

WEB PHISHING DETECTION

IBM-Project-44647-1660725800

**NALAIYA THIRAN PROJECT BASED ON LEARNING
PROFESSIONAL READLINESS FOR INNOVATION,
EMPLOYMENT AND ENTERPRENEURSHIP**

Project Report

TEAM ID: PNT2022TMID24691

TEAM MEMBERS:

1. SAADHANA. G- 210419104137
2. SREELAKSHMI. G- 210419105060
3. ANJU PRIYA. V- 210419104014
4. VENKATACHALAM. V- 210419104182

**BACHELOUR OF ENGINEERING IN COMPUTER
SCIENCE AND ENGINEERING**

**CHENNAI INSTITUTE OF TECHNOLOGY
CHENNAI- 600 069**

INDEX

1. INTRODUCTION

1. Project Overview
2. Purpose

2. LITERATURE SURVEY

1. Existing problem
2. References
3. Problem Statement Definition

3. IDEATION & PROPOSED SOLUTION

1. Empathy Map Canvas
2. Ideation & Brainstorming
3. Proposed Solution
4. Problem Solution fit

4. REQUIREMENT ANALYSIS

1. Functional requirement
2. Non-Functional requirements

5. PROJECT DESIGN

1. Data Flow Diagrams
2. Solution & Technical Architecture
3. User Stories

6. PROJECT PLANNING & SCHEDULING

1. Sprint Planning & Estimation
2. Sprint Delivery Schedule
3. Reports from JIRA

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

1. Feature 1
2. Feature 2
3. Database Schema (if Applicable)

8. TESTING

1. Test Cases
2. User Acceptance Testing

9. RESULTS

1. Performance Metrics

10. ADVANTAGES & DISADVANTAGES

11. CONCLUSION

12. FUTURE SCOPE

13. APPENDIX

Source Code

GitHub & Project Demo Link

1. INTRODUCTION

1.1 Project Overview

This project mainly focuses on applying a machine-learning algorithm to detect phishing websites. In order to detect and predict phishing websites, we proposed an intelligent, flexible, and effective system that is based on using classification algorithms. We implemented classification algorithms and techniques to extract the phishing dataset's criteria to classify their legitimacy. The phishing website can be detected based on some important characteristics, like the URL and domain identity, and security and encryption criteria in the final phishing detection rate. Once a user enters a website, our system will use a data mining algorithm to detect whether the website is a phishing website or not.

1.2 Purpose

There are a number of users who purchase products online and make payments through e-banking. Some e-banking websites ask users to provide sensitive data such as username, password, and credit card details, etc., often for malicious reasons. This type of e-banking website is known as a phishing website. Web services are one of the key communications software services for the Internet. Web phishing is one of many security threats to web services on the Internet. There are millions of incidents happening around the world in an hour. People suffer immeasurable losses due to these attacks. Therefore, protecting users from such attacks is the sole purpose of our project.

The simplest method of obtaining sensitive information from unwitting users is through phishing attacks. The goal of phishers is to obtain vital data, such as username, password, and bank account information. People working in cyber security are currently searching for reliable and consistent methods of detecting phishing websites. In this research, many properties of legal and phishing URLs are extracted and analyzed in order to detect phishing URLs. The algorithms used to identify phishing websites include decision trees, random forests, and support vector machines. By evaluating each algorithm's accuracy rate, false positive rate, and false negative rate, the study aims to identify phishing URLs as well as identify the best machine learning method.

2. LITERATURE SURVEY

2.1 Existing problem

Due to how simple it is to create a fake website that closely resembles a legitimate website, phishing has recently become a top concern for security researchers. Experts can spot fake websites, but not all users can, and those users end up falling for phishing scams. The attacker's primary goal is to steal bank account credentials. Businesses in the US lose \$2 billion annually as a result of their customers falling for phishing scams. The annual global impact of phishing was estimated to be as high as \$5 billion in the third Microsoft Computing Safer Index Report, which was published in February 2014. Because users are unaware of phishing attacks, they are becoming more successful.

Since phishing attacks take advantage of user vulnerabilities, it is highly challenging to counteract them, but it is crucial to improve phishing detection methods. The common technique, commonly referred to as the "blacklist" method, for detecting phishing websites involves adding Internet Protocol (IP) blacklisted URLs to the antivirus database. Attackers utilize clever methods to deceive people by changing the URL to seem authentic through obfuscation and many other straightforward tactics, such as fast-flux, in which proxies are automatically constructed to host the website, algorithmic production of new URLs, etc. This method's primary flaw is that it cannot identify phishing attacks that occur at zero hour.

Zero-hour phishing attacks can be detected using heuristic-based detection, which includes characteristics that have been observed to exist in phishing attacks in reality. However, the presence of these characteristics is not always guaranteed in such attacks, and the false positive rate for detection is very high.

2.2 References

S.NO	PAPER TITLE	PAPER CONCEPT	ADVANTAGE	DISADVANTAGE
1	LongfeiWu etal..., "Effective Defense Schemes for Phishing Attacks on Mobile Computing Platforms," IEEE 2016, pp.6678-6691.	In this paper, author did a comprehensive study on the Security vulnerabilities caused by mobile phishing attacks, including the web page phishing attacks.	Author propose MobiFish, a novel automated lightweight anti- phishing scheme for mobile platforms. MobiFish verifies the validity of web pages, applications, and persistent accounts by comparing the actual Identity to the claimed identity	Existing schemes designed for web phishing attacks on PCs cannot effectively address the various phishing attacks on mobile devices.
2	Surbhi Gupta etal., "A Literature Survey on Social Engineering Attacks: Phishing Attacks," in International Conference on Computing, Communication and Automation(ICCC A2016),2016, pp. 537-540.	To fool an online user into elicit personal Information. The prime objective of this review is to do literature survey on social engineering attack: Phishing attacks and techniques to detect attack.	The paper discusses various types of Phishing attacks such as Tab-napping, spoofing emails, Trojan horse, hacking and how to prevent them.	Every organization has security issues that have been of great concern to u sets, sited developers, and specialists, in order to defend the confidential data from this type of social engineering attack.

