# GOVERNMENT COLLEGE OF ENGINEERING(Formerly IRTT)
# ERODE-638 316



## BONAFIDE CERTIFICATE

Certified that this project titled **"WEB PHISHING DETECTION USING DATA SCIENCE"** is the bonafide work of **"SAGAR M(731119104040), SATHISH G(731119104042), SETHU K(731119104043), YUVAN JEYAN G(731119104054)"** who carried out the project work under my supervision.

**SIGNATURE OF HOD**

Dr.A SARADHA,M.E.,Ph.D.,
HEAD OF THE DEPARTMENT
DEPARTMENT OF CSE,
GOVERNMENT COLLEGE OF
ENGINEERING, ERODE – 638316

**SIGNATURE OF SPOC**

Dr.G.GOWRISON, M.E.,Ph.D.,
ASSISTANT PROFESSOR(SR)
DEPARTMENT OF ECE,
GOVERNMENT COLLEGE OF
ENGINEERING,ERODE - 638316

**SIGNATURE OF FACULTY MENTOR EVALUATOR**

Dr.A SARADHA,M.E.,Ph.D.,
ASSISTANT PROFESSOR(SR)
DEPARTMENT OF CSE,
GOVERNMENT COLLEGE OF
ENGINEERING, ERODE – 638316

**SIGNATURE OF FACULTY**

Dr.A KAVIDHA,M.E.,Ph.D.,
HEAD OF THE DEPARTMENT
DEPARTMENT OF CSE,
GOVERNMENT COLLEGE OF
ENGINEERING, ERODE - 638316

# Project Report Format

**1. INTRODUCTION**

      1.1 Project Overview

      1.2 Purpose

**2. LITERATURE SURVEY**

      2.1 Existing problem

      2.2 References

      2.3 Problem Statement Definition

**3. IDEATION & PROPOSED SOLUTION**

      3.1 Empathy Map Canvas

      3.2 Ideation & Brainstorming

      3.3 Proposed Solution

      3.4 Problem Solution fit

**4. REQUIREMENT ANALYSIS**

      4.1 Functional requirement

      4.2 Non-Functional requirements

**5. PROJECT DESIGN**

      5.1 Data Flow Diagrams

      5.2 Solution & Technical Architecture

      5.3 User Stories

**6. PROJECT PLANNING & SCHEDULING**

      6.1 Sprint Planning & Estimation

      6.2 Sprint Delivery Schedule

      6.3 Reports from JIRA

**7. CODING & SOLUTIONING**

      7.1 Feature 1

      7.2 Feature 2

**8. TESTING**

      8.1 Test Cases

      8.2 User Acceptance Testing

**9. RESULTS**

      9.1 Performance Metrics

**10. ADVANTAGES & DISADVANTAGES**

**11. CONCLUSION**

**12. FUTURE SCOPE**

**13. APPENDIX Source Code GitHub & Project Demo Link**

# 1.INTRODUCTION

## 1.1PROJECT OVERVIEW

Today's growing phishing websites pose significant threats due to their extremely undetectable risk. They anticipate internet users to mistake them as genuine ones in order to reveal user information and privacy, such as login ids, pass-words, credit card numbers, etc. without notice. This paper proposes a new approach to solve the anti-phishing problem. The new features of this approach can be represented by URL character sequence without phishing prior knowledge, various hyperlink information, and textual content of the webpage, which are combined and fed to train the XGBoost classifier. One of the major contributions of this paper is the selection of different new features, which are capable enough to detect 0-h attacks, and these features do not depend on any third-party services. In particular, we extract character level Term Frequency-Inverse Document Frequency (TF-IDF) features from noisy parts of HTML and plaintext of the given webpage. Moreover, our proposed hyperlink features determine the relationship between the content and the URL of a webpage. Due to the absence of publicly available large phishing data sets, we needed to create our own data set with 60,252 webpages to validate the proposed solution. This data contains 32,972 benign webpages and 27,280 phishing webpages. For evaluations, the performance of each category of the proposed feature set is evaluated, and various classification algorithms are employed. From the empirical results, it was observed that the proposed individual features are valuable for phishing detection. However, the integration of all the features improves the detection of phishing sites with significant accuracy. The proposed approach achieved an accuracy of 96.76% with only 1.39% false-positive rate on our dataset, and an accuracy of 98.48% with 2.09% false-positive rate on benchmark dataset, which outperforms the existing baseline approaches.

## 1.2 PURPOSE

The purpose of Phishing Domain Detection is detecting phishing domain names. Therefore, passive queries related to the domain name, which we want to classify as phishing or not, provide useful information to us. Phishing is a form of fraudulent attack where the attacker tries to gain sensitive information by posing as a reputable source. In a typical phishing attack, a victim opens a compromised link that poses as a credible website. The victim is then asked to enter their credentials, but since it is a "fake" website, the sensitive information is routed to the hacker and the victim gets hacked.

Phishing is popular since it is a low effort, high reward attack. Most modern web browsers, antivirus software and email clients are pretty good at detecting phishing websites at the source, helping to prevent attacks. The situation worsens when a lazy algorithm is trained and tested with a large dataset. Therefore, the performance of the research methodology used in this project may not perform so well if the wrong classifier is trained and tested with dataset size more than the classifier's capacity.

# LITERATURE SURVEY

## 2.1 EXISITING PROBLEM

- Users are easily deceived by fake websites, if they are first time consumers of the official product or service.

- User may loss very information to these sites and may face unrecoverable loss of economy or resources.

- The major issue is the lack of awareness due to which the user may loss a great range of personal information andsensitive paves way for the attack. This may also create blame on the user for unauthorised purchase, identity theft and other cybercriminal activities.

- This issue occurs when the user is in a hurry that makes them unaware of the details that indicate the fabricated site.

- Also occurs when the user is a beginner consumer of the service or website of concern.

