

# **WEB PHISHING DETECTION**

## **A PROJECT REPORT**

Submitted by

**BOOMIKA R**

**BOWTHARA G**

**MANJU S**

**VINOTHINI S**

Of

**COMPUTER SCIENCE AND ENGINEERING**

**SSM COLLEGE OF ENGINEERING**

**KOMARAPALAYAM**

# **Project Report Format**

## **1. INTRODUCTION**

- 1.1 Project Overview
- 1.2 Purpose

## **2. LITERATURE SURVEY**

- 2.1 Existing problem
- 2.2 References
- 2.3 Problem Statement Definition

## **3. IDEATION & PROPOSED SOLUTION**

- 3.1 Empathy Map Canvas
- 3.2 Ideation & Brainstorming
- 3.3 Proposed Solution
- 3.4 Problem Solution fit

## **4. REQUIREMENT ANALYSIS**

- 4.1 Functional requirement
- 4.2 Non-Functional requirements

## **5. PROJECT DESIGN**

- 5.1 Data Flow Diagrams
- 5.2 Solution & Technical Architecture
- 5.3 User Stories

## **6. PROJECT PLANNING & SCHEDULING**

- Sprint Planning & Estimation
- 6.1 Sprint Delivery Schedule
- 6.2 Reports from JIRA

## **7. CODING & SOLUTIONING (Explain the features added in the project along with code)**

- 7.1 Feature 1
- 7.2 Feature 2
- 7.3 Database Schema (if Applicable)

## **8. TESTING**

- 8.1 Test Cases
- 8.2 User Acceptance Testing

## **9. RESULTS**

- 9.1 Performance Metrics

## **10. ADVANTAGES & DISADVANTAGES**

## **11. CONCLUSION**

## **12. FUTURE SCOPE**

## **13. APPENDIX**

- Source Code
- GitHub & Project Demo Link

# INTRODUCTION

## Project Overview:

WebPhish is a website which is used to detect phishing sites to improve the customer's sense of safety whenever he/she attempts to provide any sensitive information to a site. Also, by which people won't access them which will reduce the revenue of malicious site owners. This application can be accessed online without paying instead, can be accessed via any browser of the customer's choice to detect any site with high accuracy. This system uses machine learning algorithm which implements classification algorithms and techniques to extract the phishing datasets criteria to classify their legitimacy.

The design and implementation of a comprehensive web phishing detection system instils a cyber security culture which prevents the need for the deployment of targeted anti-phishing solutions in a corporate to meet industry's compliance obligations.

## Purpose:

Web phishing is a threat in various aspects of security on the internet, which might involve scams and private information disclosure. Some of the common threats of web phishing are:

- Attempt to fraudulently solicit personal information from an individual or organization.
- Attempt to deliver malicious software by posing as a trustworthy organization or entity.
- Installing those malwares infects the data that cause a data breach or even nature's forces that takes down your company's data headquarters, disrupting access.

For this purpose, the objective of our project involves building an efficient and intelligent system to detect such websites by applying a machine-learning algorithm which implements classification algorithms and techniques to extract the phishing datasets criteria to classify their legitimacy and as a result of which whenever a user makes a transaction online and makes payment through an e- banking website our system will use a data mining algorithm to detect whether the e-banking website is a phishing website or not.

# **LITERATURE SURVEY**

## **1.Web phishing detection using a deep learning framework**

Web service is one of the key communications software services for the Internet. Web phishing is one of many security threats to web services on the Internet. Web phishing aims to steal private information, such as usernames, passwords, and credit card details, by way of impersonating a legitimate entity. It will lead to information disclosure and property damage. This paper mainly focuses on applying a deep learning framework to detect phishing websites. This paper first designs two types of features for web phishing: original features and interaction features. A detection model based on Deep Belief Networks (DBN) is then presented. The test using real IP flows from ISP (Internet Service Provider) shows that the detecting model based on DBN can achieve an approximately 90% true positive rate and 0.6% false positive rate.

## 2.Spam detection using machine learning techniques

Online reviews are often the primary factor in a customer's decision to purchase a product or service, and are a valuable source of information that can be used to determine public opinion on these products or services. Because of their impact, manufacturers and retailers are

highly concerned with customer feedback and reviews. Reliance on online reviews gives rise to the potential concern that wrongdoers may create false reviews to artificially promote or devalue products and services. This practice is known as Opinion (Review) Spam, where spammers manipulate and poison reviews (i.e., making fake, untruthful, or deceptive reviews) for profit or gain. Since not all online reviews are truthful and trustworthy, it is important to develop techniques for detecting review spam. By extracting meaningful features from the text using Natural Language Processing (NLP), it is possible to conduct review spam detection using various machine learning techniques. Additionally, reviewer information, apart from the text itself, can be used to aid in this process. In this paper, we survey the prominent machine learning techniques that have been proposed to solve the problem of review spam detection and the performance of different approaches for classification and detection of review spam. The majority of current research has focused on supervised learning methods, which require labeled data, a scarcity when it comes to online review spam. Research on methods for Big Data are of interest, since there are millions of online reviews, with many more being generated daily. To date, we have not found any papers that study the effects of Big Data analytics for review spam detection. The primary goal of this paper is to provide a strong and comprehensive comparative study of current research on detecting review spam using various machine learning techniques and to devise methodology for conducting further investigation.

### **3. A machine learning based approach for phishing detection using hyperlinks information**

This paper presents a novel approach that can detect phishing attack by analysing the hyperlinks found in the HTML source code of the website. The proposed approach incorporates various new outstanding hyperlink

specific features to detect phishing attack. The proposed approach has divided the hyperlink specific features into 12 different categories and used these features to train the machine learning algorithms. We have evaluated the performance of our proposed phishing detection approach on various classification algorithms using the phishing and non-phishing websites data set. The proposed approach is an entirely client-side solution, and does not require any services from the third party. Moreover, the proposed approach is language independent and it can detect the website written in any textual language. Compared to other methods, the proposed approach has relatively high accuracy in detection of phishing websites as it achieved more than 98.4% accuracy on logistic regression classifier

## 4. Deep Learning for Phishing Detection: Taxonomy, Current Challenges and Future Directions

Phishing has become an increasing concern and captured the attention of end-users as well as security experts. Despite decades of development and improvement, existing phishing detection techniques still suffer from the deficiency in performance accuracy and the inability to detect unknown attacks. Motivated to solve these problems, many researchers in the cybersecurity domain have shifted their attention to phishing detection that capitalizes on machine learning techniques. In recent years, deep learning has emerged as a branch of machine learning that has become a promising solution for phishing detection. As a result, this study proposes a taxonomy of deep learning algorithms for phishing detection by examining 81 selected papers using a systematic literature review approach. The paper first introduces the concept of phishing and deep learning in the context of cybersecurity. Then, phishing detection and deep learning algorithm taxonomies are provided to classify the existing literature into various categories. Next, taking the proposed taxonomy as a baseline, this study comprehensively reviews the state-of-the-art deep learning techniques and analyzes their advantages as well as disadvantages. Subsequently, the paper discusses various issues deep learning faces in phishing detection and proposes future research directions to overcome these challenges. Finally, an empirical analysis is conducted to evaluate the performance of various deep learning techniques in a practical context and highlight the related issues that motivate researchers in their future works. The results obtained from the empirical experiment showed that the common issues among most of the state-of-the-art deep learning algorithms are manual parameter-tuning, long training time, and deficient detection accuracy.

## 5. Machine learning based phishing detection from URLs

Due to the rapid growth of the Internet, users change their preference from traditional shopping to the electronic commerce. Instead of bank/shop robbery, nowadays, criminals try to find their victims in the cyberspace with some specific tricks. By using the anonymous structure of the Internet, attackers set out new techniques, such as phishing, to deceive victims with the use of false websites to collect their sensitive information such as account IDs, usernames, passwords, etc. Understanding whether a web page is legitimate or phishing is a very challenging problem, due to its semantics-based attack structure, which mainly exploits the computer users' vulnerabilities. Although software companies launch new anti-phishing products, which use blacklists, heuristics, visual and machine learning-based approaches, these products cannot prevent all of the phishing attacks. In this paper, a real-time anti-phishing system, which uses seven different classification algorithms and natural language processing (NLP) based features, is proposed. The system has the following distinguishing properties from other studies in the literature: language independence, use of a huge size of phishing and legitimate data, real-time execution, detection of new websites, independence from third-party services and use of feature-rich classifiers. For measuring the performance of the system, a new data set is constructed, and the experimental results are tested on it. According to the experimental and comparative results from the implemented classification algorithms, Random Forest algorithm with only NLP based features gives the best performance with the 97.98% accuracy rate for detection of phishing URLs.



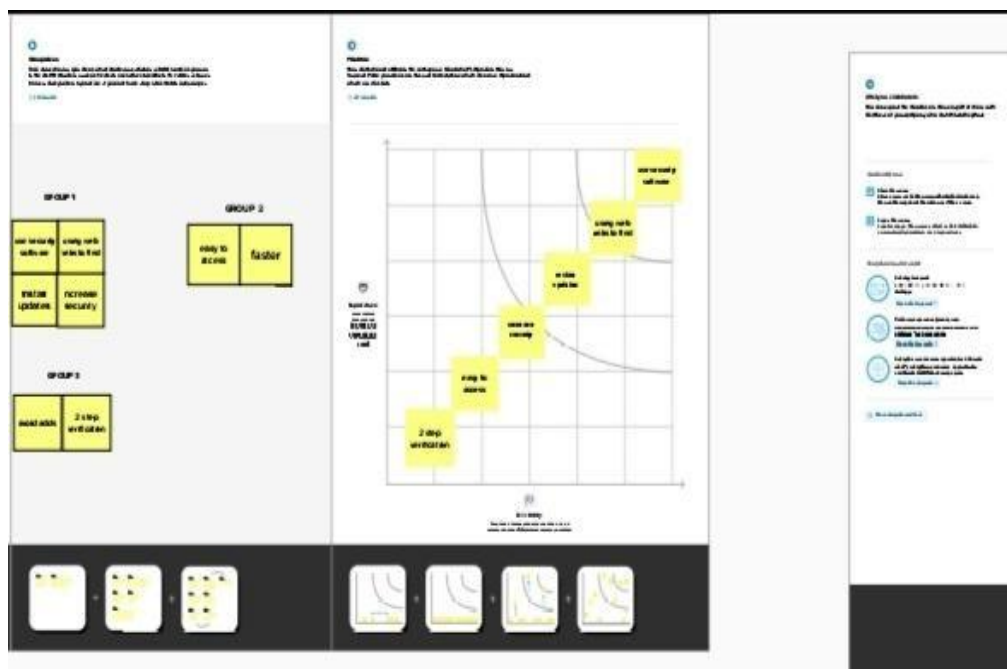
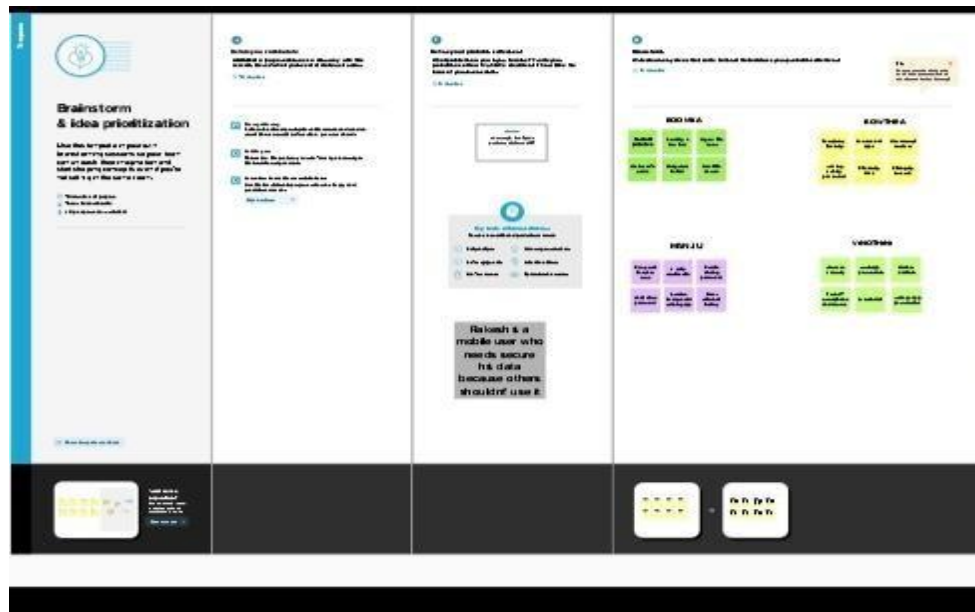
## CHAPTER 3

### IDEATION & PROPOSED SOLUTION

#### EMPATHY MAP



# IDEATION & BRAINSTROMING



### Proposed SolutionTemplate:

Project team shall fill the following information in proposed

S.No.	Parameter	Description
1.	Problem Statement (Problem to be solved)	The purpose of this web phishing detection is to help individually identify these phishing URLs and prevent the Data(Or)device from phishing websites.
2.	Idea / Solution description	If user open the phishing URLs then web phishing detection will find out the phishing websites and notified to the user with detail description of phishing websites to prevent the data (or) device of user. The user can view the status of web phishing.
3.	Novelty / Uniqueness	Occupies less storage User can view the table (date, time, phishing URL and type of phishing) of phishing websites.
4.	Social Impact / Customer Satisfaction	User satisfaction does have a positive effect on an organisation's profitability. when the users are happy with the service they receive, they are more likely to trust and be loyal to the company.
5.	Business Model (Revenue Model)	It prevents the data. It avoids web phishing URLs. It has 2-step verification method. Fast and easy to access. We can access on all devices (desktop, android). Reduce the customer problems.
6.	Scalability of the Solution	An environment where they will be able to spend less time on grunt work and more time on actually resolving critical customer issue. User does not need much to worry about web phishing..

