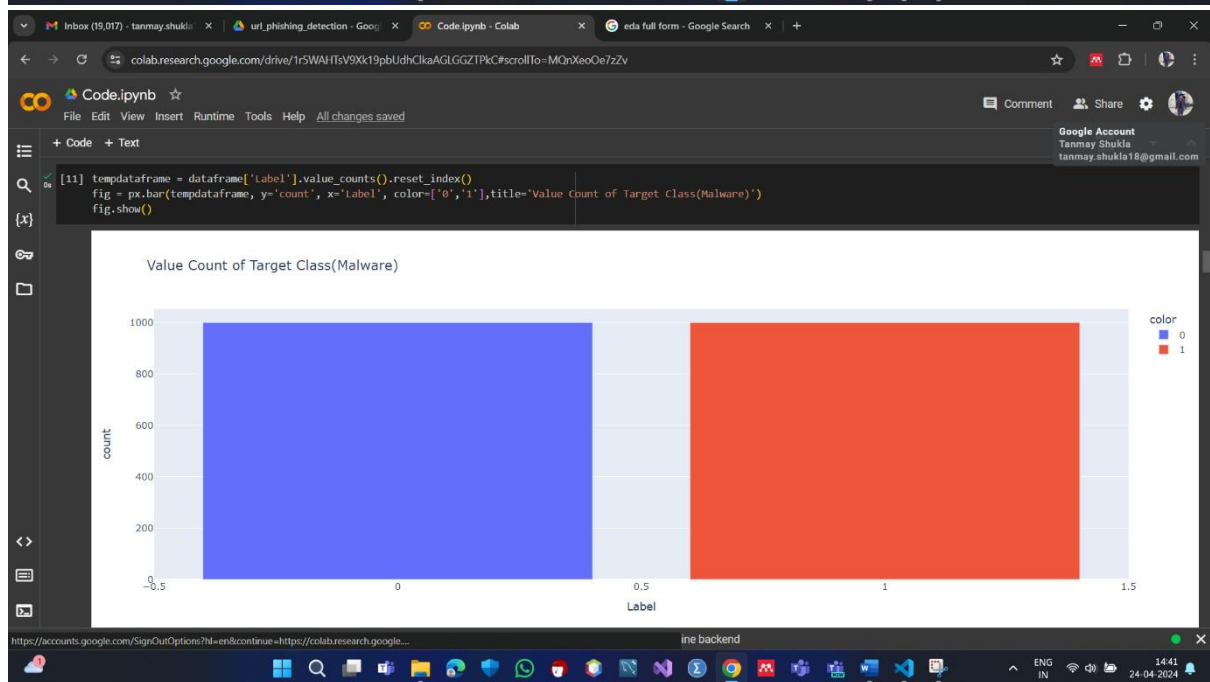
```
#statistic of dataframe
dataframe.describe()
```