- 96% of phishing attacks arrive by email. Another 3% are carried out through malicious websites and just 1% via phone.
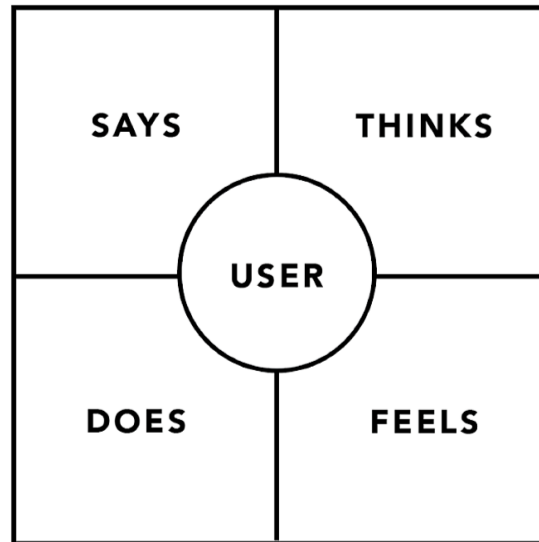
## 2.2 REFERENCES

1. Dhanalakshmi, R & Prabhu, C & Chellapan, C ' Detection of Phishing Websites and Secure Transactions Detection of Phishing Websites and Secure Transactions'. International Journal Communication & Network Security (IJCNS).

Jian Mao, Wenqian Tian, Pei Li, Tao Wei and Zhenkai Liang 'Phishing-Alarm: Robust and Efficient Phishing Detection via Page Component Similarity'. IEEE Access ( Volume: 5) 23 August 2017

2. Ramana Rao Kompella, and Minaxi Gupta. 'A machine learning based approach for phishing detection using hyperlinks information' vol.12, no.2, pp.1–27, 2007.

3. Pawan Prakash, Manish Kumar 'Phish net: predictive blacklisting to detect phishing attacks. SANS Institute, 2007. Accessed Jan 2018.

4. Dr. Gunikhan Sonowal: 'Phishing Scams Cost American Businesses Half A Billion Dollars a Year'. Forbes, 5 May 2017. Accessed Jan 2018.

## 2.3 PROBLEM STATEMENT DEFINITION

The URL of phishing websites may be very similar to real websites to the human eye, but they are different in IP. The content-based detection usually refers to the detection of phishing sites through the pages of elements, such as form information, field names, and resource reference.

# IDEATION &PROPOSED SOLUTION

## 3.1 EMPATHY MAP CANVAS



An **empathy map** is a collaborative visualization used to articulate what we know about a particular type of user. It externalizes knowledge about users in order to
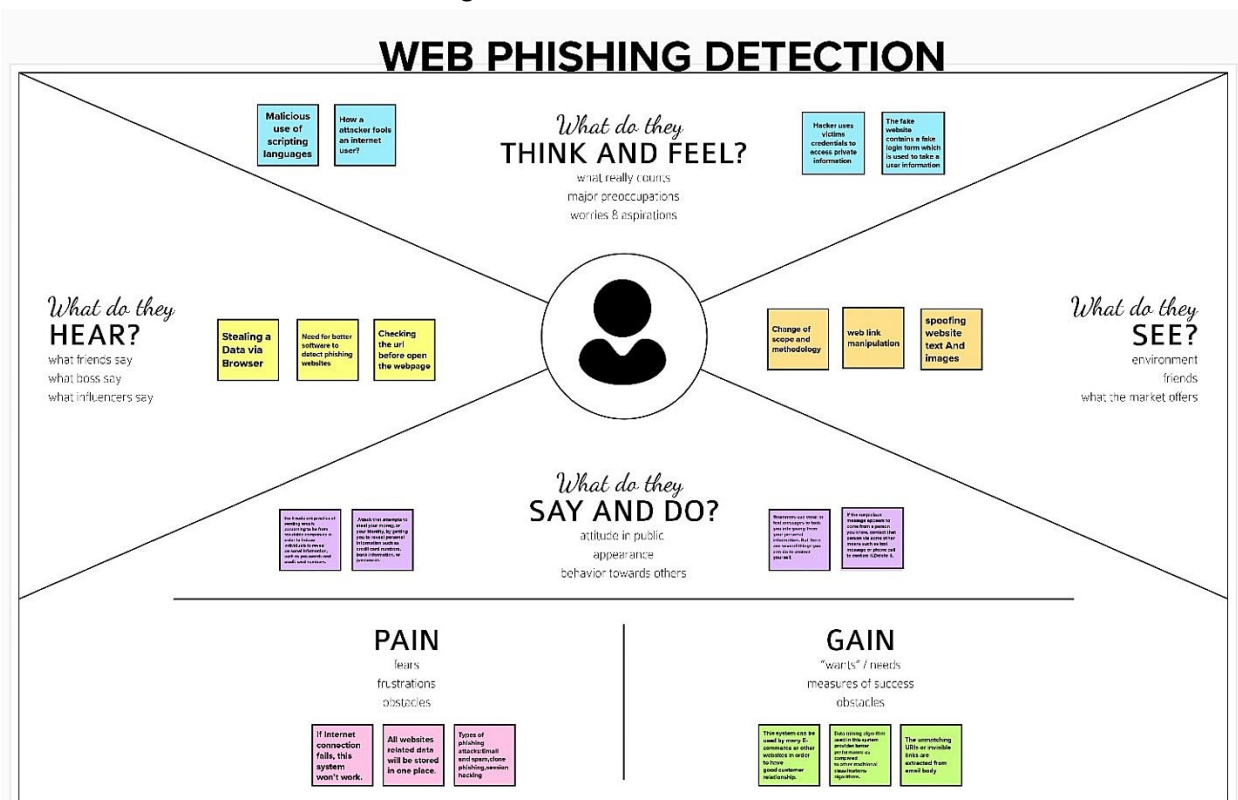➤ create a shared understanding of user needs, and
➤ aid in decision making.

Traditional empathy maps are split into 4 quadrants (*Says*, *Thinks*, *Does*, and *Feels*), with the user or persona in the middle. Empathy maps provide a glance into who a user is as a whole and are **not** chronological or sequential

➤ *Says* quadrant contains what the user says out loud in an interview or some other usability study.
➤ *Thinks* quadrant captures what the user is thinking throughout the experience.
➤ *Does* quadrant encloses the actions the user takes.
➤ *Feels* quadrant is the user's emotional state, often represented as an adjective plus a short sentence for context.