# PROBLEM SOLUTION FIT

<b>1. CUSTOMER SEGMENT(S)</b> <span>CS</span> Who is your customer? The user who wants to detect the phishing links and prevents the data.	<b>6. CUSTOMER CONSTRAINTS</b> <span>CC</span> What constraints prevent your customers from taking action or limit their choices of solutions? 1. Network connection 2. have a device	<b>5. AVAILABLE SOLUTIONS</b> <span>AS</span> Which solutions are available to the customers when they face the problem or need to get the job done? What have they tried in the past? What pros & cons do these solutions have? Its finding out the phishing links .
<b>2. JOBS-TO-BE-DONE / PROBLEMS</b> Which jobs-to-be-done (or problems) do you address for your customers? 1. People want to do the most of things with only one device. 2. They want automatically detect the phishing websites. 3. They want prevent their data from unknown person.	<b>9. PROBLEM ROOT CAUSE</b> <span>RC</span> What is the real reason that this problem exists? What is the back story behind the need to do this job? 1. Sometimes the user doesn't know about real and phishing links. 2. The unknown person send the phishing link to get the user's data	<b>7. BEHAVIOUR</b> <span>BE</span> What does your customer do to address the problem and get the job done? Detect phishing links and the user have network connection.
<b>3. TRIGGERS</b> <span>TR</span> What triggers customers to act? 1. People want to make their life easier, for connected anytime and anywhere. 2. Phishing detection act prevent the data.  <b>4. EMOTIONS: BEFORE / AFTER</b> <span>EM</span> How do customers feel when they face a problem or a job and afterwards? 1. They shouldn't want to others access their data. 2. User doesn't want much worry about web phishing	<b>10. YOUR SOLUTION</b> <span>SL</span> If you are working on an existing business, write down your current solution first, fill in the canvas, and check how much it fits reality. If you are working on a new business proposition, then keep it blank until you fill in the canvas and come up with a solution that fits within customer limitations, solves a problem and matches customer behaviour. These days data often needs to be visually presented in the form of interactive graphs or charts to be impact full and understand.	<b>8. CHANNELS of BEHAVIOUR</b> <span>CH</span> <b>8.1 ONLINE</b> What kind of actions do customers take online? online for detecting phishing websites.  <b>8.2 OFFLINE</b> What kind of actions do customers take offline? 1. Table of phishing links. 2. offline videos and photos and some data.

## CHAPTER 4

### REQUIREMENT ANALYSIS

#### Functional Requirements

Following are the functional requirements of the proposed solution.

FR NO	Functional Requirement (Epic)	Sub Requirement (Story / Sub-Task)
FR-1	User Registration	Registration through Form Registration through Gmail
FR-2	User Confirmation	Confirmation via Email  Confirmation via OTP
FR-3	User uploads	User enter the URL to know about the detail of phishing
FR-4	User accessibility	User can get the notification about the phishing web page.  User enter the URL to know about the detail of phishing.  They can access table of history which detail description of phishing URL.
FR-5	User benefits	They can prevent the their data from the unknown persons like hackers.
FR-6	notification	The notification displays whether website is a legal site or phishing site.

## Non-functional Requirements:

Following are the non-functional requirements of the proposed solution.

FR No.	Non-Functional Requirement	Description
NFR-1	<b>Usability</b>	An user friendly web application, dynamic and attractive user interface ,maintaining the table of history and easily satisfies the user needs.
NFR-2	<b>Security</b>	Users information in the web application are protected by advanced security system. Implementation of proper logging. Authentication.
NFR-3	<b>Reliability</b>	Fault less. A piece of software operating without failure. The website's load time is not more than one second for user
NFR-4	<b>Performance</b>	Fast and quick analyzation of phishing URLs, is done as a GPU used for the model is 10% more fast in analysing and uploading the user uploaded the URL. Occupation of less storage space.
NFR-5	<b>Availability</b>	Available in a google play store, in security and production.
NFR-6	<b>Scalability</b>	It works in high speed authentication of phishing URLs, quick response for queries from user, highly  reliable.

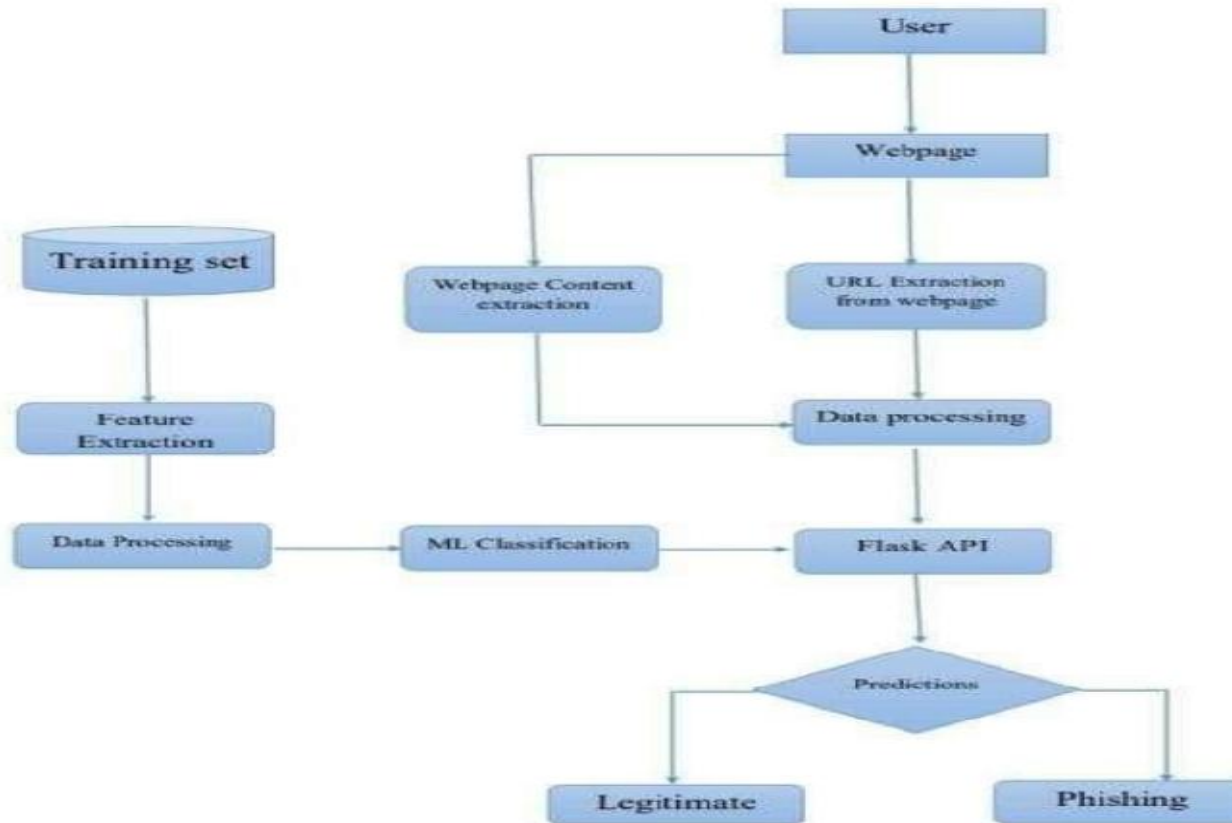
# CHAPTER 5

## PROJECT DESIGN

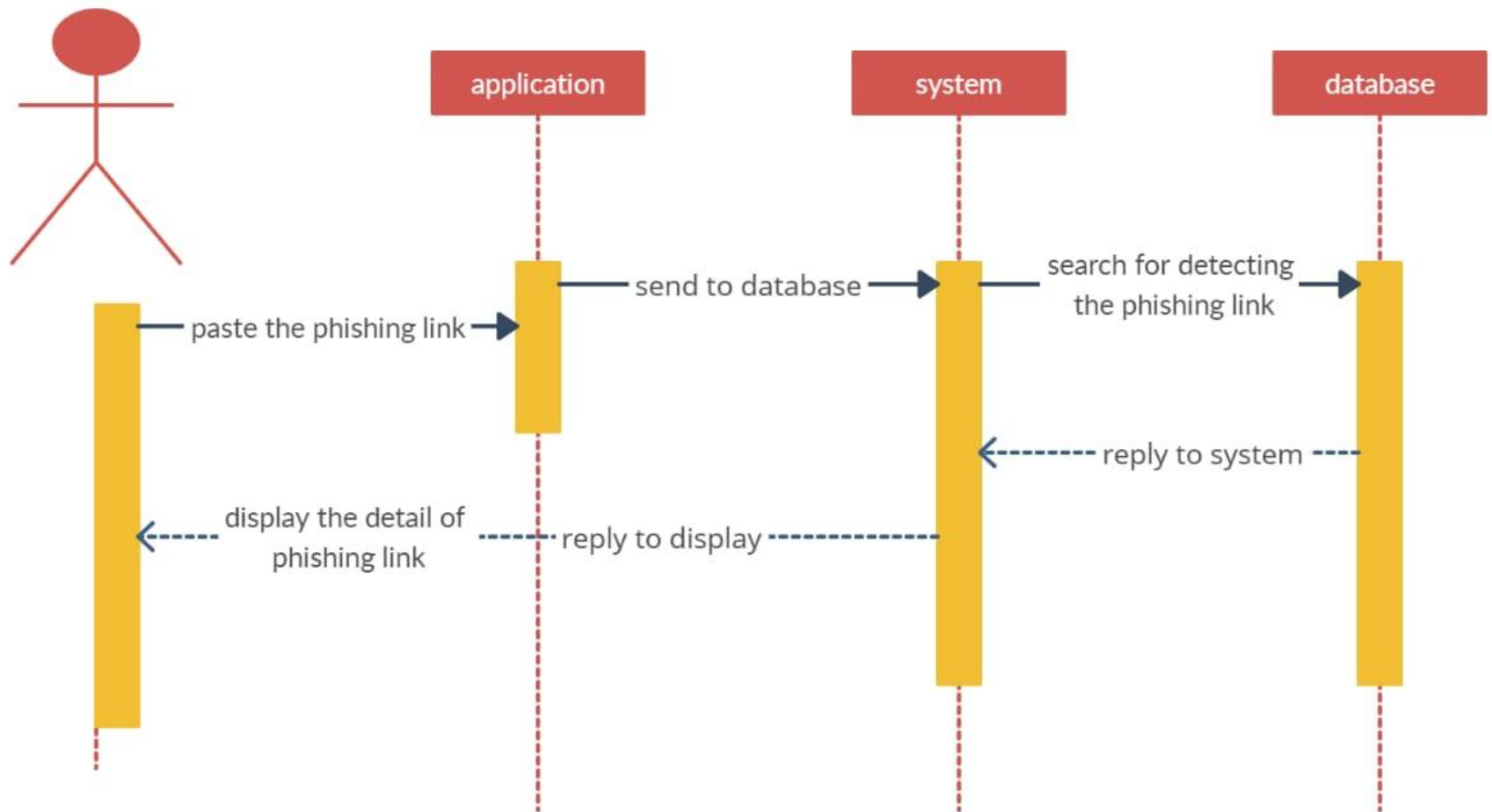
### Data Flow Diagram

A Data Flow Diagram (DFD) is a traditional visual representation of the information flows within a system. A neat and clear DFD can depict the right amount of the system requirement graphically. It shows how data enters and leaves the system, what changes the information, and where data is stored.

