



WEBPHISHINGDETECTION



Domain : Applied Data

ScienceA PROJECTREPORT

Submittedby

NAVEENKUMAR.S (920819104025)

AJAYKUMAR.M (920819104003)

SRIVATHSKARTHIC.G (920819104042)

GUNA SEAKAR.J (920819104012)

inpartialfulfillmentfortheawardofthedegree

of

BACHELOROFENGINEERING

in

COMPUTERSCIENCEANDENGINEERING

NPRCOLLEGE OFENGINEERING&TECHNOLOGY

ATHAM,DINDIGUL.

ANNAUNIVERSITY::CHENNAI600025

TABLE OF CONTENTS

CHAPTER NO.	TITLE	PAGE NO.
1.	INTRODUCTION	
	1.1 Overview	1
	1.2 Purpose	2
2.	LITERATURE SURVEY	
	2.1 Existing problem	3
	2.2 References	4
	2.3 Problem Statement Definition	5
3.	IDEATION & PROPOSED SOLUTION	
	3.1 Empathy Map Canvas	8
	3.2 Ideation & Brainstorming	9
	3.3 Proposed Solution	13
	3.4 Problem Solution Fit	14
4.	REQUIREMENT ANALYSIS	
	4.1 Functional requirement	15
	4.2 Non-Functional requirements	16

5.	PROJECTDESIGN	
	5.1 DataFlowDiagram	17
	5.2 Solution&TechnicalArchitecture	20
	5.3 UserStories	22
6.	PROJECTPLANNING&SCHEDULING	
	6.1 SprintPlanning&Estimation	23
	6.2 SprintDeliverySchedule	23
	6.3 ReportsfromJIRA	
7.	CODING&SOLUTIONING	
	7.1 Feature1	26
	7.2 Feature2	27
8.	TESTING	
	8.1 TestCases	28
	8.2 UserAcceptanceTesting	29
9.	RESULTS	
	9.1 PerformanceMetrics	29
10.	CONCLUSION	25
11.	FUTURE SCOPE	25

12.	APPENDIX	30
	SOURCECODE	30
	GITHUB&PROJECTDEMO LINK	30

CHAPTER 1

INTRODUCTION

PROJECT OVERVIEW:

Phishing is a form of fraud in which the attacker tries to learn sensitive information such as login credentials (or) account information by sending a reputable entity (or) person in email (or) other communication channels. Typically a victim receives a message that appears to have been sent by a known contact (or) organization. Phishing is popular among attackers, it is easier to trick someone into clicking a malicious link which seems legitimate than typing to break through a computer's defences, systems. So regarding this solution detecting phishing domains is a classification problem. So we need labeled data which has sample as phish. Domains & legitimate domains in the training plans. Collecting legitimate domains is another problem. For this purpose, sites reputation services are among used. This service analysis & Rank, available websites, 4 decision can be considered as an improved waste of other structures. Which segregates legitimates of phishing websites. So by implementing this. We could develop a noble Machine learning tool for phishing detection.

PURPOSE:

Our purpose is to develop a tool that is capable of detecting phishing websites.

This improves the security status of the user's machine by using the legitimate website could provide a safe user's experience. So by using our tool the user's data could be prevented from intruders.

We have included decision tree algorithm to recognize whether website is phishing or legitimate along with this random forest algorithm is used to improve the accuracy. So as a combination of such facilities we provide a secure internet experience for the end consumers.

CHAPTER 2

LITERATURE

SURVEY

- Detection of phishing websites by using machine learning based URL Analysis Helmetkrmkrmz, Banypiri, Dzguskaraysahingoz.

This approach is URL Analysis k- nearest neighbourhood (kNN), support vector machine (VMM), Decision Tree (DT) with merits comprising of effective prediction of phishing domains based on logistic regression. Using hybrid algorithm the accuracy is enhanced. The main drawbacks are time consumption of high processing power.

- Phishing website detection based on machine learning algorithm by weighbhai.

This is a phishing detection based on machine learning & logistic regression. The benefits include high threshold value due to which accuracy increases & By using logistic regression. The speed can be improvised and regarding drawbacks data processing is required.

2.2.REFERENCES:

TITLE	YEAR	AUTHORS	TECHNOLOGY ADOPTED	MERITS	DEMERITS
Real Time Detection of Phishing Websites	2016	ABDULGHANI ALI AHMED, NURUL AMIRAH ABDULLAH	Checking Uniform Resources Locators (URLs)	Increased accuracy by 32% than previous proposed systems. Identifies URL redirecting to other web pages by analyzing the URL type.	Accuracy depends on heuristic band and depends on discriminative features. Only checks validity of URLs.
Detection of Phishing Websites by Using Machine Learning-Based URL Analysis	2020	Mehmet Korkmaz, Ozgur Koray Sahingoz, Banu Diri	Machine Learning-Based URL Analysis (Logistic Regression (LR), K-Nearest Neighborhood (KNN), Support Vector Machine (SVM), Decision Tree (DT), Naive Bayes (NB), XGBoost, Random Forest (RF) and Artificial Neural Network (ANN))	Logistic Regression is the algorithm that generates effective predictions of phishing domain elements. Hybrid algorithms enhance the accuracy.	Takes time to identify the phishing website. Requires high processing power.
Phishing Website Detection Based on Machine Learning Algorithm	2020	Weiheng Bai	Machine Learning Algorithm (Logistic regression classifier)	With high threshold value, it has high accuracy rate. Logistic regression classifies being used improves the speed.	Not completely reliable because transmission of packets does not reflect the proximity of location . Data preprocessing is required.
Machine Learning Techniques for Detection of Website Phishing: A Review for Promises and Challenges	2021	Ammar Odeh, Ismail Keshta, Eman Abdelfattah	Machine learning, Deep learning (Heuristic and automated techniques)	Use of Machine learning improvised the URLs. Stacking model has improved the accuracy by detecting legitimate websites.	Large binds of datasets are difficult to handle. Low accuracy and hypertuning.

2.3 Problemstatement:

When a person asks the internet he/she might be concerned with the securitybeing provided. So when that person opens a link/website & enters their credentials ,he/shegetshacked(or)phished.

So in order to avoid this we have developed a tool backed up with machinelearning which uses certain criteria to evaluates this URL & declares a websitesasphisjhing(or)legitimateone.

IdeationPhase DefinetheProblemStatements

Date	19September2022
TeamID	PNT2022TMID48632
ProjectName	WebPhishingDetection
MaximumMarks	2Marks

CustomerProblemStatementTemplate:

Create aproblem statementto understandyour customer's pointof view. The Customer Problem Statement template helps you focus on what matterstocreateexperiencespeoplewilllove.

A well-articulated customer problem statement allows you andyourteamtofindtheidealsolutionforthechallengesyourcustomersface.Throughout the process, you'll also be able to empathize with your customers, whichhelpsyoubetterunderstandhowtheyperceiveyourproductorservice.



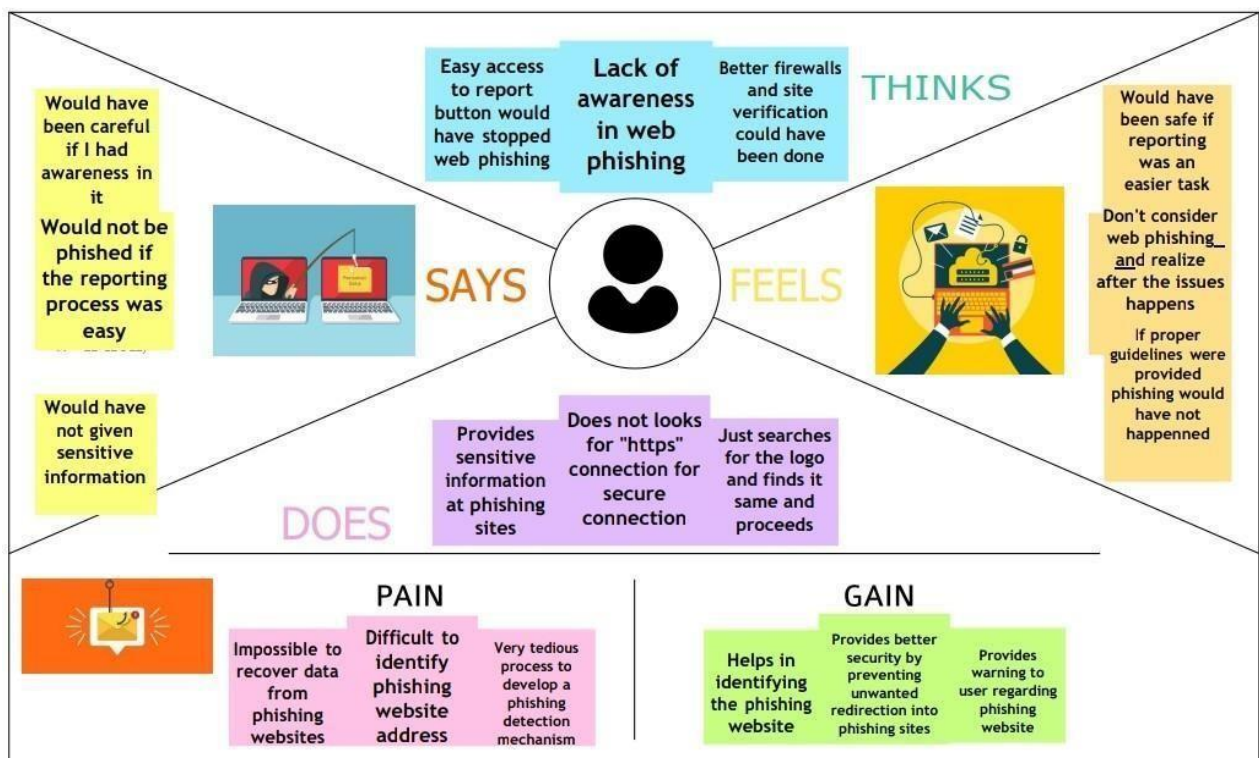
ProblemState ment(PS)	Iam(Customer)	I'mtryingto	But	Because	Whichmake smefeel
PS-1	Student	Opens acompromisedlinkan d enters thecredentials	GetsHacked	Itwasafakewebs ite	Insecure
PS-2	Shopkeeper	Opensalinkfromab ank and enterstheaccountd etails	His bankaccou ntdetails andhis creditcard detailsgetsh acked	The link wasfakeand thehacker aims forthemoneyfro mtheaccount	Untrustworthy
PS-3	Homemaker	Maketransactio ns forthe purchaseditem sfromcheckout s	Ended up asshepurcha sesfromanfa kewebsite	The website hadanadvertis ementfordiscount s	Disappointed
PS-4	Employer	Complete thetasksbyendo ftheday	Found a filethroughm ailandopensi t,andhiscom panydetails getshacked	The fileseemedto be phishingandfak eone	Annoyedanddist urbed

CHAPTER 3 IDEATION PHASE

Empathymap:

An **empathymap** is a collaborative visualization used to articulate what we know about a particular type of user. It externalizes knowledge about users in order to

- 1) Create a shared understanding of user needs, and
- 2) Aid in decision making.



