# HX 8001 – PROFESSIONAL READINESS FOR INNOVATION EMPLOYABILITY AND ENTREPRENEURSHIP

## WEB PHISHING DETECTION

### A PROJECT REPORT

**Submitted by**

Team ID : PNT2022TMID50116

**J.ABILA JESY(9513191001)**

**R.ANISHA(951319104004)**

**M.ANU PRIYA(951319104006)**

**M.CELSIYA(951319104012)**

*in partial fulfillment for the award of the degree*

*Of*

### BACHELOR OF ENGINEERING

### *In*

#### COMPUTER SCIENCE AND ENGINEERING

#### JAYARAJ ANNAPACKIAM CSI COLLEGE OF ENGINEERING NAZARETH

### ANNA UNIVERSITY: CHENNAI 600 025

**NOV-DEC 2022**

# TABLE OF CONTENTS

# CHAPTER 1

## INTRODUCTION

### 1.1 Project Overview

Now a days Phishing becomes a main area of concern for security researchers because it is not difficult to create the fake website which looks so close to legitimate website. Experts can identify fake websites but not all the users can identify the fake website and such users become the victim of phishing attack. Main aim of the attacker is to steal banks account credentials. In United States businesses, there is a loss of US$2billion per year because their clients become victim to phishing .

Phishing attacks are becoming successful because lack of user awareness. Since phishing attack exploits the weaknesses found in users, it is very difficult to mitigate them but it is very important to enhance phishing detection techniques. The general method to detect phishing websites by updating blacklisted URLs, Internet Protocol (IP) to the antivirus database which is also known as "blacklist" method. To evade blacklists attackers uses creative techniques to fool users by modifying the URL to appear legitimate via obfuscation and many other simple techniques including: fast-flux, in which proxies are automatically generated to host the web-page; algorithmic generation of new URLs; etc. Major drawback of this method is that, it cannot detect zero-hour phishing attack. Heuristic based detection which includes characteristics that are found to exist in phishing attacks in reality and can detect zero-hour phishing attack, but the characteristics are not guaranteed to always exist in such attacks and false positive rate in detection is very high.

To overcome the drawbacks of blacklist and heuristics based method, many security researchers now focused on machine learning techniques. Machine learning technology consists of a many algorithms which requires past data to make a decision or prediction on future data. Using this technique, algorithm will analyze various blacklisted and legitimate URLs and their features to accurately detect the phishing websites including zero- hour phishing websites.

### 1.2 Purpose

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. This type of e-banking website is known as a phishing website. 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.

**Common threats of web phishing:**

- 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.
- Large organizations may get trapped in different kinds of scams.

This Guided Project mainly focuses on applying a machine-learning algorithm to detect Phishing websites.

In order to detect and predict e-banking 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 datasets criteria to classify their legitimacy. The e-banking phishing website can be detected based on some important characteristics like URL and domain identity, and security and encryption criteria in the final phishing detection rate. Once a user makes a transaction online when he 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.

# CHAPTER 2

# LITERATURE SURVEY

The purpose or goal behind phishing is data, money or personal information stealing through the fake website. The best strategy for avoiding the contact with the phishing web site is to detect real time malicious URL. Phishing websites can be determined on the basis of their domains. They usually are related to URL which needs to be registered (low-level domain and upper-level domain, path, query). Recently acquired status of intra-URL relationship is used to evaluate it using distinctive properties extracted from words that compose a URL based on query data from various search engines such as Google and Yahoo. These properties are further led to the machine-learningbased classification for the identification of phishing URLs from a real dataset. This paper focus on real time URL phishing against phishing content by using phish-STORM. For this a few relationship between the register domain rest of the URL are consider also intra URL relentless is consider which help to dusting wish between phishing or non phishing URL. For detecting a phishing website certain typical blacklisted urls are used, but this technique is unproductive as the duration of phishing websites is very short. Phishing is the name of avenue. It can be defined as the manner of deception of an organization's customer to communicate with their confidential information in an unacceptable behaviour. It can also be defined as intentionally using harsh weapons such as Spasm to automatically target the victims and targeting their private information. As many of the failures being occurred in the SMTP are exploiting vectors for the phishing websites, there is a greater availability of communication for malicious message deliveries. Proposed a novel classification approach that use heuristic based feature extraction approach. In this, they have classified extracted features into different categories such as URL Obfuscation features, Hyperlink-based features. Moreover, proposed technique gives 92.5% accuracy. Also this model is purely depends on the quality and quantity of the training set and Broken links feature extraction.

## 2.1 Existing problem

Internet has been become an essential component of our everyday social and financial activities. Nevertheless, internet users may be vulnerable to different types of web threats, which may cause financial damages, identify theft,loss of private information, brand reputation damage and loss of customers confidence in e-commerce and online banking. Phishing is considered as a form of web threats that is defined as the art of impersonating a website of an honest enterprise aiming to obtain confidential information such as usernames, passwords and social security number.

So far, there is no single solution that can capture every phishing attack. In this article,we proposed an intelligent model for predicting phishing attacks based on artificial neural network particularly self- structuring neural networks. phishing is continuous problem where features significant in determining the type of web page are constantly changing. Thus,we need to constantly improve the network structure in order to cope with these changes.

Our model solves this problem by automating the process of structuring the network and shows high acceptance for noisy data , fault tolerance and high prediction accuracy. several experiments were conducted in our research,and the number of epochs differs in each experiment

## 2.2 References

[1] Gunter Ollmann, "The Phishing Guide Understanding &  Preventing Phishing Attacks", IBMInternet Security  Systems, 2007.

[2] https://resources.infosecinstitute.com/category/enterprise /phishing/the-phishing-landscape/phishing-data-attack statistics/#gref

[3] Mahmoud Khonji, Youssef Iraqi, "Phishing Detection: A  Literature Survey IEEE, and Andrew Jones, 2013

[4] Mohammad R., Thabtah F. McCluskey L., (2015)  Phishing websites dataset. Available:  https://archive.ics.uci.edu/ml/datasets/Phishing+Websites  Accessed January 2016

[5] http://dataaspirant.com/2017/01/30/how-decision-tree algorithm-works/

[6] http://dataaspirant.com/2017/05/22/random-forest algorithm-machine-learing/

[7] https://www.kdnuggets.com/2016/07/support-vector machines-simple-explanation.html

[8] www.alexa.com

[9] www.phishtank.com

## 2.3 Problem Statement Definition

Web Phishing is a form of cyber fraud, which implies that fraudsters use various means to impersonate the URL address and page content of a real website or use vulnerabilities in the server program of a real website to insert dangerous HTML code in certain pages of the site.

It 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:

> ➢ Obtaining personal information from an individual or organization.
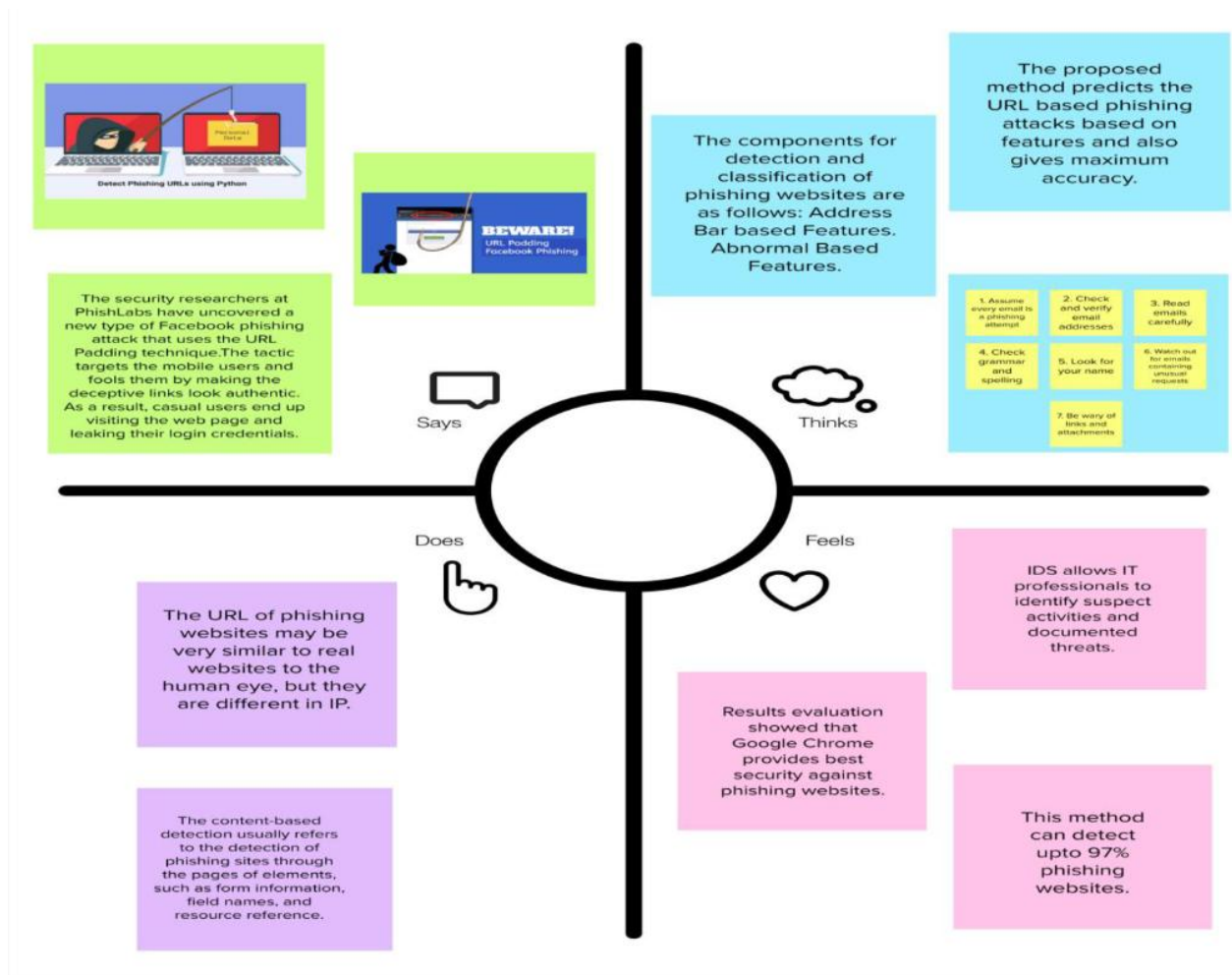> ➢ Impersonating as a trustworthy organization to deliver malicious websites.

To avoid these threats, we build an efficient and intelligent system to detect such websites using machine-learning algorithms which implements classification algorithms and techniques to extract the phishing datasets criteria to classify their legitimacy.

# CHAPTER 3

## IDEATION & PROPOSED SOLUTION

### 3.1 Empathy Map Canvas

An empathy map is a collaborative tool teams can use to gain a deeper insight into their customers. Much like a user persona, an empathy map can represent a group of users, such as a customer segment. Empathy maps should be used throughout any UX process to establish common ground among team members and to understand and prioritize user needs. In user-centered design, empathy maps are best used from the very beginning of the design process.



Reference:

### 3.2 Ideation & Brainstorming

Ideation essentially refers to the whole creative process of coming up with and communicating new ideas. Ideation is innovative thinking, typically aimed at solving a problem or providing a more efficient means of doing or accomplishing something.

Ideation is often closely related to the practice of brainstorming, a specific technique that is utilized to generate new ideas. A principal difference between ideation and brainstorming is that ideation is commonly more thought of as being an individual pursuit, while brainstorming is almost always a group activity.

Reference:-
https://app.mural.co/t/webphishingdetectionteamjaci1182/m/webphishingdetectionteamjaci1182/1664368513873/e5e6f585ef90541a2811cb1e1be418fe739ec9dd?sender=udb86796b365a1dd6434c0159

## 3.3 Proposed Solution

Project team shall fill the following information in proposed solution template.

| S.No | Parameter | Description |
|------|-----------|-------------|
| 1. | Problem Statement (Problem tobesolved) | **PROBLEM**<br>➢ Phishing detection techniques suffer low detectionaccuracy &high false alarm speciallyin phishing<br>**SOLUTION**<br>➢ Use anti-phishing protection & anti-spam software to protect yourself when maliciousmessage slip through to your computer<br>Anti-malware is included to prevent other types ofthreats<br>➢ Spam software<br>➢ Anti-malware software<br>These are programmed by security researches to spoteven the stealthiest +malware. |
| 2. | Idea / Solution description | **IDEA**<br>➢ Phishing assaults are usually detected by experiencedusers however, security is a primaryconcern for system users who are unaware & such situations<br>**SOLUTOIN**<br>➢ Anti-phishing technology is designed to identityand block phishing emails using a variety of methods certain anti-phishing solution scan the content of inbound & internal emails for any<br>sign of language that suggest a potentialphishing<br>or impersonation attack. |
| 3. | Novelty / Uniqueness | **NOVELTY/UNIQNESS**<br>➢ Since many phishing sites examines stayed live for atleast 48 hours, we monitored all sites for atleast 2 days. Based on Cyrene's analysis, googlechrome & fire fox did the best job dectection & blocking knows phishing sites with chrome blocking 74% phis site with in 6 hours 20 min<br>on average. |
| 4. | Social Impact / Customer Satisfaction | ➢ As a result of this, the organization as well consumers are facing enormous social effectsphishing is causing 2-way damage<br>➢ In the process of phishing, hackers mainly focuson extraction the details like bank account details, redidcard passwords and Aadhar card<br>verification |
| 5. | Business Model (Revenue Model) | **BUSINESS MODEL**<br>➢ Linking Aadhar card<br>➢ Frame work |

| 6. | Scalability of the Solution | **SCALABILITY OF THE SOLUTION** |
|---|---|---|
| | | ➢ Avoid falling victims to cyber-attacks including phishing, spear phishing and whaling phishing attempts account for the majority of malware &ransomware attacks. <br> **Highlights:** <br> • CRN Magazines available <br> • Newsletters available |

## 3.4 Problem Solution fit

The Problem-Solution Fit simply means that you have found a problem with your customer andthat the solution you have realized for it actually solves the customer's problem. It helps entrepreneurs, marketers and corporate innovators identify behavioral patterns and recognize what would workand why

**Purpose:**

❑ Solve complex problems in a way that fits the state of your customers.

❑ Succeed faster and increase your solution adoption by tapping into existing mediums andchannels of behaviour.

❑ Sharpen your communication and marketing strategy with the right triggers and messaging.

❑ Understand the existing situation in order to improve it for your target group.

**Template:**

| 3. TRIGGERS TR | 10. YOUR SOLUTION SL | 8. CHANNELS of BEHAVIOUR |
|---|---|---|
| Phishing is a type of social engineering attack often used to steal user data, including login credentials and credit card numbers. It occurs when an attacker, masquerading as a trusted entity, dupes a victim into opening an email, instant message, or text message. | 1. Don't click on that link<br>2. Get free anti-phishing add-ons<br>3. Don't give your information to an unsecured site<br>4. Rotate passwords regularly<br>5. Install firewalls | Types of phishing attacks range from classic email phishing schemes to more inventive approaches such as spear phishing and smishing. All have the same purpose – to steal your personal details. |
| **4. EMOTIONS: BEFORE / AFTER** EM<br><br>The findings from the study show that participants, generally, found it difficult to detect modern phishing email attacks. Saying that, participants were alert to the spelling mistakes of the older phishing email attacks, sensitive information being requested from them and any slight change to what they were normally used to from an email. | | |

**References:**

1. https://www.ideahackers.network/problem-solution-fit-canvas/

2. https://medium.com/@epicantus/problem-solution-fit-canvas-aa3dd59cb4fe

# CHAPTER 4

## REQUIREMENT ANALYSIS

### 4.1 FUNCTIONAL REQUIREMENTS

A function of software system is defined in functional requirement and the behavior of the system is evaluated when presented with specific inputs or conditions which may include calculations, data manipulation and processing and other specific functionality.