3	Guardian Analytics, " A Practical Guide to Anomaly Detection Implications of meeting new FFIEC minimum expectations for layered security ". [Accessed : 08 Jan 2015]	Commercial and retail account holders at financial institutions of all sizes are under attacks by sophisticated, Organized, Well-funded cyber criminals.	Anomaly detection solutions are readily available, are deployed quickly and immediately and automatically protect all account holders against all types of fraud attack with minimal Disruption to legitimate online banking activity.	Implementing anomaly detection will not only meet FFIEC expectations, it will decrease the total cost of fraud, and will increase customer loyalty and trust.
4	SANS Institute, " Phishing : Analysis of a Growing Problem ",2007. 1417[Accessed : 23 May 2017]	This paper gives an in depth analysis of phishing: what it is, the technologies and security. Weaknesses it takes advantage of the dangers it poses to end users.	In this analysis author explain the concepts and technology behind phishing, show how the threat is much more than just a nuisance or passing trend, and discuss how gangs of criminals are using these scams to make a great deal of money.	Unfortunately, a growing number of cyber-thieves are using these same systems to manipulate us and steal our private information.

5	<p>J. Phys.: Conf. Ser.</p> <p>"A literature survey on Retraction: Phishing website detection using machine learning and deep learning techniques" 1916 (2021) 012407.</p>	<p>Nowadays, website phishing is more damaging. It is becoming a big threat to people's daily life and networking environment. In these attacks, the intruder puts on an act as if it is a trusted organization with an intention to purloin liable and essential information.</p> <p>The methodology we discovered is a powerful technique to detect the phished websites and can provide more effective defenses for phishing attacks of the future.</p>	<p>The association between independent variables as well as dependent variables can be formed without any presumptions about the statistical depiction of the aspect. It contributes positive gains on regression algorithm which includes its competence to act with noisy data.</p>	<p>The ANN's are not suitable for infrequent or utmost events where data is inadequate in order to train it.</p> <p>ANNs do not permit the embodiment of human mastery to be substitutive for perceptible proof.</p>
6	<p>"Phishing Website Detection Based on Deep Convolutional Neural Network and Random</p>	<p>This paper proposes an integrated phishing website detection method based on convolutional neural networks</p>	<p>A 99.35% correct classification rate of phishing websites was obtained on the dataset. Experiments were conducted on the test set and training set, and the</p>	<p>It takes longer to train. However, the trained model is better than the others in terms of accuracy of phishing website detection. Another disadvantage is</p>

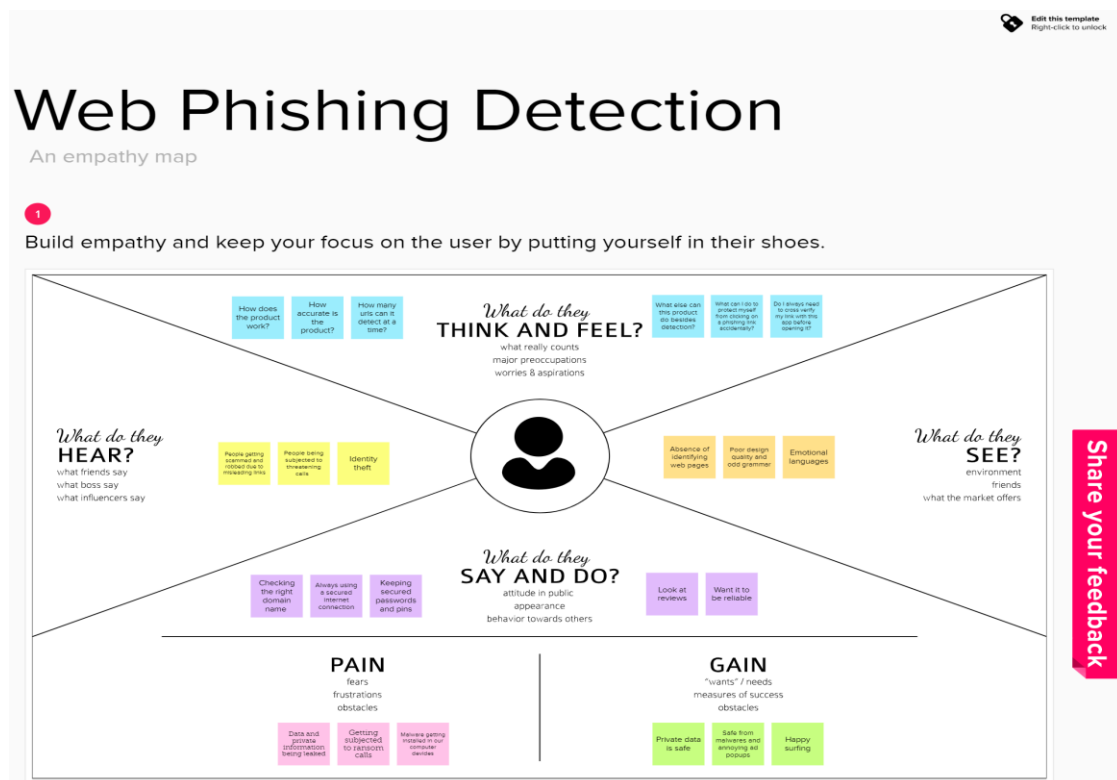
	<p>Forest Ensemble Learning"</p> <p>,This research was funded by the National Key R & D Program of China Grant Numbers 2017YFB0802800 and Beijing Natural Science Foundation (4202002)</p>	<p>(CNN) and random forest (RF).</p> <p>The method can predict the legitimacy of URLs without accessing the web content or using third-party services.</p> <p>The proposed technique uses character embedding techniques to convert URLs into fixed-size matrices, extract features at different levels using CNN models, classify multi-level features using multiple RF classifiers, and, finally, output prediction results using a winner-take-all approach.</p>	<p>experimental results proved that the proposed method has good generalization ability and is useful in practical applications.</p>	<p>that the model cannot determine whether the URL is active or not, so it is necessary to test whether the URL is active or not before detection to ensure the effectiveness of detection. In addition, some attackers use URLs that are not imitations of other websites, and such URLs will not be detected.</p>
--	---	--	--	---

2.3 Problem Statement Definition

Human users' inability to recognize phishing sites allows phishing attacks to succeed. Past work in anti-phishing can be broadly divided into four categories: studies to understand why people fall for phishing attacks, strategies for teaching people not to fall for phishing attacks, user interfaces for assisting people in making better decisions about trusting email and websites, and automated tools to detect phishing. Our research outlines a method for automatically identifying phishing. Most end users typically base their decisions only on how they feel and how they look. When a user accesses the internet, all they see is a browser's screen. After that, he or she works on a web page's command. Most phishing efforts take use of this sort of unintended chance provided by the user and trick them since the user is unconcerned with the back end procedure.