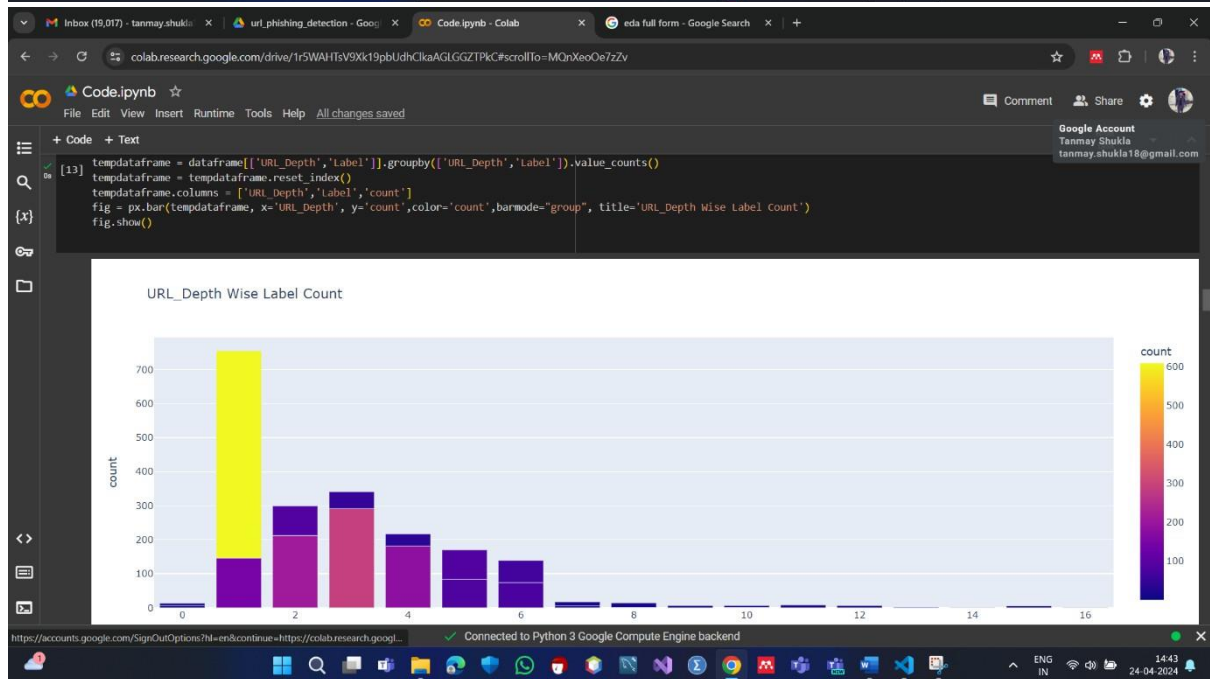
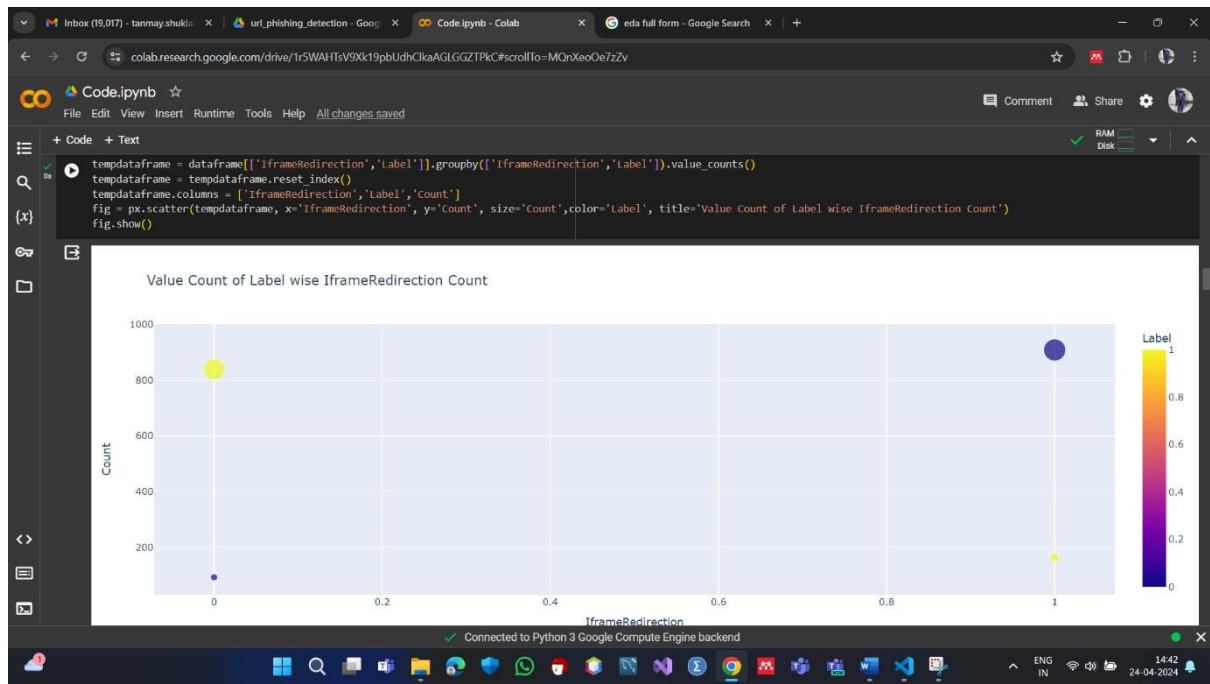
refix/Suffix	DNS_Record	Web_Traffic	Domain_Age	Domain_End	IframeRedirection	StatusBarCust	DisablerightClick	WebsiteForwarding	LinksPointingToPage	GoogleIndex	Label
2000.000000	2000.0	2000.000000	2000.0	2000.0	2000.000000	2000.000000	2000.0	2000.000000	2000.000000	2000.0	2000.000000
0.337000	1.0	0.596500	1.0	1.0	0.535000	0.459000	1.0	0.508500	-0.679500	1.0	0.500000
0.472803	0.0	0.490722	0.0	0.0	0.498898	0.488441	0.0	0.500053	0.709247	0.0	0.500125
0.000000	1.0	0.000000	1.0	1.0	0.000000	0.000000	1.0	0.000000	-1.000000	1.0	0.000000
0.000000	1.0	0.000000	1.0	1.0	0.000000	0.000000	1.0	0.000000	-1.000000	1.0	0.000000
0.000000	1.0	1.000000	1.0	1.0	1.000000	0.000000	1.0	1.000000	-1.000000	1.0	0.500000
1.000000	1.0	1.000000	1.0	1.0	1.000000	1.000000	1.0	1.000000	-1.000000	1.0	1.000000
1.000000	1.0	1.000000	1.0	1.0	1.000000	1.000000	1.0	1.000000	1.000000	1.0	1.000000