| FR NO. | Functional Requirement (Epic) | Sub Requirement (Story / Sub-Task) |
|--------|-------------------------------|------------------------------------|
| FR-1 | User Input | User inputs an URL in required field to check its validation. |
| FR-2 | Website Comparison | Model compares the websites using Blacklist andWhite list approach. |
| FR-3 | Feature extraction | After comparing, if none found on comparison then it extracts feature using heuristic and visual similarity approach. |
| FR-4 | Prediction | Model predicts the URL using Machine Learning algorithms such as Logistic Regression, KNN |
| FR-5 | Classifier | Model sends all output to classifier and producesfinal result. |
| FR-6 | Announcement | Model then displays whether website is a legal siteor a phishing site. |
| FR-7 | Events | This model needs the capability of retrieving anddisplaying accurate result for a website |

## 4.2 NON-FUNCTIONAL REQUIREMENTS

| FR No. | Non-Functional Requirement | Description |
|--------|----------------------------|-------------|
| NFR-1 | Usability | It is an easy to use and access interface which results in greater efficiency. |
| NFR-2 | Security | It is a secure website which protects the sensitive information of the user and prevents malicious attacks. |
| NFR-3 | Reliability | The system can detect phishing websites with greater accuracy using ML algorithms. |
| NFR-4 | Performance | The system produces responses within seconds and execution is faster. |
| NFR-5 | Availability | Users can access the website via any browser from anywhere at any time. |
| NFR-6 | Scalability | This application can be accessed online without paying. It can detect any web site with high accuracy. |

# CHAPTER 5

## PROJECT DESIGN

### 5.1 Data Flow Diagrams

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



### 5.2 Solution & Technical Architecture

Our solution is to build an efficient and intelligent system to detect phishing sites by applying a machine learning algorithm which implements classification algorithms and techniques to extract the phishing datasets criteria to classify their legitimacy by carefully analysing and identifying various factors that could be used to detect a phishing site. These factors fall under the categories of address bar-based features, domain-based features, HTML & JavaScript based features. Using these features, we can identify a phishing site with high accuracy.

### Technical Architecture

Technical architecture which is also often referred to as application architecture includes the major components of the system, their relationships, and the contracts that define the interactions

between the components. The goal of technical architects is to achieve all the business needs with an application that is optimized for both performance and security.



## 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 | I can receive confirmation email & click confirm | High | Sprint-1 |

| | | | application | | | |
|---|---|---|---|---|---|---|
| | | USN-3 | As a user, I can register for the applicationthrough Facebook | I can register & access the dashboard with FacebookLogin | 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 Exec: utive | 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 b etween websites for security. | High | Sprint-1 |
| Administ rator | Prediction | USN-1 | Here the Model will predict the URL websites using Machine Learning algorithms such as LogisticRegression, KNN | In this i can have correct prediction on th | High | Sprint-1 |
| | Classifier | USN-2 | Here i will send all the model output to classifier inorder to produce final result. | I this i will find the correct classifier f | Medium | Sprint-2 |

# CHAPTER 6

## PROJECT PLANNING & SCHEDULING

### 6.1 Sprint Planning & Estimation

| Sprint | Functional Requirement (Epic) | User Story Number | User Story / Task | Story Points | Priority | Team Members |
|--------|-------------------------------|-------------------|-------------------|--------------|----------|--------------|
| Sprint-1 | User input | USN-1 | User inputs an URL in the required field to check its validation. | 1 | Medium | Abila Jesy.J |
| Sprint-1 | Website Comparison | USN-2 | Model compares the websites using Blacklistand Whitelist approach. | 1 | High | Anisha.S |
| Sprint-2 | Feature Extraction | USN-3 | After comparison, if none found on comparison thenit extract feature using heuristic and visual similarity. | 2 | high | Anu Priya.M |
| Sprint-2 | Prediction | USN-4 | Model predicts the URL using Machine learningalgorithms such as logistic Regression, KNN. | 1 | Medium | Celsiya.M |
| Sprint-3 | Classifier | USN-5 | Model sends all the output to the classifier andproduces the final result. | 1 | Medium | Abila Jesy.J |
| Sprint-4 | Announcement | USN-6 | Model then displays whether the website islegal site or a phishing site. | 1 | High | Anisha.S |
| Sprint-4 | Events | USN-7 | This model needs the capability of retrieving anddisplaying accurate result for a website. | 1 | High | Celsiya.M |

## 6.2 Sprint Delivery Schedule

| Sprint | Total Story Points | Duration | Sprint Start Date | Sprint End Date (Planned) | Story Points Completed (as on Planned End Date) | Sprint Release Date (Actual) |
|---|---|---|---|---|---|---|
| Sprint-1 | 20 | 6 Days | 24 Oct 2022 | 29 Oct 2022 | 20 | 29 Oct 2022 |
| Sprint-2 | 20 | 6 Days | 31 Oct 2022 | 05 Nov 2022 | 20 | 05 Nov 2022 |
| Sprint-3 | 20 | 6 Days | 07 Nov 2022 | 12 Nov 2022 | 20 | 12 Nov 2022 |
| Sprint-4 | 20 | 6 Days | 14 Nov 2022 | 19 Nov 2022 | 20 | 19 Nov 2022 |

## 6.3 Reports from JIRA

**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)

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

**Burndown Chart:**

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.

https://www.visual-paradigm.com/scrum/scrum-burndown-chart/
https://www.atlassian.com/agile/tutorials/burndown-charts

**Reference:**

https://www.atlassian.com/agile/project-management

https://www.atlassian.com/agile/tutorials/how-to-do-scrum-with-jira-software
https://www.atlassian.com/agile/tutorials/epics

https://www.atlassian.com/agile/tutorials/sprints

https://www.atlassian.com/agile/project-management/estimation
https://www.atlassian.com/agile/tutorials/burndown-charts

# CHAPTER 7

## CODING & SOLUTIONING

### 7.1 Feature 1

We define two kinds of features to detect web phishing, and they are an original feature and interactive feature.

**Original Feature**

There are some features in the phishing URL, such as special characters. We definite these features in URL as an original feature as follows:

- ➢ O1:there are special characters in URL, such as @, Unicode, and so on. Those special characters are not allowed in a normal URL.
- ➢ O2: there are too many dots or less than four dots in normal URL.
- ➢ O3: the age of the domain is too short. For example, the age of the normal domain is more than 3 months.

In order to quantify the above characteristics, all the characteristic values are binary, that is, one of 0 or 1. Intuitively, the more of the 1 appear in the feature, the higher the likelihood that the site will be a phishing site.

**Interaction Feature**

There are some features in graph G=(V,E), such as access frequency. We define these features through a node relationship as interaction feature as follows:

- ➢ I1:  in-degree of URL node from REF is very small. In general, the normal websites do not link to phishing sites. The phishing sites are directly accessed.
- ➢ I2: out-degree of URL node is very small. In order to get personal private information, the phishing sites are usually terminal websites and do not link to the other sites.
- ➢ I3: the frequency of URL from AD is one. In general, one user accesses the phishing site only one time and the user cannot access the phishing site more than one time.
- ➢ I4: when AD accesses URL , user browser type UA is not the main browser. Well-known browser vendors often have a built-in filtering phishing site plug-in. A user who uses unknown browsers is more likely to access the phishing sites.
- ➢ I5: there is no cookie in user. The phishing site does not leave its cookie in user.

### 7.2 Feature 2

We have implemented python program to extract featuresfrom URL. Below are the features that we have extracted fordetection of phishing URLs.

**1) Presence of IP address in URL**: If IP address present in URL then the feature is set to 1 else set to 0. Most of the benign sites do not use IP address as an URL to download a webpage. Use of IP address in URL indicates that attacker is trying to steal sensitive information.

**2) Presence of @ symbol in URL**: If @ symbol present in URL then the feature is set to 1 else set to 0. Phishers add special symbol @ in the URL leads the browser to ignore everything preceding the "@" symbol and the real address often follows the "@" symbol.