3. IDEATION & PROPOSED SOLUTION

3.1 Empathy Map Canvas



3.2 Ideation & Brainstorming

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

- Name gathering**
Define who should participate in the session and send an invite. Share relevant information or pre-work ahead.
- Set the goal**
Think about the problem you'll be focusing on solving in the brainstorming session.
- Learn how to use the facilitation tools**
Join the facilitator. Supplement to it as a helpful and productive session.

Open article

Define your problem statement

What problem are you trying to solve? Frame your problem as a how might the statement. This will be the focus of your brainstorm.

5 minutes

Problem

How might we [your problem statement]?

As online transactions grew in popularity, cybercrimes also grew quickly. Because of the anonymity offered by the internet, hackers try to trick users by using techniques like phishing, SQL injection, malware, man-in-the-middle attacks, domain name system tunnelling, ransomware, web trojans, and other forms of attack. Phishing is said to be the most misleading of all of these tactics.

Brainstorm

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

10 minutes

Saadhana, G(Team Leader)

- Difficult to block "redirecting pages", as these URLs can be disguised
- To keep track of all potential phishing techniques, use an intelligent online security gateway.
- Analytical comparison of trustworthy and fraudulent websites

Sreelakshmi, G(Team member 1)

- Web phishing detection with advanced deployment
- Maintaining web phishing records with a spam repository
- Observation of resource loading times to identify fraudulent websites' tendencies to display just costs rather than their functionality

Anju Priya, V(Team member 2)

- Link pathways are being traced after detection of tampering
- By contrasting it with real online domain names, homograph spoofing can be identified.
- To check for extended links and prevent link shortening because most fraudulent websites use the "bit.ly" extension

Venkatachalam V(Team member 3)

- Incorrectly worded keywords and odd website changes
- Requesting confirmation of unrelated qualifications
- Verify secure web protocols like "https" visually

3.3 Proposed Solution

S.No.	Parameter	Description
1.	Problem Statement (Problem to be solved)	As opposed to software vulnerabilities, "phishing sites" are a particular kind of internet security problems that primarily target human vulnerabilities. Phishing sites are harmful websites that pretend to be trustworthy websites or web pages in order to steal users' personal information, including their user name, password, and credit card number. Since phishing is mostly a semantics-based attack that focuses on human vulnerabilities, identifying these phishing websites can be difficult. The main goal of this project is to classify phishing websites using a

		variety of machine learning approaches in order to produce a model with the highest level of accuracy and simplicity.
2.	Idea / Solution description	<ul style="list-style-type: none">• The method includes the extraction of lexical features from collected webpages as well as host- and page-based feature extraction. The first stage is gathering phishing and legitimate websites. In the host-based technique, attribute extractions based on admiration and lexical bases are carried out to create a database of attribute value. This database contains knowledge that has been extracted using various machine learning methods. A selective classifier is chosen after comparing the methods, and it is put into practice in Python.
		<ul style="list-style-type: none">• The suggested approach gathered URLs of safe websites from sites like www.alex.com, www.dmoz.org, and browsing history. We gathered the phishing URLs from www.phishtak.com. 20000 benign URLs and 17000 phishing URLs make up the data collection.

3.	Novelty / Uniqueness	<p>The dataset provided by UCI Machine Learning repository⁴ and compiled by Mohammad et al³ was used by the suggested system. The dataset contains 6157 legal URLs and 4898 phishing URLs across 11055 data points.</p> <p>Each data point had 30 features that were sorted into the three categories below:</p> <ul style="list-style-type: none">• Features extracted from the URL• Features based on the page's source code, such as URLs that are incorporated into the webpage and HTML and JavaScript-based features.• Features based on domains.
4.	Social Impact / Customer Satisfaction	<p>The majority of the public (users) were assisted by the project in determining if a website was a phishing website or not. It assisted them in classifying the hazardous locations. Machine learning methods were employed in this research. The URL is entered, and it will recognize it and provide users with precise results.</p>

5.	Business Model (Revenue Model)	<p>In the literature, a number of methods for phishing attack detection and filtering have been suggested. Researchers are still looking for a solution that can protect consumers from phishing attacks and produce better outcomes. It might be easier to spot phishing websites if we can recognize the specific traits and patterns they exhibit.</p> <p>The classification problem of identifying such traits can be resolved using machine learning approaches.</p>
6.	Scalability of the Solution	<p>This project offers an effective method for phishing detection that pulls features from the URL and HTML source code of websites. In particular, we suggested a hybrid feature set that included features for the HTML source code's plaintext and noisy HTML data, different hyperlink information, and URL character sequence characteristics without the knowledge of experts. The suggested anti-phishing technique has demonstrated competitive performance on actual datasets in terms of several assessment statistics, according to extensive trials.</p> <p>The following criteria have been established for our anti-phishing strategy.</p> <ul style="list-style-type: none">• Target independent• Real-time detection• High detection efficiency• Third-party independent

3.4 Problem Solution fit

Define CS fit into CC	1. CUSTOMER SEGMENT(S) <p>Internet users who frequent millions of websites</p> <p>especially those who utilise websites for e-banking and e-commerce.</p>	6. CUSTOMER CONSTRAINTS <ul style="list-style-type: none"> Phishing attempts frequently result in the loss of a customer's credentials and valuable personal information. 	5. AVAILABLE SOLUTIONS <ul style="list-style-type: none"> Manual self-analysis using address features as a basis for confirmation. Double checking the link with a phishing database. 	Explore AS differentiate

Focus on I&P fit into RC understand RC	2. JOBS-TO-BE-DONE / PROBLEMS <p>Obtaining the URLs of websites from customers,</p>	9. PROBLEM ROOT CAUSE <p>Developments in technology that encourage hacking and phishing.</p> <p>Low effectiveness of algorithms.</p> <p>Credential access that</p>	7. BEHAVIOUR <ul style="list-style-type: none"> Making use of a unique extension that examines the current link The user can access the extension that offers results. 	Focus on I&P fit into RC understand RC

Identify the	3. TRIGGERS <ul style="list-style-type: none"> As alerted with the urge or temptation to commit to a task. 	10. YOUR SOLUTION <p>Making a website in Python where a user may enter a URL and the system classifies it as a phishing website or not using machine learning algorithms and then provides the user with feedback</p>	8. CHANNELS of BEHAVIOUR <p>ONLINE</p> <p>Using the website link to examine the phishing website's behaviour and receiving feedback from the build</p>	Identify the
	4. EMOTIONS: BEFORE / AFTER <ul style="list-style-type: none"> Before: Fear of Uncertainty, Vulnerability. After: Relief of maintaining privacy and confidence in website access. 			