Explanation:

✓ What do they think & feel?

- Feels safe to use internet.
- Avoid stress due to data loss.
- Eliminates forgery.

✓ What do they see ?

- Machine learning tool that detects the URL
- Popups saying safe/danger to use.

✓ What do they see & do ?

- Safe browsing environment.
- Stress free surfing on the internet.

✓ What do they hear?

- They might be a basic software pre-built on all devices.

Ideation and Brainstorming:

Brainstorming provides a free and open environment that encourages everyone within a team to participate in the creative thinking process that leads to problem solving. Prioritizing volume over value, out-of-the-box ideas are welcome and built upon, and all participants are encouraged to collaborate, helping each other develop a rich amount of creative solutions.

Step-1:TeamGathering,CollaborationandSelectthe ProblemStatement

Template



Brainstorm & idea prioritization

Use this template in your own brainstorming sessions so your team can unleash their imagination and start shaping concepts even if you're not sitting in the same room.

🕒 10 minutes to prepare

🕒 1 hour to collaborate

👤 2-8 people recommended

Share template feedback

➔

Before you collaborate

A little bit of preparation goes a long way with this session. Here's what you need to do to get going.

🕒 10 minutes

A

Team gathering

All the team members are invited to participate in this session so as to gather more knowledge about our project.

B

Set the goal

There are much more road accidents happened in the last year because of uncontrolled speed and road conditions.

C

Learn how to use the facilitation tools

Use the Facilitation Superpowers to run a happy and productive session.

Open article ➔

1

Define your problem statement

The road conditions on our city might be lead to accident and there is some physical sign boards to intimate the speed limit.

🕒 5 minutes

PROBLEM

How might we replace physical sign boards with digital sign boards and help the drivers to avoid accidents due to road conditions?

➔

Key rules of brainstorming

To run an smooth and productive session

Stay in topic.

Encourage wild ideas.

Defer judgment.

Listen to others.

Go for volume.

If possible, be visual.



Need some inspiration?

See a finished version of this template by kickingstart your work.

Open example ➔

Step-2:Brainstorm,IdeaListingandGrouping:

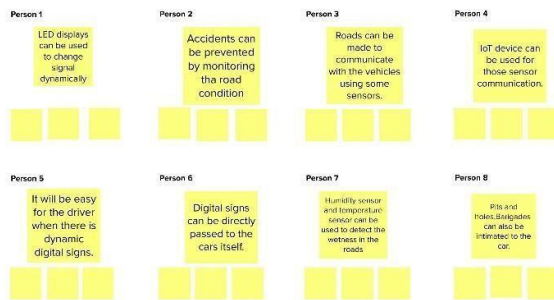
2

Brainstorm

Write down any ideas that come to mind that address your problem statement.

10 minutes

TIP
You can select a sticky note and the first point (arrow to sketch) to start drawing!

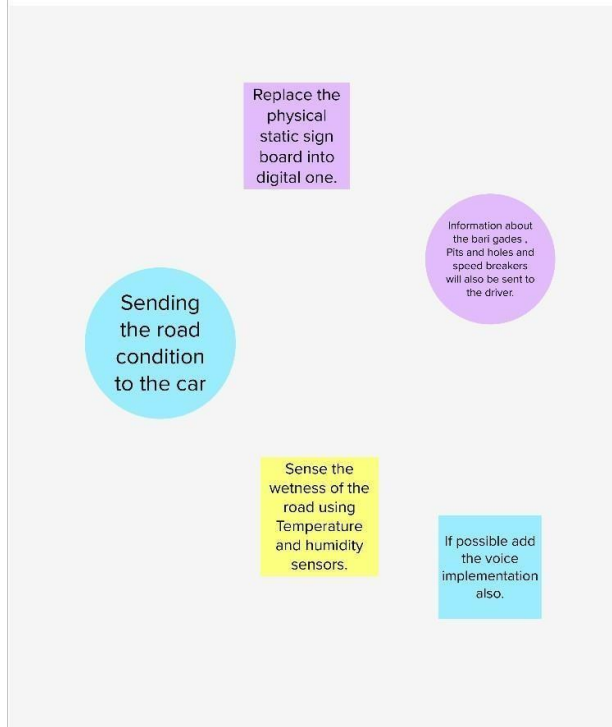


3

Group ideas

Take turns sharing your ideas while clustering similar or related notes as you go. Once all sticky notes have been grouped, give each cluster a sentence-like label. If a cluster is bigger than six sticky notes, try and see if you can break it up into smaller sub-groups.

20 minutes



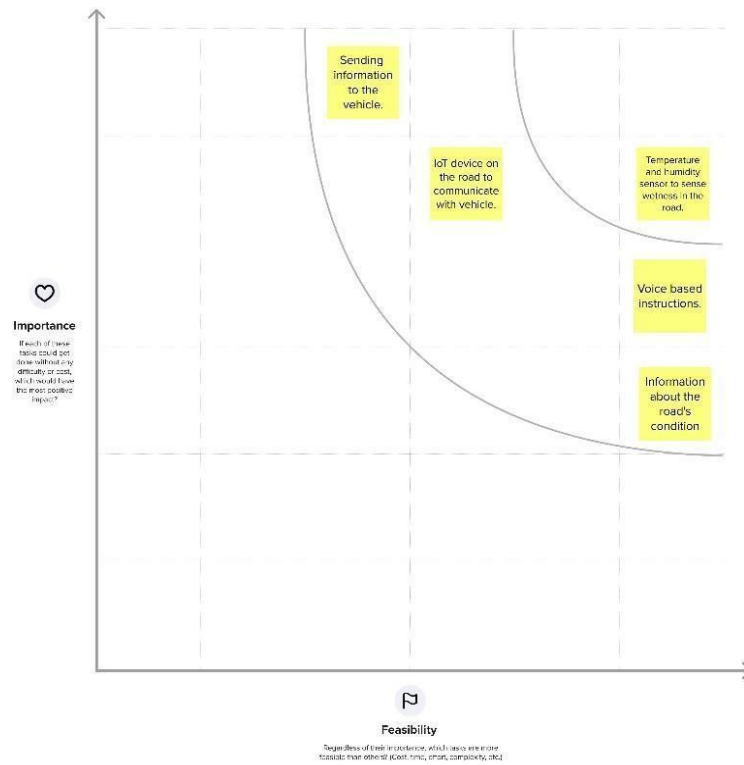
Step-3:IdeaPrioritization

4

Prioritize

Your team should all be on the same page about what's important moving forward. Place your ideas on this grid to determine which ideas are important and which are feasible.

20 minutes



ProposedSolution:

S.No.	Parameter	Description
1.	ProblemStatement (Problem to be solved)	<p>There are a number of users who purchase products online and make payments through e-banking. There are e-banking websites that ask users to provide sensitive data such as username, password & credit card details, etc., often for malicious reasons. Web phishing is one of many security threats to web services on the Internet.</p> <p>Common threats of web phishing:</p> <ul style="list-style-type: none"> • Web phishing emphasizes on steal private information • It will lead to information disclosure and property damage. • Large organizations may get trapped in different kinds of scams.
2.	Idea/Solution description	<p>The proposed system handles the dataset and classifies the dataset to be tested as genuine or not by verifying certain criteria that are necessary to validate it. The datasets are now evaluated where the number of criteria that has been fulfilled plays a vital role. So based on the results that we get from the tools that we used (Decision tree algorithm), the website is declared genuine or not.</p>
3.	Novelty/ Uniqueness	<p>Implementation of Random Forest algorithm along with the decision tree improves the accuracy of the detection.</p>
4.	Social Impact / Customer Satisfaction	<p>The proposed system prevents the user or customer from falling prey to the phishing and scam websites by detecting the phishing websites.</p>
5.	Business Model (Revenue Model)	<p>By including premium subscription facility to provide enhanced features to the customer we generate revenue to the development team.</p>
6.	Scalability of the Solution	<p>In this emerging world of technology, phishing URLs with new scamming methods could be identified using our system.</p>

ProblemSolutionFit:

Project Title: Web Phishing Detection

Project Design Phase-I - Solution Fit Template

Team ID: PNT2022TMID48632

Define CS, fit into CC	1. CUSTOMER SEGMENT(S) CS <p>All internet users</p>	6. CUSTOMER CONSTRAINTS CC <ul style="list-style-type: none"> Difficult to detect exact replica of the actual site. Difficult to scan for malwares in files downloaded from that website. 	5. AVAILABLE SOLUTIONS AS <ul style="list-style-type: none"> Alerts the user that the website redirects to another website. Alerts the user if the website automatically downloads files into the system. 	Explore AS, differentiate

Focus on J&P, up into BE, understand RC	2. JOBS-TO-BE-DONE / PROBLEMS J&P <ul style="list-style-type: none"> Improper firewall protection leads to security threats. Difficult to retrieve data from phishing websites Tedious process to develop an phishing detection software. 	9. PROBLEM ROOT CAUSE RC <p>Certain systems doesn't support the software as it's requirements are based on a baseline.</p>	7. BEHAVIOUR BE <ul style="list-style-type: none"> User pretends to be precautious while surfing on the internet. 	Focus on J&P, up into BE, understand RC

Identify strong TR & EM	3. TRIGGERS TR <ul style="list-style-type: none"> Seamless browsing facilitated with safety Protection from virus and malware 	10. YOUR SOLUTION SI <ul style="list-style-type: none"> Linear regression technique and Decision tree algorithm are used to classify the website based on the dataset that was used to train the AI model Random forest algorithm is used to improvise the accuracy of the detection 	8. CHANNELS of BEHAVIOUR CH <p>ONLINE</p> <ul style="list-style-type: none"> The customer shares his/her feedback on the product through the online portals <p>OFFLINE</p> <ul style="list-style-type: none"> The customer also shares his/her experience on the product to their peer group which seeks popularity for the product 	Identify strong TR & EM
	4. EMOTIONS: BEFORE / AFTER EM <ul style="list-style-type: none"> Before implementation of our projects people weren't aware whether they have been phished, but after which they started browsing without any stress 			

CHAPTER4

REQUIREMENT ANALYSIS

FUNCTIONAL REQUIREMENTS

Functional requirements may involve calculations, technical details, data manipulation and processing, and other specific functionality that define what a system is supposed to accomplish. Behavioral requirements describe all the cases where the system uses the functional requirements, these are captured in use cases.