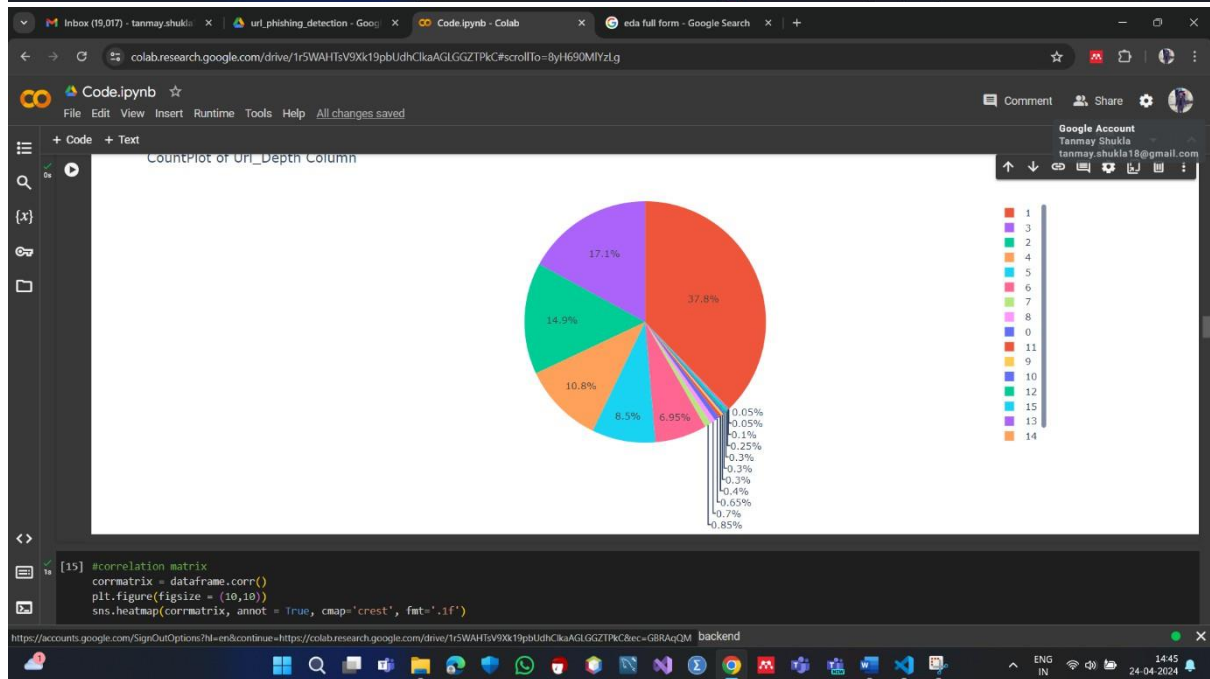
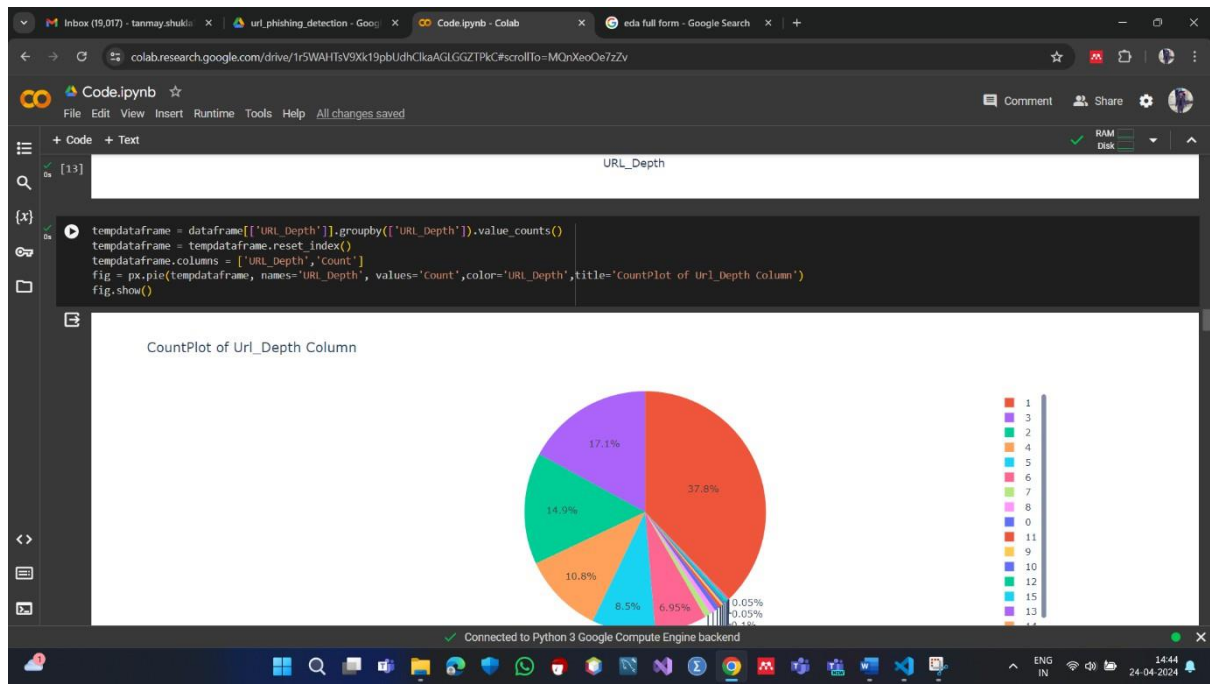
Exploratory Data Analysis (EDA)