As their name suggests, empathy maps simply help us build empathy with our end users. When based on real data and when combined with other mapping method, they can:

➤ Remove bias from our designs and align the team on a single, shared understanding of the user
➤ Discover weaknesses in our research

➤ Uncover user needs that the user themselves may not even be aware of
➤ Understand what drives users' behaviours
➤ Guide us towards meaningful innovation



## WEB PHISHING DETECTION

**What do they THINK AND FEEL?**
what really counts
major preoccupations
worries & aspirations

Malicious use of scripting languages

How a attacker fools an internet user?

Hacker uses victim credentials to access private information

The fake website contains a fake login form which is used to take a user information

**What do they HEAR?**
what friends say
what boss say
what influencers say

Stealing a Data via Browser

Need for better software to detect phishing websites

Checking the url before open the webpage

Change of scope and methodology

web link manipulation

spoofing website text And images

**What do they SEE?**
environment
friends
what the market offers

**What do they SAY AND DO?**
attitude in public
appearance
behavior towards others

**PAIN**
fears
frustrations
obstacles

If Internet connection fails, this system won't work.

All websites related data will be stored in one place.

Types of phishing attacks:Email and spam,close phishing,session hacking

**GAIN**
"wants" / needs
measures of success
obstacles

This system can be used by every E-commerce or other websites in order to have good customer relationship.

Data mining algorithm used in this system provide better performance as compared to other machine/ classifications algorithms.

The unmatching URIs or invisible links are extracted from email body

## 3.2 IDEATION & BRAINSTORMING

Brainstorming is a great way to generate a lot of ideas that you would not be able to generate by just sitting down with a pen and paper. The intention of brainstorming is to leverage the collective thinking of the group, by engaging with each other, listening, and building on other ideas. Conducting a brainstorm also creates a distinct segment of time when you intentionally turn up the generative part of your brain and turn down the evaluative part. You can use brainstorming throughout any design or work process, of course, to generate ideas for design solutions, but also any time you are trying to generate ideas, such as planning where to do empathy work, or thinking about product and services related to our project.

Brainstorming has remained the cornerstone of the creative industry for decades and has evolved over the years as experience and learning developed from its variety of applications. Brainstorming essentially relies on a group of people coming together with their prior knowledge and research in order to gather ideas for solving the stated problem. It evokes images of exploration, experimental thinking, and wild ideas.

However, all too often it takes the form of controlled sessions where dominant figures assert themselves over others and creativity ends up getting stifled. Or, in other cases, the facilitator does not succeed in helping the team steer towards the goal by keeping the user, the user's need and the team's insights about them in a problem statement – also called Point-of-View at the front of the team's minds.

Ideation, is an art form, which is dependent on appointing an experienced facilitator and having an experienced team. However, we're all here to learn, and here is how you can start learning to become a successful facilitator of brainstorming sessions. Brainstorming is about setting a safe, creative space for people to feel like they can say anything and be wild—and know that they will not be judged for doing so—so that new ideas can be born.

- Set a time limit

- Start with a problem statement, Point of View, How Might We questions, a plan or a goal – and stay focused on the topic

- Defer judgement or criticism, including non-verbal

- Encourage weird, wacky, and wild ideas

- Aim for quantity

- Build on each other's ideas

- Be visual

- One conversation at a time

# BRAINSTORM AND IDEATION FOR

# WEB PHISHING DETECTION

## 3.3 PROPOSED SOLUTIONS

| S.No | Parameter | Description |
|------|-----------|-------------|
| 1. | Problem Statement (Problem to be solved) | A phishing website is a fake website, with domain similar in name and appearance to an official website. They're made in order to fool someone into believing it is legitimate. And then extract login credentials or confidential information such as credit card details from victims to perform malicious activities |
| 2. | Idea/ Solution description | Using data visualization, machine learning algorithm, the user will find the Legitimate websites for their transaction from this project. Keep all systems current with the latest security patches and updates. We implemented classification algorithm and techniques to extract the phishing data sets criteria (URL and Domain Identity, security and encryption) to classify their legitimacy. Then detect whether the website is a phishing site or not |
| 3. | Novelty/ Uniqueness | One of the major contributions of this project is the selection of different new features, which are capable enough to detect 0-h attacks, and these features do not depend on any third-party services. In particular, we extract character level Term Frequency Inverse Document Frequency (TF-IDF) features from noisy parts of HTML and plaintext of the given webpage |
| 4. | Social Impact/ Customer Satisfaction | Since Data mining algorithm used in this system, it provides better performance as compared to other traditional classifications algorithms. With the help of this system user can also purchase products online without any hesitation. The accuracy of phishing site identification is around 89%. |
| 5. | Business Model (Revenue Model) | A free web application system that can be used by many any person, E-commerce website or other websites in order to perform protected transaction. It's quick , free and trustworthy. |
| 6. | Scalability of the Solution | Additional features such as reporting the phishing website can be made. Also this web application can be made as a chrome extension so that users can securely surf through websites without any hesitation. |

## 3.4 PROBLEM SOLUTION FIT

The Problem-Solution Fit simply means that you have found a problem with your customer and that the solution you have realized for it actually solves the customer's problem. When we enter the stage of Problem-Solution-Fit assessment we do not engage in testing some specific product ideas, — it is yet too early for a product. It's all about testing whether the kind of solution itself looks attractive. It's enough to describe the value proposition you came up within a few sentences without features and technical details. Keep it as simple as possible.