Following are the functional requirements of the proposed solution.

FRNo.	Functional Requirement(Epic)	SubRequirement(Story/Sub-Task)
FR-1	HomePage	The user finds the homepage easy to navigate and feels comfortable with the user interface
FR-2	Sign up	The user will get authentication over their account by security measures
FR-3	Login	The user can register using either their google account or their mobile number
FR-4	Dashboard	The user can go through the facilities provided by the product
FR-5	Prediction	User would be prompted with a popup indicating the trustfulness of the website
FR-6	Results page	The user would be able to analyze the website whether it's genuine or not
FR-7	Reporting	The user can report for any bugs or ask any queries on the product

NON FUNCTIONAL REQUIREMENTS:

Non-Functional Requirements are the constraints or the requirements imposed on the system. They specify the quality attribute of the software. Non-Functional Requirements deal with issues like scalability, maintainability, performance, portability, security, reliability, and many more.

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

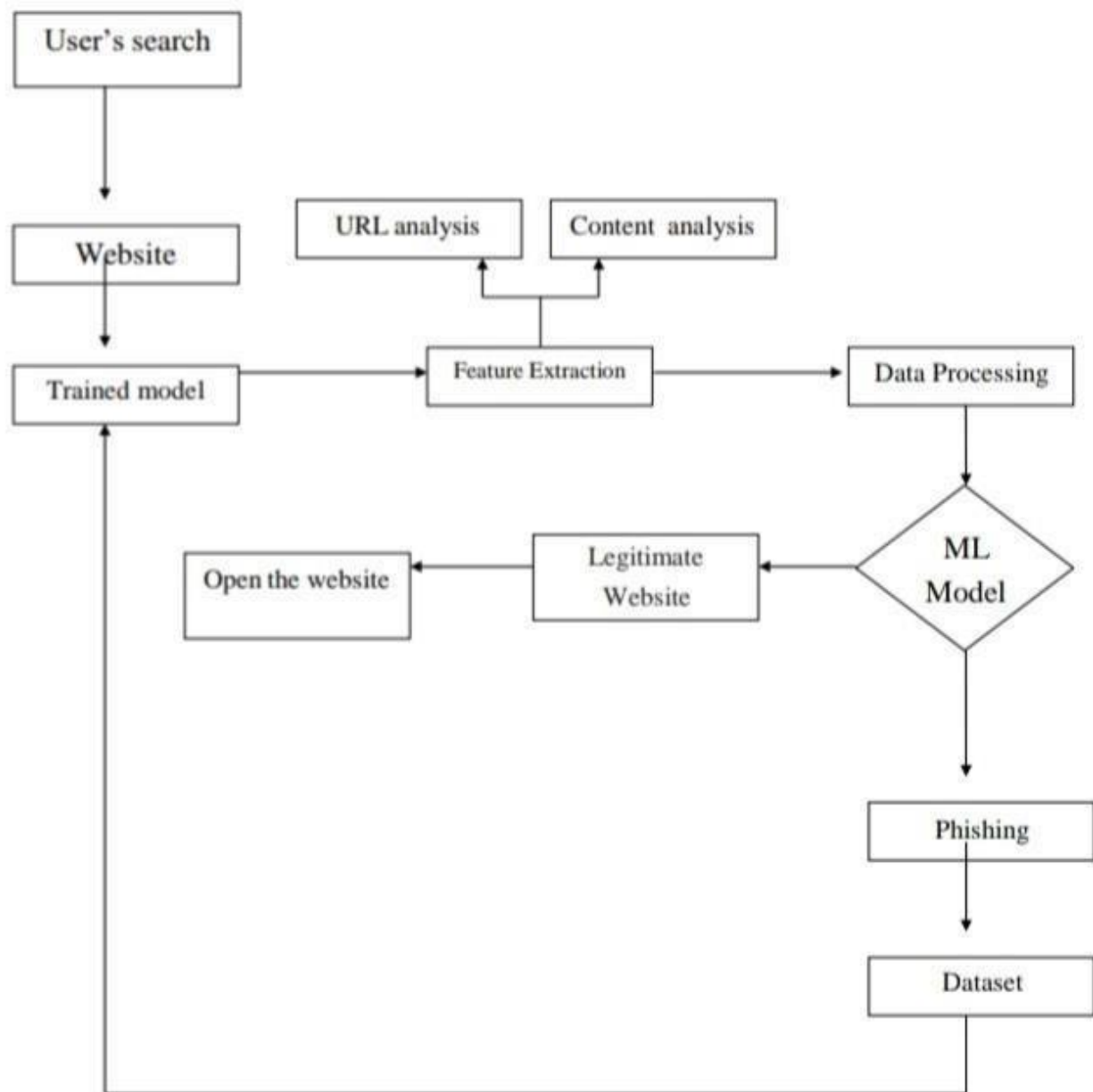
FRNo.	Non-Functional Requirement	Description
NFR-1	Usability	The proposed system is efficient and easy to configure in detecting the phishing websites.
NFR-2	Security	The system is secured as it prevents the user from the unauthorized access and prevents them from the threat.
NFR-3	Reliability	The system will perform the tasks it was supposed to do.
NFR-4	Performance	The system will perform the task efficiently and with a good accuracy rate.
NFR-5	Availability	The proposed system is always available whenever it is required to be executed.
NFR-6	Scalability	The system is scalable to handle the increasing and decreasing workloads.

CHAPTER 5

PROJECT DESIGN

Data Flow Diagrams:

A Data Flow Diagram (DFD) is a traditional visual representation of the information flow 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. A data flow diagram (DFD) maps out the flow of information for any process or system. It uses defined symbols like rectangles, circles and arrows, plus short text labels, to show data inputs, outputs, storage points and the routes between each destination. Data flowcharts can range from simple, even hand-drawn process overviews, to in-depth, multi-level DFDs that dig progressively deeper into how the data is handled. They can be used to analyze an existing system or model a new one. Like all the best diagrams and charts, a DFD can often visually “say” things that would be hard to explain in words, and they work for both technical and non-technical audiences, from developer to CEO.



User Stories:

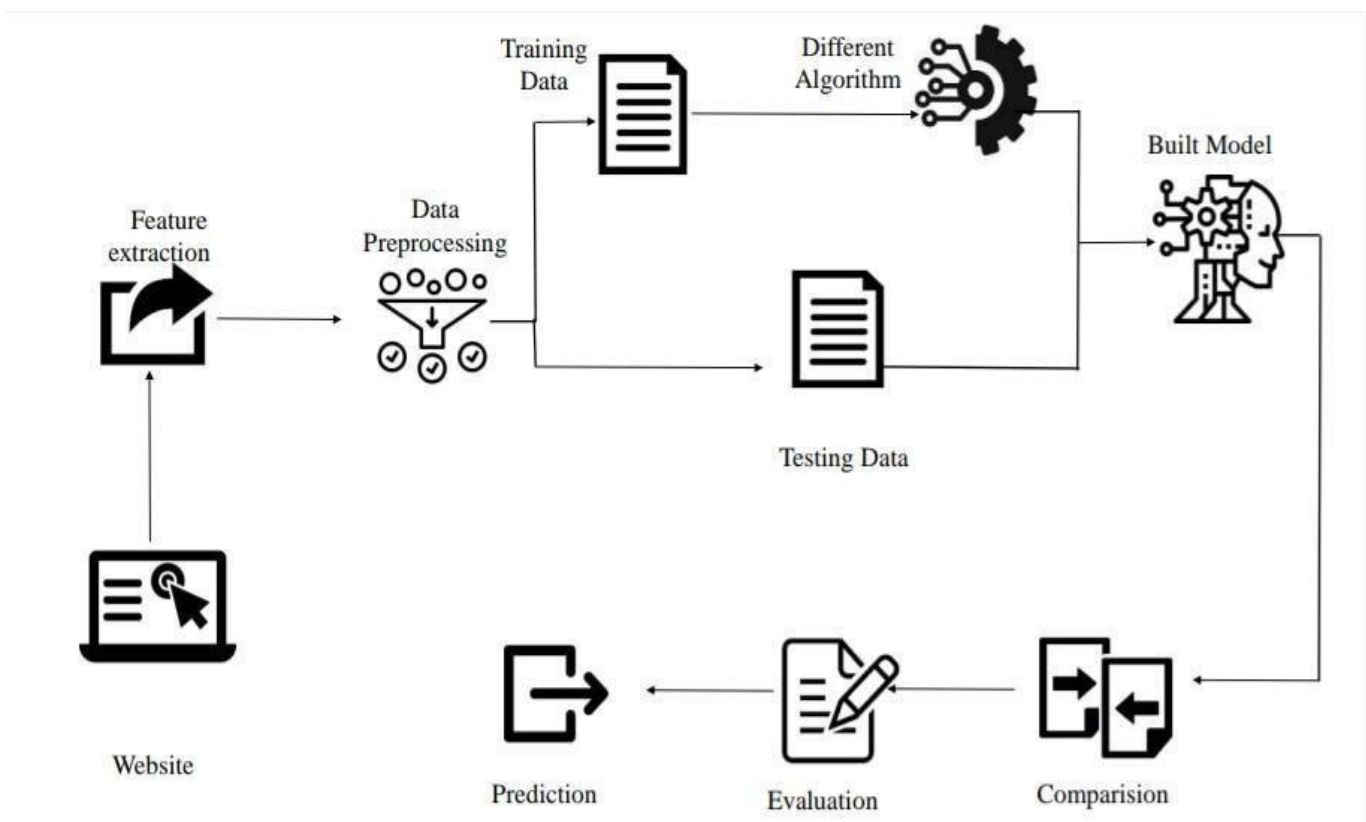
User Type	Functional Requirement (Epic)	User Story Number	User Story / Task	Acceptance criteria	Priority	Release
General Users	Home Page	USN-1	The user finds the home page easy to navigate and feels comfortable with the user interface	The every user can access and visit the Homepage	High	Sprint1
Client user(Customer)	Sign up	USN-2	The user will get authentication over their account by security measures	The client customer would register for the application using their email and setting up their own passwords for accessing the application.	High	Sprint 2
		USN-3	The user will be able to authorise their account only if they remember their authentication key (biometrics if suitable hardware available, password)	The client user can access the application if their registration is confirmed by verifying with their email.	Medium	Sprint 2
	Login	USN-4	The user can register using either their google account or their mobile number	The client user can login after getting registered	Medium	Sprint 2
	Dashboard	USN-5	The user can go through the facilities provided by the product	The user can view their profile and status in their dashboard	Low	Sprint 3
Administrator	Prediction	USN-6	User would be prompted with a pop up indicating the trustfulness of the website	The user can accurately forecast about the algorithms used.	High	Sprint 3
	Results page	USN-7	The user would be able to analyse website whether it's genuine or not	The results of the website on whether the web pages is genuine or not.	Medium	Sprint 4
Customer Care Executive	Reporting	USN-8	The user can report for any bugs or ask any queries on the product	Bugs or queries can be enquired by the user.	High	Sprint 4