---



## SOLUTION & TECHNICAL ARCHITECTURE





# TECHNICAL ARCHITECTURE

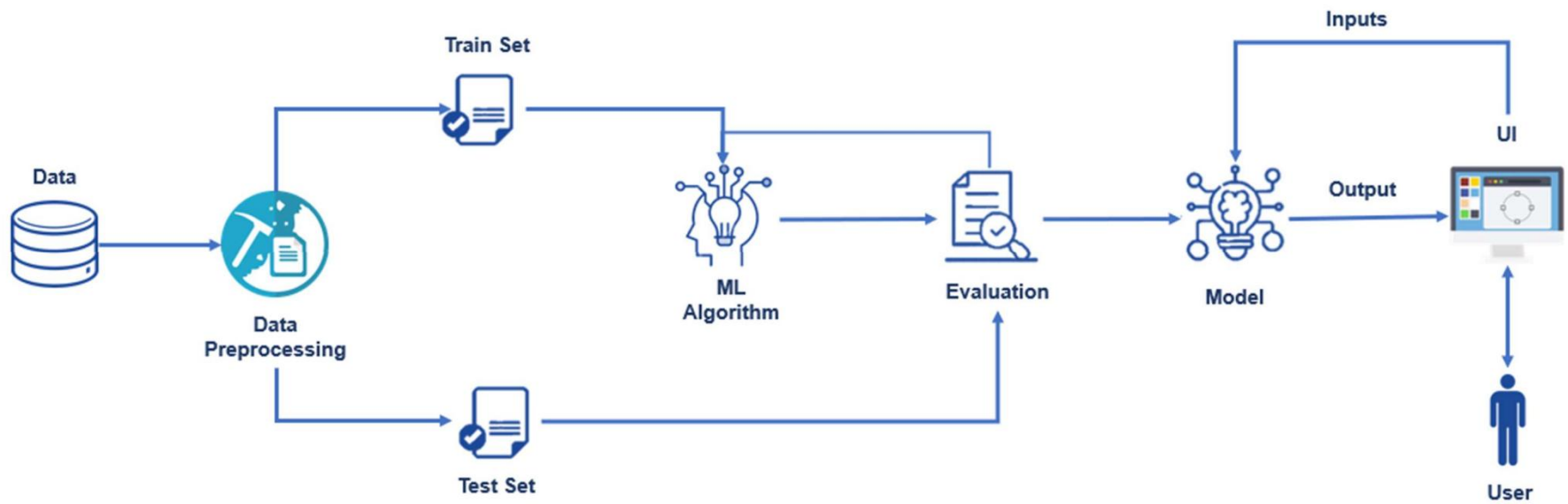


Table-1 : Components & Technologies:

S. No	Component	Description	Technology
1.	User Interface	Mobile application	HTML, CSS, JavaScript ,Python.
2.	Application Logic-1	Data pre-processing	Python.
3.	Application Logic-2	Model creating.	Keras, numpy, pandas.
4.	Application Logic-3	Web application(UI)	IBM Watson Assistant
5.	Database	data	MySQL, No SQL, etc.
6.	File Storage	File storage requirements	IBM Block Storage or Other Storage Service or Local File system
7.	External API	Purpose of External API used in the application	Data processing API.
8.	Machine Learning Model	Binary, multi class, regression classification.	Object Recognition Model.
9.	Infrastructure (Server / Cloud)	Application Deployment on web server.	Flask – python HTTP server.

Table-2: Application Characteristics:

S. No	Characteristics	Description	Technology
1.	Open-Source Frameworks	Flask	Technology of Opensource framework
2.	Security Implementations	Secure flag for cookies.	Flask WTF Session_cookie_secure.
3.	Scalable Architecture	Micro-services	Micro web applications.
4.	Availability	Integrated supporting for unit testing.	SQL, Sinatra ruby, framework, jinja 2
5.	Performance	HTTP request handling High flexibility.	SQL, Sinatra ruby, framework,jinja 2.

# User Stories

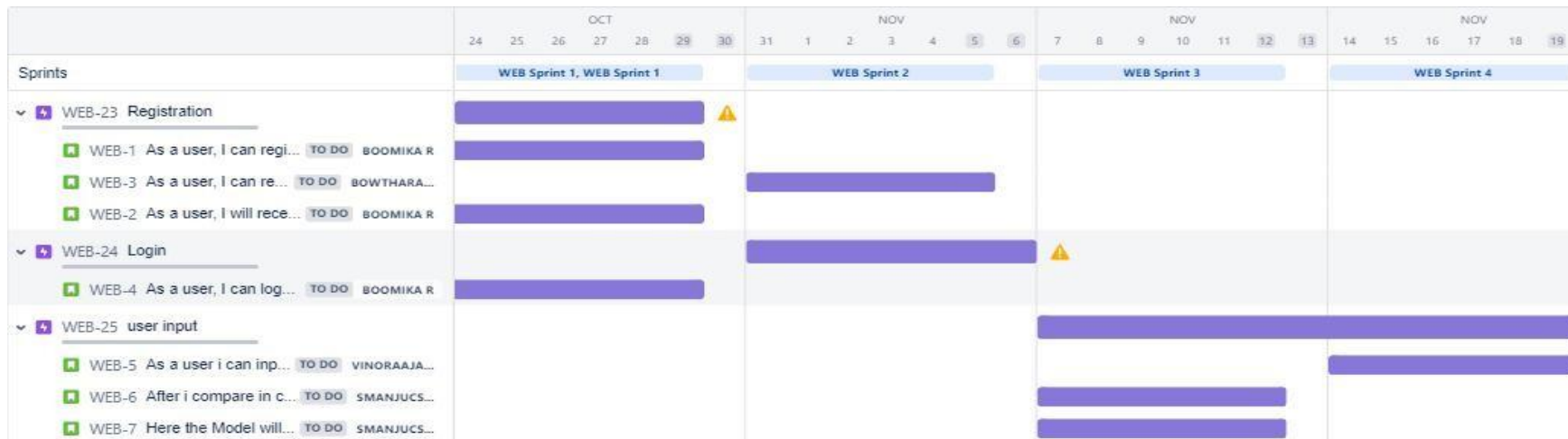
Use the below template to list all the user stories for the product.

User Type	Functional Requirement (Epic)	User Story Number	User Story / Task	Acceptance criteria	Priority	Release
Customer (Mobile user)	Registration	USN-1	As a user, I can register for the application By entering my email, password, and confirming my password.	I can access my account / dashboard	High	Sprint-1
		USN-2	As a user, I will receive confirmation email once I have registered for the application	I can receive confirmation email & click confirm	High	Sprint-1
		USN-3	As a user, I can register for the application through Facebook	I can register & access The dashboard with Facebook Login	Low	Sprint-2
		USN-4	As a user, I can register for the application through Gmail		Medium	Sprint-1
	Login	USN-5	As a user, I can log into the application by entering email & password		High	Sprint-1
	Dashboard					
Customer (Web user)	User input	USN-1	As a user i can input the particular URL in the required field and waiting for validation.	I can go access the website without any problem	High	Sprint-1
Customer Care Executive	Feature extraction	USN-1	After i compare in case if none found on comparison then we can extract feature using heuristic and visual similarity approach.	As a User i can have comparison between websites for security.	High	Sprint-1
Administrator	Prediction	USN-1	Here the Model will predict the URL websites using Machine Learning algorithms such as Logistic Regression, KNN	In this i can have correct prediction on the particular algorithms	High	Sprint-1
	Classifier	USN-2	Here i will send all the model output to classifier in order to produce final result.	I this i will find the correct classifier for producing the result	Medium	Sprint-2

## SPRINT PLANNING & ESTIMATION

Title	Description	Date
<b>Literature Survey and Information Gathering</b>	Gathering Information by referring the technical papers, research publications etc	10 September 2022
<b>Prepare Empathy Map</b>	To capture user pain and gains Prepare List of Problem Statement	17 September 2022
<b>Ideation</b>	Prioritise a top 3 ideas based on feasibility and Importance	18 September 2022
<b>Proposed Solution</b>	Solution include novelty, feasibility, business model, social impact and scalability of solution	1 October 2022
<b>Problem Solution Fit</b>	Solution fit document	1 October 2022
<b>Solution Architecture</b>	Solution Architecture	1 October 2022
<b>Customer Journey</b>	To Understand User Interactions and experiences with application	8 October 2022
<b>Functional Requirement</b>	Prepare functional Requirement	9 October 2022
<b>Data flow Diagrams</b>	Data flow diagram	11 October 2022
<b>Technology Architecture</b>	Technology Architecture diagram	15 October 2022
<b>Milestone sprint &amp; delivery plan</b>	Activity what we done &further plans	21 October 2022
<b>Project Development-Delivery of sprint 1,2,3 &amp;4</b>	Develop and submit the developed code by testing it	24 October 2022 – 19 November 2022

# REPORTS FROM JIRA



# CODING & SOLUTIONING

## FEATURE 1

### Phishing URL Detection

```
Phishing URL Detection (2).ipynb
File Edit View Insert Run/Tools Help All changes saved
+ Code + Text
[ ] #Information about the dataset
data.info()
> 'pandas.core.frame.DataFrame'>
Index: 11054 entries, 0 to 11053
Columns (total 32 columns):
Column Non-Null Count
-----
Index 11054 non-null
UsingIP 11054 non-null
LongURL 11054 non-null
ShortURL 11054 non-null
Symbol@ 11054 non-null
Redirecting// 11054 non-null
PrefixSuffix- 11054 non-null
SubDomains 11054 non-null
HTTPS 11054 non-null
DomainRegLen 11054 non-null
Favicon 11054 non-null
NonStdPort 11054 non-null
HTTPSDomainURL 11054 non-null
RequestURL 11054 non-null

[ ] # unique value in columns
data.nunique()
Index 11054
UsingIP 2
LongURL 3
ShortURL 2
Symbol@ 2
Redirecting// 2
PrefixSuffix- 2
SubDomains 3
HTTPS 1
```

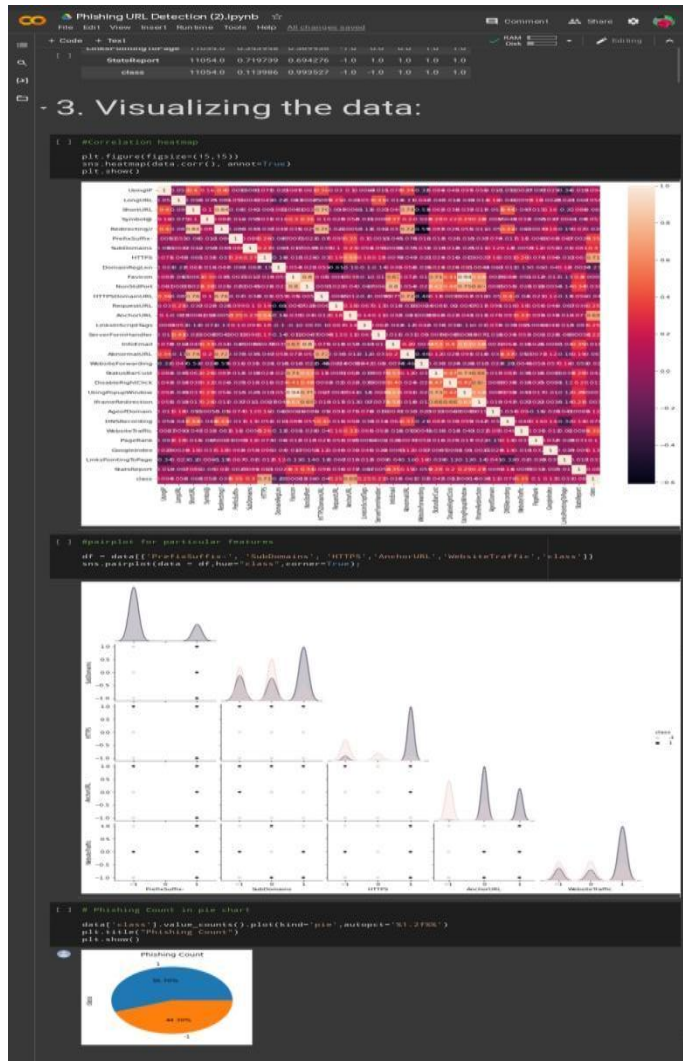
```
Phishing URL Detection (2).ipynb
File Edit View Insert Run/Tools Help All changes saved
+ Code + Text
[ ] #Information about the dataset
data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11054 entries, 0 to 11053
Data columns (total 32 columns):
# Column Non-Null Count
---
0 Index 11054 non-null
1 UsingIP 11054 non-null
2 LongURL 11054 non-null
3 ShortURL 11054 non-null
4 Symbol@ 11054 non-null
5 Redirecting// 11054 non-null
6 PrefixSuffix- 11054 non-null
7 SubDomains 11054 non-null
8 HTTPS 11054 non-null
9 DomainRegLen 11054 non-null
10 Favicon 11054 non-null
11 NonStdPort 11054 non-null
12 HTTPSDomainURL 11054 non-null
13 RequestURL 11054 non-null

[ ] # unique value in columns
data.nunique()
Index 11054
UsingIP 2
LongURL 3
ShortURL 2
Symbol@ 2
Redirecting// 2
```