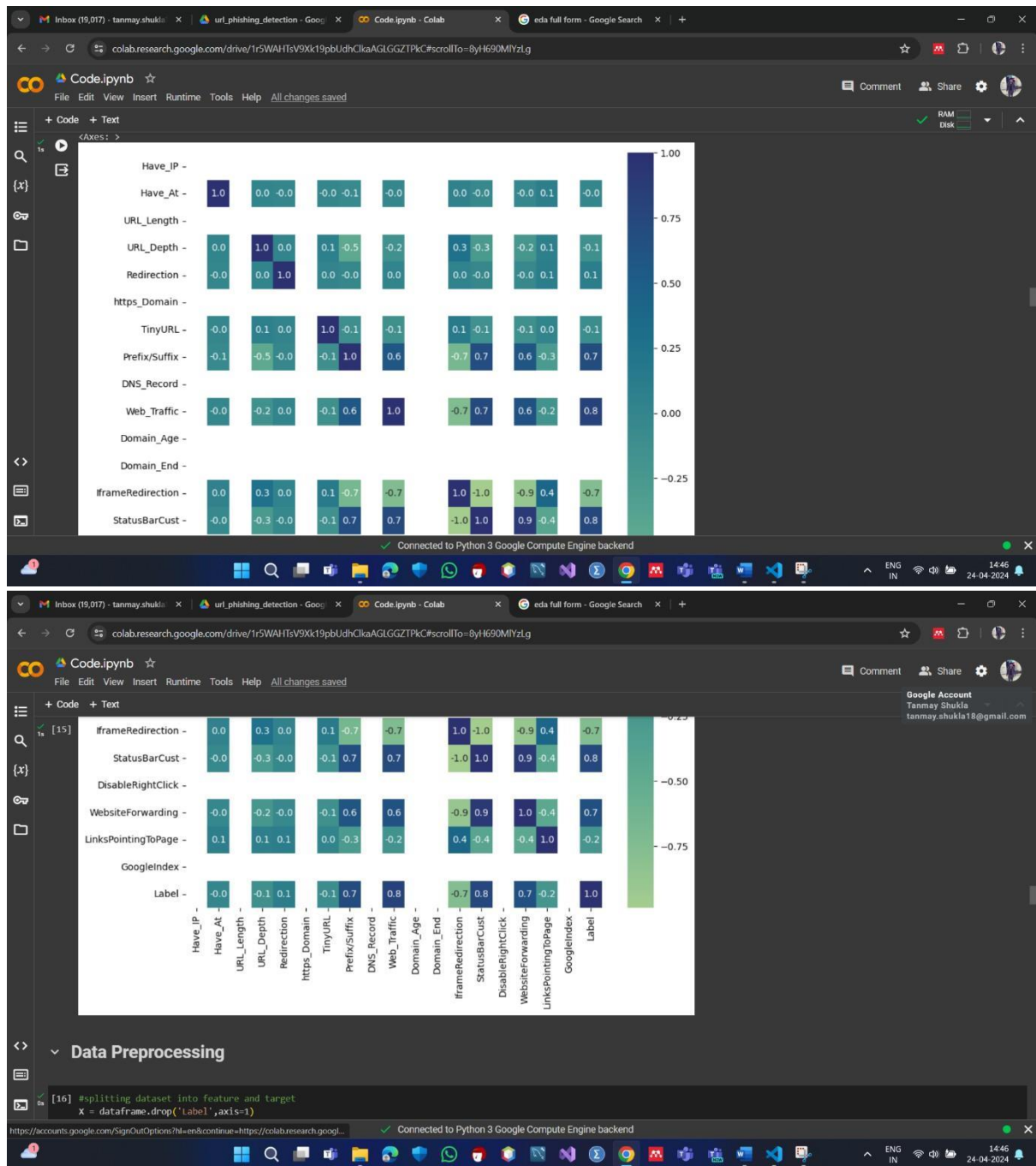
```
[11] tempdataframe = dataframe['Label'].value_counts().reset_index()
fig = px.bar(tempdataframe, y='count', x='Label', color=['0','1'],title='Value Count of Target Class(Malware)')
fig.show()
```

Value Count of Target Class(Malware)