SOLUTION AND TECHNICAL ARCHITECTURE: SOLUTION ARCHITECTURE:

SOLUTION ARCHITECTURE:

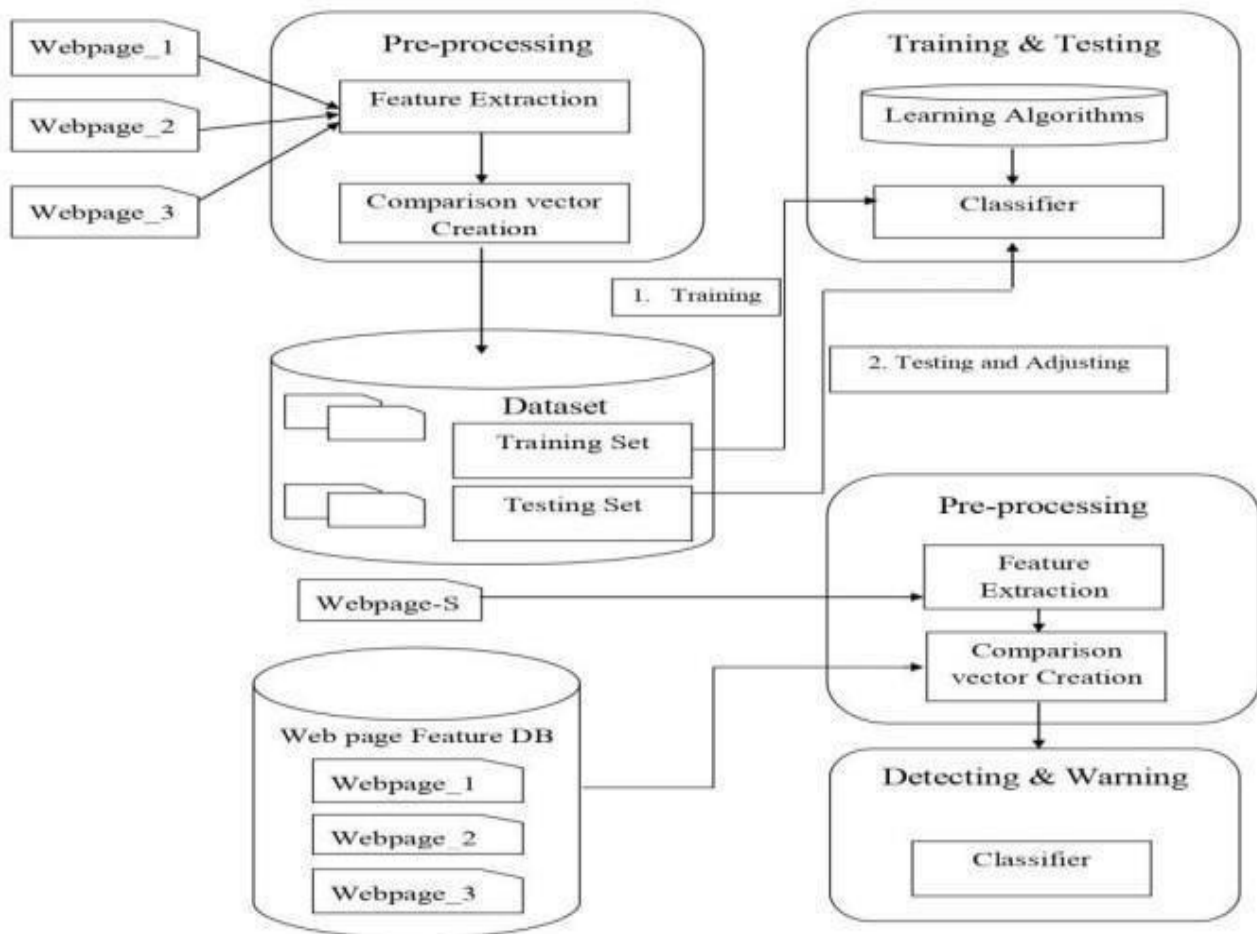
Solution architecture is a complex process – with many sub-processes – that bridges the gap between business problems and technology solutions. Its goals are to:

- Find the best tech solution to solve existing business problems.
- Describe the structure, characteristics, behavior, and other aspects of the software to project stakeholders.
- Define features, development phases, and solution requirements.
- Provide specifications according to which the solution is defined, managed, and delivered.



TECHNICAL ARCHITECTURE:

Technology architecture deals with the deployment of application components on technology components. A standard set of predefined technology components is provided in order to represent servers, network, workstations, and so on.



USERSTORIES:

Table-1:Components&Technologies:

S.No	Component	Description	Technology
1.	User Interface	How user interacts with application e.g. Web UI, Mobile App, Chatbot etc.	HTML, CSS, JavaScript
2.	Application Logic	Logic for a process in the application	Flask (Python)
3.	Database	Data Type, Configurations etc.	MySQL
4.	Cloud Database	Database Service on Cloud	IBM Watson.
5.	File Storage	File storage requirements	IBM Block Storage ,MongoDB
6.	Machine Learning Model	Purpose of Machine Learning Model	Decision tree algorithm
7.	Infrastructure (Server / Cloud)	Application Deployment on Local System / Cloud Local Server Configuration: Cloud Server Configuration :	Local, IBM Cloud

Table-2:ApplicationCharacteristics:

S.No	Characteristics	Description	Technology
1.	Open-Source Frameworks	The package Sckit Learn in Python is used to handle Machine Learning Algorithms	Machine Learning
2.	Security Implementations	Typosquatting, Cybersquatting	Cyber security
3.	Scalable Architecture	The system will be able to detect maximum of the recently updated phishing websites and is highly scalable to use.	Technology used
4.	Availability	The system is always available whenever it is required to be executed by balancing the load traffic among the servers.	IBM Cloud Load Balancers
5.	Performance	The system would have efficiency and good accuracy rate in detecting the phishing websites.	Machine Learning algorithm(Decision tree algorithm)

CHAPTER 6

PROJECT PLANNING AND SCHEDULING

6.1 SPRINT PLANNING AND ESTIMATION:

Sprint	Functional Requirement (Epic)	User Story Number	User Story / Task	Story Points	Priority	Team Members
Sprint-1	Home Page	USN-1	The user finds the home page easy to navigate and feels comfortable with the user interface	10	High	G. Srivaths Karthic
Sprint-2	Sign up	USN-2	The user will get authentication over their account by security measures	10	High	G. Srivaths Karthic S. Naveenkumar
Sprint-2		USN-3	The user will be able to authorise their account only if they remember their authentication key (biometrics if suitable hardware available, password)		Medium	S. Naveenkumar M. Ajaykumar J. Guna Seakar
Sprint-2	Login	USN-4	The user can register using either their google account or their mobile number	10	Medium	G. Srivaths Karthic
Sprint-3	Dashboard	USN-5	The user can go through the facilities provided by the product	5	Low	S. Naveenkumar
Sprint-3	Prediction	USN-6	User would be prompted with a pop up indicating the trustfulness of the website	15	High	M. Ajaykumar J. Guna Seakar
Sprint-4	Results page	USN-7	The user would be able to analyse website whether it's genuine or not	5	Medium	G. Srivaths Karthic S. Naveenkumar
Sprint-4	Reporting	USN-8	The user can report for any bugs or ask any queries on the product	15	High	S. Naveenkumar

Project Tracker, Velocity & Burndown Chart:

S.No	Characteristics	Description	Technology
1.	Open-Source Frameworks	The package Scikit Learn in Python is used to handle Machine Learning Algorithms	Machine Learning
2.	Security Implementations	Typosquatting, Cybersquatting	Cyber security
3.	Scalable Architecture	The system will be able to detect maximum of the recently updated phishing websites and is highly scalable to use.	Technology used
4.	Availability	The system is always available whenever it is required to be executed by balancing the load traffic among the servers.	IBM Cloud Load Balancers
5.	Performance	The system would have efficiency and good accuracy rate in detecting the phishing websites.	Machine Learning algorithm(Decision tree algorithm)