**3) Number of dots in Hostname:** Phishing URLs have many dots in URL. For example http://shop.fun.amazon.phishing.com, in this URL phishing.com is an actual domain name, whereas use of "amazon" word is to trick users to click on it. Average number of dots in benign URLs is 3. If the number of dots in URLs is more than 3 then the feature is set to 1 else to 0.

**4) Prefix or Suffix separated by (-) to domain:** If domain name separated by dash (-) symbol then feature is set to 1 else to 0. The dash symbol is rarely used I legitimate URLs. Phishers add dash symbol (-) to the domain name so that users feel that they are dealing with a legitimate webpage. For example Actual site is http://www.onlineamazon.com but phisher can create another fake website like http://www.online-amazon.com to confuse the innocent users.

**5) URL redirection:** If "//" present in URL path then feature is set to 1 else to 0. The existence of "//" within the URL path means that the user will be redirected to another website.

**6) HTTPS token in URL:** If HTTPS token present in URL then the feature is set to 1 else to 0. Phishers may add the "HTTPS" token to the domain part of a URL in order to trick users. For example, http://https-www-paypal-it-mpp-home.soft-hair.com.

**7) Information submission to Email**: Phisher might use "mail()" or "mailto:" functions to redirect the user's information to his personal email[4]. If such functions are present in the URL then feature is set to 1 else to 0.

**8) URL Shortening Services "TinyURL":** TinyURL service allows phisher to hide long phishing URL by making it short. The goal is to redirect user to phishing websites. If the URL is crafted using shortening services (like bit.ly) then feature is set to 1 else 0

**9) Length of Host name:** Average length of the benign URLs is found to be a 25, If URL's length is greater than25 then the feature is set to 1 else to 0

**10) Presence of sensitive words in URL**: Phishing sites use sensitive words in its URL so that users feel that they are dealing with a legitimate webpage. Below are the words that found in many phishing URLs :- 'confirm', 'account', 'banking', 'secure', 'ebyisapi', 'webscr', 'signin', 'mail', 'install', 'toolbar', 'backup', 'paypal', 'password', 'username', etc;

**11) Number of slash in URL**: The number of slashes in benign URLs is found to be a 5; if number of slashes in URL is greater than 5 then the feature is set to 1 else to 0.

**12) Presence of Unicode in URL:** Phishers can make a use of Unicode characters in URL to trick users to click on it. For example the domain "xn--80ak6aa92e.com" is equivalent to "apple.com". Visible URL to user is "apple.com" but after clicking on this URL, user will visit to "xn--80ak6aa92e.com" which is a phishing site.

**13) Age of SSL Certificate:** The existence of HTTPS is very important in giving the impression of website legitimacy. But minimum age of the SSL certificate ofbenign website is between 1 year to 2 year.

**14) URL of Anchor:** We have extracted this feature by crawling the source code oh the URL. URL of the anchor is defined by <a> tag. If the <a> tag has a maximum number of hyperlinks which are from the other domain then the feature is set to 1 else to 0.

**15) IFRAME:** We have extracted this feature by crawling the source code of the URL. This tag is used to add another web page into existing main webpage. Phishers can make use of the "iframe" tag and make it invisible i.e. without frame borders. Since border of inserted webpage is invisible, user seems that the inserted web page is also the part of the main web page and can enter sensitive information.

**16) Website Rank:** We extracted the rank of websites and compare it with the first One hundred thousand websites of Alexa database. If rank of the website is greater than 10,0000 then feature

## 7.3 Database Scheme

## MACHINE LEARNING ALGORITHM

## Decision Tree Algorithm

One of the most widely used algorithm in machine learning technology. Decision tree algorithm is easy to understand and also easy to implement. Decision tree begins its work by choosing best splitter from the available attributes for classification which is considered as a root of the tree. Algorithm continues to build tree until it finds the leaf node. Decision tree creates training model which is used to predict target value or class in tree representation each internal node of the tree belongs to attribute and each leaf node of the tree belongs to class label. In decision tree algorithm, gini index and information gain methods are used to calculate these nodes.

## Random Forest Algorithm

Random forest algorithm is one of the most powerful algorithms in machine learning technology and it is based on concept of decision tree algorithm. Random forest algorithm creates the forest with number of decision trees. High number of tree gives high detection accuracy. Creation of trees are based on bootstrap method. In bootstrap method features and samples of dataset are randomly selected with replacement to construct single tree. Among randomly selected features,random forest algorithm will choose best splitter for the classification and like decision tree algorithm; Random forest algorithm also uses gini index and information gain methods to find the best splitter. This process will get continue until random forest creates n number of trees. Each tree in forest predicts the target value and then algorithm will calculate the votes for each

predicted target. Finally random forest algorithm considers high voted predicted target as a final prediction.

**<u>Support Vector Machine Algorithm</u>**

Support vector machine is another powerful algorithm in machine learning technology. In support vector machine algorithm each data item is plotted as a point in n-dimensional space and support vector machine algorithm constructs separating line for classification of two classes, this separating line is well known as hyper plane. Support vector machine seeks for the closest points called as support vectors and once it finds the closest point it draws a line connecting to them. Support vector machine then construct separating line which bisects and perpendicular to the connecting line. In order to classify data perfectly the margin should be maximum. Here the margin is a distance between hyperplane and support vectors. In real scenario it is not possible to separate complex and non linear data, to solve this problem support vector machine uses kernel trick which transforms lower dimensional space to higher dimensional space.

# CHAPTER 8

## TESTING

## 8.1 Test Cases

In machine learning systems, however, data and desired behavior are the inputs and the models learn the logic as the outcome of the training and optimization processes. In this case, testing involves validating the consistency of the model's logic and our desired behavior.

We usually write two different classes of tests for Machine Learning systems:

- ➢ Pre-train tests
- ➢ Post-train tests

**Pre-train tests:** The intention is to write such tests which can be run without trained parameters so that we can catch implementation errors early on. This helps in avoiding the extra time and effort spent in a wasted training job.

We can test the following in the pre-train test:

- ➢ the model predicted output shape is proper or not
- ➢ test dataset leakage i.e. checking whether the data in training and testing datasets have no duplication
- ➢ temporal data leakage which involves checking whether the dependencies between training and test data do not lead to unrealistic situations in the time domain like training on a future data point and testing on a past data point
- ➢ check for the output ranges. In the cases where we are predicting outputs in a certain range (for example when predicting probabilities), we need to ensure the final prediction is not outside the expected range of values.
- ➢ Ensuring a gradient step training on a batch of data leads to a decrease in the loss
- ➢ data profiling assertions

**Post-train tests:** Post-train tests are aimed at testing the model's behavior. We want to test the learned logic and it could be tested on the following points and more:

- ➢ invariance tests which involve testing the model by tweaking only one feature in a data point and checking for consistency in model predictions. For example, if we are working with a loan prediction dataset then change in sex should not affect an individual's eligibility for the loan given all other features are the same or in the case of titanic survivor probability prediction data, change in the passenger's name should not affect their chances of survival.
- ➢ Directional expectations wherein we test for a direct relation between feature values and predictions. For example, in the case of a loan prediction problem, having a higher credit score should definitely increase a person's eligibility for a loan.
- ➢ Apart from this, you can also write tests for any other failure modes identified for your model.

- ➢ Now, let's try a hands-on approach and write tests for the Medical Cost Personal Datasets. Here, we are given a bunch of features and we have to predict the insurance costs

## 8.2 User Acceptance Testing

**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 |

# 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 |

# CHAPTER 9

## RESULTS

### 9.1 Performance Metrics

The median efficiency is used to assess each categorization model's effectiveness. The final item will appear in the way it was envisioned. Graphical representations are used to depict information during classification. The percentage of predictions made using the testing dataset is used to gauge accuracy. By dividing the entire number of forecasts even by properly predicted estimates, it is simple to calculate. The difference between actual and anticipated output is used to calculate accuracy.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

Where TP = True Positives, TN = True Negatives, FN = False Negatives and FP = False Positives.