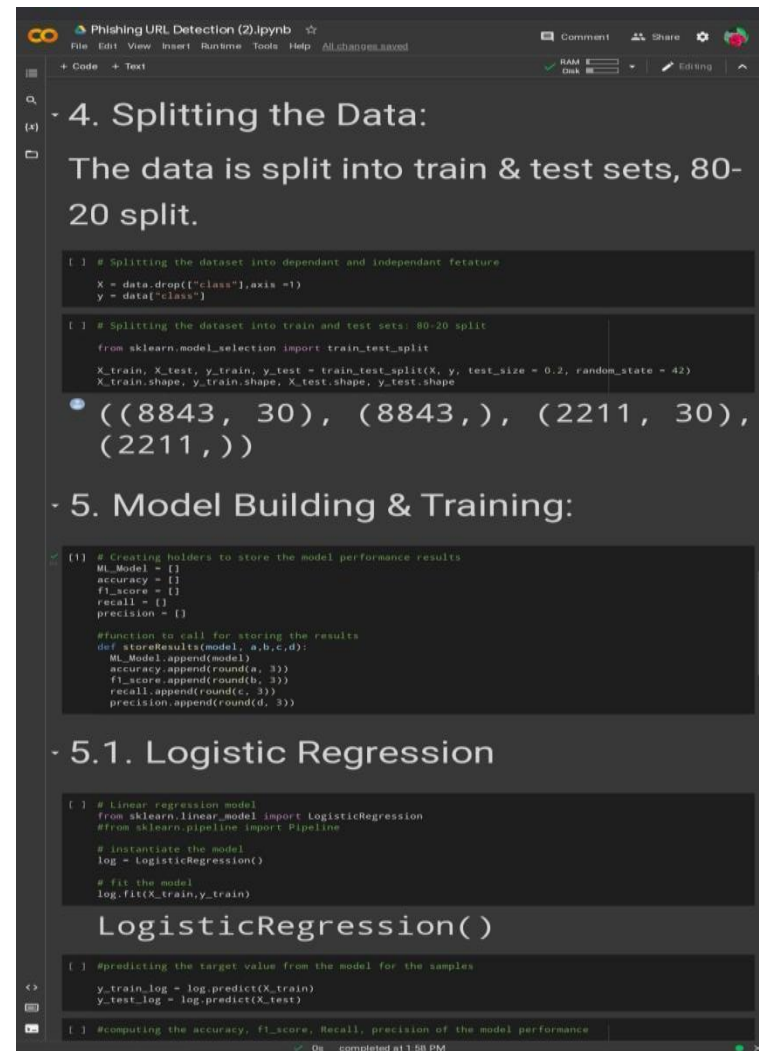




## Visualizing the data



## Splitting the data



## 5. Model Building & Training

```
Phishing URL Detection (2).ipynb
File Edit View Insert Runtime Tools Help All changes saved
+ Code + Text
RAM 100% Disk 100%
5. Model Building & Training:

[ ] # Creating holders to store the model performance results
ML_Model = []
accuracy = []
f1_score = []
recall = []
precision = []

#function to call for storing the results
def storeResults(model, a,b,c,d):
    ML_Model.append(model)
    accuracy.append(round(a, 3))
    f1_score.append(round(b, 3))
    recall.append(round(c, 3))
    precision.append(round(d, 3))

- 5.1. Logistic Regression

[ ] # Linear regression model
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import Pipeline

# instantiate the model
log = LogisticRegression()

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

LogisticRegression()

[ ] #predicting the target value from the model for the samples
y_train_log = log.predict(X_train)
y_test_log = log.predict(X_test)

[ ] #computing the accuracy, f1_score, Recall, precision of the model performance
acc_train_log = metrics.accuracy_score(y_train,y_train_log)
acc_test_log = metrics.accuracy_score(y_test,y_test_log)
print("Logistic Regression : Accuracy on training Data: {:.3f}".format(acc_train_log))
print("Logistic Regression : Accuracy on test Data: {:.3f}".format(acc_test_log))
print()

f1_score_train_log = metrics.f1_score(y_train,y_train_log)
f1_score_test_log = metrics.f1_score(y_test,y_test_log)
print("Logistic Regression : f1_score on training Data: {:.3f}".format(f1_score_train_log))
print("Logistic Regression : f1_score on test Data: {:.3f}".format(f1_score_test_log))
print()

recall_score_train_log = metrics.recall_score(y_train,y_train_log)
recall_score_test_log = metrics.recall_score(y_test,y_test_log)
print("Logistic Regression : Recall on training Data: {:.3f}".format(recall_score_train_log))
print("Logistic Regression : Recall on test Data: {:.3f}".format(recall_score_test_log))
print()

precision_score_train_log = metrics.precision_score(y_train,y_train_log)
precision_score_test_log = metrics.precision_score(y_test,y_test_log)
print("Logistic Regression : precision on training Data: {:.3f}".format(precision_score_train_log))
print("Logistic Regression : precision on test Data: {:.3f}".format(precision_score_test_log))

sion : Accuracy on training Data:
sion : Accuracy on test Data: 0.9:
sion : f1_score on training Data:
```

## K-Nearest Neighbors: Classifier

```
Phishing URL Detection (2).ipynb
File Edit View Insert Runtime Tools Help All changes saved
+ Code + Text
RAM 100% Disk 100%
5.2. K-Nearest Neighbors : Classifier

[ ] # K-Nearest Neighbors Classifier model
from sklearn.neighbors import KNeighborsClassifier

# instantiate the model
knn = KNeighborsClassifier(n_neighbors=1)

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

KNeighborsClassifier(n_neighbors=

[ ] #predicting the target value from the model for the samples
y_train_knn = knn.predict(X_train)
y_test_knn = knn.predict(X_test)

[ ] #computing the accuracy, f1_score, Recall, precision of the model performance
acc_train_knn = metrics.accuracy_score(y_train,y_train_knn)
acc_test_knn = metrics.accuracy_score(y_test,y_test_knn)
print("K-Nearest Neighbors : Accuracy on training Data: {:.3f}".format(acc_train_knn))
print("K-Nearest Neighbors : Accuracy on test Data: {:.3f}".format(acc_test_knn))
print()

f1_score_train_knn = metrics.f1_score(y_train,y_train_knn)
f1_score_test_knn = metrics.f1_score(y_test,y_test_knn)
print("K-Nearest Neighbors : f1_score on training Data: {:.3f}".format(f1_score_train_knn))
print("K-Nearest Neighbors : f1_score on test Data: {:.3f}".format(f1_score_test_knn))
print()

recall_score_train_knn = metrics.recall_score(y_train,y_train_knn)
recall_score_test_knn = metrics.recall_score(y_test,y_test_knn)
print("K-Nearest Neighbors : Recall on training Data: {:.3f}".format(recall_score_train_knn))
print("K-Nearest Neighbors : Recall on test Data: {:.3f}".format(recall_score_test_knn))
print()

precision_score_train_knn = metrics.precision_score(y_train,y_train_knn)
precision_score_test_knn = metrics.precision_score(y_test,y_test_knn)
print("K-Nearest Neighbors : precision on training Data: {:.3f}".format(precision_score_train_knn))
print("K-Nearest Neighbors : precision on test Data: {:.3f}".format(precision_score_test_knn))

• K-Nearest Neighbors : Accuracy on
K-Nearest Neighbors : Accuracy on
K-Nearest Neighbors : f1_score on
K-Nearest Neighbors : f1_score on
K-Nearest Neighborsn : Recall on
Logistic Regression : Recall on t
K-Nearest Neighbors : precision o
K-Nearest Neighbors : precision o

[ ] #computing the classification report of the model
print(metrics.classification_report(y_test, y_test_knn))
```

# Decision Trees: Classifier

```
recall_score_train_nb,precision_score_train_nb)

- 5.5. Decision Trees : Classifier

[ ] # Decision Tree Classifier model
from sklearn.tree import DecisionTreeClassifier

# instantiate the model
tree = DecisionTreeClassifier(max_depth=30)

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

DecisionTreeClassifier(max_depth=

[ ] #predicting the target value from the model for the samples
y_train_tree = tree.predict(X_train)
y_test_tree = tree.predict(X_test)

[ ] #computing the accuracy, f1_score, Recall, precision of the model performance

acc_train_tree = metrics.accuracy_score(y_train,y_train_tree)
acc_test_tree = metrics.accuracy_score(y_test,y_test_tree)
print("Decision Tree : Accuracy on training Data: {:.3f}".format(acc_train_tree))
print("Decision Tree : Accuracy on test Data: {:.3f}".format(acc_test_tree))
print()

f1_score_train_tree = metrics.f1_score(y_train,y_train_tree)
f1_score_test_tree = metrics.f1_score(y_test,y_test_tree)
print("Decision Tree : f1_score on training Data: {:.3f}".format(f1_score_train_tree))
print("Decision Tree : f1_score on test Data: {:.3f}".format(f1_score_test_tree))
print()

recall_score_train_tree = metrics.recall_score(y_train,y_train_tree)
recall_score_test_tree = metrics.recall_score(y_test,y_test_tree)
print("Decision Tree : Recall on training Data: {:.3f}".format(recall_score_train_tree))
print("Decision Tree : Recall on test Data: {:.3f}".format(recall_score_test_tree))
print()

precision_score_train_tree = metrics.precision_score(y_train,y_train_tree)
precision_score_test_tree = metrics.precision_score(y_test,y_test_tree)
print("Decision Tree : precision on training Data: {:.3f}".format(precision_score_train_tree))
print("Decision Tree : precision on test Data: {:.3f}".format(precision_score_test_tree))

Decision Tree : Accuracy on train
Decision Tree : Accuracy on test
Decision Tree : f1_score on train
```

# Random Forest : Classifier

```
Phishing URL Detection (2).ipynb
File Edit View Insert Runtime Tools Help Last saved at 2:50 PM

- 5.6. Random Forest : Classifier

[ ] # Random Forest Classifier Model
from sklearn.ensemble import RandomForestClassifier

# instantiate the model
forest = RandomForestClassifier(n_estimators=10)

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

RandomForestClassifier(n_estimators=

[ ] #predicting the target value from the model for the samples
y_train_forest = forest.predict(X_train)
y_test_forest = forest.predict(X_test)

[ ] #computing the accuracy, f1_score, Recall, precision of the model performance

acc_train_forest = metrics.accuracy_score(y_train,y_train_forest)
acc_test_forest = metrics.accuracy_score(y_test,y_test_forest)
print("Random Forest : Accuracy on training Data: {:.3f}".format(acc_train_forest))
print("Random Forest : Accuracy on test Data: {:.3f}".format(acc_test_forest))
print()

f1_score_train_forest = metrics.f1_score(y_train,y_train_forest)
f1_score_test_forest = metrics.f1_score(y_test,y_test_forest)
print("Random Forest : f1_score on training Data: {:.3f}".format(f1_score_train_forest))
print("Random Forest : f1_score on test Data: {:.3f}".format(f1_score_test_forest))
print()

recall_score_train_forest = metrics.recall_score(y_train,y_train_forest)
recall_score_test_forest = metrics.recall_score(y_test,y_test_forest)
print("Random Forest : Recall on training Data: {:.3f}".format(recall_score_train_forest))
print("Random Forest : Recall on test Data: {:.3f}".format(recall_score_test_forest))
print()

precision_score_train_forest = metrics.precision_score(y_train,y_train_forest)
precision_score_test_forest = metrics.precision_score(y_test,y_test_forest)
print("Random Forest : precision on training Data: {:.3f}".format(precision_score_train_forest))
print("Random Forest : precision on test Data: {:.3f}".format(precision_score_test_forest))

Random Forest : Accuracy on train
Random Forest : Accuracy on test

Random Forest : f1_score on train
Random Forest : f1_score on test

Random Forest : Recall on trainin
Random Forest : Recall on test Da

Random Forest : precision on trai
Random Forest : precision on test

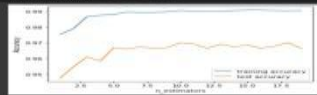
[ ] #printing the classification report of the model
print(metrics.classification_report(y_test, y_test_forest))

              precision    recall
-1              0.96      0.96
1              0.97      0.97

accuracy
macro avg      0.97      0.97
weighted avg    0.97      0.97

[ ] #training accuracy = {}
train_accuracy = {}
for n_estimators in range(1, 20):
    forest = RandomForestClassifier(n_estimators=n_estimators)
    forest.fit(X_train, y_train)
    training_accuracy.append(forest.score(X_train, y_train))
    # predict accuracy on test set
    test_accuracy.append(forest.score(X_test, y_test))

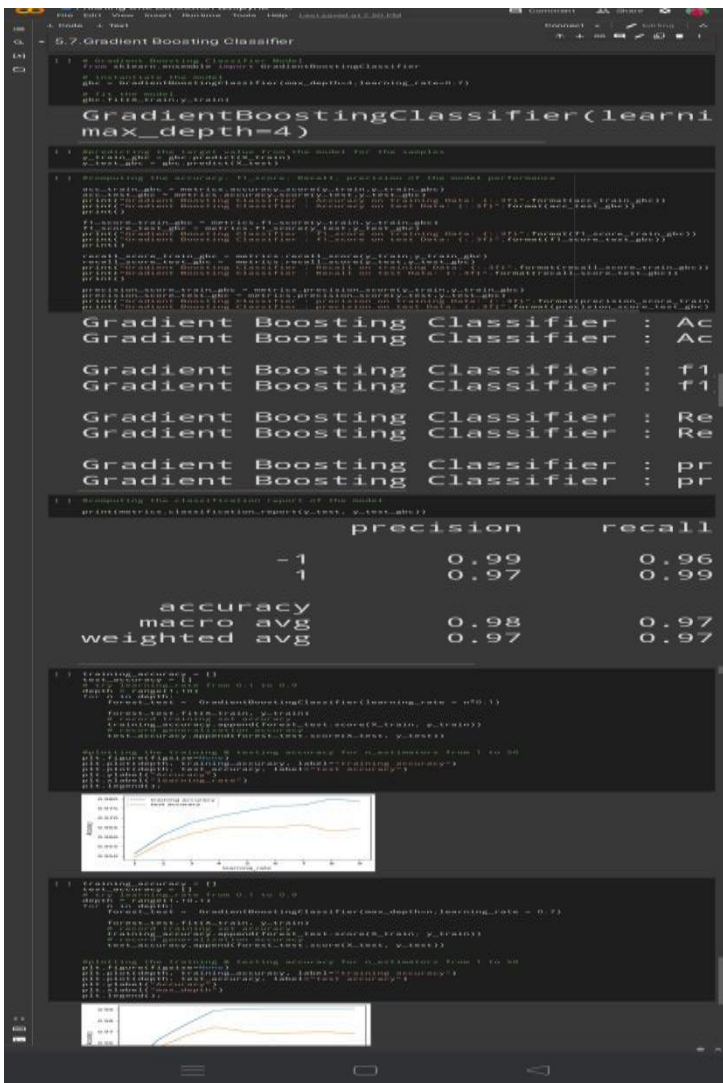
#plotting the training & testing accuracy for n_estimators from 1 to 20
plt.figure(figsize=(10,5))
plt.plot(range(1,20), training_accuracy, label="training accuracy")
plt.plot(range(1,20), test_accuracy, label="test accuracy")
plt.xlabel("n_estimators")
plt.ylabel("accuracy")



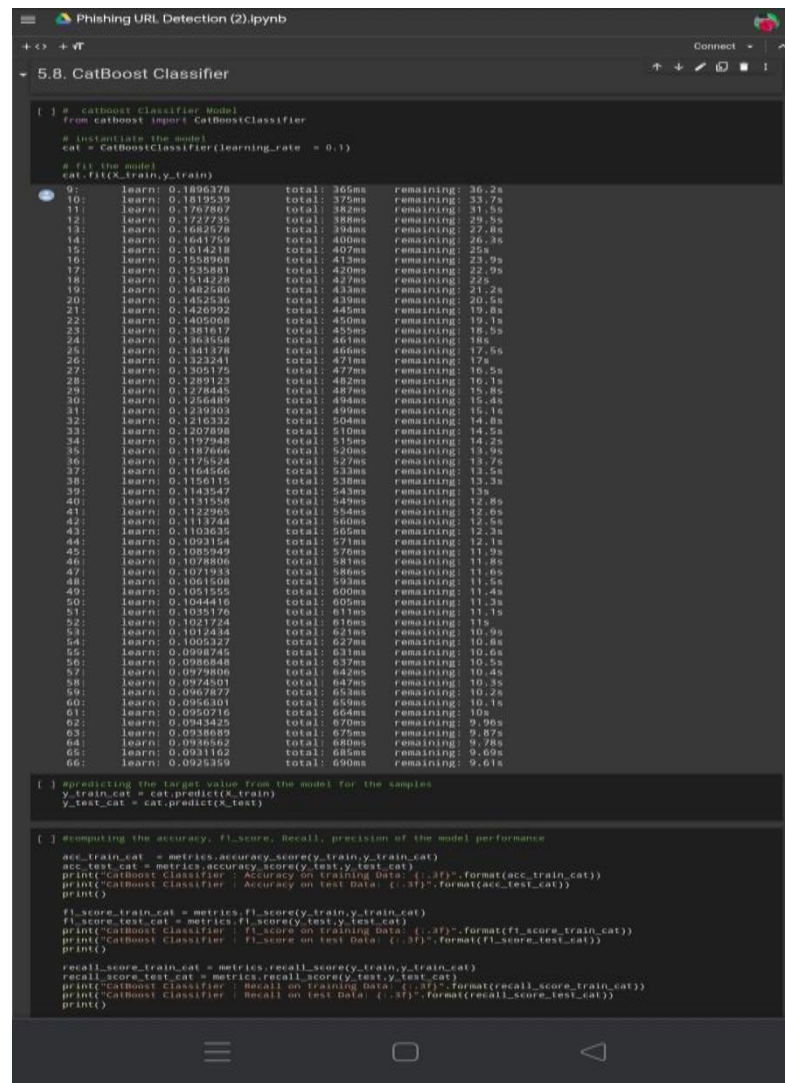
[ ] #printing the results, the below mentioned order of parameter passing is important.
store_result['Random Forest'] = acc_test_forest, f1_score_test_forest,
recall_score_train_forest, precision_score_train_forest
```