| Define CS, fit into CC | **1. CUSTOMER SEGMENT(S)** **CS**<br>Who is your customer?<br>i.e. working parents of 0-5 y.o. kids<br><br>• User who uses online shopping websites.<br>• Individual who handle sensitive data and online transactions. | **6. CUSTOMER CONSTRAINTS** **CC**<br>What constraints prevent your customers from taking action or limit their choices of solutions? i.e. spending power, budget, no cash, network connection, available devices.<br><br>• the customer don't know where to report the issue.<br>• Lack of budget to improve the security system.<br>• they were not aware of the person behind these attacks | **5. AVAILABLE SOLUTIONS** **AS**<br>Which solutions are available to the customers when they face the problem<br>or need to get the job done? What have they tried in the past? What pros & cons do these solutions have? i.e. pen and paper is an alternative to digital notetaking<br><br>• Change the passwords on all accounts that use the same credentials.<br>• the website can be scanned so that the virus is prevented in user's mobile and computer | Explore AS, differentiate |
|---|---|---|---|---|
| Focus on J&P, tap into BE, understand RC | **2. JOBS-TO-BE-DONE / PROBLEMS** **J&P**<br>Which jobs-to-be-done (or problems) do you address for your customers? There could be more than one; explore different sides.<br><br>• Help to identify between fake and original websites.<br>• The user while visiting the website can be warned prior while they get into it. | **9. PROBLEM ROOT CAUSE** **RC**<br>What is the real reason that this problem exists?<br>What is the back story behind the need to do this job?<br>i.e. customers have to do it because of the change in regulations.<br><br>• Low security configurations and poor authentication.<br>• Not having prior knowledge to the users<br>• The ML prediction accuracy is less.<br>• There were not that much research are carried out in this field | **7. BEHAVIOUR** **BE**<br>What does your customer do to address the problem and get the job done?<br>i.e. directly related: find the right solar panel installer, calculate usage and benefits; indirectly associated: customers spend free time on volunteering work (i.e. Greenpeace)<br><br>• Report the phishing incident to cyber cell, turn off internet, scan the whole device to clear the virus.<br>• If the user has these kind of experience then they give a warning to the one who doesn't have prior knowledge about the pioblems while using the website | Focus on J&P, tap into BE, understand RC |
| Identify strong TR & EM | **3. TRIGGERS** **TR**<br>What triggers customers to act? i.e. seeing their neighboring installing solar panels, leading about a more efficient solution in the news.<br><br>• When a user is tricked into clicking a bad link.<br>• They might have no prior knowledge about the kind of attacks done while clicking the websites<br><br>**4. EMOTIONS: BEFORE / AFTER** **EM**<br>How do customers feel when they face a problem or a job and afterwards?<br><br>• They may feel insecure while using the website.<br>• They lose all their details and credit card information and because of that they feel frustrated.<br><br>i.e. lost, insecure > confident, in control - use it in your communication strategy & design. | **10.** **YOUR SL**<br>**SOLUTION**<br>If you are working on an existing business, write down you client solution first, fill in the canvas, and check how much it fits reality.<br>If you are working on a new business proposition, then keep it blank until you fill in the canvas and come up with a solution that fits within customer limitations, solves a problem and matches customer behavior.<br><br>• Allows the customer to check whether the attachment or the link received is legitimate in a more user-friendly manner.<br>• We can give prior alert box while using the website to predict that the website we are using is secure or not<br>• User must be aware of the phishing websites and they can prevent the loss of their personal information | **8. CHANNELS of BEHAVIOUR** **CH**<br>**8.1 ONLINE**<br>What kind of actions do customers take online? Extract online channels from 7<br><br>• They provide all their personal details including credit card information to some websites<br><br>**8.2 OFFLINE**<br>What kind of actions do customers take offline? Extract offline channels from 7 and use them for customer development.<br><br>• They try to research more information regarding attacks through books or from public<br>• Know what's a phishing scam. | Identify strong TR & EM |

# REQUIREMENT ANALYSIS

## 4.1 FUNCTIONAL REQUIREMENT

**Functional requirements** are product features or functions that developers must implement to enable users to accomplish their tasks. So, it's important to make them clear both for the development team and the stakeholders. Generally, functional requirements describe system behaviour under specific conditions. For example: The system sends an approval request after the user enters personal information. A search feature allows a user to hunt among various invoices if they want to credit an issued invoice. The system sends a confirmation email when a new user account is created.

| FR No | Functional Requirement (Epic) | Sub Requirement (Story / Sub-Task) |
|---|---|---|
| FR-1 | User Registration | Registration through Form. |
| FR-2 | User Confirmation | Confirmation via Email. |
| FR-3 | User Authentication | Authentication via Password. |
| FR-4 | User Input | User input an URL to check it is legal or phishing site. |
| FR-5 | Website Comparison | Model comparing the entered URL with the help of Blacklist and Whitelist. |
| FR-6 | Feature extraction | After comparing, if none found on comparison the it extracts feature using heuristic and visual similarity approach. |
| FR-7 | Prediction | Model Predicts the URL using Machine Learning algorithm such as Logistic Regression, KNN. |
| FR-8 | Classifier | Model sends output to classifier and it produce final result. |
| FR-9 | Announcement | Model the displays whether the website is a legal or phishing site. |
| FR-10 | Events | Model needs the capability of retrieving and displaying accurate resultfor a website. |

## 4.2 NON-FUNCTIONAL REQUIREMENT

**Non-functional Requirements** (NFRs) define 'how' systems do what they do. This includes characteristics such as their performance, security, maintainability scalability, and ease of use. Essentially, they provide the proper checks and balances to the functional requirements. NFRs are essential. But they're also a matter of balance and compromise. While you don't want to forget about them until it's too late, you also don't want to focus too much on NFRs at the expense of your costs, schedule, and resources.

| FR No | Non-Functional Requirement | Description |
|-------|----------------------------|-------------|
| NFR-1 | Usability | A set of specifications that describe the system's operation capabilities and constraints and attempt to improve its functionality. |
| NFR-2 | Security | Assuring all data inside the system or its part will be protected against malware attacks or unauthorized access. |
| NFR-3 | Reliability | This approach gives more accuracy then existing system. |
| NFR-4 | Performance | Performance Parameters for the proposed system gives accurate predicted value which is compared to the existing system. |
| NFR-5 | Availability | The system is accessible by user at any time using web browser. |
| NFR-6 | Scalability | The design will be suitable and performs with full efficiency according to rising demands. |

# PROJECT DESIGN

## 5.1 DATA FLOW DIAGRAMS

A Data Flow Diagram (DFD) is a traditional visual representation of the information flows within a system. A neat and clear DFD can depict the right amount of the system requirement graphically. It shows how data enters and leaves the system, what changes the information, and where data is stored. A data flow diagram (DFD) maps out the flow of information for any process or system. It uses defined symbols like rectangles, circles and arrows, plus short text labels, to show data inputs, outputs, storage points and the routes between each destination. Data flowcharts can range from simple, even hand-drawn process overviews, to in-depth, multi-level DFDs that dig progressively deeper into how the data is handled.