Thus, accuracy for all the four used models were calculated and ranked. XGBoost performed better than other models.

|   | ML Model | Train Accuracy | Test Accuracy |
|---|---|---|---|
| 3 | XGBoost | 0.913 | 0.905 |
| 0 | Decision Tree | 0.898 | 0.894 |
| 1 | Random Forest | 0.893 | 0.886 |
| 2 | SVM | 0.886 | 0.883 |

# CHAPTER 10

## 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.

# CHAPTER 12

## 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.

# CHAPTER 13

## APPENDIX

### 13.1 Source Code

### App.py

#importing required libraries


```python
from flask import Flask, request, render_template

import numpy as np

import pandas as pd

from sklearn import metrics

import warnings

import pickle

warnings.filterwarnings('ignore')

from feature import FeatureExtraction


file = open("pickle/model.pkl","rb")

gbc = pickle.load(file)

file.close()

app = Flask(_name_)

@app.route("/", methods=["GET", "POST"])

def index():

    if request.method == "POST":

        url = request.form["url"]

        obj = FeatureExtraction(url)

        x = np.array(obj.getFeaturesList()).reshape(1,30)
```

```python
    y_pred =gbc.predict(x)[0]

    #1 is safe

    #-1 is unsafe

    y_pro_phishing = gbc.predict_proba(x)[0,0]

    y_pro_non_phishing = gbc.predict_proba(x)[0,1]

    # if(y_pred ==1 ):

    pred = "It is {0:.2f} % safe to go ".format(y_pro_phishing*100)

    return render_template('index.html',xx =round(y_pro_non_phishing,2),url=url )

  return render_template("index.html", xx =-1)

if _name_ == "_main_":

  app.run(debug=True)
```