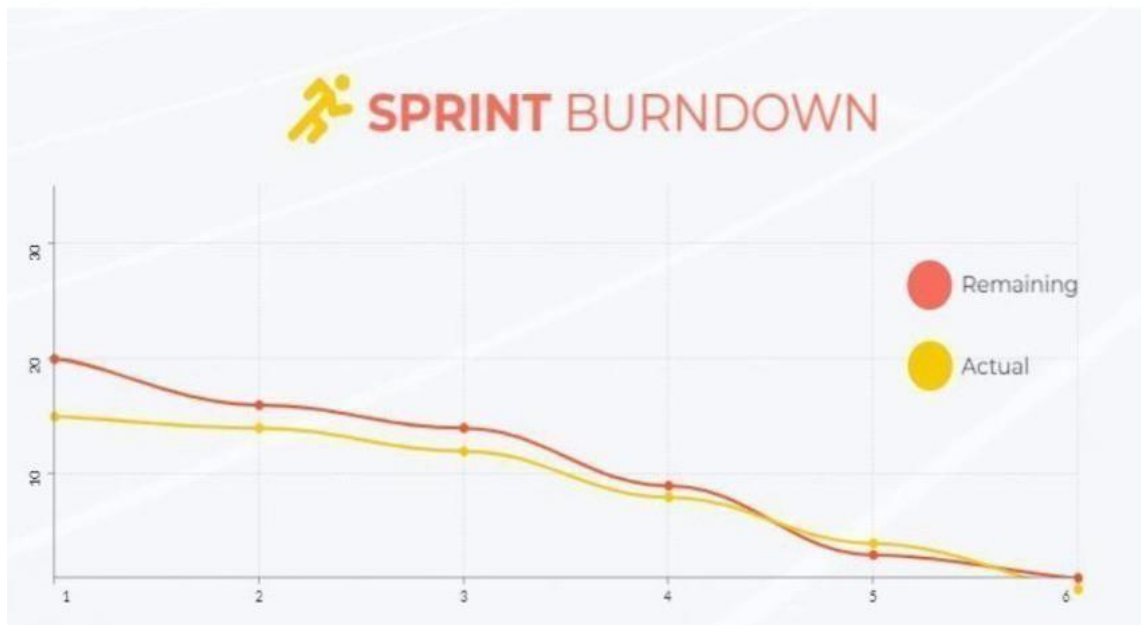
Velocity:

Imagine we have a 10-day sprint duration, and the velocity of the team is 20 (points per sprint). Let's calculate the team's average velocity (AV) per iteration unit (story points per day).

$$AV = \frac{\text{sprint duration}}{\text{velocity}} = \frac{20}{10} = 2$$

BurndownChart:

A burn down chart is a graphical representation of work left to do versus time. It is often used in agile software development methodologies such as Scrum. However, burn down charts can be applied to any project containing measurable progress over time.



CHAPTER 7

CONCLUSION

Internet is something that could not be eliminated from the daily routine. So having so much importance in it, we must also ensure the safety of ourselves while using it. So the model developed by us would certainly be a great change in this society. We also make use of large data set to train this model. So by doing this model is robust enough to detect the phishing website.

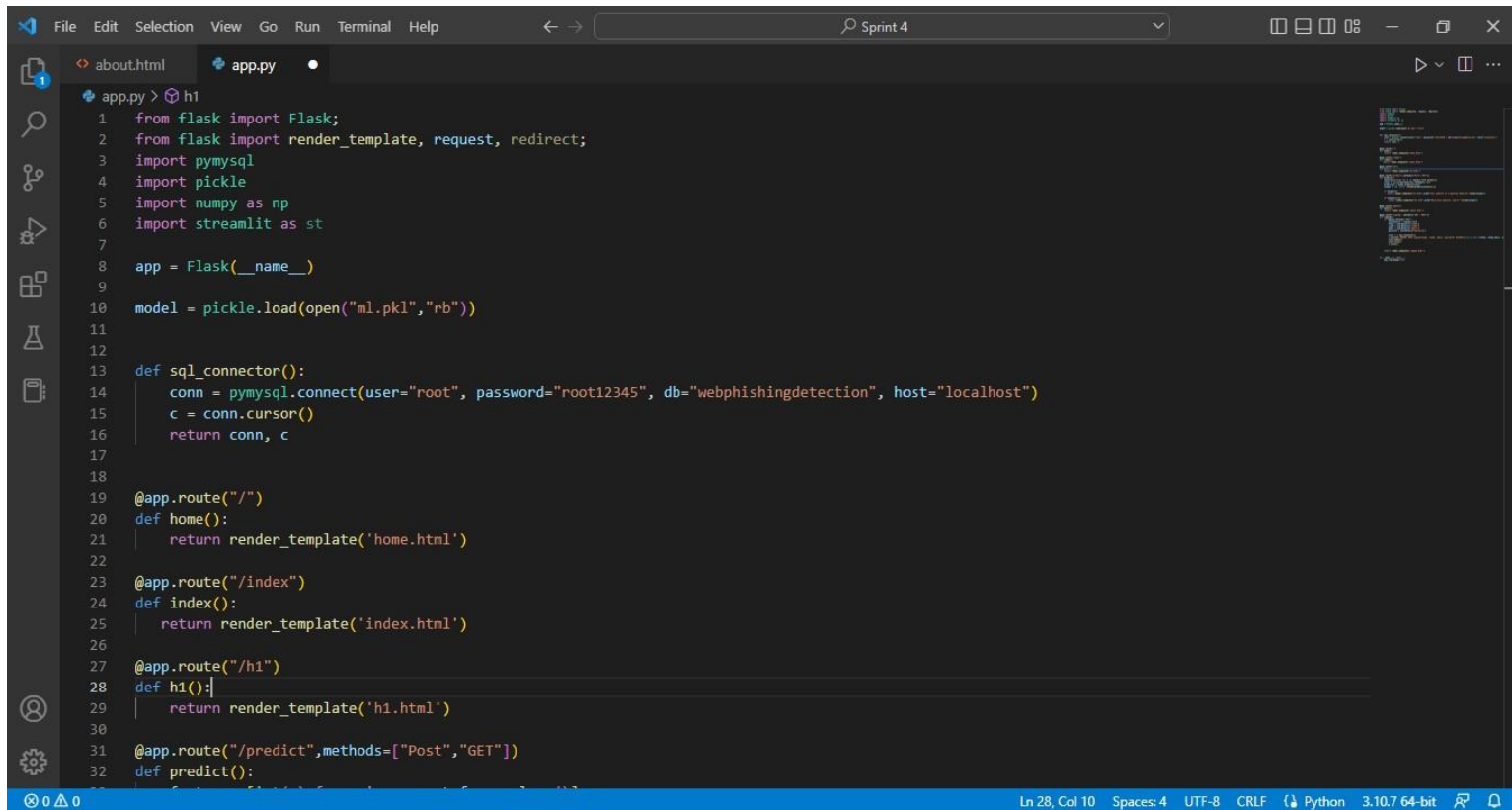
So along with this reporting facility is provided to help users to address their issues while browsing. So this tool is highly purposeful to be used by all internet users. This enhances the data integrity and browsing experience.

CHAPTER 8

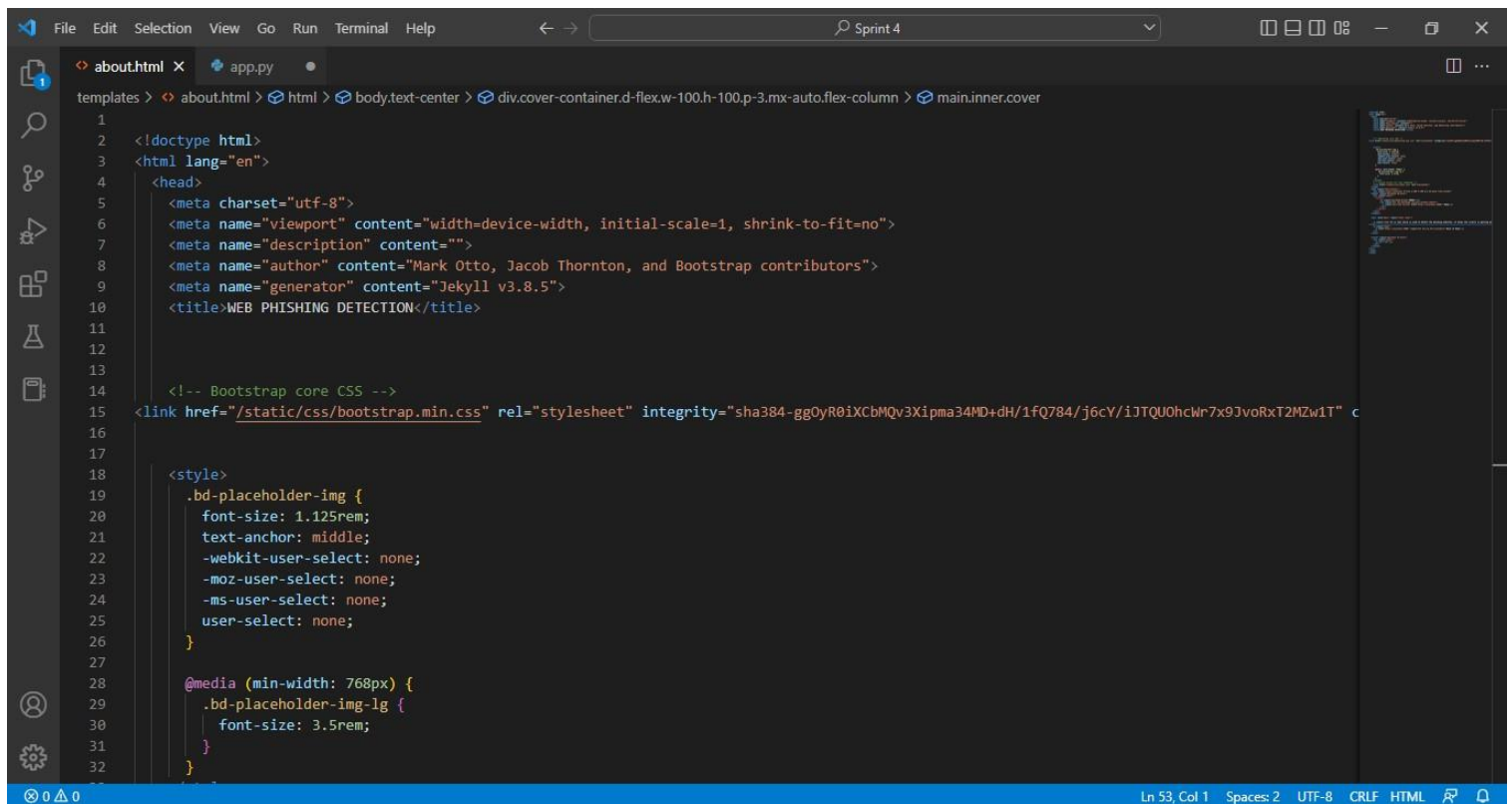
FUTURE SCOPE

In near future we have also planned to update the tool with phishing techniques so that it facilitates the detection process. It also improves the accuracy of detection.