## 5.2 SOLUTIONS & TECHNICAL ARCHITECTURE

Solution architecture is comprised of several complex processes and sub-processes. It plays a central role in an organization's efforts to introduce and successfully implement new technology solutions. In the first step, solution architecture specialists closely look at how the different elements of business, information, and technology can be applied to solve a specific problem. Next, they propose a combination of building blocks that provides the best possible fix. This process is very detail-oriented and serves as a connecting piece between enterprise architecture and technical architecture. After solution architects have designed a solution for an existing problem, it is their job to manage the tasks and activities that are involved with its successful implementation.

The technical architect's main task is to realize particular technical implementation processes. As this requires a high level of in-depth expertise, technical architects usually specialize in one single technology. Technical architects can be in charge of leading large teams of developers and technical professionals. They act as technical project managers who define the structure of a specific system and oversee the related IT assignments. Out of all IT architects, they are the closest to an organization's end-user. Thus, they have to ensure that the technology is not only delivered in a timely manner but fully functional for the end-user.

## MODEL FOR WEB PHISHING DETECTION

## 5.3 USER STORIES

A user story is the smallest unit of work in an agile framework. It's an end goal, not a feature, expressed from the software user's perspective. A user story is an informal, general explanation of a software feature written from the perspective of the end user or customer. The purpose of a user story is to articulate how a piece of work will deliver a particular value back to the customer. Note that "customers" don't have to be external end users in the traditional sense, they can also be internal customers or colleagues within your organization who depend on your team. User stories are a few sentences in simple language that outline the desired outcome. They don't go into detail. Requirements are added later, once agreed upon by the team.

| Sprint | Functioal Requirement (Epic) | User Story Number | User Story/ Task | Story Points | Priority | Team Members |
|--------|------------------------------|-------------------|------------------|--------------|----------|--------------|
| Sprint-1 | Registration | USN-1 | As a user, I can registerfor the application by entering my Phone number/Gmail, Username password, and confirming my password. | 2 | High | Sagar M |
| Sprint-1 | Registration | USN-2 | As a user, I will receive confirmation in phone or Gmail once I have registered for the application | 1 | High | Sethu K |
| Sprint-1 | Login | USN-3 | As a user, I can log into the application by entering email & password | 1 | High | Sathish G |
| Sprint-1 | Verification | USN-4 | As a user, I can verify the registration and login of the user | 1 | High | Yuvan Jeyan G |
| Sprint-2 | Dataset | USN-5 | Collect number of datasets and get | 2 | Medium | Sagar M |

| | Collect | | accuracy | | | |
|---|---|---|---|---|---|---|
| Sprint-2 | Pre-processing | USN-5 | The dataset is extracted | 2 | High | Sathish G |
| Sprint-2 | Train the model | USN-6 | Train the model. | 4 | High | Sethu K |
| Sprint-2 | Test the model | USN-7 | Test the model | 6 | High | Yuvan Jeyan G |
| Sprint-3 | Detection | USN-8 | Load the trained model. | 3 | High | Sagar M |
| Sprint-3 | Detection | USN-9 | Prediction of legitimate site | 5 | Medium | Sathish G |
| Sprint-3 | Detection | USN-10 | classify it by using a trained model to predict the output | 8 | High | Sethu K |
| Sprint-4 | Detection | USN-11 | Alerts the user about the legitimate site | 7 | High | Yuvan Jeyan G |
| Sprint-4 | Detection | USN-12 | As a User, I can detect the phished site | 3 | Medium | Sagar M |
| Sprint-4 | Logout | USN-13 | As a User, I can logout the application. | 2 | Low | Sathish G |

# PROJECT PLANNING & SCHEDULING

## 6.1 SPRINT PLANNING & ESTIMATION

Sprint planning is an event in scrum that kicks off the sprint. The purpose of sprint planning is to define what can be delivered in the sprint and how that work will be achieved. Sprint planning is done in collaboration with the whole scrum team.

In scrum, the sprint is a set period of time where all the work is done. However, before you can leap into action you have to set up the sprint. You need to decide on how long the time box is going to be, the sprint goal, and where you're going to start. The sprint planning session kicks off the sprint by setting the agenda and focus. If done correctly, it also creates an environment where the team is motivated, challenged, and can be successful. Bad sprint plans can derail the team by setting unrealistic expectations.

- **The What** – The product owner describes the objective(or goal) of the sprint and what backlog items contribute to that goal. The scrum team decides what can be done in the coming sprint and what they will do during the sprint to make that happen.
- **The How** – The development team plans the work necessary to deliver the sprint goal. Ultimately, the resulting sprint plan is a negotiation between the development team and product owner based on value and effort.
- **The Who** – You cannot do sprint planning without the product owner or the development team. The product owner defines the goal based on the value that they seek. The development team needs to understand how they can or cannot deliver that goal. If either is missing from this event it makes planning the sprint almost impossible.
- **The Inputs** – A great starting point for the sprint plan is the product backlog as it provides a list of 'stuff' that could potentially be part of the current sprint. The team should also look at the existing work done in the increment and have a view to capacity
- **The Outputs** – The most important outcome for the sprint planning meeting is that the team can describe the goal of the sprint and how it will start working toward that goal. This is made visible in the sprint backlog.

Sprint planning requires some level of estimation. The team needs to define what can or cannot be done in the sprint:  estimated effort vs capacity. Estimation is often confused with commitments. Estimates are by their very nature forecasts based on the knowledge at hand. Techniques such as story points or t-shirt sizing add value to the process by giving the team a different way of looking at the problem. They are not,

however, magical tools that can find out the truth when there is none to be found. The more unknowns, the less likely the estimate will be correct.

Good estimation requires a trust-based environment where information is given freely, and assumptions are discussed in the pursuit of learning and improvement. If estimates are used in a negative, confrontational way after the work is completed, then it's likely that future estimates will be either be much bigger to ensure they never are wrong again or the time taken to create them will be much longer as the team second guesses itself worrying about the implications of getting them wrong.

## 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 REPORT FROM JIRA

Jira is a software platform for agile development and customer support. It's one of the most popular platforms for agile development — which makes it a go-to for many software development teams.