**feature.py**

```python
import ipaddress

import re

import urllib.request

from bs4 import BeautifulSoup

import socket

import requests

from googlesearch import search

import whois

from datetime import date, datetime

import time

from dateutil.parser import parse as date_parse

from urllib.parse import urlparse
```

```python
class FeatureExtraction:

    features = []

    def __init__(self,url):

        self.features = []

        self.url = url

        self.domain = ""

        self.whois_response = ""

        self.urlparse = ""

        self.response = ""

        self.soup = ""


        try:

            self.response = requests.get(url)

            self.soup = BeautifulSoup(response.text, 'html.parser')

        except:

            pass


        try:

            self.urlparse = urlparse(url)

            self.domain = self.urlparse.netloc

        except:

            pass


        try:

            self.whois_response = whois.whois(self.domain)
```

```python
    except:
        pass

self.features.append(self.UsingIp())

self.features.append(self.longUrl())

self.features.append(self.shortUrl())

self.features.append(self.symbol())

self.features.append(self.redirecting())

self.features.append(self.prefixSuffix())

self.features.append(self.SubDomains())

self.features.append(self.Hppts())

self.features.append(self.DomainRegLen())

self.features.append(self.Favicon())


self.features.append(self.NonStdPort())

self.features.append(self.HTTPSDomainURL())

self.features.append(self.RequestURL())

self.features.append(self.AnchorURL())

self.features.append(self.LinksInScriptTags())

self.features.append(self.ServerFormHandler())

self.features.append(self.InfoEmail())

self.features.append(self.AbnormalURL())

self.features.append(self.WebsiteForwarding())

self.features.append(self.StatusBarCust())
```

```python
        self.features.append(self.DisableRightClick())

        self.features.append(self.UsingPopupWindow())

        self.features.append(self.IframeRedirection())

        self.features.append(self.AgeofDomain())

        self.features.append(self.DNSRecording())

        self.features.append(self.WebsiteTraffic())

        self.features.append(self.PageRank())

        self.features.append(self.GoogleIndex())

        self.features.append(self.LinksPointingToPage())

        self.features.append(self.StatsReport())



    # 1.UsingIp

    def UsingIp(self):

        try:

            ipaddress.ip_address(self.url)

            return -1

        except:

            return 1



    # 2.longUrl

    def longUrl(self):

        if len(self.url) < 54:

            return 1

        if len(self.url) >= 54 and len(self.url) <= 75:
```

```python
            return 0
        return -1



    # 3.shortUrl
    def shortUrl(self):
        match = re.search('bit\.ly|goo\.gl|shorte\.st|go2l\.ink|x\.co|ow\.ly|t\.co|tinyurl|tr\.im|is\.gd|cli\.gs|'
                          'yfrog\.com|migre\.me|ff\.im|tiny\.cc|url4\.eu|twit\.ac|su\.pr|twurl\.nl|snipurl\.com|'
                          'short\.to|BudURL\.com|ping\.fm|post\.ly|Just\.as|bkite\.com|snipr\.com|fic\.kr|loopt\.us|'
                          'doiop\.com|short\.ie|kl\.am|wp\.me|rubyurl\.com|om\.ly|to\.ly|bit\.do|t\.co|lnkd\.in|'
                          'db\.tt|qr\.ae|adf\.ly|goo\.gl|bitly\.com|cur\.lv|tinyurl\.com|ow\.ly|bit\.ly|ity\.im|'
                          'q\.gs|is\.gd|po\.st|bc\.vc|twitthis\.com|u\.to|j\.mp|buzurl\.com|cutt\.us|u\.bb|yourls\.org|'

'x\.co|prettylinkpro\.com|scrnch\.me|filoops\.info|vzturl\.com|qr\.net|1url\.com|tweez\.me|v\.gd|tr\.im|link\.zip\.net',
self.url)
        if match:
            return -1
        return 1



    # 4.Symbol@
    def symbol(self):
        if re.findall("@",self.url):
            return -1
        return 1



    # 5.Redirecting//
    def redirecting(self):
        if self.url.rfind('//')>6:
```

```python
            return -1

        return 1


# 6.prefixSuffix

def prefixSuffix(self):

    try:

        match = re.findall('\-', self.domain)

        if match:

            return -1

        return 1

    except:

        return -1


# 7.SubDomains

def SubDomains(self):

    dot_count = len(re.findall("\.", self.url))

    if dot_count == 1:

        return 1

    elif dot_count == 2:

        return 0

    return -1


# 8.HTTPS

def Hppts(self):

    try:
```

```python
        https = self.urlparse.scheme

        if 'https' in https:

            return 1

        return -1

    except:

        return 1


# 9.DomainRegLen

def DomainRegLen(self):

    try:

        expiration_date = self.whois_response.expiration_date

        creation_date = self.whois_response.creation_date

        try:

            if(len(expiration_date)):

                expiration_date = expiration_date[0]

        except:

            pass

        try:

            if(len(creation_date)):

                creation_date = creation_date[0]

        except:

            pass


        age = (expiration_date.year-creation_date.year)*12+ (expiration_date.month-creation_date.month)

        if age >=12:
```

```python
            return 1

        return -1

    except:

        return -1


# 10. Favicon

def Favicon(self):

    try:

        for head in self.soup.find_all('head'):

            for head.link in self.soup.find_all('link', href=True):

                dots = [x.start(0) for x in re.finditer('\.', head.link['href'])]

                if self.url in head.link['href'] or len(dots) == 1 or domain in head.link['href']:

                    return 1

        return -1

    except:

        return -1


# 11. NonStdPort

def NonStdPort(self):

    try:

        port = self.domain.split(":")

        if len(port)>1:

            return -1

        return 1

    except:
```

```python
        return -1


# 12. HTTPSDomainURL

def HTTPSDomainURL(self):

    try:

        if 'https' in self.domain:

            return -1

        return 1

    except:

        return -1


# 13. RequestURL

def RequestURL(self):

    try:

        for img in self.soup.find_all('img', src=True):

            dots = [x.start(0) for x in re.finditer('\.', img['src'])]

            if self.url in img['src'] or self.domain in img['src'] or len(dots) == 1:

                success = success + 1

            i = i+1


        for audio in self.soup.find_all('audio', src=True):

            dots = [x.start(0) for x in re.finditer('\.', audio['src'])]

            if self.url in audio['src'] or self.domain in audio['src'] or len(dots) == 1:

                success = success + 1

            i = i+1
```

```python
    for embed in self.soup.find_all('embed', src=True):

        dots = [x.start(0) for x in re.finditer('\.', embed['src'])]

        if self.url in embed['src'] or self.domain in embed['src'] or len(dots) == 1:

            success = success + 1

        i = i+1


    for iframe in self.soup.find_all('iframe', src=True):

        dots = [x.start(0) for x in re.finditer('\.', iframe['src'])]

        if self.url in iframe['src'] or self.domain in iframe['src'] or len(dots) == 1:

            success = success + 1

        i = i+1


    try:

        percentage = success/float(i) * 100

        if percentage < 22.0:

            return 1

        elif((percentage >= 22.0) and (percentage < 61.0)):

            return 0

        else:

            return -1

    except:

        return 0

except:

    return -1
```

```python
# 14. AnchorURL

def AnchorURL(self):

    try:

        i,unsafe = 0,0

        for a in self.soup.find_all('a', href=True):

            if "#" in a['href'] or "javascript" in a['href'].lower() or "mailto" in a['href'].lower() or not (url in a['href'] or self.domain in a['href']):

                unsafe = unsafe + 1

            i = i + 1


        try:

            percentage = unsafe / float(i) * 100

            if percentage < 31.0:

                return 1

            elif ((percentage >= 31.0) and (percentage < 67.0)):

                return 0

            else:

                return -1

        except:

            return -1


    except:

        return -1


# 15. LinksInScriptTags
```

```python
def LinksInScriptTags(self):

    try:

        i,success = 0,0


        for link in self.soup.find_all('link', href=True):

            dots = [x.start(0) for x in re.finditer('\.', link['href'])]

            if self.url in link['href'] or self.domain in link['href'] or len(dots) == 1:

                success = success + 1

            i = i+1


        for script in self.soup.find_all('script', src=True):

            dots = [x.start(0) for x in re.finditer('\.', script['src'])]

            if self.url in script['src'] or self.domain in script['src'] or len(dots) == 1:

                success = success + 1

            i = i+1


        try:

            percentage = success / float(i) * 100

            if percentage < 17.0:

                return 1

            elif((percentage >= 17.0) and (percentage < 81.0)):

                return 0

            else:

                return -1

        except:
```

```python
            return 0

        except:

            return -1



# 16. ServerFormHandler

def ServerFormHandler(self):

    try:

        if len(self.soup.find_all('form', action=True))==0:

            return 1

        else :

            for form in self.soup.find_all('form', action=True):

                if form['action'] == "" or form['action'] == "about:blank":

                    return -1

                elif self.url not in form['action'] and self.domain not in form['action']:

                    return 0

                else:

                    return 1

    except:

        return -1



# 17. InfoEmail

def InfoEmail(self):

    try:

        if re.findall(r"[mail\(\)|mailto:?]", self.soap):

            return -1
```

```python
        else:

            return 1

    except:

        return -1


# 18. AbnormalURL

def AbnormalURL(self):

    try:

        if self.response.text == self.whois_response:

            return 1

        else:

            return -1

    except:

        return -1


# 19. WebsiteForwarding

def WebsiteForwarding(self):

    try:

        if len(self.response.history) <= 1:

            return 1

        elif len(self.response.history) <= 4:

            return 0

        else:

            return -1

    except:
```

```python
        return -1



# 20. StatusBarCust

def StatusBarCust(self):

    try:

        if re.findall("<script>.+onmouseover.+</script>", self.response.text):

            return 1

        else:

            return -1

    except:

        return -1



# 21. DisableRightClick

def DisableRightClick(self):

    try:

        if re.findall(r"event.button ?== ?2", self.response.text):

            return 1

        else:

            return -1

    except:

        return -1



# 22. UsingPopupWindow

def UsingPopupWindow(self):

    try:
```

```python
        if re.findall(r"alert\(", self.response.text):

            return 1

        else:

            return -1

    except:

        return -1


# 23. IframeRedirection

def IframeRedirection(self):

    try:

        if re.findall(r"[<iframe>|<frameBorder>]", self.response.text):

            return 1

        else:

            return -1

    except:

        return -1


# 24. AgeofDomain

def AgeofDomain(self):

    try:

        creation_date = self.whois_response.creation_date

        try:

            if(len(creation_date)):

                creation_date = creation_date[0]

        except:
```

```python
            pass

        today  = date.today()

        age = (today.year-creation_date.year)*12+(today.month-creation_date.month)

        if age >=6:

            return 1

        return -1

    except:

        return -1

# 25. DNSRecording

def DNSRecording(self):

    try:

        creation_date = self.whois_response.creation_date

        try:

            if(len(creation_date)):

                creation_date = creation_date[0]

        except:

            pass

        today  = date.today()

        age = (today.year-creation_date.year)*12+(today.month-creation_date.month)

        if age >=6:

            return 1

        return -1

    except:

        return -1

# 26. WebsiteTraffic
```

```python
def WebsiteTraffic(self):

    try:

        rank = BeautifulSoup(urllib.request.urlopen("http://data.alexa.com/data?cli=10&dat=s&url=" + url).read(), "xml").find("REACH")['RANK']

        if (int(rank) < 100000):

            return 1

        return 0

    except :

        return -1

# 27. PageRank

def PageRank(self):

    try:

        prank_checker_response = requests.post("https://www.checkpagerank.net/index.php", {"name": self.domain})


        global_rank = int(re.findall(r"Global Rank: ([0-9]+)", rank_checker_response.text)[0])

        if global_rank > 0 and global_rank < 100000:

            return 1

        return -1

    except:

        return -1




# 28. GoogleIndex

def GoogleIndex(self):

    try:

        site = search(self.url, 5)
```

```python
        if site:

            return 1

        else:

            return -1

    except:

        return 1


# 29. LinksPointingToPage

def LinksPointingToPage(self):

    try:

        number_of_links = len(re.findall(r"<a href=", self.response.text))

        if number_of_links == 0:

            return 1

        elif number_of_links <= 2:

            return 0

        else:

            return -1

    except:

        return -1


# 30. StatsReport

def StatsReport(self):

    try:

        url_match = re.search(

'at\.ua|usa\.cc|baltazarpresentes\.com\.br|pe\.hu|esy\.es|hol\.es|sweddy\.com|myjino\.ru|96\.lt|ow\.ly', url)
```

```python
        ip_address = socket.gethostbyname(self.domain)

        ip_match =
re.search('146\.112\.61\.108|213\.174\.157\.151|121\.50\.168\.88|192\.185\.217\.116|78\.46\.211\.158|181\.174\.165\.
13|46\.242\.145\.103|121\.50\.168\.40|83\.125\.22\.219|46\.242\.145\.98|'


'107\.151\.148\.44|107\.151\.148\.107|64\.70\.19\.203|199\.184\.144\.27|107\.151\.148\.108|107\.151\.148\.109|119\.
28\.52\.61|54\.83\.43\.69|52\.69\.166\.231|216\.58\.192\.225|'


'118\.184\.25\.86|67\.208\.74\.71|23\.253\.126\.58|104\.239\.157\.210|175\.126\.123\.219|141\.8\.224\.221|10\.10\.10\.
10|43\.229\.108\.32|103\.232\.215\.140|69\.172\.201\.153|'


'216\.218\.185\.162|54\.225\.104\.146|103\.243\.24\.98|199\.59\.243\.120|31\.170\.160\.61|213\.19\.128\.77|62\.113\.
226\.131|208\.100\.26\.234|195\.16\.127\.102|195\.16\.127\.157|'


'34\.196\.13\.28|103\.224\.212\.222|172\.217\.4\.225|54\.72\.9\.51|192\.64\.147\.141|198\.200\.56\.183|23\.253\.164\.
103|52\.48\.191\.26|52\.214\.197\.72|87\.98\.255\.18|209\.99\.17\.27|'


'216\.38\.62\.18|104\.130\.124\.96|47\.89\.58\.141|78\.46\.211\.158|54\.86\.225\.156|54\.82\.156\.19|37\.157\.192\.10
2|204\.11\.56\.48|110\.34\.231\.42', ip_address)

        if url_match:

            return -1

        elif ip_match:

            return -1

        return 1

    except:

        return 1

 def getFeaturesList(self):

        return self.features
```