SOURCECODE



```
1 from flask import Flask;
2 from flask import render_template, request, redirect;
3 import pymysql
4 import pickle
5 import numpy as np
6 import streamlit as st
7
8 app = Flask(__name__)
9
10 model = pickle.load(open("ml.pkl", "rb"))
11
12
13 def sql_connector():
14     conn = pymysql.connect(user="root", password="root12345", db="webphishingdetection", host="localhost")
15     c = conn.cursor()
16     return conn, c
17
18
19 @app.route("/")
20 def home():
21     return render_template('home.html')
22
23 @app.route("/index")
24 def index():
25     return render_template('index.html')
26
27 @app.route("/h1")
28 def h1():
29     return render_template('h1.html')
30
31 @app.route("/predict", methods=["Post", "GET"])
32 def predict():
```

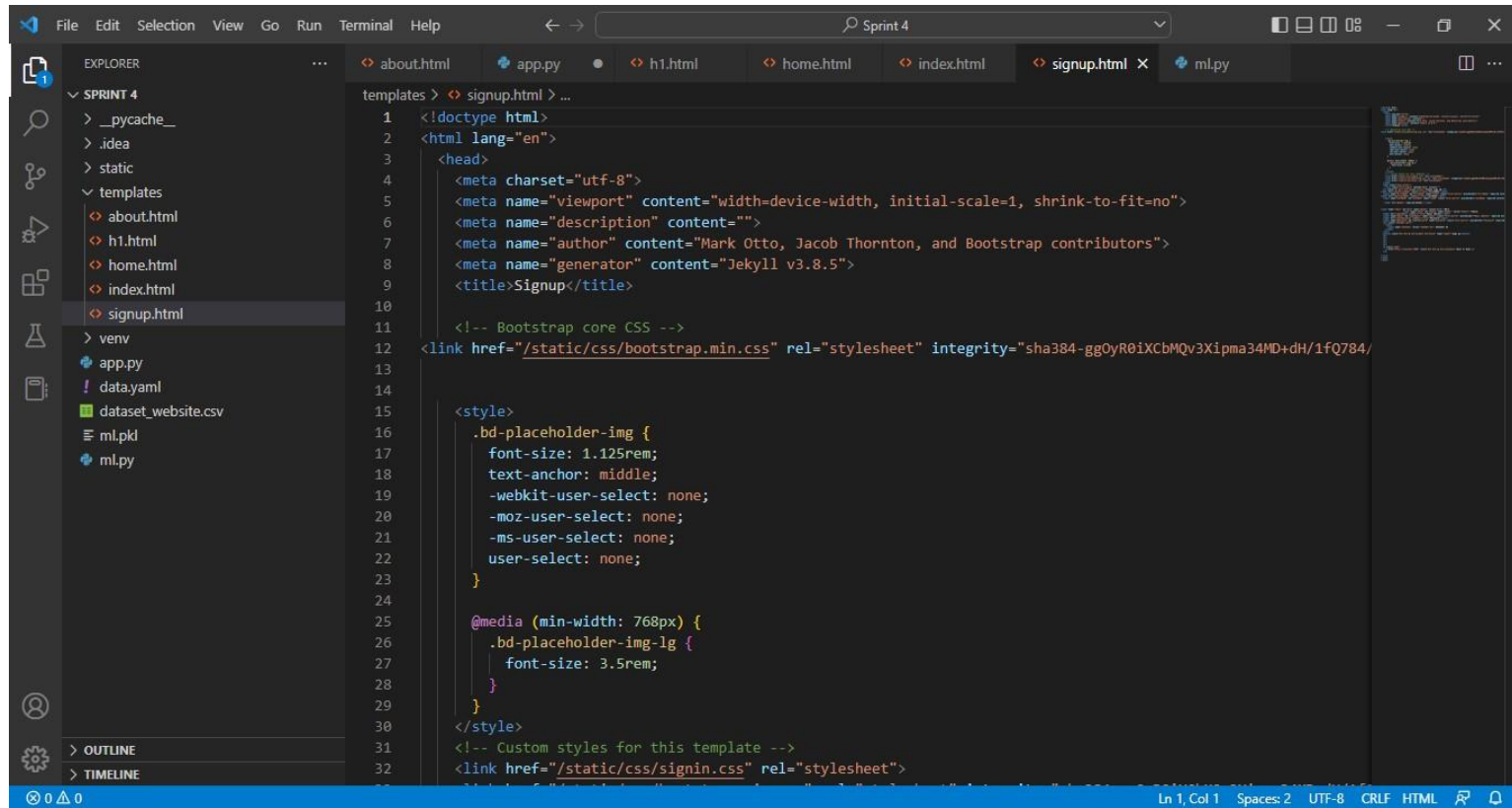


```
1 <!doctype html>
2 <html lang="en">
3 <head>
4     <meta charset="utf-8">
5     <meta name="viewport" content="width=device-width, initial-scale=1, shrink-to-fit=no">
6     <meta name="description" content="">
7     <meta name="author" content="Mark Otto, Jacob Thornton, and Bootstrap contributors">
8     <meta name="generator" content="Jekyll v3.8.5">
9     <title>WEB PHISHING DETECTION</title>
10
11 <!-- Bootstrap core CSS -->
12 <link href="/static/css/bootstrap.min.css" rel="stylesheet" integrity="sha384-ggOyR0iXCbMQv3Xipma34ND+dH/1fQ784/j6cY/iJTQUOhcWr7x9JvoRxT2MZw1T" c
13
14
15 <style>
16     .bd-placeholder-img {
17         font-size: 1.125rem;
18         text-align: middle;
19         -webkit-user-select: none;
20         -moz-user-select: none;
21         -ms-user-select: none;
22         user-select: none;
23     }
24
25     @media (min-width: 768px) {
26         .bd-placeholder-img-lg {
27             font-size: 3.5rem;
28         }
29     }
30 </style>
```

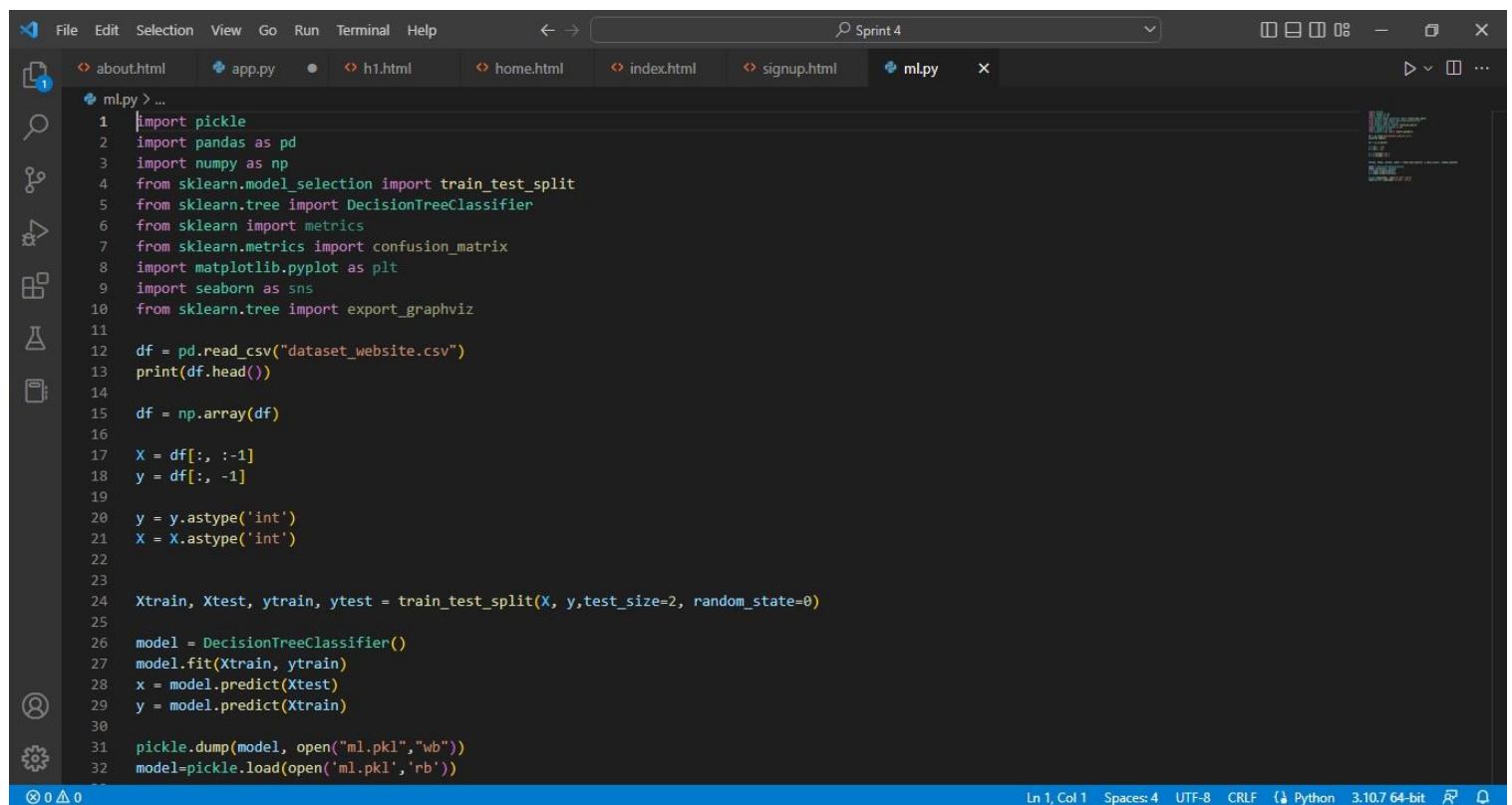


```
1 <!doctype html>
2 <html lang="en">
3 <head>
4   <meta charset="utf-8">
5   <meta name="viewport" content="width=device-width, initial-scale=1, shrink-to-fit=no">
6   <meta name="description" content="">
7   <meta name="author" content="Mark Otto, Jacob Thornton, and Bootstrap contributors">
8   <meta name="generator" content="Jekyll v3.8.5">
9   <title>WEB PHISHING DETECTION</title>
10
11 <!-- Bootstrap core CSS -->
12 <link href="/static/css/bootstrap.min.css" rel="stylesheet" integrity="sha384-ggOyR0iXCbMQv3Xipma34MD+dH/1fQ784/
13
14 <style>
15   .bd-placeholder-img {
16     font-size: 1.125rem;
17     text-align: middle;
18     -webkit-user-select: none;
19     -moz-user-select: none;
20     -ms-user-select: none;
21     user-select: none;
22   }
23
24   @media (min-width: 768px) {
25     .bd-placeholder-img-lg {
26       font-size: 3.5rem;
27     }
28   }
29 </style>
```