# CODING & SOLUTIONING

## 7.1 FEATURE 1

**home.html**

```html
1   <!doctype html>

2   <html lang="en">

3   <head>

4     <link rel="stylesheet" type="text/css"
    href="{{url_for('static',filename='css/style.css')}}">

5     <meta charset="utf-8">

6     <meta name="viewport" content="width=device-width, initial-scale=1">

7     <title>Home</title>

8
```

```html
9    <style>

10

11

12

13    </style>
14  <link href="https://cdn.jsdelivr.net/npm/bootstrap@5.2.2/dist/css/bootstrap.min.css"
    rel="stylesheet" integrity="sha384-
    Zenh87qX5JnK2Jl0vWa8Ck2rdkQ2Bzep5IDxbcnCeuOxjzrPF/et3URy9Bv1WTRi" crossorigin="anonymous">
15  </head>

16

17

18  <body class="bg-co"  >

19    <div class="bg-nav text-light d-flex flex-column flex-md-row align-items-center pb-3 mb-
    4 border-bottom">

20      <h5 class="my-0 mr-md-auto font-weight-bold mt-3" style="font-size:20px;opacity: 0.5;
    font-family: Georgia, serif; font-weight: bold;  padding-left: 50px;">URL Prediction</h5>

21      <nav class="d-inline-flex mt-2 mt-md-0 ms-md-auto ">

22        <a style="font-family: Georgia, serif;font-weight: bold;margin-right:20px;
    color:black;"  href="/predicturl">Predict URL</a>

23        <a style="font-family: Georgia, serif;font-weight: bold;margin-right:20px;
    color:black;" href="/addurl">Add url</a>

24        <a style="font-family: Georgia, serif;font-weight: bold; margin-right:20px;
    color:black" href="/about">About</a>

25      </nav>

26    </div>

27    <div class="container bg-co">

28      <div class="row">

29        <div class= "col-md-6">

30          <h1 style="font-size:60px; font-weight: bold;font-family: Georgia, serif;
    color:black;">

31            THE MOST RELAIABLE WAY TO PREDICT THE FUTURE IS TO CREATE IT

32          </h1>

33          <b><h5 style="text-align: right;font-family: Georgia, serif;color:black; font-
    weight: bold;color: rgba(0, 0, 0, 0.705);">- ABRAHAM LINCOLN</h5></b>
```

```
34        </div>

35        <div class="col-md-6 " >

36        <div  style="margin-top: 120px;margin-left: 190px;">

37              <form class="form"  action="/predicturl">

38                <center><img src="{{url_for('static', filename='Code.jpg')}}"
   style="height: 40%; width: 40%; opacity: 0.5; " class="img img-responsive img-circle  mx-
   auto d-block" /><br>

39                  <input type="submit" style="background-color: rgba(0, 0, 0,
   0.801);color: white; font-weight: bolder;"class="btn  btn-lg btn-block mx-auto "
   value="PREDICT YOUR URL">

40                  </center>

41              </form>

42          </div>

43          </div>

44        </div>

45    </div>

46  </div>

47

48      <script
   src="https://cdn.jsdelivr.net/npm/@popperjs/core@2.11.6/dist/umd/popper.min.js"
   integrity="sha384-oBqDVmMz9ATKxIep9tiCxS/Z9fNfEXiDAYTujMAeBAsjFuCZSmKbSSUnQlmh/jp3"
   crossorigin="anonymous"></script>

49      <script src="https://cdn.jsdelivr.net/npm/bootstrap@5.2.2/dist/js/bootstrap.min.js"
   integrity="sha384-IDwe1+LCz02ROU9k972gdyvl+AESN10+x7tBKgc9I5HFtuNz0wWnPclzo6p9vxnk"
   crossorigin="anonymous"></script>

50      <script
   src="https://cdn.jsdelivr.net/npm/bootstrap@5.2.2/dist/js/bootstrap.bundle.min.js"
   integrity="sha384-OERcA2EqjJCMA+/3y+gxIOqMEjwtxJY7qPCqsdltbNJuaOe923+mo//f6V8Qbsw3"
   crossorigin="anonymous"></script>

51    </body>

52  </html>
```

**EXPLANATION**

URL is the first thing to analyse a website to decide whether it is a phishing or not. As we mentioned before, URLs of phishing domains have some distinctive points. Features which are related to these points are obtained when the URL is processed. Some of URL-Based Features are given below.

- Digit count in the URL
- Total length of URL
- Checking whether the URL is Typo squatted or not. (google.com → goggle.com)
- Checking whether it includes a legitimate brand name or not (apple-icloud-login.com)
- Number of subdomains in URL
- Is Top Level Domain (TLD) one of the commonly used one?