# Gradient Boosting Classifier



## CatBoost Classifier



# XGBoost Classifier

```
Phishing URL Detection (2).ipynb
File Edit View Insert Runtime Tools Help Last saved at 2:50 PM
+ Code + Text
Connect + Editing
+ + + + +

5.9. XGBoost Classifier

1 # XGBoost Classifier Model
from xgboost import XGBClassifier

# instantiate the model
xgb = XGBClassifier()

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

interaction_constraints='',
learning_rate=0.300000012,
    max_delta_step=0,
max_depth=6, min_child_weight=1,
missing=nan,

monotone_constraints='()',
n_estimators=100, n_jobs=4,

num_parallel_tree=1,
predictor='auto', random_state=0,
    reg_alpha=0,
reg_lambda=1, scale_pos_weight=1,
subsample=1,

tree_method='exact',
validate_parameters=1,
verbosity=None)

1 #predicting the target value from the model for the samples
y_train_xgb = xgb.predict(X_train)
y_test_xgb = xgb.predict(X_test)

1 #computing the accuracy, f1_score, Recall, precision of the model performance
acc_train_xgb = metrics.accuracy_score(y_train,y_train_xgb)
acc_test_xgb = metrics.accuracy_score(y_test,y_test_xgb)
print(XGBClassifier : Accuracy on training Data: {:.3f}".format(acc_train_xgb))
print(XGBClassifier : Accuracy on test Data: {:.3f}".format(acc_test_xgb))
print()

f1_score_train_xgb = metrics.f1_score(y_train,y_train_xgb)
f1_score_test_xgb = metrics.f1_score(y_test,y_test_xgb)
print(XGBClassifier : f1_score on training Data: {:.3f}".format(f1_score_train_xgb))
print(XGBClassifier : f1_score on test Data: {:.3f}".format(f1_score_test_xgb))
print()

recall_score_train_xgb = metrics.recall_score(y_train,y_train_xgb)
recall_score_test_xgb = metrics.recall_score(y_test,y_test_xgb)
print(XGBClassifier : Recall on training Data: {:.3f}".format(recall_score_train_xgb))
print(XGBClassifier : Recall on test Data: {:.3f}".format(recall_score_test_xgb))
print()

precision_score_train_xgb = metrics.precision_score(y_train,y_train_xgb)
precision_score_test_xgb = metrics.precision_score(y_test,y_test_xgb)
print(XGBClassifier : precision on training Data: {:.3f}".format(precision_score_train_xgb))
print(XGBClassifier : precision on test Data: {:.3f}".format(precision_score_test_xgb))

XGBoost Classifier : Accuracy on
XGBoost Classifier : Accuracy on

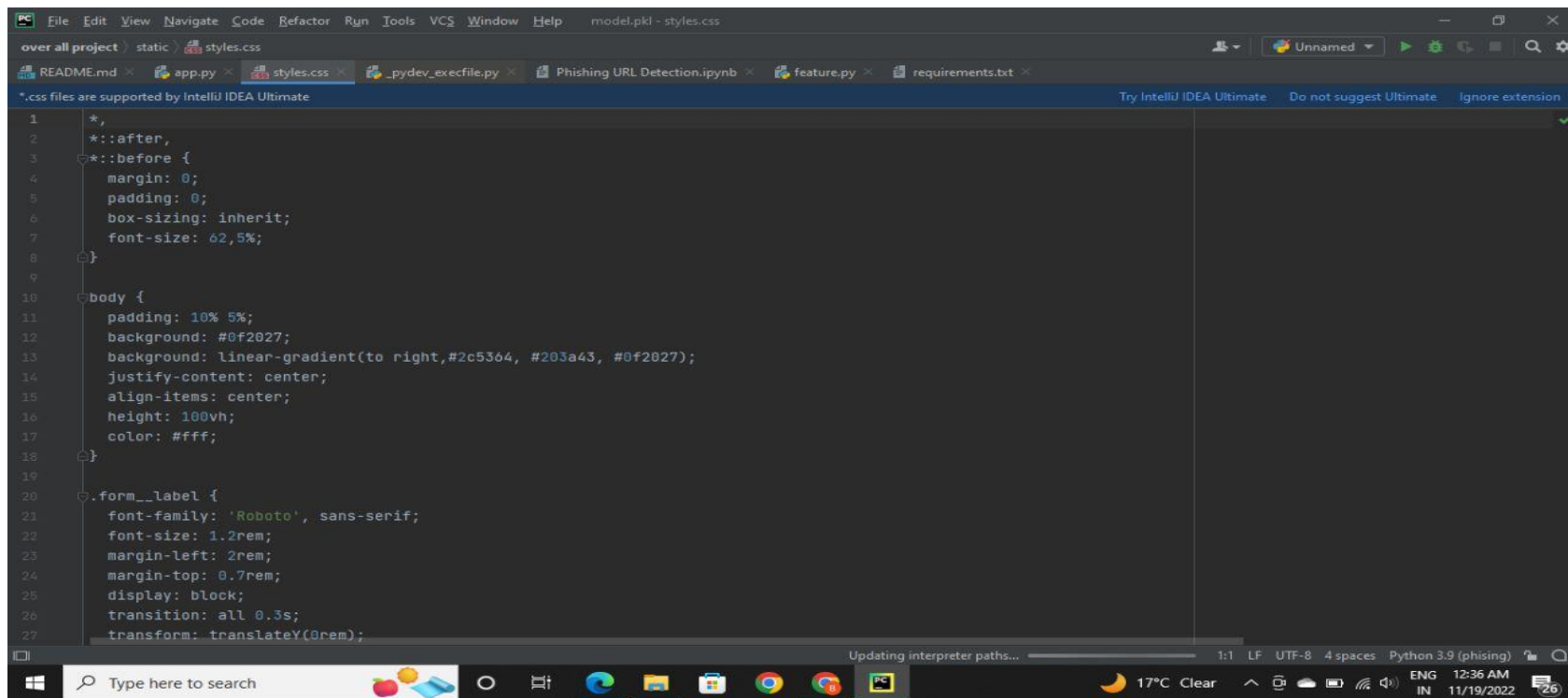
XGBoost Classifier : f1_score on
XGBoost Classifier : f1_score on

XGBoost Classifier : Recall on tr
XGBoost Classifier : Recall on te

XGBoost Classifier : precision on
XGBoost Classifier : precision on
```

## FEATURE 2

### Styles.css



```
1  *,
2  *::after,
3  *::before {
4      margin: 0;
5      padding: 0;
6      box-sizing: inherit;
7      font-size: 62.5%;
8  }
9
10 body {
11     padding: 10% 5%;
12     background: #0f2027;
13     background: linear-gradient(to right, #2c5364, #203a43, #0f2027);
14     justify-content: center;
15     align-items: center;
16     height: 100vh;
17     color: #fff;
18 }
19
20 .form__label {
21     font-family: 'Roboto', sans-serif;
22     font-size: 1.2rem;
23     margin-left: 2rem;
24     margin-top: 0.7rem;
25     display: block;
26     transition: all 0.3s;
27     transform: translateY(0rem);
```

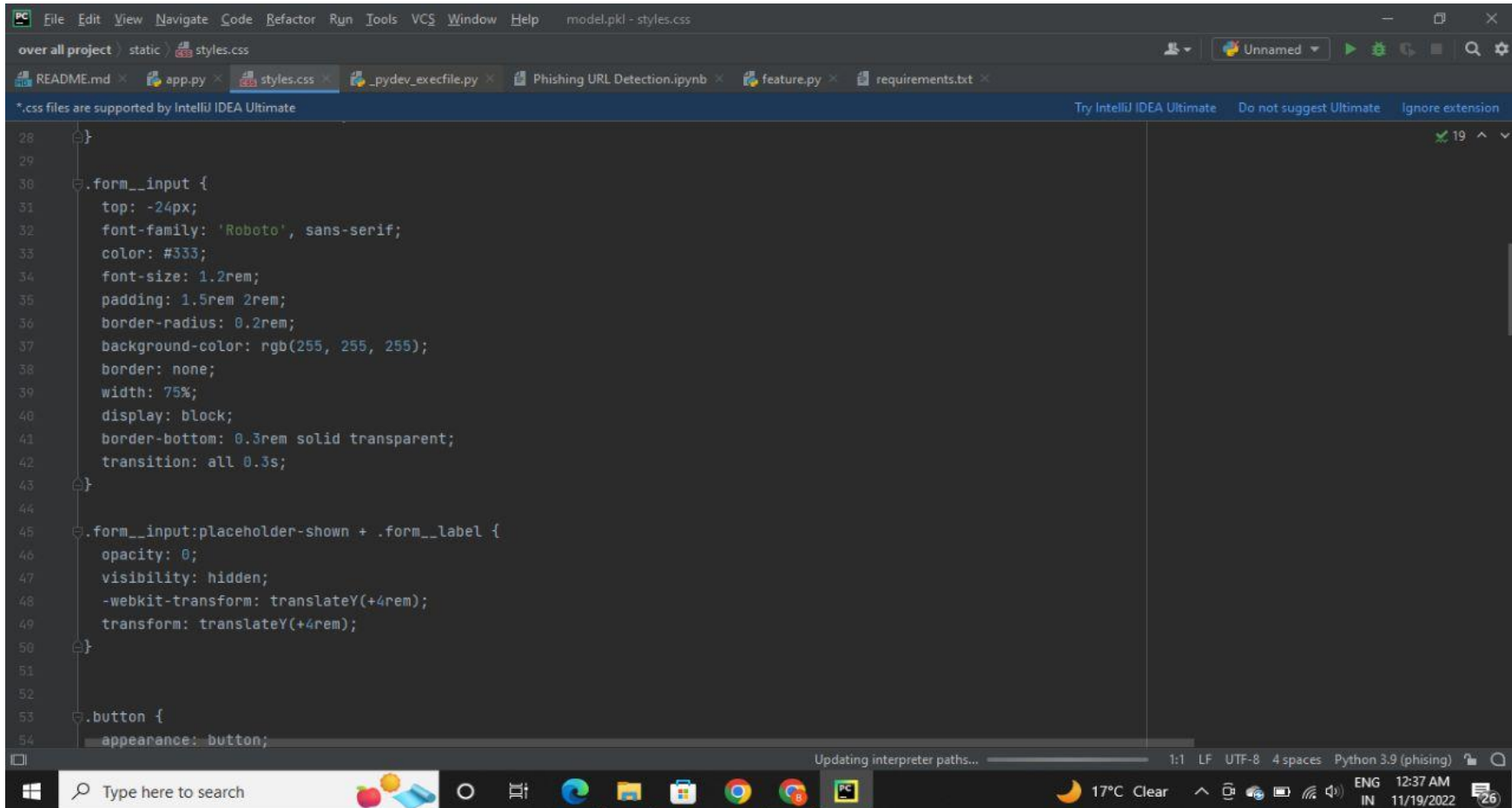
The screenshot shows an IDE window with the file 'styles.css' open. The code defines a reset for all elements, a body with a gradient background and centered content, and a class for form labels. The IDE interface includes a menu bar, a toolbar, and a status bar at the bottom showing the current interpreter as Python 3.9 (phising).