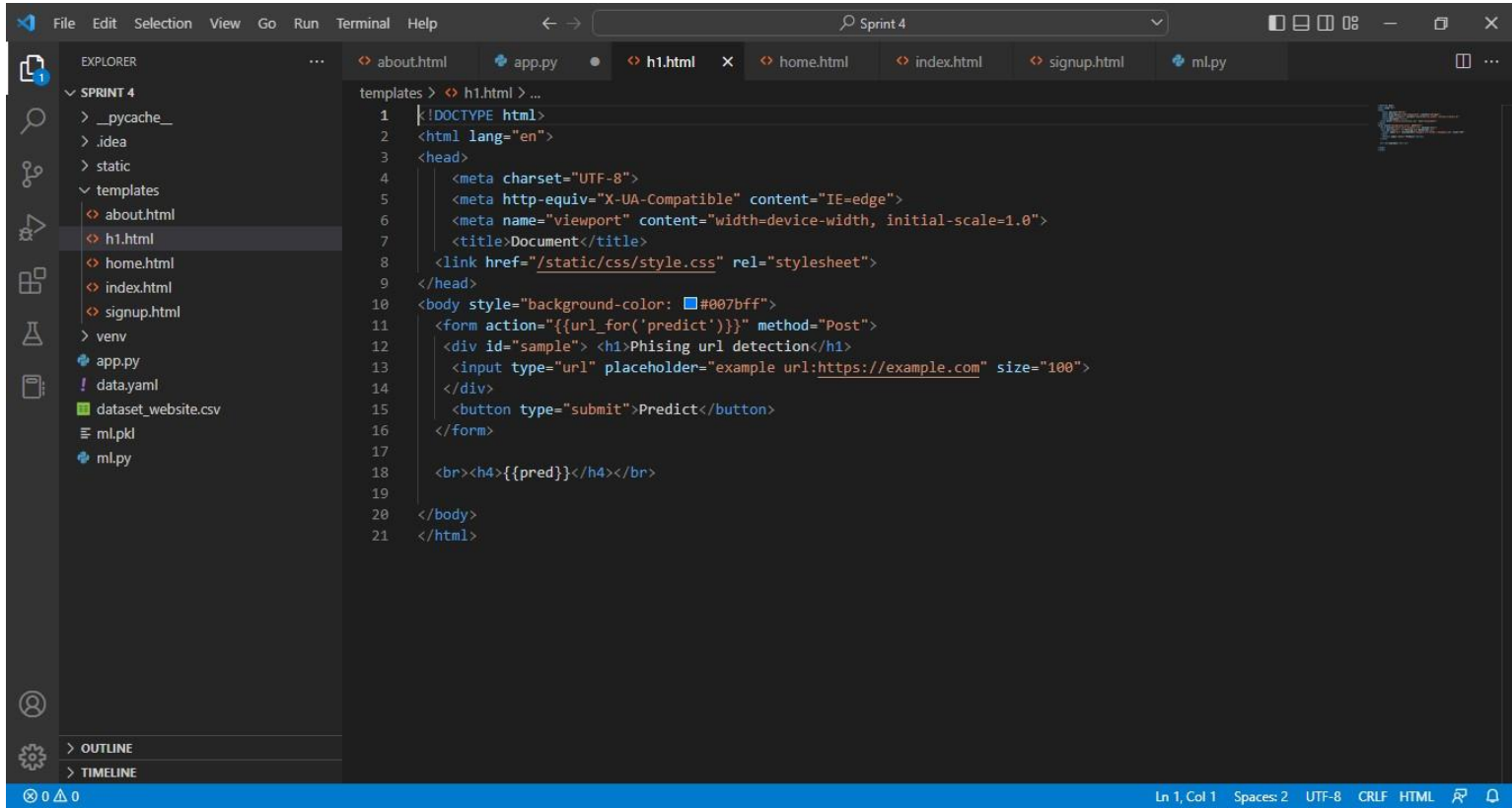
```
1 <!doctype html>
2 <html lang="en">
3 <head>
4   <meta charset="utf-8">
5   <meta name="viewport" content="width=device-width, initial-scale=1, shrink-to-fit=no">
6   <meta name="description" content="">
7   <meta name="author" content="Mark Otto, Jacob Thornton, and Bootstrap contributors">
8   <meta name="generator" content="Jekyll v3.8.5">
9   <title>Signin</title>
10
11 <!-- Bootstrap core CSS -->
12 <link href="/static/css/bootstrap.min.css" rel="stylesheet" integrity="sha384-ggOyR0iXCbMQv3Xipma34MD+dH/1fQ784/
13
14 <style>
15   .bd-placeholder-img {
16     font-size: 1.125rem;
17     text-align: middle;
18     -webkit-user-select: none;
19     -moz-user-select: none;
20     -ms-user-select: none;
21     user-select: none;
22   }
23
24   @media (min-width: 768px) {
25     .bd-placeholder-img-lg {
26       font-size: 3.5rem;
27     }
28   }
29 </style>
30 <!-- Custom styles for this template -->
31 <link href="/static/css/signin.css" rel="stylesheet">
```



```
1 <!doctype html>
2 <html lang="en">
3   <head>
4     <meta charset="utf-8">
5     <meta name="viewport" content="width=device-width, initial-scale=1, shrink-to-fit=no">
6     <meta name="description" content="">
7     <meta name="author" content="Mark Otto, Jacob Thornton, and Bootstrap contributors">
8     <meta name="generator" content="Jekyll v3.8.5">
9     <title>Signup</title>
10
11   <!-- Bootstrap core CSS -->
12   <link href="/static/css/bootstrap.min.css" rel="stylesheet" integrity="sha384-ggOyR0iXCbMQv3Xipma34MD+dH/1fQ784/
13
14
15   <style>
16     .bd-placeholder-img {
17       font-size: 1.125rem;
18       text-align: middle;
19       -webkit-user-select: none;
20       -moz-user-select: none;
21       -ms-user-select: none;
22       user-select: none;
23     }
24
25     @media (min-width: 768px) {
26       .bd-placeholder-img-lg {
27         font-size: 3.5rem;
28       }
29     }
30   </style>
31
32   <!-- Custom styles for this template -->
33   <link href="/static/css/signin.css" rel="stylesheet">
```