4. REQUIREMENT ANALYSIS

4.1 Functional requirement

Following are the functional requirements of the proposed solution.

FR NO	Functional Requirements	Classification
FR-1	Fetch Electronic Mail Messages	Core
FR-2	Extract URLS	Core
FR-3	Extract Header Information	Core
FR-4	Classify Email	Core
FR-5	Static or Dynamic (Inbox)	Core
FR-6	Provide User Feedback	Core

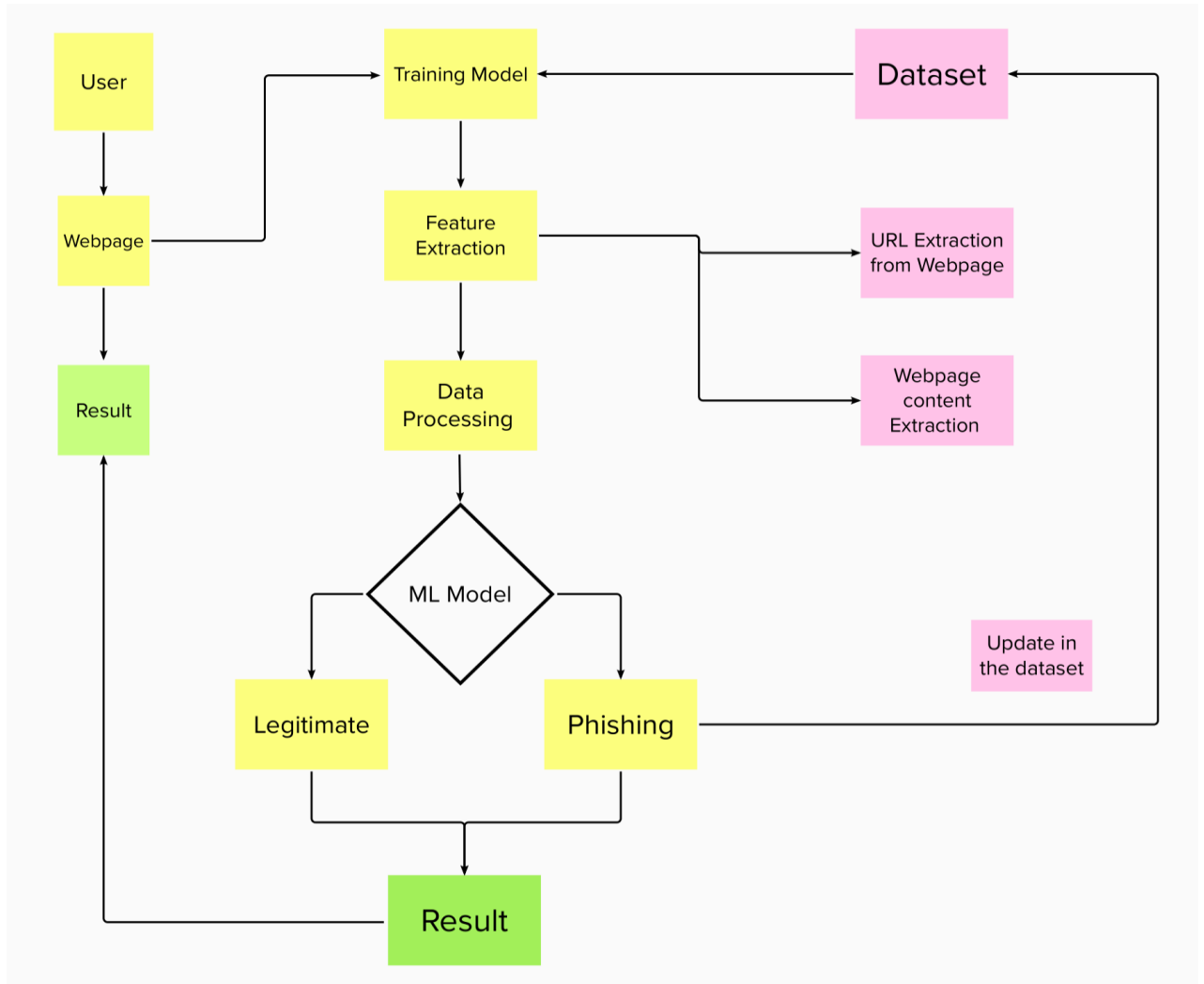
4.2 Non-Functional requirements

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

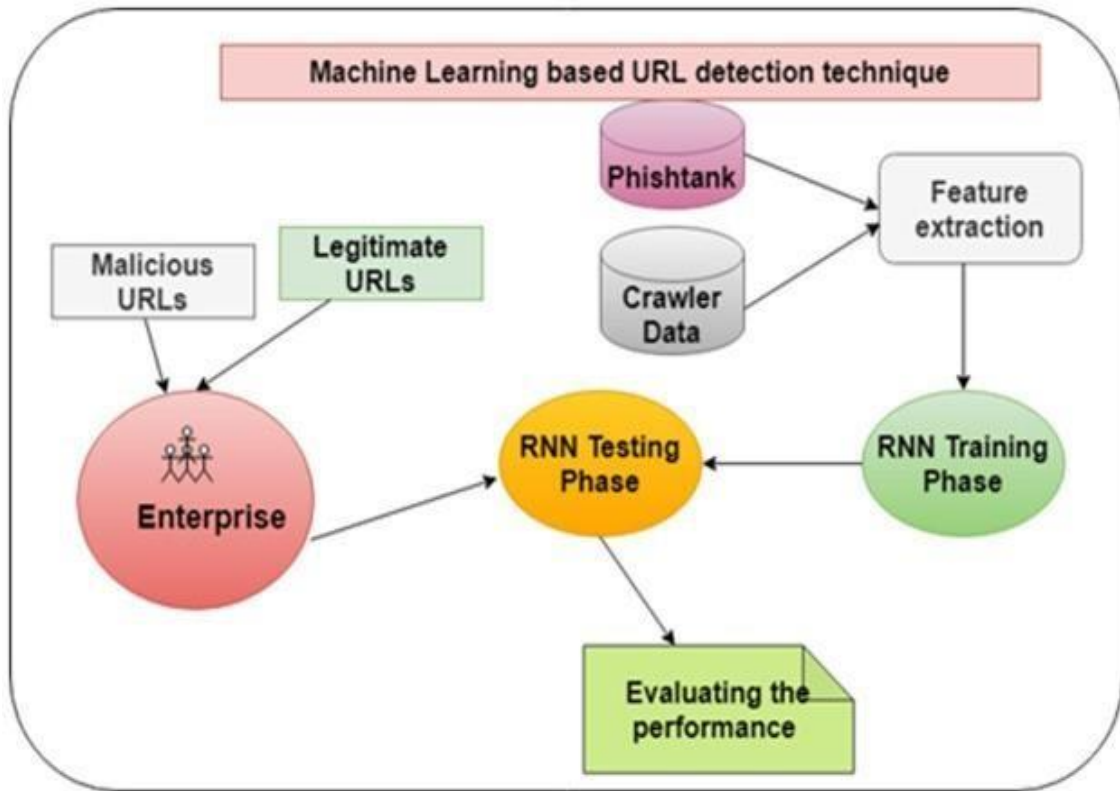
FR NO	Non-Functional Requirements	Description
NFR-1	Usability	System is easy to configure and is efficient in carrying out user tasks.
NFR-2	Availability	System is available to work asrequired when it is required.
NFR-3	Reliability	System will perform the tasks it was designed to do.
NFR-4	Performance	System will perform tasks in a fashion that complies with predetermined criteria.
NFR-5	Security	System will protect all data manipulated internally from unauthorized access and threats.
NFR-6	Scalability	System will appropriately handle increasing and decreasing workloads.