Code.ipynb

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

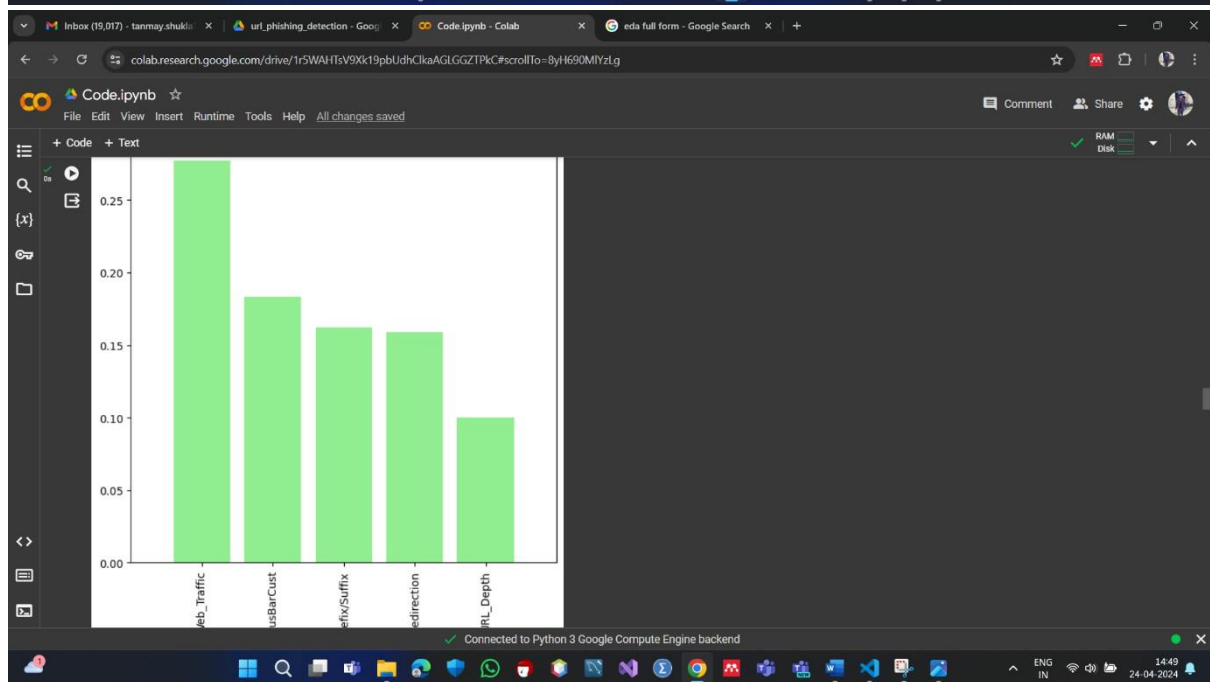
Data Preprocessing

```
[16] #splitting dataset into feature and target
X = dataframe.drop('label',axis=1)
Y = dataframe['label']

[17] #splitting dataset into train and test with ratio of 90:10
X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.1, stratify=Y, shuffle=True)

[18] #model feature importance
feature_name = X_train.columns.values
model = ensemble.ExtraTreesClassifier()
model.fit(X_train,y_train)
#plot imp
importance = model.feature_importances_
std = np.std([tree.feature_importances_ for tree in model.estimators_],axis=0)
indices = np.argsort(importance)[::-1][:5]
plt.figure(figsize=(7,7))
plt.title("Top 5 Important Features")
plt.bar(range(len(indices)), importance[indices], color="lightgreen")
plt.xticks(range(len(indices)), feature_name[indices], rotation='vertical')
plt.xlim([-1, len(indices)])
plt.show()
```

Connected to Python 3 Google Compute Engine backend



Code.ipynb

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

Web_Traffic
StatusBarCust
Prefix/Suffix
Iframe/redirect
URL_Depth

ML Models

Random Forest Classifier

```
[19] rfc = RandomForestClassifier(n_estimators=1, max_depth=2, criterion='log_loss', max_leaf_nodes=2, max_features='log2')
      rfc.fit(X_train,y_train)
      ypred = rfc.predict(X_test)
      print("Accuracy Score: ", accuracy_score(y_test,ypred))
```

Accuracy Score: 0.825

```
[20] matrix = confusion_matrix(y_test, ypred)
      sns.heatmap(pd.DataFrame(matrix), annot=True, cmap="crest", fsize=.1f)
      plt.title('Confusion matrix', y=1.1)
      plt.ylabel('Actual label')
      plt.xlabel('Predicted label')
      plt.show()
```

Connected to Python 3 Google Compute Engine backend

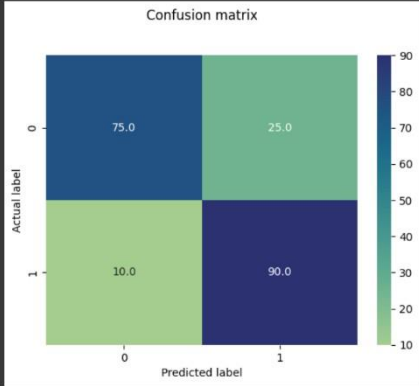
Code.ipynb

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

Google Account
Tanmay Shukla
tanmay.shukla18@gmail.com

[20]



Confusion matrix

Actual label

Predicted label

```
[21] print(classification_report(y_test, ypred))
```

precision recall f1-score support

Connected to Python 3 Google Compute Engine backend

Code.ipynb

```
print(classification_report(y_test, ypred))
```

	precision	recall	f1-score	support
0	0.88	0.75	0.81	100
1	0.78	0.90	0.84	100
accuracy			0.82	200
macro avg	0.83	0.82	0.82	200
weighted avg	0.83	0.82	0.82	200

Extra Tree Classifier

```
[22] etc = ExtraTreesClassifier(n_estimators=1, max_depth=2, criterion='log_loss', max_leaf_nodes=2, max_features='log2')
etc.fit(X_train, y_train)
ypred = etc.predict(X_test)
print("Accuracy Score: ", accuracy_score(y_test, ypred))
```

Accuracy Score: 0.885

```
[23] matrix = confusion_matrix(y_test, ypred)
sns.heatmap(pd.DataFrame(matrix), annot=True, cmap='crest', fsize=10)
plt.title('Confusion matrix', y=1.1)
plt.ylabel('Actual label')
plt.xlabel('Predicted label')
plt.show()
```