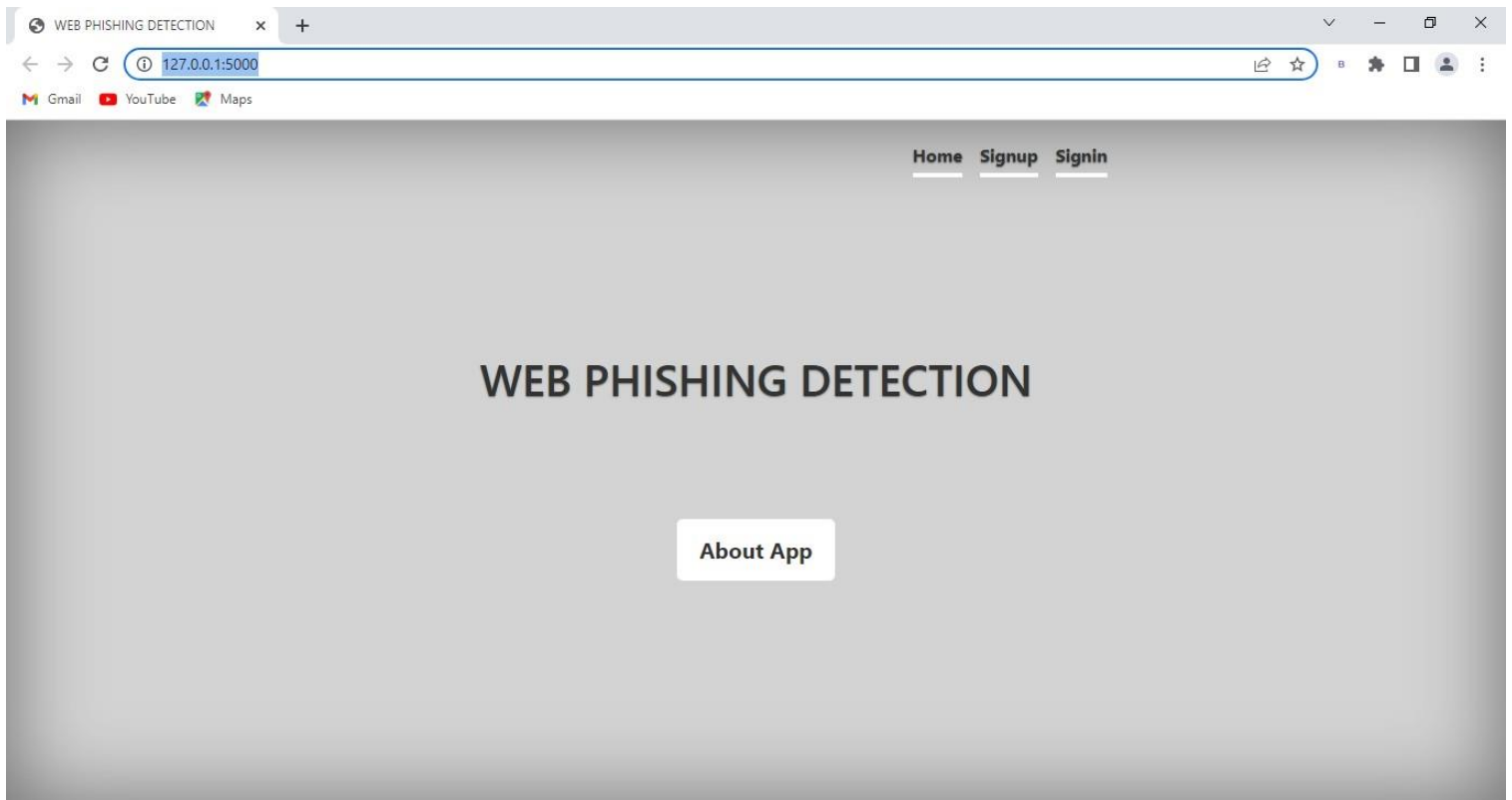


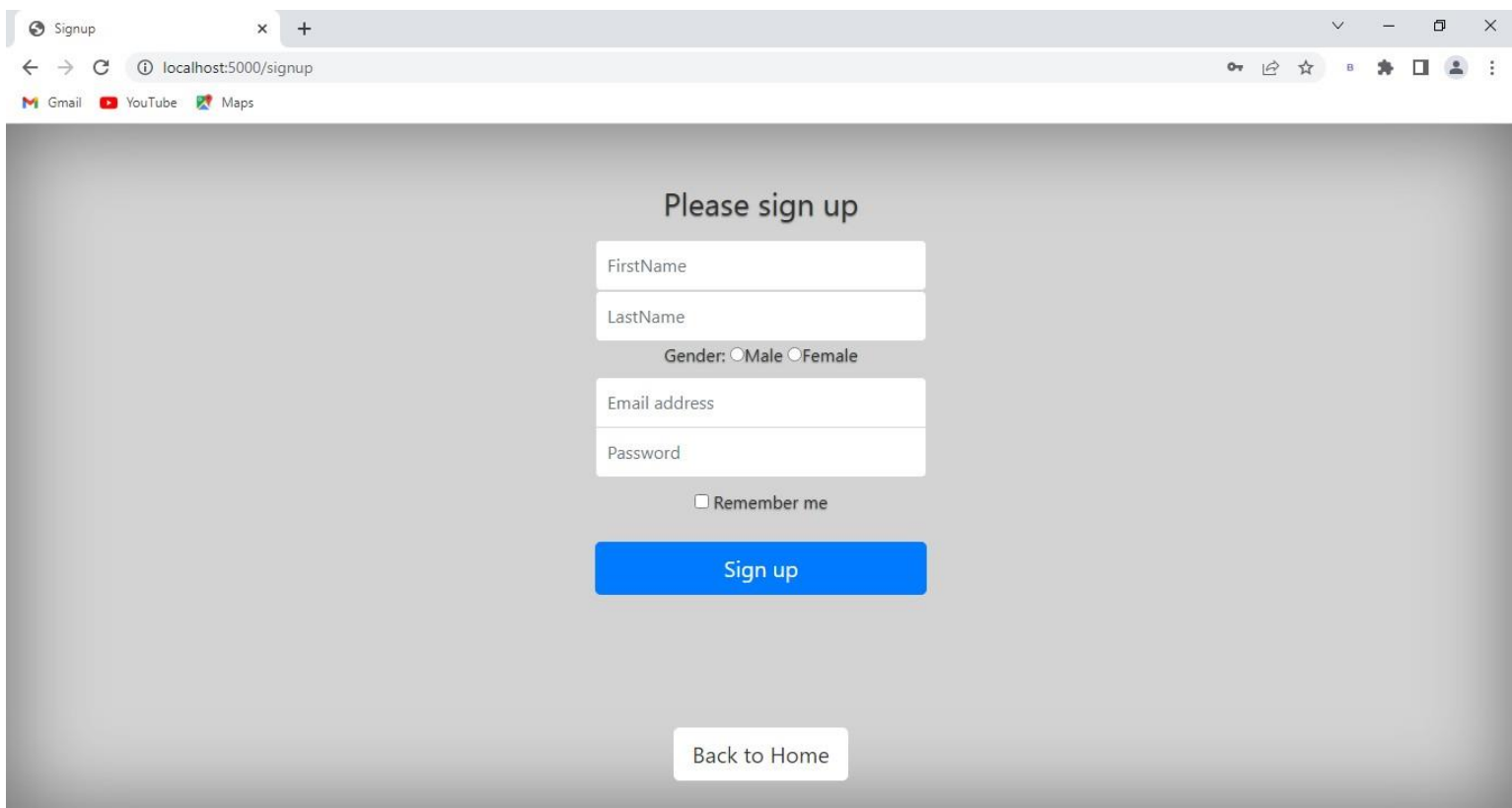
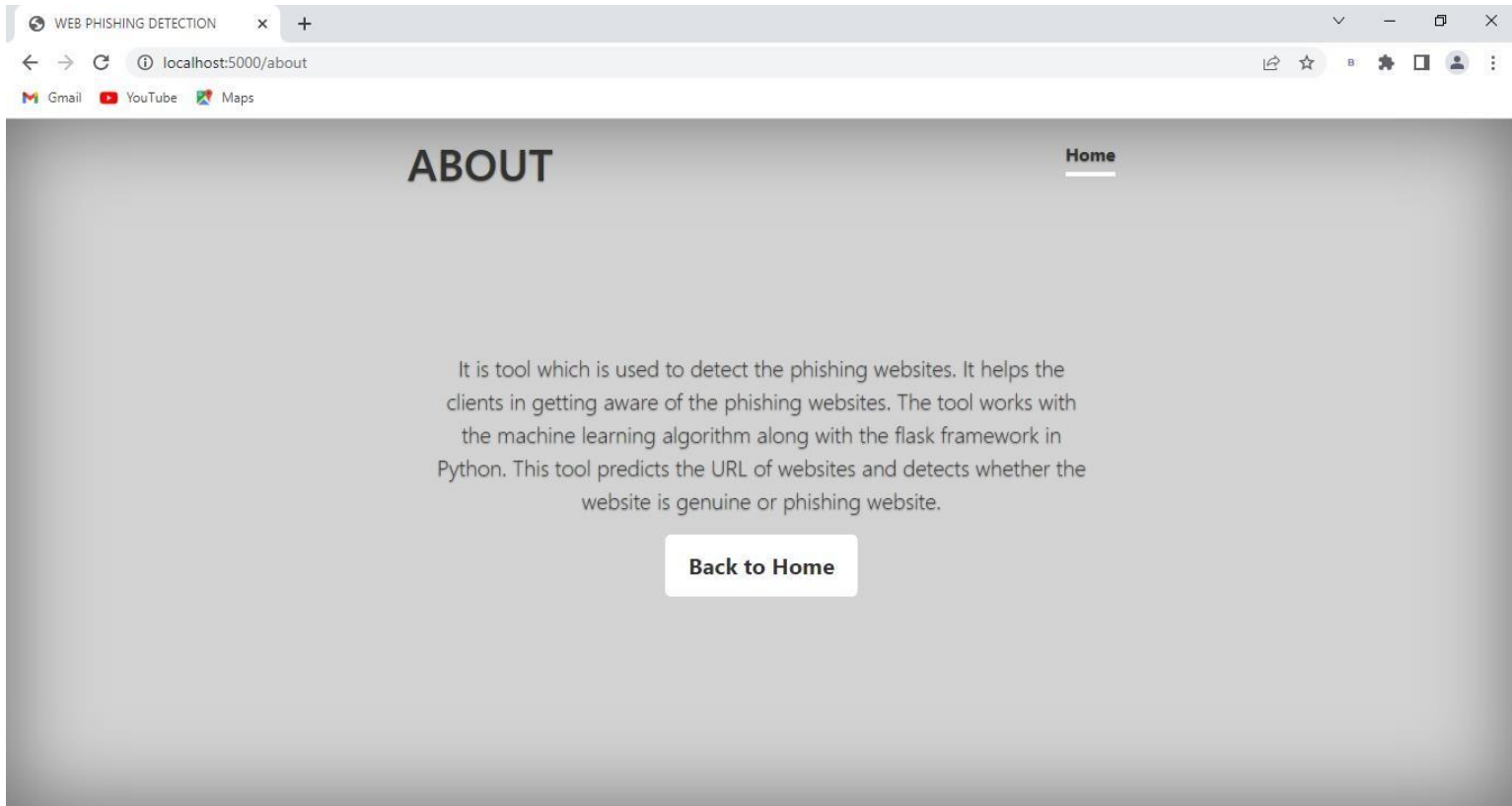
```
1 import pickle
2 import pandas as pd
3 import numpy as np
4 from sklearn.model_selection import train_test_split
5 from sklearn.tree import DecisionTreeClassifier
6 from sklearn import metrics
7 from sklearn.metrics import confusion_matrix
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10 from sklearn.tree import export_graphviz
11
12 df = pd.read_csv("dataset_website.csv")
13 print(df.head())
14
15 df = np.array(df)
16
17 X = df[:, :-1]
18 y = df[:, -1]
19
20 y = y.astype('int')
21 X = X.astype('int')
22
23
24 Xtrain, Xtest, ytrain, ytest = train_test_split(X, y, test_size=2, random_state=0)
25
26 model = DecisionTreeClassifier()
27 model.fit(Xtrain, ytrain)
28 x = model.predict(Xtest)
29 y = model.predict(Xtrain)
30
31 pickle.dump(model, open("ml.pkl", "wb"))
32 model=pickle.load(open('ml.pkl', 'rb'))
```

```
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   <title>Document</title>
8   <link href="/static/css/style.css" rel="stylesheet">
9 </head>
10 <body style="background-color: #007bff">
11   <form action="{{url_for('predict')}}" method="Post">
12     <div id="sample"> <h1>Phising url detection</h1>
13     <input type="url" placeholder="example url:https://example.com" size="100">
14   </div>
15   <button type="submit">Predict</button>
16 </form>
17
18   <br><h4>{{pred}}</h4></br>
19
20 </body>
21 </html>
```

OUTPUT





GITHUB LINK –<https://github.com/IBM-EPBL/IBM-Project-42668-1660703736DEMOLINK> – https://drive.google.com/file/d/1EynOE7Tfe1rGl8-QPzNKIUYN_GTRqm0p/view?usp=share_link