5. PROJECT DESIGN

5.1 Data Flow Diagrams:



5.2 Solution & Technical Architecture



5.3 User Stories:

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

TEAM ID: PNT2022TMID24691

		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
Customer (Web user)	User Input	USN-1	As a user, I can enter the required URL in the box while awaiting validation.	I can access the website without any problem	High	Sprint-1
Customer Care Executive	Feature Extraction	USN-1	In the event that nothing is discovered during comparison, we can extract features using a heuristic and a visual similarity technique.	As a user I can have comparison between websites for security	High	Sprint-1
Administrator	Prediction	USN-1	The model will use machine learning algorithms like a logistics regression and KNN to forecast the URLs of the websites.	I can accurately forecast the specific algorithms in this way.	High	Sprint-1
	Classifier	USN-2	To create the final product, I will now feed all of the model output to classifier.	I'll use this to identify the appropriate classifier for generating the outcome.	Medium	Sprint-2

6. PROJECT PLANNING & SCHEDULING

6.1 Sprint Planning & Estimation

Sprints	User Type	Functional Requirement (Epic)	User Story No	User Story / Task	Story points	Team members	Priority
Sprint-1	Dataset collection and preprocessing	Fetch electronic mail messages	USN-1	As a new user, I will register first.	35	Saadhana G Sree Lakshmi G Anju Priya V Venkatachalam M	High
Sprint-2	Model and application building	Extract URLs	USN-2	As a user, I will provide specific URL for checking	15	Saadhana G Sree Lakshmi G Anju Priya V Venkatachalam M	High
Sprint-3	Feature addition for prediction page	Extract Header Information	USN-3	As a user, I wait for the application to classify it based on certain criteria.	25	Saadhana G Sree Lakshmi G Anju Priya V Venkatachalam M	High
Sprint-4	User acceptance testing, performance testing, migration from mongo DB to DB2	Classify the website	USN-4	As a user, I will be informed whether the link is suspicious or safe to use	25	Saadhana G Sree Lakshmi G Anju Priya V Venkatachalam M	High

6.2 Sprint Delivery Schedule

Sprint	Total Story Points	Duration	Sprint StartDate	Sprint End Date (Planned)	Story Points Completed (as on Planned End Date)	Sprint Release Date (Actual)
Sprint-1	35	7Days	29-10-2022	5-11-2022	35	4-11-2022
Sprint-2	15	8 Days	7-11-2022	14-11-2022	15	13-11-2022
Sprint-3	25	8 Days	16-11-2022	23-11-2022	25	23-11-2022
Sprint-4	25	8 Days	23-11-2022	30-11-2022	25	30-11-2022

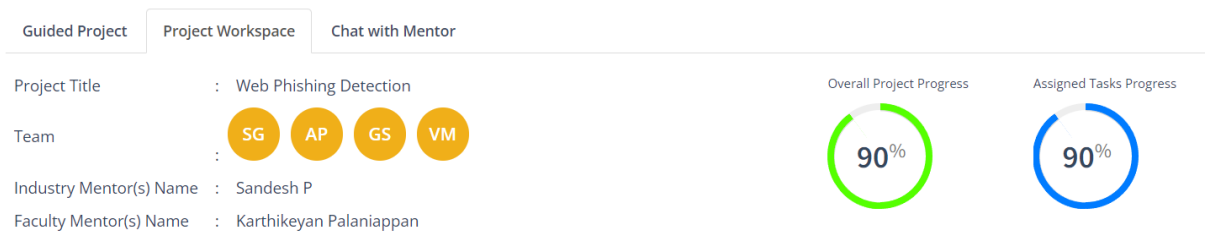
Velocity:

$$AV = \text{Velocity} / \text{Duration} = 35 / 7 = 5$$

$$AV = \text{Velocity} / \text{Duration} = 15 / 8 = 1.875A$$

$$V = \text{Velocity} / \text{Duration} = 25 / 8 = 3.125$$

6.3 Reports from JIRA



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

7.1 Feature 1

LOGIN

```
@app.route('/login/', methods=['POST'])
def login():
    if request.method=="POST":
        email=request.form.get("email")
        password=request.form.get("password")
        if(account.find_one({"email":email})):
            user=account.find_one({"email":email})
            if(user and pbkdf2_sha256.verify(password,user['password'])):
                return start_session(user)
            else:
                flash("Password is incorrect","loginError")
                return redirect(url_for('index',loginError=True))
        flash("Sorry, user with this email id does not exist","loginError")
        return redirect(url_for('index',loginError=True))
```

SIGNUP

```
@app.route('/signup/',methods=['POST'])
def signup():
    if request.method=="POST":
        userInfo={
            "fullName":request.form.get('fullName'),
            "email":request.form.get('email'),
            "phoneNumber":request.form.get('phoneNumber'),
            "password":request.form.get('password'),
        }
        userInfo['password']=pbkdf2_sha256.encrypt(userInfo['password'])
        if(account.find_one({"email":userInfo['email']})):
            flash("Sorry,user with this email already exist","signupError")
            return redirect(url_for('index',signupError=True))
        if(account.insert_one(userInfo)):
            return start_session(userInfo)
        flash("Signup failed","signupError")
        return redirect(url_for('index',signupError=True))
```