## 7.2 FEATURE 2

**predict1.html**

```
1   <!doctype html>
2   <html lang="en">
3   <head>
4     <link        rel="stylesheet"        type="text/css"        href=
      "{{url_for('static',filename='css/style.css')}}">
5     <meta charset="utf-8">
6     <meta name="viewport" content="width=device-width, initial-scale=1">
7     <title>URL Prediction</title>
8   <script>
9   function clearInput() {
10    document.getElementById("Form").reset();
11  }
12  </script>
13  <link
    href="https://cdn.jsdelivr.net/npm/bootstrap@5.2.2/dist/css/bootstrap.min.css"
    rel="stylesheet"                                integrity="sha384-
    Zenh87qX5JnK2Jl0vWa8Ck2rdkQ2Bzep5IDxbcnCeuOxjzrPF/et3URy9Bv1WTRi"
    crossorigin="anonymous">
14  </head>
15  <body class="bg-co">
16    <div class="bg-nav text-light d-flex flex-column flex-md-row align-items-center
    pb-3 mb-4 border-bottom">
17          <h5    class="my-0    mr-md-auto    font-weight-bold    mt-3"   style="font-
    size:20px;opacity: 0.5; font-family: Georgia, serif; font-weight: bold;  padding-
    left: 50px;">URL Prediction</h5>
18      <nav class="d-inline-flex mt-2 mt-md-0 ms-md-auto ">
19          <a  class="me-3  py-2  text-light  text-decoration-none  mt-3"  style="font-
    family: Georgia, serif;font-weight: bold;margin-right:20px; "  href="/">Home</a>
20          <a  class="me-3  py-2  text-light  text-decoration-none  mt-3"  style="font-
```

```html
        family: Georgia, serif;font-weight: bold;margin-right:20px; " href="/addurl">Add
    url</a>
21          <a class="py-2 text-light text-decoration-none mt-3" style="font-family:
    Georgia, serif;font-weight: bold; margin-right:20px;" href="/about">About</a>
22      </nav>
23  </div>
24    <div class="bg-co" style="margin-top: 120px;  margin-left: 400px;margin-right:
    400px;">
25      <div class="card-body">
26        <form id="Form" action="/predict" method='post'class="form">
27              <img src="{{url_for('static', filename='code.png')}}" style="height:
    20%; width: 20%;" class="img img-responsive img-circle  mx-auto d-block" /><br>
28          <label><b>Enter URL to predict</b></label><br>
29                <input type="text" name="url" id= "myText" placeholder="Ex :
    https://abcde.com/" class="form-control" required><br>
30              <div class="w3-bar "><center>
31                    <input type="submit"  style="background-color: black; font-
    weight: bold; color: white;" class="btn "  value="Predict URL" >
32                    <input type= "button" style="background-color: black; font-
    weight:  bold;  color:  white;"  class="  btn  "  value=  "Clear"  onclick=
    "clearInput()"></center>
33              </div>
34        </form>
35      </div>
36      </div>
37                                                                    <script
    src="https://cdn.jsdelivr.net/npm/@popperjs/core@2.11.6/dist/umd/popper.min.js"
    integrity="sha384-
    oBqDVmMz9ATKxIep9tiCxS/Z9fNfEXiDAYTujMAeBAsjFuCZSmKbSSUnQlmh/jp3"
    crossorigin="anonymous"></script>
38                                                                    <script
    src="https://cdn.jsdelivr.net/npm/bootstrap@5.2.2/dist/js/bootstrap.min.js"
    integrity="sha384-
    IDwe1+LCz02ROU9k972gdyvl+AESN10+x7tBKgc9I5HFtuNz0wWnPclzo6p9vxnk"
    crossorigin="anonymous"></script>
39                                                                    <script
    src="https://cdn.jsdelivr.net/npm/bootstrap@5.2.2/dist/js/bootstrap.bundle.min.j
    s"                                                           integrity="sha384-
    OERcA2EqjJCMA+/3y+gxIOqMEjwtxJY7qPCqsdltbNJuaOe923+mo//f6V8Qbsw3"
    crossorigin="anonymous"></script>
40    </body>
41 </html>
42
```

**EXPLANATION**

Page-Based Features are using information about pages which are calculated reputation ranking services. Some of these features give information about how much reliable a web site is. Some of Page-Based Features are given below.

- Global PageRank
- Country PageRank
- Position at the Alexa Top 1 Million Site
- Some Page-Based Features give us information about user activity on target site. Some of these features are given below. Obtaining these types of features is not easy. There are some paid services for obtaining these types of features.
- Estimated Number of Visits for the domain on a daily, weekly, or monthly basis
- Average Pageviews per visit
- Average Visit Duration
- Web traffic share per country
- Count of reference from Social Networks to the given domain
- Category of the domain
- Similar websites etc.

# TESTING

## 8.1 TEST CASE

| Test Scenario | Expected Result |
|---|---|
| Verify user is able to enter the URL in the form | Result of classification will be displayed |
| Verify the UI elements in the form | Application should show below UI elements:<br> a.input form box<br> c.submit button<br> d.services offered<br> e.team |
| Verify user is able to see an alert when nothing is entered in the textbox | Alert of incomplete input |
| Verify user is able to see the result when URL is entered in the textbox | Result of classification will be displayed |
| Verify user is able to enter their name, email and query message in the form | Details are stored in the database |

## 8.2 USER ACCEPTANCE TESTING

1. Purpose of Document

The purpose of this documentis to briefly explain the test coverageand open issues ofthe Web PhishingDetection project at the time of the release to User Acceptance Testing (UAT).

2. Defect Analysis

This reportshows the numberof resolved or closed bugs at each severity level,and how they wereresolvedalysis

| 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 |

3. 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 |

# RESULTS

## 9.1 PERFORMANCE METRICS

# ADVANTAGE & DISADVANTAGE

**ADVANTAGES :**

- ➤ Measure the degrees of corporate and employee vulnerability
- ➤ Eliminate the cyber threat risk level
- ➤ Increase user alertness to phishing risks
- ➤ Instill a cyber security culture and create cyber security heroes
- ➤ Change behavior to eliminate the automatic trust response
- ➤ Deploy targeted anti-phishing solutions
- ➤ Protect valuable corporate and personal data
- ➤ Meet industry compliance obligations
- ➤ Assess the impacts of cyber security awareness training
- ➤ Segment phishing simulation

**DISADVANTAGES:**

- ➤ Loss of data
- ➤ Damaged Reputation
- ➤ Direct Monetry Loss
- ➤ Loss Of Productivity
- ➤ Loss of Customers
- ➤ Financial Penalties

# 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 that perform phishing detection by classification of websites using trained machine learning models. URL based analysis increases the speed of detection. Furthermore, by applying feature selection algorithms and dimensionality reduction techniques, we can reduce the number of features and remove irrelevant data. There are many machine learning algorithms that perform classification with good performance measures. In this paper, we have done a study of the process of phishing detection and the phishing detection schemes in the recent research literature. This will serve as a guide for new researchers to understand the process and to develop more accurate phishing detection systems.

# FUTURE SCOPE

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

The widely publicized Gmail phishing scam earlier this year is just one example of a modern threat that affected users on a large scale. In this case, users were sent an email that appeared legitimate and directed them to an actual Google page. While most phishing scams rely on pushing users to a malicious domain, this particular attack simply led unsuspecting individuals to granting broad permissions to a malicious application. Hackers could then see victims' contacts, read their emails, have insight into the users' locations, and see files created in G Suite.

The Gmail phishing attack shows us just how advanced these techniques have become – it was difficult to detect and difficult to prevent. A critical takeaway is that the attack was able to clear the psychological trust hurdle. Users were tricked into giving permissions to a third party application because they trusted it; they believed the

application to be a Google-approved service. A minute change in how the application domain was disguised successfully convinced users that the application was trustworthy.