Connected to Python 3 Google Compute Engine backend

Code.ipynb

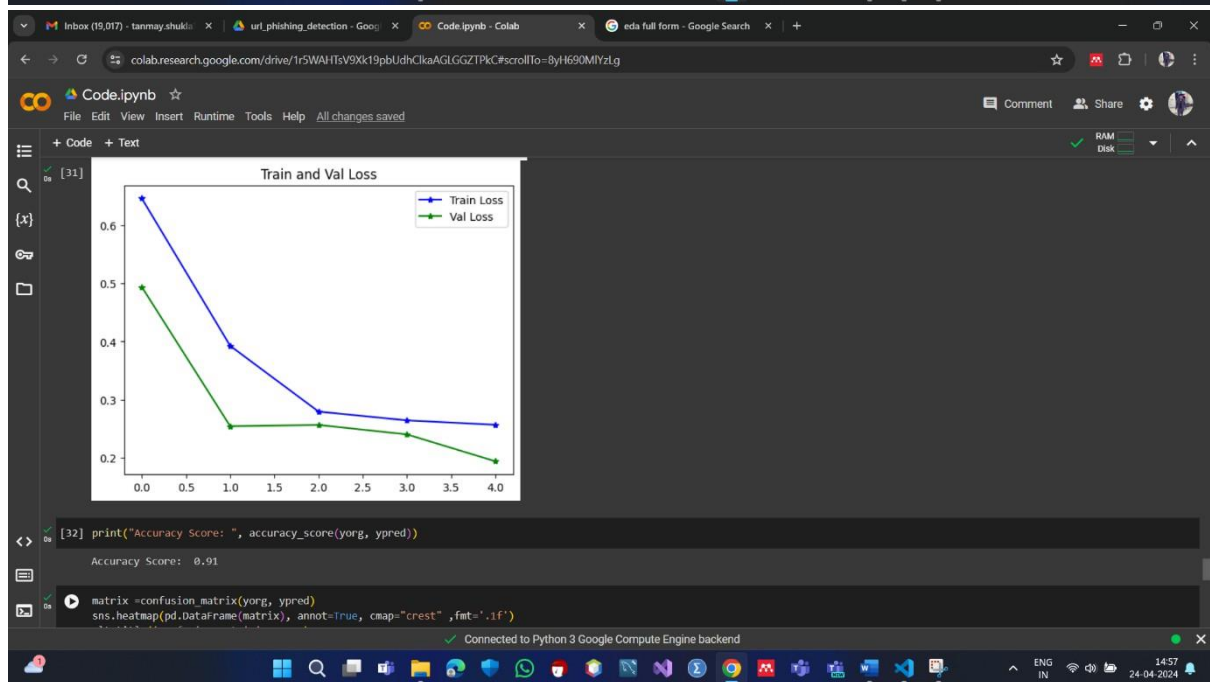
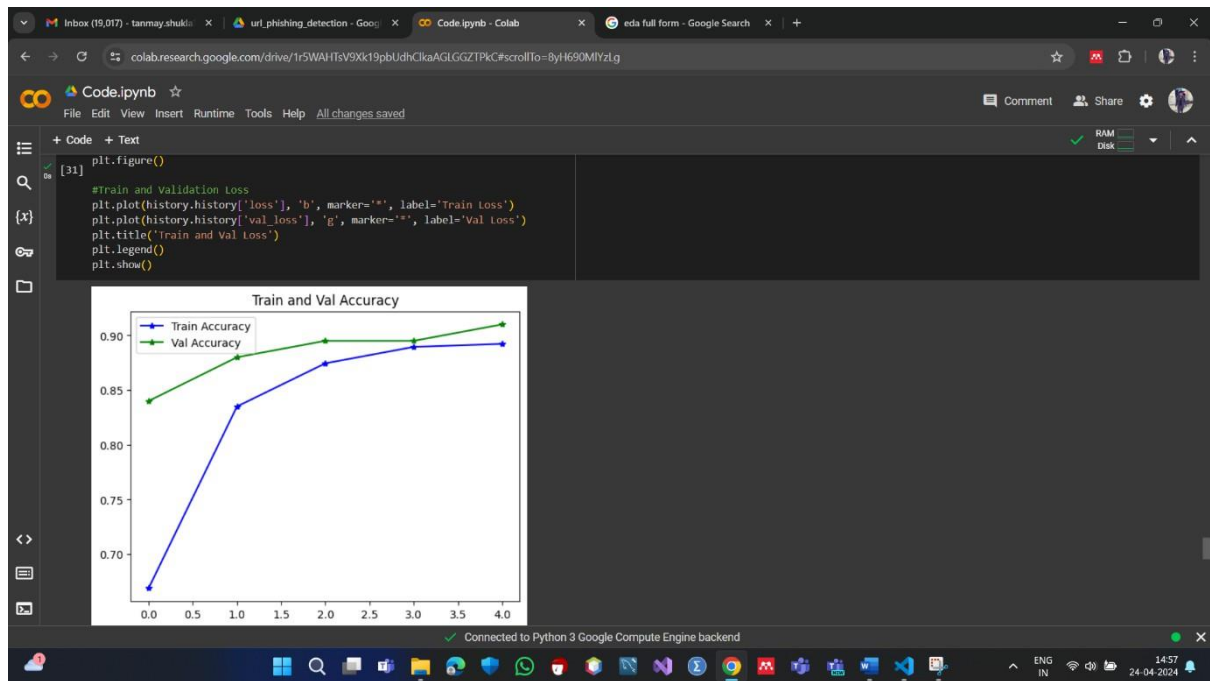
Confusion matrix