This screenshot shows the IntelliJ IDEA interface with the 'styles.css' file open. The code defines the default style for a button with the selector `.button1`. The CSS includes properties for appearance, background, border, color, font, padding, and transitions. The background features a linear gradient from transparent to a light blue color. The font is set to a sans-serif typeface, and the padding is 1rem 1.5rem. A transition is applied to all properties with a duration of 0.1s and a cubic-bezier curve.

```
97
98 .button1{
99     appearance: button;
100     background-color: transparent;
101     background-image: linear-gradient(to bottom, rgb(160, 245, 174), #37ee65);
102     border: 0 solid #e5e7eb;
103     border-radius: .5rem;
104     box-sizing: border-box;
105     color: #482307;
106     column-gap: 1rem;
107     cursor: pointer;
108     display: flex;
109     font-family: ui-sans-serif,system-ui,-apple-system,system-ui,"Segoe UI",Roboto,"Helvetica Neue",Arial,"Noto Sans",sans-serif,"Apple Color Emoji","Segoe UI Emoji";
110     font-size: 100%;
111     font-weight: 700;
112     line-height: 24px;
113     margin: 0;
114     outline: 2px solid transparent;
115     padding: 1rem 1.5rem;
116     text-align: center;
117     text-transform: none;
118     transition: all .1s cubic-bezier(.4, 0, .2, 1);
119     user-select: none;
120     -webkit-user-select: none;
121     touch-action: manipulation;
122     box-shadow: -6px 8px 10px rgba(81,41,10,0.1),0px 2px 2px rgba(81,41,10,0.2);
123     display: none;
```

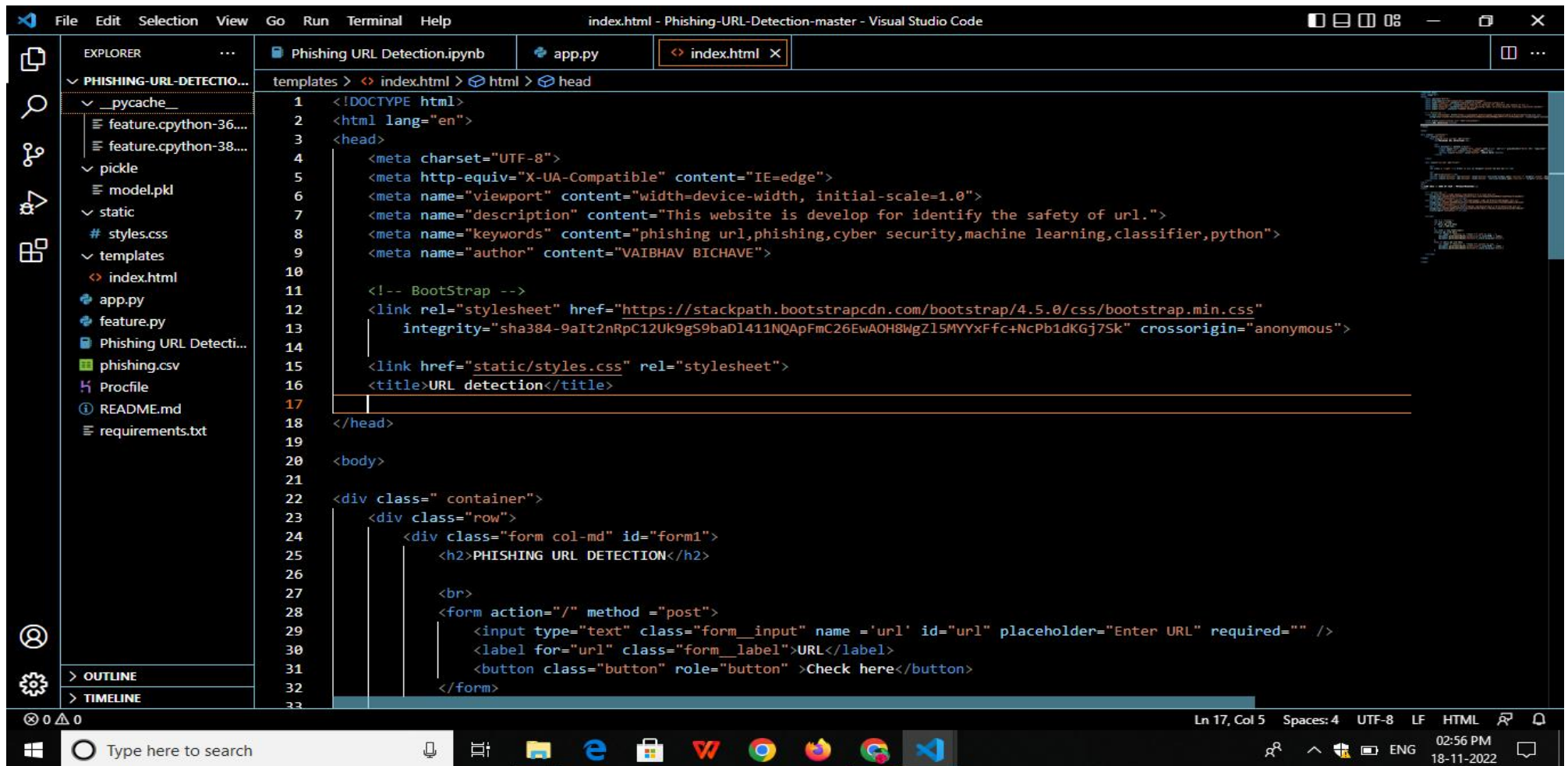
This screenshot shows the continuation of the CSS code in the 'styles.css' file. It defines the styles for the button when it is active or focused. The `.button:active` state has a solid background color and a box shadow. The `.button:focus` state has a box shadow with a blue border. The `.main-body` is also defined with flex properties. The `.button1` selector is repeated at the bottom of the visible code.

```
78
79
80 .button:active {
81     background-color: #f3f4f6;
82     box-shadow: -1px 2px 5px rgba(81,41,10,0.15),0px 1px 1px rgba(81,41,10,0.15);
83     transform: translateY(0.125rem);
84 }
85
86 .button:focus {
87     box-shadow: rgba(72, 35, 7, .46) 0 0 0 4px, -6px 8px 10px rgba(81,41,10,0.1), 0px 2px 2px rgba(81,41,10,0.2);
88 }
89
90
91 .main-body{
92     display: flex;
93     flex-direction: row;
94     width: 75%;
95     justify-content:space-around;
96 }
97
98 .button1{
99     appearance: button;
100     background-color: transparent;
101     background-image: linear-gradient(to bottom, rgb(160, 245, 174), #37ee65);
102     border: 0 solid #e5e7eb;
103     border-radius: .5rem;
```





# Index.html



```
1 <!DOCTYPE html>
2 <html lang="en">
3 <head>
4   <meta charset="UTF-8">
5   <meta http-equiv="X-UA-Compatible" content="IE=edge">
6   <meta name="viewport" content="width=device-width, initial-scale=1.0">
7   <meta name="description" content="This website is develop for identify the safety of url.">
8   <meta name="keywords" content="phishing url,phishing,cyber security,machine learning,classifier,python">
9   <meta name="author" content="VAIBHAV BICHAVE">
10
11   <!-- Bootstrap -->
12   <link rel="stylesheet" href="https://stackpath.bootstrapcdn.com/bootstrap/4.5.0/css/bootstrap.min.css"
13     integrity="sha384-9aIt2nRpC12Uk9gS9baDl411NQApFmC26EwAOH8WgZl5MYYYxFfc+NcPb1dKGj7Sk" crossorigin="anonymous">
14
15   <link href="static/styles.css" rel="stylesheet">
16   <title>URL detection</title>
17
18 </head>
19
20 <body>
21
22   <div class="container">
23     <div class="row">
24       <div class="form col-md" id="form1">
25         <h2>PHISHING URL DETECTION</h2>
26
27         <br>
28         <form action="/" method="post">
29           <input type="text" class="form__input" name='url' id="url" placeholder="Enter URL" required="" />
30           <label for="url" class="form__label">URL</label>
31           <button class="button" role="button">Check here</button>
32         </form>
33       </div>
34     </div>
35   </div>
36 </body>
37 </html>
```

FileEditSelectionViewGoRunTerminalHelpindex.html - Phishing-URL-Detection-master - Visual Studio Code

EXPLORER

PHISHING-URL-DETECTIO...

- \_pycache\_
  - feature.cpython-36...
  - feature.cpython-38...
- pickle
  - model.pkl
- static
  - styles.css
- templates
  - index.html
- app.py
- feature.py
- Phishing URL Detecti...
- phishing.csv
- Procfile
- README.md
- requirements.txt

OUTLINE

TIMELINE

Phishing URL Detection.ipynbapp.pyindex.html

templates > index.html > html > body > div..container

36<div class="col-md" id="form2">

37

38<br>

39<h6 class = "right "><a href= {{ url }} target="\_blank">{{ url }}</a></h6>

40

41<br>

42<h3 id="prediction"></h3>

43<button class="button2" id="button2" role="button" onclick="window.open('{{url}}')" target="\_blank" >Still want to

44<button class="button1" id="button1" role="button" onclick="window.open('{{url}}')" target="\_blank">Continue</butt

45</div>

46</div>

47<br>

48<p>©IBM 2022 -> DONE BY TEAM : PNT2022TMID44601</p>

49</div>

50

51<!-- JavaScript -->

52<script src="https://code.jquery.com/jquery-3.5.1.slim.min.js"

53integrity="sha384-DfXdz2htPH0lsSSs5nCTpuj/zy4C+OGpamoFVy38MVBnE+IbbVYUew+OrCXaRkfj"

54crossorigin="anonymous"></script>

55<script src="https://cdn.jsdelivr.net/npm/popper.js@1.16.0/dist/umd/popper.min.js"

56integrity="sha384-Q6E9RHvbIyZFJoft+2mJbHaEWldlvI9IOYy5n3zV9zzTtmI3UksdQRVvoxMfooAo"

57crossorigin="anonymous"></script>

58<script src="https://stackpath.bootstrapcdn.com/bootstrap/4.5.0/js/bootstrap.min.js"

59integrity="sha384-OgVRvuATP1z7JjHLku0U7Xw704+h835Lr+6QL9UvYjZE3Ipu6Tp75j7Bh/kR0JKI"

60crossorigin="anonymous"></script>

61

62

63<script>

64

65let x = '{{xx}}';

66let num = x\*100;