ABOUT US

```
@app.route('/about/')
def about():
    if(session and session['logged_in']):
        if(session['logged_in']==True):
            return render_template('./templates/about.html',userInfo=session['user'],aboutContents=aboutData['aboutContents'])
        else:
            return render_template('./templates/about.html',aboutContents=aboutData['aboutContents'])
    else:
        return render_template('./templates/about.html',aboutContents=aboutData['aboutContents'])
```

7.2 Feature 2

HISTORY PAGE

```
@app.route('/detection-history/')
@login_required
def detectionHistory():
    if(session and session['logged_in']):
        if(session['logged_in']==True):
            get_detection_history_stmt = "SELECT title,url,status FROM detectionHistory where email=?"
            get_detection_history = ibm_db.prepare(conn, get_detection_history_stmt)
            ibm_db.bind_param(get_detection_history,1,session['user']['email'])
            ibm_db.execute(get_detection_history)
            fetch_detection_history = ibm_db.fetch_assoc(get_detection_history)
            detection_history = []
            ind = 0
            while fetch_detection_history != False:
                detection_history.append(fetch_detection_history)
                ind += 1
                fetch_detection_history = ibm_db.fetch_assoc(get_detection_history)
            detection_history= detection_history[::-1]
            return render_template('./templates/detection-history.html',userInfo=session['user'],detectionHistory=detection_history)
```

CONTACT US PAGE

```
@app.route('/contact/')
def contact():
    if(session and session['logged_in']):
        if(session['logged_in']==True):
            return render_template('./templates/contact.html',userInfo=session['user'])
        else:
            return render_template('./templates/contact.html')
    else:
        return render_template('./templates/contact.html')
```

FAQ

```
<h1 class="faq-title">FAQs about phishing URL</h1>
<div>
    <ul class="faq-list">
        <li>
            <h4 class="faq-heading">How can I identify a Phishing scam? </h4>
            <p class="read faq-text">
                The first rule to remember is to never give out any personal information in an email. No institution, bank or other
            </p>
        </li>
        <li>
            <h4 class="faq-heading">Do I only need to worry about Phishing attacks via email? </h4>
            <p class="read faq-text">
                No. Phishing attacks can also occur through phone calls, texts, instant messaging, or malware on your computer wh
            </p>
        </li>
        <li>
            <h4 class="faq-heading">What kind of information should I protect? </h4>
            <p class="read faq-text">
                You should protect all sensitive and confidential data. For information on what is considered sensitive and confide
            </p>
        </li>
        <li>
            <h4 class="faq-heading">
                Why Is Phishing Dangerous?
            </h4>
            <p class="read faq-text">
                Phishing is dangerous for anyone who is even remotely touched by technology because it puts them under the risk of
            </p>
        </li>
        <li>
            <h4 class="faq-heading">
                What Do You Do If You Suspect Phishing?
            </h4>
            <p class="read faq-text">
                Cybersecurity experts recommend users to treat every email they receive as a phishing email so that they are extra
            </p>
        </li>
    </ul>
</div>
```


7.3 Database Schema (if Applicable)




Tables			New table +				
<input type="checkbox"/> Name ▼	Schema	Properties					
<input type="checkbox"/> ACCOUNT	YSX70667	...					
<input type="checkbox"/> DETECTIONHISTORY	YSX70667	...					

Table definition

: ✕

ACCOUNT

No statistics available.





Name	Data type	Nullable	Length	Scale	
FULLNAME	VARCHAR	N	100	0	
EMAIL	VARCHAR	N	100	0	
PHONENUMBER	LONG VARCHAR	N	32700	0	
PASSWORD	VARCHAR	N	100	0	

Table definition

: ✕

DETECTIONHISTORY

No statistics available.

Name	Data type	Nullable	Length	Scale	
EMAIL	VARCHAR	N	100	0	
TITLE	VARCHAR	N	100	0	
URL	VARCHAR	N	100	0	
STATUS	VARCHAR	N	100	0	

8. TESTING

8.1 Test Cases

Test case ID	Feature Type	Component	Test Scenario	Pre-Requisite	Steps To Execute	Test Data	Expected Result	Actual Result	Status	Comments
HomePage_TC_OO1	Functional	Home Page	Verify user is able to enter the URL in the form	Run the flask app in local host	1.Open our phishing website 2.Login to use the phishing services 3.Enter the link to be detected and click on predict button	https://google.com/	Result of classification will be displayed	Working as expected	Pass	Since www.google.com is a safe link, the output would display and say it is a safe link
ResultPage_TC_OO1	UI	Contact us page	Verify the UI elements in the form	Run the flask app in local host	1.Enter name, email and message 2.Press submit	-	An email received stating that the message has been forwarded to the team	Working as expected	Pass	Email JS is used to send automatic email
ResultPage_TC_OO2	Functional	Prediction result page	Verify user is able to see an alert when	Run the flask app in local host	1.Enter URL and click go		Alert of incomplete input	Working as expected	Pass	

TEAM ID: PNT2022TMID24691

			nothing is entered in the textbox		2.Enter nothing and click submit 3.An alert is displayed to provide proper input					
Prediction Page_TC_OO1	Functional	Prediction form page	Verify user is able to see the result when URL is entered in the textbox	Run the flask app in local host	1.Enter URL and click go 2. Enter any URL and click submit 3. The result of the classification is displayed in a new page.	https://google.com/	Result of classification will be displayed with a corresponding emoticon	Working as expected	Pass	

8.2 User Acceptance Testing

1. Defect Analysis:

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

Resolution	Severity 1	Severity 2	Severity 3	Severity 4	Subtotal
By Design	10	4	2	3	20
Duplicate	1	0	3	0	4
External	2	3	0	1	6
Fixed	11	2	4	20	37
Not Reproduced	0	0	1	0	1
Skipped	0	0	1	1	2
Won't Fix	0	5	2	1	8
Totals	24	14	13	26	77


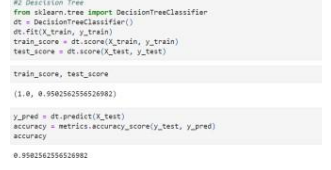
2. Test Case Analysis:

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

Section	Total Cases	Not Tested	Fail	Pass
Print Engine	5	0	0	5-
Client Application	51	0	0	51
Security	2	0	0	2
Outsource Shipping	3	0	0	3
Exception Reporting	9	0	0	9
Final Report Output	4	0	0	4
Version Control	2	0	0	2