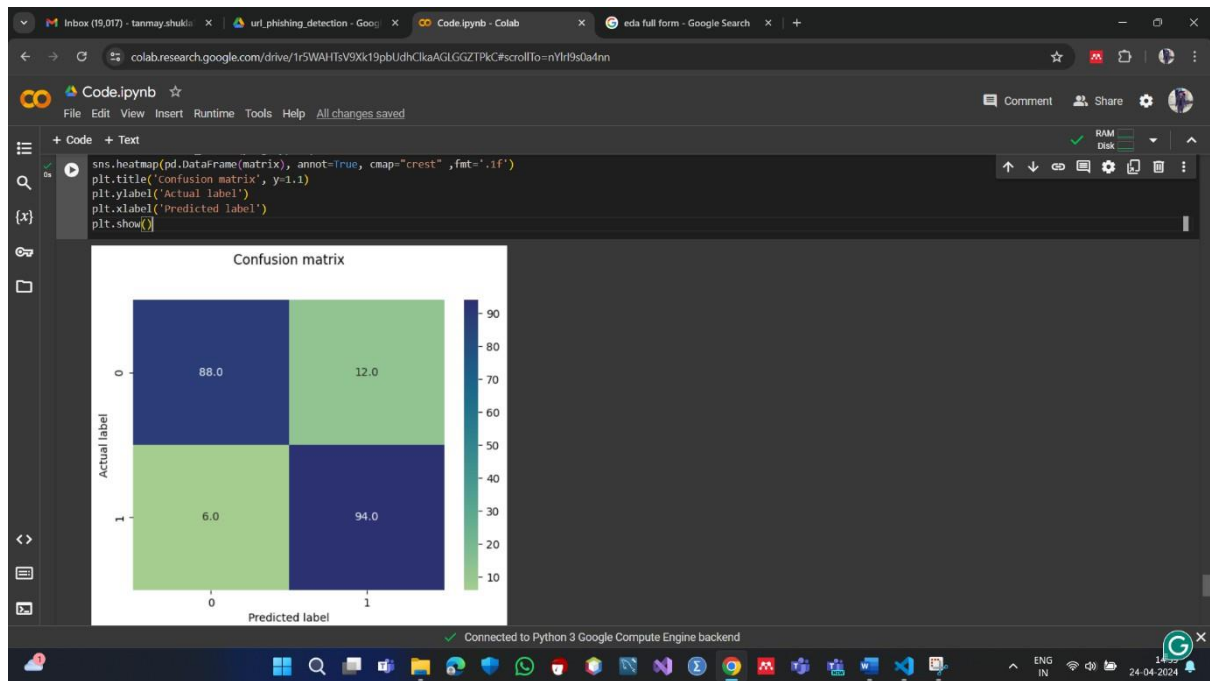
	0	1
0	77.0	23.0
1	0.0	100.0

```
[24] print(classification_report(y_test, ypred))
```

	precision	recall	f1-score	support
0	0.88	0.75	0.81	100
1	0.78	0.90	0.84	100
accuracy			0.82	200
macro avg	0.83	0.82	0.82	200
weighted avg	0.83	0.82	0.82	200

Connected to Python 3 Google Compute Engine backend





Code.ipynb

```

print(classification_report(yorg, ypred))

```

```

[[
      precision    recall  f1-score   support

     0.          0.94      0.88      0.91         100
     1.          0.89      0.94      0.91         100

 accuracy          0.91      0.91      0.91         200
 macro avg          0.91      0.91      0.91         200
 weighted avg          0.91      0.91      0.91         200
]]

```

```

[35] """model_json = model.to_json()
with open("/content/drive/MyDrive/url_phishing_detection/models/bestmodel.json", "w") as json_file:
    json_file.write(model_json)
# serialize weights to HDF5
model.save_weights("/content/drive/MyDrive/url_phishing_detection/models/bestmodel.h5")
print("Saved model to disk")"""

```

```

'model_json = model.to_json()\nwith open("/content/drive/MyDrive/url_phishing_detection/models/bestmodel.json", "w") as json_file:\n    json_file.write(model_json)\n# serialize weights to HDF5\nmodel.save_weights("/content/drive/MyDrive/url_phishing_detection/models/bestmodel.h5")\nprint("Saved model to disk")'

```

Connected to Python 3 Google Compute Engine backend

```

46
47 phishing_dataframe = pd.DataFrame(phishing_feat, columns= feature_names)
48 print(phishing_dataframe.head())
49
50
51 frames = [legit_dataframe, phishing_dataframe]
52 df = pd.concat(frames)
53
54 df.to_csv('final_dataframe.csv', index= False)
55
56
57