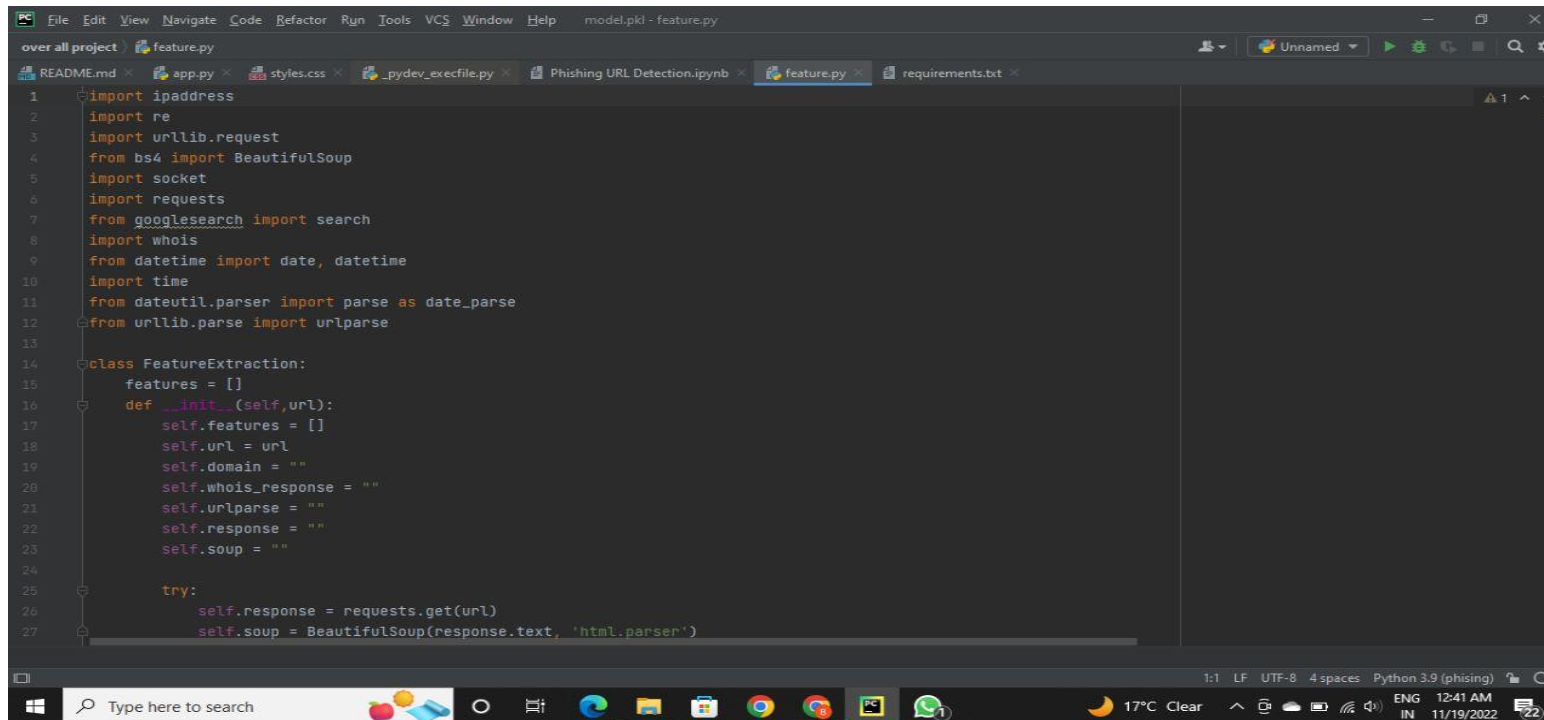
67if (0<=x && x<0.50){

68num = 100-num;

Ln 49, Col 7Spaces: 4UTF-8LFHTML

03:02 PM18-11-2022

# Feature.py



```
File Edit View Navigate Code Refactor Run Tools VCS Window Help model.pkl - feature.py
over all project feature.py
README.md app.py styles.css _pydev_execfile.py Phishing URL Detection.ipynb feature.py requirements.txt
1 import ipaddress
2 import re
3 import urllib.request
4 from bs4 import BeautifulSoup
5 import socket
6 import requests
7 from googlesearch import search
8 import whois
9 from datetime import date, datetime
10 import time
11 from dateutil.parser import parse as date_parse
12 from urllib.parse import urlparse
13
14 class FeatureExtraction:
15     features = []
16     def __init__(self,url):
17         self.features = []
18         self.url = url
19         self.domain = ""
20         self.whois_response = ""
21         self.urlparse = ""
22         self.response = ""
23         self.soup = ""
24
25     try:
26         self.response = requests.get(url)
27         self.soup = BeautifulSoup(response.text, 'html.parser')
```

The screenshot shows an IDE window titled 'model.pkl - feature.py'. The editor displays the first 64 lines of a Python script. The script starts with a GET request to a URL, followed by a BeautifulSoup parser. It then attempts to parse the URL and perform a WHOIS lookup. Finally, it appends several features to a list named 'self.features'.

```
26 self.response = requests.get(url)
27 self.soup = BeautifulSoup(response.text, 'html.parser')
28 except:
29     pass
30
31 try:
32     self.urlparse = urlparse(url)
33     self.domain = self.urlparse.netloc
34 except:
35     pass
36
37 try:
38     self.whois_response = whois.whois(self.domain)
39 except:
40     pass
41
42
43
44
45 self.features.append(self.UsingIp())
46 self.features.append(self.longUrl())
47 self.features.append(self.shortUrl())
48 self.features.append(self.symbol())
49 self.features.append(self.redirecting())
50 self.features.append(self.prefixSuffix())
51 self.features.append(self.SubDomains())
52 self.features.append(self.Hppts())
53 self.features.append(self.DomainRegLen())
54 self.features.append(self.Favicon())
55
56
57 self.features.append(self.NonStdPort())
58 self.features.append(self.HTTPSDomainURL())
59 self.features.append(self.RequestURL())
60 self.features.append(self.AnchorURL())
61 self.features.append(self.LinksInScriptTags())
62 self.features.append(self.ServerFormHandler())
63 self.features.append(self.InfoEmail())
64 self.features.append(self.AbnormalURL())
65 self.features.append(self.WebsiteForwarding())
66 self.features.append(self.StatusBarCust())
67
68 self.features.append(self.DisableRightClick())
69 self.features.append(self.UsingPopupWindow())
70 self.features.append(self.IframeRedirection())
71 self.features.append(self.AgeofDomain())
```

The screenshot shows the same IDE window, but the editor now displays lines 45 through 71 of the Python script. This section continues the list of features being appended to 'self.features'.

```
45 self.features.append(self.UsingIp())
46 self.features.append(self.longUrl())
47 self.features.append(self.shortUrl())
48 self.features.append(self.symbol())
49 self.features.append(self.redirecting())
50 self.features.append(self.prefixSuffix())
51 self.features.append(self.SubDomains())
52 self.features.append(self.Hppts())
53 self.features.append(self.DomainRegLen())
54 self.features.append(self.Favicon())
55
56
57 self.features.append(self.NonStdPort())
58 self.features.append(self.HTTPSDomainURL())
59 self.features.append(self.RequestURL())
60 self.features.append(self.AnchorURL())
61 self.features.append(self.LinksInScriptTags())
62 self.features.append(self.ServerFormHandler())
63 self.features.append(self.InfoEmail())
64 self.features.append(self.AbnormalURL())
65 self.features.append(self.WebsiteForwarding())
66 self.features.append(self.StatusBarCust())
67
68 self.features.append(self.DisableRightClick())
69 self.features.append(self.UsingPopupWindow())
70 self.features.append(self.IframeRedirection())
71 self.features.append(self.AgeofDomain())
```



The screenshot shows a code editor with the following Python code:

```
69 self.features.append(self.UsingPopupWindow())
70 self.features.append(self.IframeRedirection())
71 self.features.append(self.AgeofDomain())
72 self.features.append(self.DNSRecording())
73 self.features.append(self.WebsiteTraffic())
74 self.features.append(self.PageRank())
75 self.features.append(self.GoogleIndex())
76 self.features.append(self.LinksPointingToPage())
77 self.features.append(self.StatsReport())
78
79
80 # 1.UsingIp
81 def UsingIp(self):
82     try:
83         ipaddress.ip_address(self.url)
84         return -1
85     except:
86         return 1
87
88 # 2.longUrl
89 def longUrl(self):
90     if len(self.url) < 54:
91         return 1
92     if len(self.url) >= 54 and len(self.url) <= 75:
93         return 0
94     return -1
95
```

The IDE interface includes a menu bar (File, Edit, View, etc.), a toolbar with icons for running and saving, and a Windows taskbar at the bottom with various application icons and system information (17°C, 12:41 AM, 11/19/2022).

The screenshot shows the same code editor with the following Python code:

```
44 self.features.append(self.UsingIp())
45 self.features.append(self.longUrl())
46 self.features.append(self.shortUrl())
47 self.features.append(self.symbol())
48 self.features.append(self.redirecting())
49 self.features.append(self.prefixSuffix())
50 self.features.append(self.SubDomains())
51 self.features.append(self.Hppts())
52 self.features.append(self.DomainRegLen())
53 self.features.append(self.Favicon())
54
55
56 self.features.append(self.NonStdPort())
57 self.features.append(self.HTTPSDomainURL())
58 self.features.append(self.RequestURL())
59 self.features.append(self.AnchorURL())
60 self.features.append(self.LinksInScriptTags())
61 self.features.append(self.ServerFormHandler())
62 self.features.append(self.InfoEmail())
63 self.features.append(self.AbnormalURL())
64 self.features.append(self.WebsiteForwarding())
65 self.features.append(self.StatusBarCust())
66
67
68 self.features.append(self.DisableRightClick())
69 self.features.append(self.UsingPopupWindow())

```

A yellow error banner at the top of the editor area states: "Package requirements 'numpy==1.21.4', 'pandas==1.3.4', 'scikit\_learn==1.0.1' are not satisfied". To the right of the banner are buttons for "Install requirements" and "Ignore requirements". The IDE interface and Windows taskbar are consistent with the first screenshot.

The image shows a Visual Studio Code editor window with a Python file named 'feature.py' open. The editor is displaying a script for 'Phishing URL Detection'. The script includes three methods: 'redirecting', 'prefixSuffix', and 'SubDomains'. The 'redirecting' method checks if a URL contains a redirect sequence. The 'prefixSuffix' method checks if a URL has a valid domain prefix. The 'SubDomains' method checks if a URL contains subdomains. The interface includes a file explorer on the left, a search bar, and a taskbar at the bottom. The taskbar shows the Windows Start button, a search bar, and several application icons including VS Code, Chrome, and a terminal. The system tray at the bottom right shows the date and time as 11/19/2022, 12:42 AM.

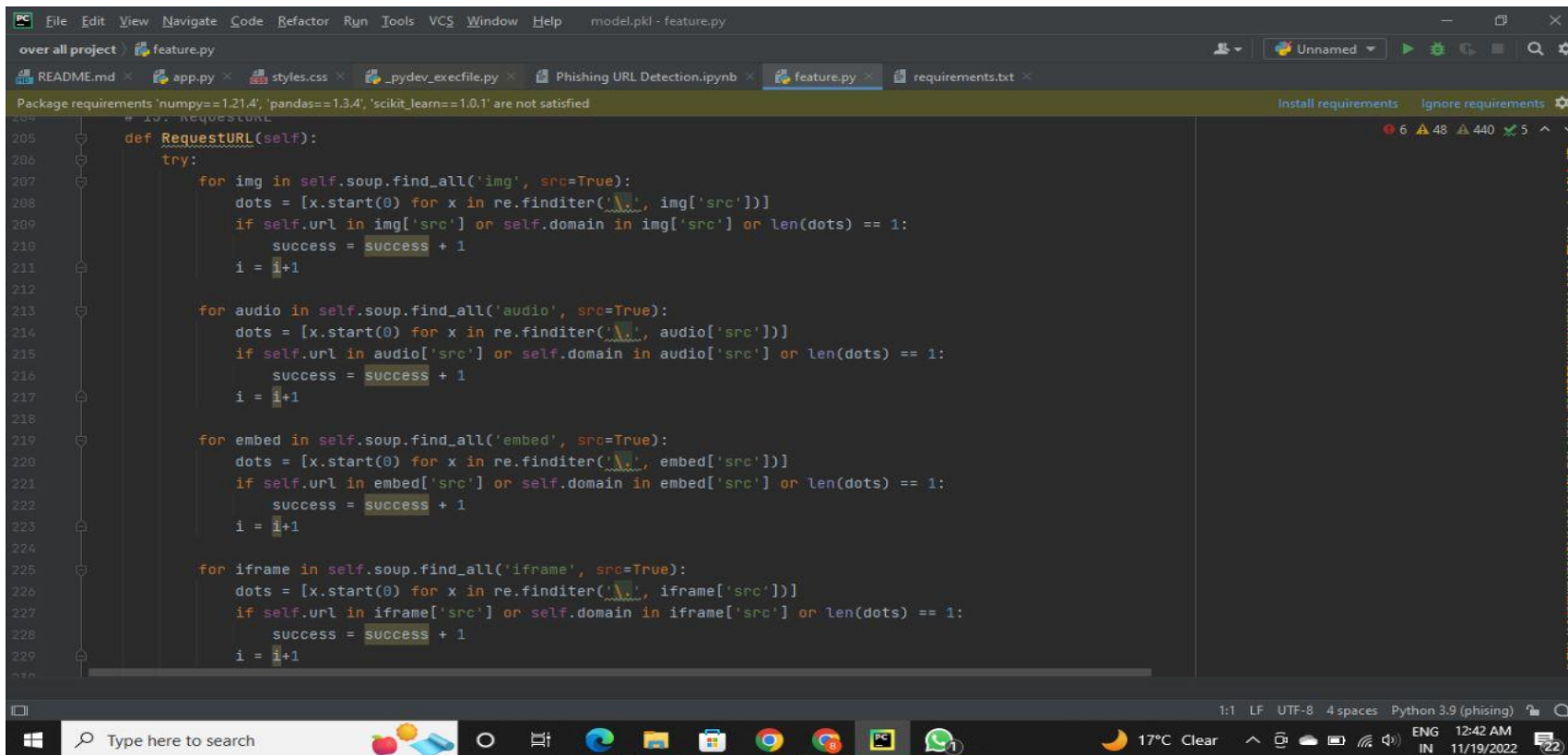
The image shows a screenshot of a Visual Studio Code (VS Code) editor window. The top menu bar includes File, Edit, View, Navigate, Code, Refactor, Tools, VCS, Window, Help, and model.pkl - feature.py. The file explorer on the left shows a project named 'feature.py' with files README.md, app.py, styles.css, \_pydev\_execfile.py, Phishing URL Detection.ipynb, feature.py, and requirements.txt. The search bar is empty. The main editor area displays the content of 'feature.py', which is a Python script for URL validation. The script includes a package requirements section at the top, followed by three methods: UsingIp, longUrl, and shortUrl. The UsingIp method checks if the URL contains an IP address. The longUrl method checks if the URL length is greater than 54 and less than or equal to 75. The shortUrl method checks if the URL contains any of a list of common short URL domain suffixes. The taskbar at the bottom shows the Windows Start button, a search bar, and several application icons including VS Code, Chrome, and a terminal. The system tray on the right shows the date and time as 12:42 AM on 11/19/2022.

```
PC File Edit View Navigate Code Refactor Run Tools VCS Window Help model.pkl - feature.py
over all project feature.py
README.md app.py styles.css _pydev_execfile.py Phishing URL Detection.ipynb feature.py requirements.txt
Package requirements 'numpy==1.21.4', 'pandas==1.3.4', 'scikit_learn==1.0.1' are not satisfied Install requirements ignore requirements
139
140 # 8.HTTPS
141 def Hhttps(self):
142     try:
143         https = self.urlparse.scheme
144         if 'https' in https:
145             return 1
146         return -1
147     except:
148         return 1
149
150 # 9.DomainRegLen
151 def DomainRegLen(self):
152     try:
153         expiration_date = self.whois_response.expiration_date
154         creation_date = self.whois_response.creation_date
155         try:
156             if(len(expiration_date)):
157                 expiration_date = expiration_date[0]
158             except:
159                 pass
160         try:
161             if(len(creation_date)):
162                 creation_date = creation_date[0]
163             except:
164                 pass
```