9. RESULTS

9.1 Performance Metrics

S.No.	Parameter	Values	Screenshot
1.	Model Summary	Decision Tree ModelAccuracy – 97%	 <pre> #0 Decision Tree from sklearn.tree import DecisionTreeClassifier dt = DecisionTreeClassifier() dt.fit(X_train, y_train) train_score = dt.score(X_train, y_train) test_score = dt.score(X_test, y_test) train_score, test_score (1.0, 0.9582562556526982) y_pred = dt.predict(X_test) accuracy = metrics.accuracy_score(y_test, y_pred) accuracy 0.9582562556526982 </pre>
2.	Accuracy	Training Accuracy - Test	 <pre> #0 Decision Tree from sklearn.tree import DecisionTreeClassifier dt = DecisionTreeClassifier() dt.fit(X_train, y_train) train_score = dt.score(X_train, y_train) test_score = dt.score(X_test, y_test) train_score, test_score (1.0, 0.9582562556526982) y_pred = dt.predict(X_test) accuracy = metrics.accuracy_score(y_test, y_pred) accuracy 0.9582562556526982 </pre>

Model Performance Comparison:

https://dataplatform.cloud.ibm.com/analytics/notebooks/v2/8188710d-09bc-4dd6-824d-ec5519a26fea/view?access_token=c66af8d9cede342710725423df04821dc21cc53af0d4eae2f7496d9e1d3f5a7f

10. ADVANTAGES & DISADVANTAGES:

Phishing is the attempt to obtain a user's financial and personal information, such as credit card numbers and passwords, through electronic communication such as email and other messaging services. Attackers pose as representatives of a company and direct users to a fake website that looks like a phishing website, which is then used to gather personal data about users. A link embedded in the email can be used by attackers to trick users into downloading malware or malicious software.

To protect users from phishing attacks, numerous studies have been conducted. Firewalls, the blocking of specific domains and IP addresses, spam filtering methods, the detection of phoney websites, client-side toolbars, and user education are some of them. Both benefits and drawbacks may be seen in any of these methods now in use. The requirement to automatically identify phishing targets is a significant issue for anti-phishing initiatives. Knowing the website that is thought to be the target website allows us to identify which specific pages are phishing attempts. The owners may benefit from being able to recognize phishing attempts and take the appropriate countermeasures right away.

11. CONCLUSION

Using machine learning technologies, this initiative seeks to improve the detection process for phishing websites. Using the random forest approach, we had the lowest percentage of false positives and 97.14% detection accuracy. The outcome further demonstrates that classifiers perform better when more data is utilized as training data. Future phishing website detection will be more accurate thanks to the implementation of hybrid technology, which combines the blacklist approach with the random forest algorithm of machine learning.

12. FUTURE SCOPE:

Future study will evaluate the effectiveness of the current finding with the use of a different method, such as deep learning, for phishing web page identification. Additionally, a web browser plug-in that can identify phishing websites and shield consumers in real time will be created based on an effective algorithm.

For simple access to human life, service providers provide a variety of the quickest instruments online. Additionally, online crime such as phishing is disseminated similarly to real-world crime. However, there is no online security team protecting users from these crimes. All types of internet users can benefit greatly from an anti-phishing program. These security tools are more necessary for beginners or people with limited internet or e-commerce knowledge. Phishing's primary targets are online banking or payments. The ideal method for identifying cybercrime or e-marketing fraud is thus an automated anti-phishing technique.

13. APPENDIX

Source Code

```
1 import datetime
2 import os
3 from os.path import join, dirname
4 from dotenv import load_dotenv
5 from functools import wraps
6 from http.client import HTTPException
7 import numpy as np
8 from flask import Flask, request, render_template, session,
9     url_for, redirect, flash
10 import json
11 import pickle
12 from passlib.hash import pbkdf2_sha256
13 import json
14 import inputScript
15 import ibm_db
16 app = Flask(__name__, template_folder='../Flask')
17 model = pickle.load(open('../Flask/Phishing_Website.pkl', 'rb'
18 ))
19
20 dotenv_path = join(dirname(__file__), '.env')
21 load_dotenv(dotenv_path)
22 conn = ibm_db.connect(os.environ.get('IBMDB_URL'), '', '')
23 SECRET_KEY = os.environ.get("SECRET_KEY")
24 app.secret_key = SECRET_KEY
25 carouselDataFile = open('./static/json/carouselData.json')
26 carouselData = json.load(carouselDataFile)
27 aboutDataFile = open('./static/json/aboutData.json')
28 aboutData = json.load(aboutDataFile)
```



```

1  def login_required(f):
2      @wraps(f)
3      def wrap(*args, **kwargs):
4          if('logged_in' in session):
5              return f(*args, **kwargs)
6          else:
7              return redirect('/')
8      return wrap
9
10
11 def start_session(userInfo):
12     del userInfo['password']
13     session['logged_in']=True
14     session['user']=userInfo
15     session['predicted']=False
16     return redirect(url_for('index'))
17
18
19 @app.route('/login/',methods=['POST'])
20 def login():
21     if request.method=="POST":
22         email=request.form.get("email")
23         password=request.form.get("password")
24         verify_account = "SELECT * FROM account WHERE email =?"
25         stmt = ibm_db.prepare(conn, verify_account)
26         ibm_db.bind_param(stmt,1,email)
27         ibm_db.execute(stmt)
28         fetch_account = ibm_db.fetch_assoc(stmt)
29         if(fetch_account):
30             if(pbkdf2_sha256.verify(password,fetch_account['PASSWORD'])):
31                 userInfo={
32                     "fullName":fetch_account['FULLNAME'],
33                     "email":fetch_account['EMAIL'],
34                     "phoneNumber":fetch_account['PHONENUMBER'],
35                     "password":fetch_account['PASSWORD'],
36                 }
37                 return start_session(userInfo)
38             else:
39                 flash("Password is incorrect","loginError")
40                 return redirect(url_for('index',loginError=True))
41         flash("Sorry, user with this email id does not exist","loginError")
42         return redirect(url_for('index',loginError=True))
43