```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

Python + -

Ln 40, Col 19 Spaces: 2 UTF-8 CRLF Python 3.10.11 64-bit

14:24 24-04-2024

```

1 create_dataset.py > ...
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import sklearn
6 import seaborn as sns
7 import re
8 from urllib.parse import urlparse,urllencode
9 import ipaddress
10 from urlfeatureextraction import *
11
12 feature_names = ['Domain', 'Have_IP', 'Have_At', 'URL_Length', 'URL_Depth','Redirection',
13                 'https_Domain', 'TinyURL', 'Prefix/Suffix', 'DNS_Record', 'Web_Traffic',
14                 'Domain_Age', 'Domain_End', 'IframeRedirection', 'StatusBarCust','DisableRightClick','WebsiteForwarding',
15                 'LinksPointingToPage', 'GoogleIndex', 'Label']
16
17 data1 = pd.read_csv('./dataset/Benign_list_big_final.csv',nrows=1000,header=None)
18 data1.columns = ['URL']
19

```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

Python + -

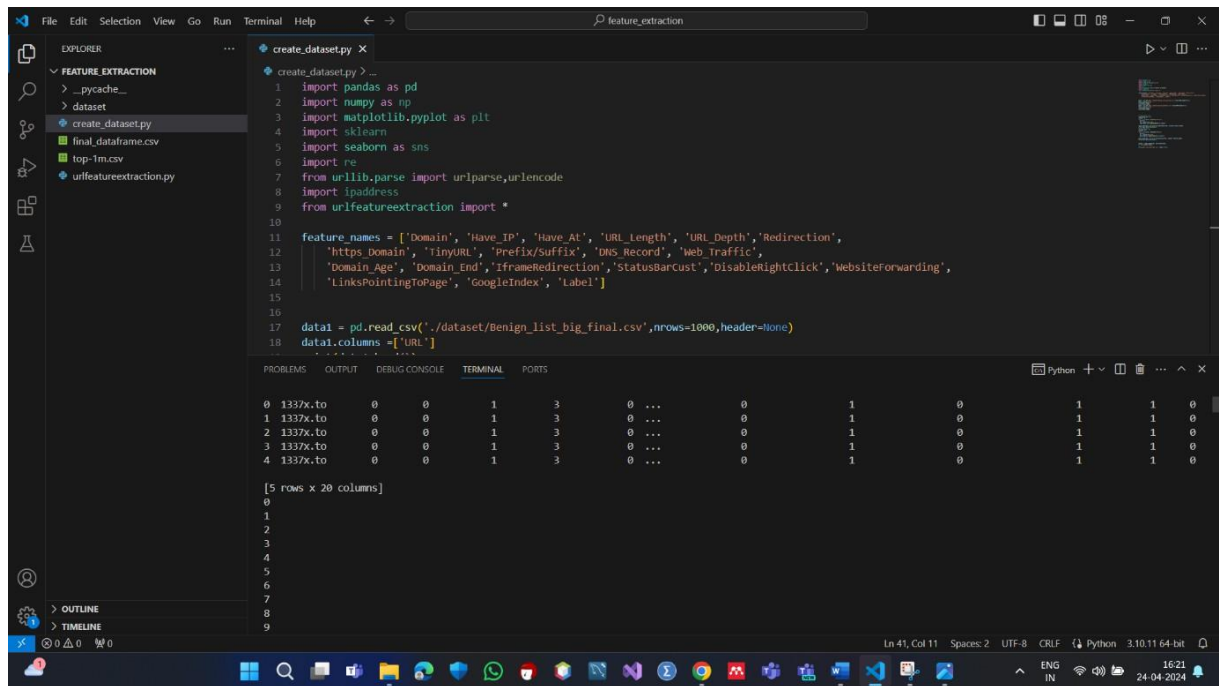
Ln 41, Col 11 Spaces: 2 UTF-8 CRLF Python 3.10.11 64-bit

16:20 24-04-2024

	Domain	Have_IP	Have_At	URL_Length	URL_Depth	...	DisableRightClick	WebsiteForwarding	LinksPointingToPage	GoogleIndex	Label
0	v2.email-marketing.adminsimple.com	0	0	1	2	...	1	1	1	1	1
1	bid.openrx.net	0	0	1	1	...	1	0	1	1	1
2	webmail2.centurytel.net	0	0	1	3	...	1	0	1	1	1
3	google.com.ng	0	0	1	1	...	1	0	-1	1	1
4	webmail2.centurytel.net	0	0	1	3	...	1	0	1	1	1

[5 rows x 20 columns]

c:\Users\IANWAY\Desktop\url_phishing_detection\feature_extraction>



References

Best Python libraries for Machine Learning - GeeksforGeeks (2023). Available at: <https://www.geeksforgeeks.org/best-python-libraries-for-machine-learning/> (Accessed: 25 April 2024).

colab.google (n.d). Available at: <https://colab.google/> (Accessed: 25 April 2024).

URL 2016 | Datasets | Research | Canadian Institute for Cybersecurity | UNB (n.d). Available at: <https://www.unb.ca/cic/datasets/url-2016.html> (Accessed: 25 April 2024).

Visual Studio Code - Code Editing. Redefined (n.d). Available at: <https://code.visualstudio.com/> (Accessed: 25 April 2024).

Welcome to Flask — Flask Documentation (3.0.x) (no date). Available at: <https://flask.palletsprojects.com/en/3.0.x/> (Accessed: 25 April 2024).