17°C Clear 12:42 AM 11/19/2022

```
180         return -1
181     except:
182         return -1
183
184 # 11. NonStdPort
185 def NonStdPort(self):
186     try:
187         port = self.domain.split(":")
188         if len(port)>1:
189             return -1
190         return 1
191     except:
192         return -1
193
194 # 12. HTTPSDomainURL
195 def HTTPSDomainURL(self):
196     try:
197         if 'https' in self.domain:
198             return -1
199         return 1
200     except:
201         return -1
202
203 # 13. RequestURL
204 def RequestURL(self):
```

17°C Clear 12:42 AM 11/19/2022





## TEST CASES

[illegible]

# User Acceptance Testing

## 1.Purpose of Document

The purpose of this document is to briefly explain the test coverage and open issues of the [ProductName] project at the time of the release to User Acceptance Testing (UAT).

## 2.Defect Analysis

This report shows the number of resolved or closed bugs at each severity level, and how they were resolved

Resolution	Severit y 1	Severit y 2	Severit y 3	Severit y 4	Subtota l
By Design	14	7	4	3	28
Duplicate	1	0	3	0	4
External	2	3	0	1	6
Fixed	15	5	3	20	43
Not Reproduced	0	0	1	0	1
Skipped	0	2	1	1	4
Won't Fix	0	1	1	2	4
Totals	32	18	13	27	90

3.Test Case Analysis

This report shows the number of test cases that have passed, failed, and untested

Section	Total Cases	Not Tested	Fail	Pas s
Print Engine	6	0	0	6
Client Application	40	0	2	38
Security	4	0	0	4
Outsource Shipping	3	1	0	2
Exception Reporting	7	1	0	6
Final Report Output	5	0	0	5
Version Control	3	0	1	2

# RESULTS

## Performance Matrics

### Model Performance Testing:

Project team shall fill the following information in model performance testing template.

S.No.	Parameter	Value s
1.	Metrics	<b>Random Forest Classifier Accuracy score-96.653</b>
2.	Tune the Model	<b>Hyperparameter Tuning - Validation Method -</b>

#### 14. METRICS

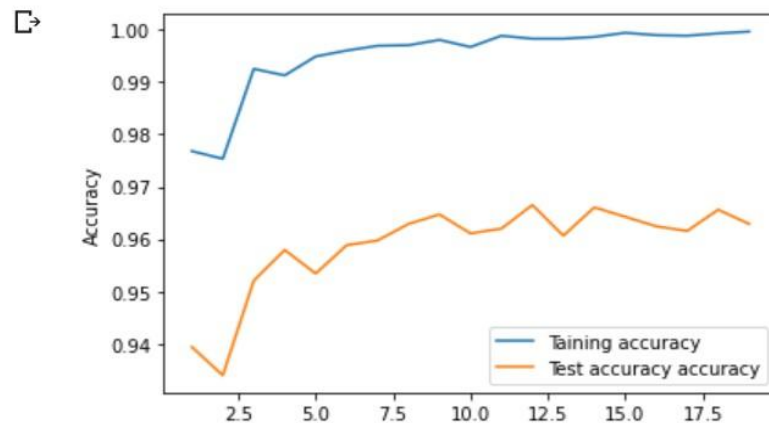
##### Classification Report:

```
[ ] #classification report of Randomm Forest model  
print(metrics.classification_report(y_test,y_test_rf))
```

	precision	recall	f1-score	support
-1	0.98	0.95	0.96	1014
1	0.96	0.98	0.97	1197
accuracy			0.97	2211
macro avg	0.97	0.97	0.97	2211
weighted avg	0.97	0.97	0.97	2211

# Peformance

```
▶ training_accuracy=[]  
test_accuracy=[]  
depth=range(1,20)  
for n in depth:  
    rf_test=RandomForestClassifier(n_estimators=n)  
    rf_test.fit(x_train,y_train)  
    training_accuracy.append(rf_test.score(x_train,y_train))  
    test_accuracy.append(rf_test.score(x_test,y_test))  
plt.figure(figsize=None)  
plt.plot(depth,training_accuracy,label="Taining accuracy")  
plt.plot(depth,test_accuracy,label="Test accuracy accuracy")  
plt.ylabel("Accuracy")  
plt.xlabel("max_depth")  
plt.legend();
```



✓ 1s completed at 1:37 PM

## 2.Tune the model

```
[ ]
```

	ML Model	Accuracy	f1_score	Recall	Precision
0	Logistic Regression	91.814	92.567	94.496	94.496
1	Random Forest	96.653	96.942	100.000	100.000
2	XgbClassifier	94.754	95.207	96.714	96.714
3	Decision tree	95.206	95.605	100.000	100.000

```
[ ] sorted_result=result.sort_values(by=['Accuracy', 'f1_score'],ascending=False).reset_index(drop=True)
sorted_result
```

	ML Model	Accuracy	f1_score	Recall	Precision
0	Random Forest	96.653	96.942	100.000	100.000
1	Decision tree	95.206	95.605	100.000	100.000
2	XgbClassifier	94.754	95.207	96.714	96.714
3	Logistic Regression	91.814	92.567	94.496	94.496

# ADVANTAGES & DISADVANTAGES

## ADVANTAGES:

- **Increases user alertness to phishing risks** Whenever the user navigates into the website and provide the URL of the website that needs to be verified for legitimacy, the system detects phishing sites by applying a machine learning algorithm which implements classification algorithms and techniques to extract the phishing datasets criteria to classify their legitimacy which in turn helps the customers to eliminate the risks of cyber threat and protect their valuable corporate or personal data.
- **Users will also be able to pose any query to the admin through the report page designed** Our system is also provided with an option for the clients to report to the administrator which helps them to ask their questions significantly improving their experience on our site.

## DISADVANTAGES:

- Not a generalized model
- Huge number of rules
- Needs feed continuously

## CHAPTER 11

### CONCLUSION

Phishing detection is now an area of great interest among the researchers due to its significance in protecting privacy and providing security. There are many methods to perform phishing detection. Our system aims to enhance the detection method to detect phishing websites using machine learning technology. We achieved a high detection accuracy, and the results show that the classifiers give better performance when we use more data as training data.

In future, hybrid technology will be implemented to detect phishing websites more accurately.

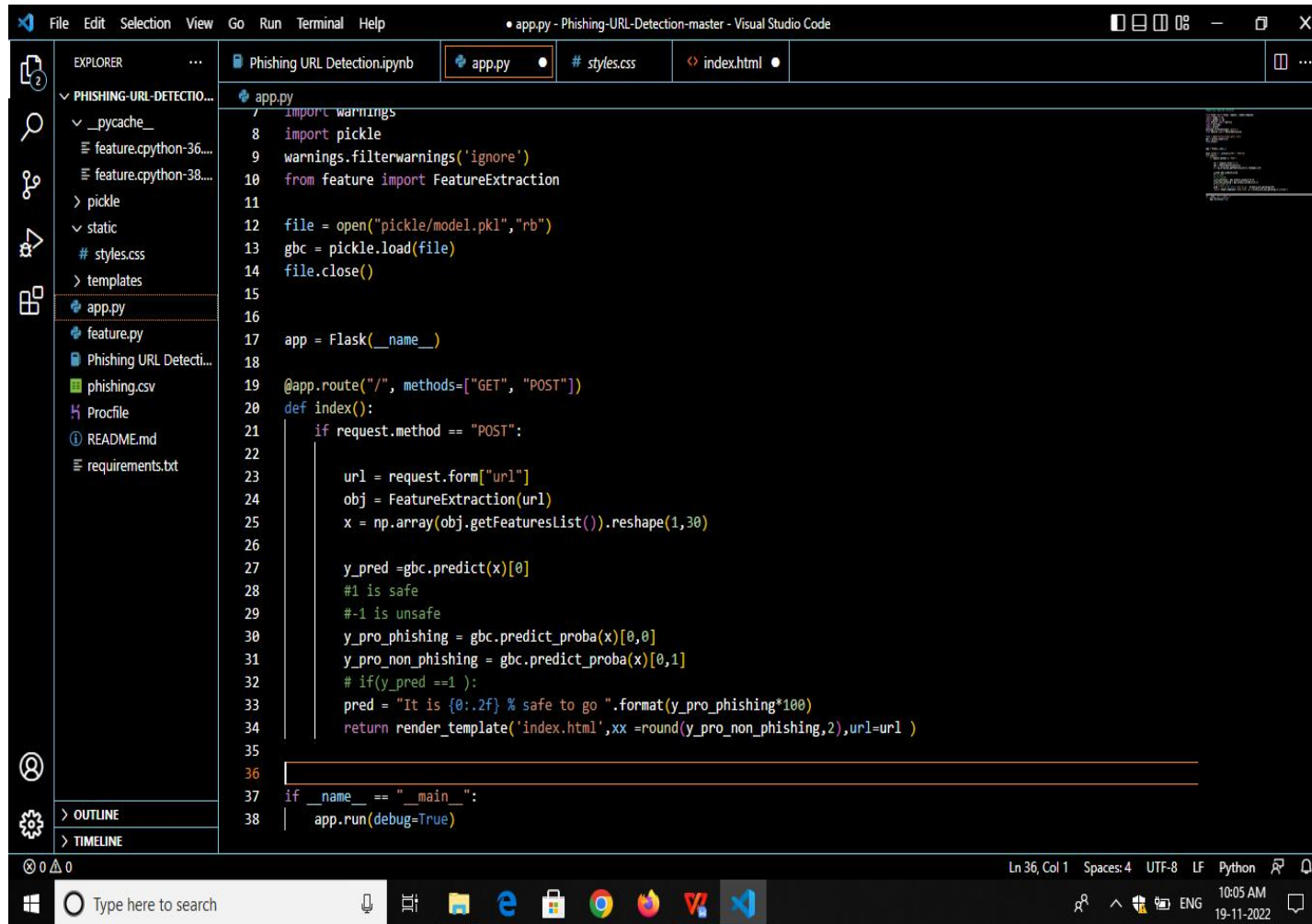


## FUTURE SCOPE

In future we intend to build an add-ons for our system and if we get a structured dataset of phishing, we can perform phishing detection much faster than any other technique. We can also use a combination of any two or more classifiers to get maximum accuracy. We plan to explore various phishing techniques which use Network based features, Content based features, Webpage based features and HTML and JavaScript features of web pages which will improve the performance of the system. In particular, we extract features from URLs and pass it through the various classifiers.

# APPENDIX

## Source Code



```
File Edit Selection View Go Run Terminal Help • app.py - Phishing-URL-Detection-master - Visual Studio Code

EXPLORER
PHISHING-URL-DETECTIO...
  _pycache_
  feature.cpython-36...
  feature.cpython-38...
  pickle
  static
  # styles.css
  > templates
  app.py
  feature.py
  Phishing URL Detecti...
  phishing.csv
  Profile
  README.md
  requirements.txt

OUTLINE
TIMELINE

app.py
7 import warnings
8 import pickle
9 warnings.filterwarnings('ignore')
10 from feature import FeatureExtraction
11
12 file = open("pickle/model.pkl", "rb")
13 gbc = pickle.load(file)
14 file.close()
15
16 app = Flask(__name__)
17
18 @app.route("/", methods=["GET", "POST"])
19 def index():
20     if request.method == "POST":
21
22         url = request.form["url"]
23         obj = FeatureExtraction(url)
24         x = np.array(obj.getFeaturesList()).reshape(1,30)
25
26         y_pred = gbc.predict(x)[0]
27         #1 is safe
28         #-1 is unsafe
29         y_pro_phishing = gbc.predict_proba(x)[0,0]
30         y_pro_non_phishing = gbc.predict_proba(x)[0,1]
31         # if(y_pred ==1 ):
32         pred = "It is {0:.2f} % safe to go ".format(y_pro_phishing*100)
33         return render_template("index.html",xx =round(y_pro_non_phishing,2),url=url )
34
35
36
37 if __name__ == "__main__":
38     app.run(debug=True)
```

Ln 36, Col 1 Spaces: 4 UTF-8 LF Python 10:05 AM 19-11-2022