```

```

1  @app.route('/signup/',methods=['POST'])
2  def signup():
3      if request.method=="POST":
4          userInfo={
5              "fullName":request.form.get('fullName'),
6              "email":request.form.get('email'),
7              "phoneNumber":request.form.get('phoneNumber'),
8              "password":request.form.get('password'),
9          }
10         userInfo['password']=pbkdf2_sha256.encrypt(userInfo['
password'])
11         sql = "SELECT * FROM account WHERE email =?"
12         stmt = ibm_db.prepare(conn, sql)
13         ibm_db.bind_param(stmt,1,userInfo['email'])
14         ibm_db.execute(stmt)
15         account = ibm_db.fetch_assoc(stmt)
16         if account:
17             flash("Sorry,user with this email already exist","
signupError")
18             return redirect(url_for('index',signupError=True))
19         else:
20             insert_sql = "
INSERT INTO  account(fullName, email, phoneNumber, password) VALUES
(?, ?, ?, ?)
"
21             prep_stmt = ibm_db.prepare(conn, insert_sql)
22             ibm_db.bind_param(prepare_stmt, 1, userInfo['fullName'])
23             ibm_db.bind_param(prepare_stmt, 2, userInfo['email'])
24             ibm_db.bind_param(prepare_stmt, 3, userInfo['phoneNumber'
])
25             ibm_db.bind_param(prepare_stmt, 4, userInfo['password'])
26             ibm_db.execute(prepare_stmt)
27             return start_session(userInfo)
28         flash("Signup failed","signupError")
29         return redirect(url_for('index',signupError=True))
30
31
32 @app.route('/logout/',methods=["GET"])
33 def logout():
34     if request.method=="GET":
35         session.clear()
36         return redirect(url_for('index'))
37

```

```

1  @app.route('/')
2  def index():
3      if(session and '_flashes' in dict(session)):
4          loginError=request.args.get('loginError')
5          signupError=request.args.get('signupError')
6          if(loginError):
7              return render_template('./index.html',loginError=loginError,
8 carousel_content=carouselData['carousel_content'],currentYear=datetime.date.today().
9 year)
10         if(signupError):
11             return render_template('./index.html',signupError=signupError,
12 carousel_content=carouselData['carousel_content'],currentYear=datetime.date.today().
13 year)
14         if(session and '_flashes' not in dict(session)):
15             if(session['logged_in']==True):
16                 return render_template('./index.html',userInfo=session['user'],
17 carousel_content=carouselData['carousel_content'],currentYear=datetime.date.today().
18 year)
19             else:
20                 return render_template('./index.html',carousel_content=carouselData['
21 carousel_content'],currentYear=datetime.date.today().year)
22             else:
23                 return render_template('./index.html',carousel_content=carouselData['
24 carousel_content'],currentYear=datetime.date.today().year)
25
26 @app.route('/detect/', methods=['GET','POST'])
27 @login_required
28 def predict():
29     if request.method == 'POST':
30         title=request.form['title']
31         url = request.form['url']
32         checkprediction = inputScript.main(url)
33         prediction = model.predict(checkprediction)
34         output=prediction[0]
35         session['predicted']=True
36         print(output)
37         if(output==1):
38             pred = "Wohoo! You are good to go."
39             session['status']='safe'
40             session['pred'] = pred
41         else:
42             pred = "Oh no! This is a Malicious URL"
43             session['status']='unsafe'
44             session['pred'] = pred
45         session['title']=title
46         session['url']=url
47         insert_detection_info_stmt="
48 INSERT INTO DETECTIONHISTORY(email,title,url,status) VALUES(?,?,?,?)"
49         insert_detection_info = ibm_db.prepare(conn, insert_detection_info_stmt)
50         ibm_db.bind_param(insert_detection_info,1,session['user']['email'])
51         ibm_db.bind_param(insert_detection_info,2,session['title'])
52         ibm_db.bind_param(insert_detection_info,3,session['url'])
53         ibm_db.bind_param(insert_detection_info,4,session['status'])
54         ibm_db.execute(insert_detection_info)
55         if(session and session['logged_in']):
56             if(session['logged_in']==True):
57                 return redirect(url_for('predictionResult'))
58     if request.method == 'GET':
59         return render_template('./templates/predict-form.html',userInfo=session['user'
60 ])
61

```

```

1 @app.route('/detection-result/')
2 @login_required
3 def predictionResult():
4     if(session['predicted']==True):
5         urlInfo={
6             'message':session['pred'] ,
7             'title':session['title'],
8             'url':session['url'],
9             'status':session['status']
10        }
11        return render_template("./templates/prediction-result.html", urlInfo
=urlInfo,userInfo=session['user'])
12    else:
13        return redirect(url_for('predict'))
14
15
16 @app.route('/detection-history/')
17 @login_required
18 def detectionHistory():
19     if(session and session['logged_in']):
20         if(session['logged_in']==True):
21             get_detection_history_stmt = "
22 SELECT title,url,status FROM detectionHistory where email=?"
23             get_detection_history = ibm_db.prepare(conn,
24 get_detection_history_stmt)
25             ibm_db.bind_param(get_detection_history,1,session['user']['email
26 '])
27             ibm_db.execute(get_detection_history)
28             fetch_detection_history = ibm_db.fetch_assoc(
29 get_detection_history)
30             detection_history = []
31             ind = 0
32             while fetch_detection_history != False:
33                 detection_history.append(fetch_detection_history)
34                 ind += 1
35                 fetch_detection_history = ibm_db.fetch_assoc(
36 get_detection_history)
37             detection_history= detection_history[::-1]
38             return render_template('./templates/detection-history.html',
39 userInfo=session['user'],detectionHistory=detection_history)
40
41
42 @app.route('/about/')
43 def about():
44     if(session and session['logged_in']):
45         if(session['logged_in']==True):
46             return render_template('./templates/about.html',userInfo=session
47 ['user'],aboutContents=aboutData['aboutContents'])
48         else:
49             return render_template('./templates/about.html',aboutContents=
50 aboutData['aboutContents'])
51     else:
52         return render_template('./templates/about.html',aboutContents=
53 aboutData['aboutContents'])
54
55

```



```
1 @app.route('/contact/')
2 def contact():
3     if(session and session['logged_in']):
4         if(session['logged_in']==True):
5             return render_template('
6             ./templates/contact.html',userInfo=session['user'])
7         else:
8             return render_template('
9             ./templates/contact.html')
10
11
12 if __name__ == '__main__':
13     app.run(host='127.0.0.1', debug=True)
```

GitHub & Project Demo Link

GitHub: <https://github.com/IBM-EPBL/IBM-Project-44647-1660725800>

Demo Link: <https://drive.google.com/file/d/1pqEgfEm-q96Oa9qVDKYxfQSVFLwsgKi0/view?usp=sharing>