**index.html**

```html
<!DOCTYPE html>

<html lang="en">

<head>
```

```html
<meta charset="UTF-8">

<meta http-equiv="X-UA-Compatible" content="IE=edge">

<meta name="viewport" content="width=device-width, initial-scale=1.0">

<meta name="description" content="This website is develop for identify the safety of url.">

<meta name="keywords" content="phishing url,phishing,cyber security,machine learning,classifier,python">

<meta name="author" content="VAIBHAV BICHAVE">


<!-- BootStrap -->

<link rel="stylesheet" href="https://stackpath.bootstrapcdn.com/bootstrap/4.5.0/css/bootstrap.min.css"

    integrity="sha384-9aIt2nRpC12Uk9gS9baDl411NQApFmC26EwAOH8WgZl5MYYxFfc+NcPb1dKGj7Sk" crossorigin="anonymous">


<link href="static/styles.css" rel="stylesheet">

<title>URL detection</title>


</head>


<body>


<div class=" container">

  <div class="row">

    <div class="form col-md" id="form1">

      <h2>PHISHING URL DETECTION</h2>


      <br>

      <form action="/" method ="post">
```

```html
<input type="text" class="form__input" name ='url' id="url" placeholder="Enter URL" required="" />

<label for="url" class="form__label">URL</label>

<button class="button" role="button" >Check here</button>

</form>


</div>


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


<br>

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


<br>

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

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

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

</div>
</div>

<br>

<h1>GitHub Team ID : PNT2022TMID10819</h1>


</div>


<!-- JavaScript -->

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

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

    crossorigin="anonymous"></script>

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

    integrity="sha384-Q6E9RHvbIyZFJoft+2mJbHaEWldlvI9IOYy5n3zV9zzTtmI3UksdQRVvoxMfooAo"

    crossorigin="anonymous"></script>

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

    integrity="sha384-OgVRvuATP1z7JjHLkuOU7Xw704+h835Lr+6QL9UvYjZE3Ipu6Tp75j7Bh/kR0JKI"

    crossorigin="anonymous"></script>

<script>


    let x = '{{xx}}';

    let num = x*100;

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

        num = 100-num;

    }

    let txtx = num.toString();

    if(x<=1 && x>=0.50){

        var label = "Website is "+txtx +"% safe to use...";

        document.getElementById("prediction").innerHTML = label;

        document.getElementById("button1").style.display="block";

    }

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

        var label = "Website is "+txtx +"% unsafe to use..."

        document.getElementById("prediction").innerHTML = label ;

        document.getElementById("button2").style.display="block";
```

```
        }


    </script>


</body>


</html>
```

**Styles.css**

```css
*,

*::after,

*::before {

  margin: 0;

  padding: 0;

  box-sizing: inherit;

  font-size: 62,5%;

}


body {

  padding: 10% 5%;

  background: #0f2027;

  background: linear-gradient(to right,#2c5364, #203a43, #0f2027);

  justify-content: center;

  align-items: center;

  height: 100vh;

  color: #fff;
```

```css
}


.form__label {

  font-family: 'Roboto', sans-serif;

  font-size: 1.2rem;

  margin-left: 2rem;

  margin-top: 0.7rem;

  display: block;

  transition: all 0.3s;

  transform: translateY(0rem);

}


.form__input {

  top: -24px;

  font-family: 'Roboto', sans-serif;

  color: #333;

  font-size: 1.2rem;

  padding: 1.5rem 2rem;

  border-radius: 0.2rem;

  background-color: rgb(255, 255, 255);

  border: none;

  width: 75%;

  display: block;

  border-bottom: 0.3rem solid transparent;

  transition: all 0.3s;
```

```css
}


.form__input:placeholder-shown + .form__label {

  opacity: 0;

  visibility: hidden;

  -webkit-transform: translateY(+4rem);

  transform: translateY(+4rem);

}




.button {

  appearance: button;

  background-color: transparent;

  background-image: linear-gradient(to bottom, #fff, #f8eedb);

  border: 0 solid #e5e7eb;

  border-radius: .5rem;

  box-sizing: border-box;

  color: #482307;

  column-gap: 1rem;

  cursor: pointer;

  display: flex;

  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","Segoe UI Symbol","Noto Color Emoji";

  font-size: 100%;

  font-weight: 700;

  line-height: 24px;
```

```css
  margin: 0;

  outline: 2px solid transparent;

  padding: 1rem 1.5rem;

  text-align: center;

  text-transform: none;

  transition: all .1s cubic-bezier(.4, 0, .2, 1);

  user-select: none;

  -webkit-user-select: none;

  touch-action: manipulation;

  box-shadow: -6px 8px 10px rgba(81,41,10,0.1),0px 2px 2px rgba(81,41,10,0.2);

}


.button:active {

  background-color: #f3f4f6;

  box-shadow: -1px 2px 5px rgba(81,41,10,0.15),0px 1px 1px rgba(81,41,10,0.15);

  transform: translateY(0.125rem);

}


.button:focus {

  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);

}


.main-body{

  display: flex;
```

```css
  flex-direction: row;

  width: 75%;

  justify-content:space-around;

}


.button1{

  appearance: button;

  background-color: transparent;

  background-image: linear-gradient(to bottom, rgb(160, 245, 174), #37ee65);

  border: 0 solid #e5e7eb;

  border-radius: .5rem;

  box-sizing: border-box;

  color: #482307;

  column-gap: 1rem;

  cursor: pointer;

  display: flex;

  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","Segoe UI Symbol","Noto Color Emoji";

  font-size: 100%;

  font-weight: 700;

  line-height: 24px;

  margin: 0;

  outline: 2px solid transparent;

  padding: 1rem 1.5rem;

  text-align: center;

  text-transform: none;
```

```css
  transition: all .1s cubic-bezier(.4, 0, .2, 1);

  user-select: none;

  -webkit-user-select: none;

  touch-action: manipulation;

  box-shadow: -6px 8px 10px rgba(81,41,10,0.1),0px 2px 2px rgba(81,41,10,0.2);

  display: none;

}

.button2{

  appearance: button;

  background-color: transparent;

  background-image: linear-gradient(to bottom, rgb(252, 162, 162), #ee3737);

  border: 0 solid #e5e7eb;

  border-radius: .5rem;

  box-sizing: border-box;

  color: #482307;

  column-gap: 1rem;

  cursor: pointer;

  display: flex;

  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","Segoe UI Symbol","Noto Color Emoji";

  font-size: 100%;

  font-weight: 700;

  line-height: 24px;

  margin: 0;

  outline: 2px solid transparent;

  padding: 1rem 1.5rem;
```

```css
  text-align: center;

  text-transform: none;

  transition: all .1s cubic-bezier(.4, 0, .2, 1);

  user-select: none;

  -webkit-user-select: none;

  touch-action: manipulation;

  box-shadow: -6px 8px 10px rgba(81,41,10,0.1),0px 2px 2px rgba(81,41,10,0.2);

  display: none;

}

.right {

  right: 0px;

  width: 300px;

}

@media (max-width: 576px) {

  .form {

    width: 100%;

  }

}

.abc{

  width: 50%;

}
```
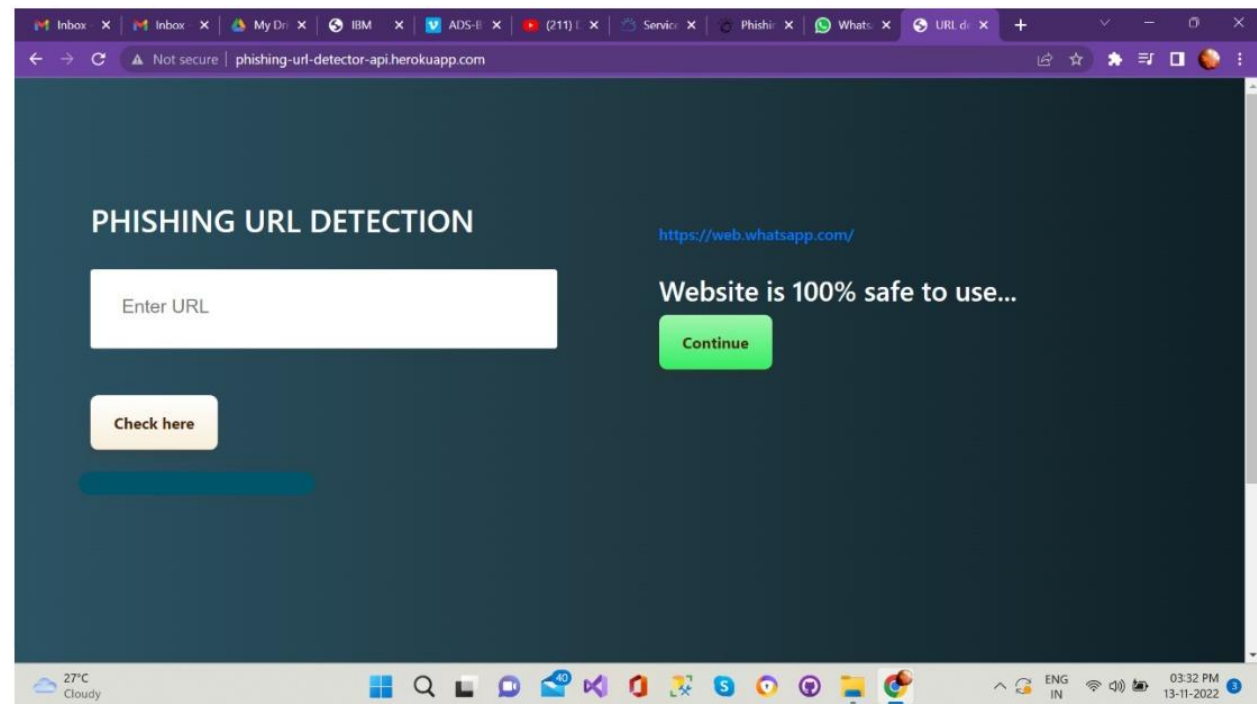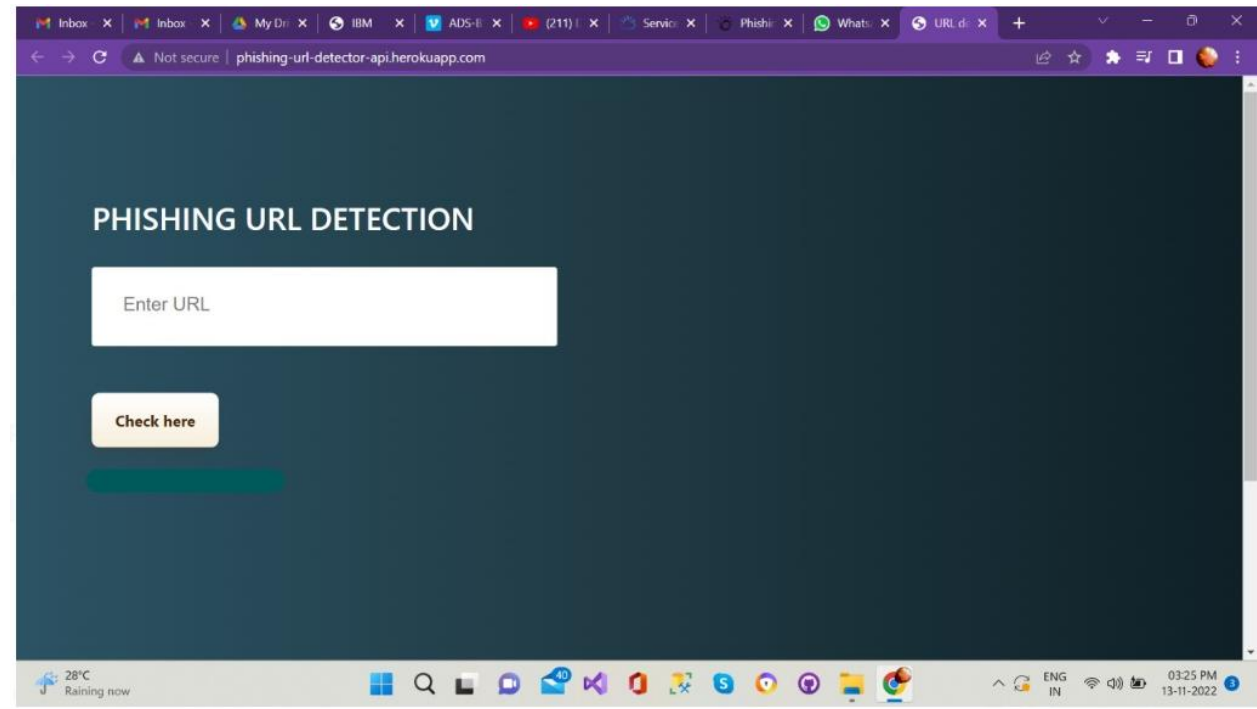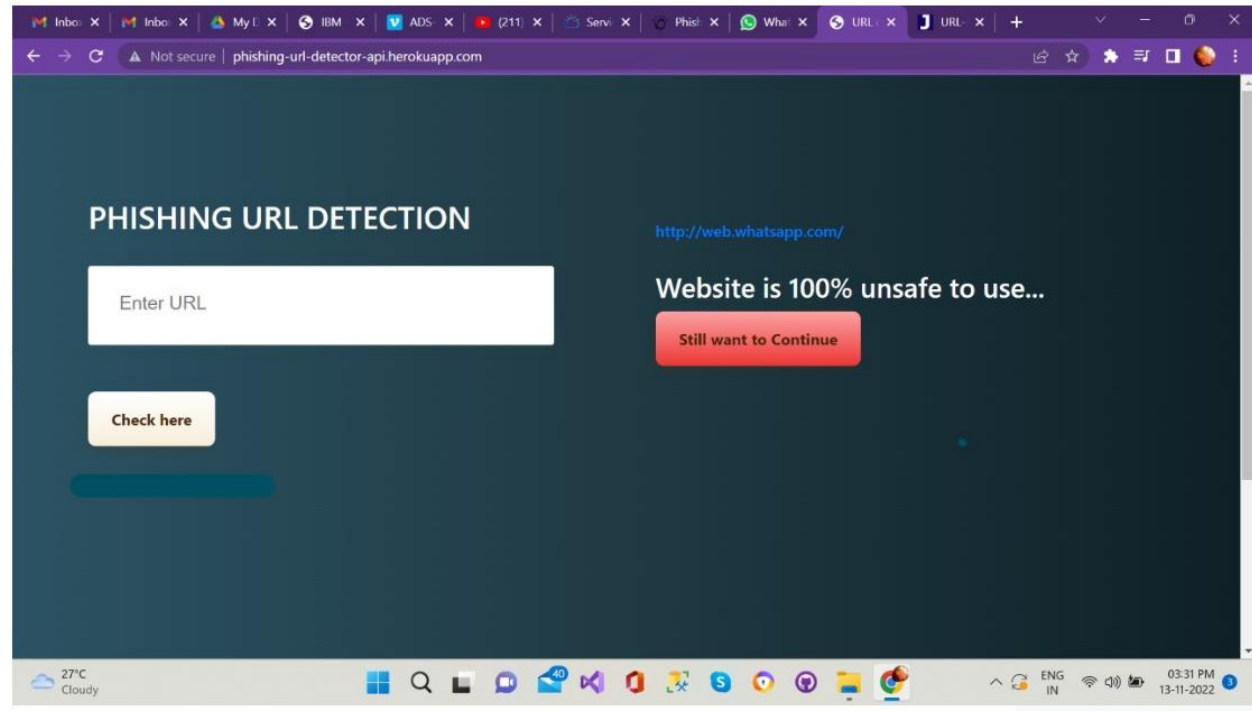
## Screenshots

## 13.2 GitHub & Project Demo Link

**GitHub Link:-** https://github.com/IBM-EPBL/IBM-Project-32984-1660213361

**Demo Link:-** https://youtu.be/pHqOrgJYKCg