This is the future of phishing. The ability to spoof cloud apps while masking the true identity of the sender in order to steal personal information – an alarming trend given the rapid increase of cloud adoption in verticals around the world.

## Traditional Phishing

Traditional phishing was rather simplistic in execution and relied on the user's lack of knowledge. For example, social engineering driven by phone calls and emails wherein malicious actors would pose as government agents or corporate customer service representatives. Many targets of these attacks – elderly and young internet users alike – would readily provide any and all information to avoid the threat of legal action, penalties, and account shutdowns.

There are two key reasons why traditional attacks have become less effective: advances in detection and the increase in awareness among the average user. Major email providers, for example, alert users when a message is deemed suspicious or the source domain is not as it seems. What's more, users are more computer savvy than ever and know not to trust inbound messages that request personal information. That said, while each major breach prompts reaction from users to update their privacy settings and login credentials, targeted attacks are and will remain relentless.

# APPENDIX

## SOURCE CODE

### app.py

```python
1   from flask import Flask, request, render_template,session, url_for,redirect,flash

2   import pickle

3   import inputScript

4   from passlib.hash import pbkdf2_sha256

5   import json

6   app = Flask(__name__,template_folder='templates')

7   model = pickle.load(open('Website_dt.pkl','rb'))

8

9   @app.route("/")

10  def helloworld():

11      return render_template("/home.html")

12

13  @app.route("/predicturl")

14  def predicturl():

15      return render_template("/predict1.html")

16

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

18  def predict():

19      if request.method == 'POST':

20          url = request.form['url']

21          checkprediction = inputScript.main(url)

22          print(url)

23          print(checkprediction)

24          prediction = model.predict(checkprediction)

25          print(prediction)

26          output=prediction[0]
```

```python
27              if output==1 :
28                  return render_template("/output1.html")
29          elif output==-1 :
30                  return render_template("/output.html")
31
32  @app.route("/project_details")
33  def support():
34      return render_template("/project_details.html")
35
36  @app.route("/addurl")
37  def addurl():
38      return render_template("/addurl.html")
39
40  @app.route("/about")
41  def about():
42      return render_template("/about.html")
43
44
45  if __name__ =="__main__":
46      app.run(debug=True,host='0.0.0.0',port=2000)
47
```

## inputScript.py

```python
1  import regex
2  from tldextract import extract
3  import ssl
4  import socket
5  from bs4 import BeautifulSoup
6  import urllib.request
```

```python
7   import whois
8   import datetime
9
10
11  def url_having_ip(url):
12      symbol =
    regex.findall(r'(http((s)?)://)(((( \d)+).)*)(( \w)+)(/(( \w)+))?',url)
13      if(len(symbol)!=0):
14          having_ip = 1 #phishing
15      else:
16          having_ip = -1 #legitimate
17      return(having_ip)
18      return 0
19
20
21  def url_length(url):
22      length=len(url)
23      if(length<54):
24          return -1
25      elif(54<=length<=75):
26          return 0
27      else:
28          return 1
29
30
31  def url_short(url):
32      #ongoing
33      return 0
```

```python
34
35 def having_at_symbol(url):
36     symbol=regex.findall(r'@',url)
37     if(len(symbol)==0):
38         return -1
39     else:
40         return 1
41 def doubleSlash(url):
42     #ongoing
43     return 0
44
45 def prefix_suffix(url):
46     subDomain, domain, suffix = extract(url)
47     if(domain.count('-')):
48         return 1
49     else:
50         return -1
51
52 def sub_domain(url):
53     subDomain, domain, suffix = extract(url)
54     if(subDomain.count('.')==0):
55         return -1
56     elif(subDomain.count('.')==1):
57         return 0
58     else:
59         return 1
60
61 def SSLfinal_State(url):
```

```python
62     try:
63 #check wheather contains https
64         if(regex.search('^https',url)):
65             usehttps = 1
66         else:
67             usehttps = 0
68 #getting the certificate issuer to later compare with trusted issuer
69         #getting host name
70         subDomain, domain, suffix = extract(url)
71         host_name = domain + "." + suffix
72         context = ssl.create_default_context()
73         sct = context.wrap_socket(socket.socket(), server_hostname = host_name)
74         sct.connect((host_name, 443))
75         certificate = sct.getpeercert()
76         issuer = dict(x[0] for x in certificate['issuer'])
77         certificate_Auth = str(issuer['commonName'])
78         certificate_Auth = certificate_Auth.split()
79         if(certificate_Auth[0] == "Network" or certificate_Auth == "Deutsche"):
80             certificate_Auth = certificate_Auth[0] + " " + certificate_Auth[1]
81         else:
82             certificate_Auth = certificate_Auth[0]
83         trusted_Auth =
   ['Comodo','Symantec','GoDaddy','GlobalSign','DigiCert','StartCom','Entrust','Ve
   rizon','Trustwave','Unizeto','Buypass','QuoVadis','Deutsche Telekom','Network
   Solutions','SwissSign','IdenTrust','Secom','TWCA','GeoTrust','Thawte','Doster',
   'VeriSign']
84 #getting age of certificate
85         startingDate = str(certificate['notBefore'])
86         endingDate = str(certificate['notAfter'])
```

```python
87              startingYear = int(startingDate.split()[3])
88              endingYear = int(endingDate.split()[3])
89              Age_of_certificate = endingYear-startingYear
90  #checking final conditions
91              if((usehttps==1) and (certificate_Auth in trusted_Auth) and
    (Age_of_certificate>=1) ):
92                  return -1 #legitimate
93              elif((usehttps==1) and (certificate_Auth not in trusted_Auth)):
94                  return 0 #suspicious
95              else:
96                  return 1 #phishing
97      except Exception as e:
98          return 1
99
100 def domain_registration(url):
101     try:
102         w = whois.whois(url)
103         updated = w.updated_date
104         exp = w.expiration_date
105         length = (exp[0]-updated[0]).days
106         if(length<=365):
107             return 1
108         else:
109             return -1
110     except:
111         return 0
112
113 def favicon(url):
```

```python
114    #ongoing
115    return 0
116
117 def port(url):
118    #ongoing
119    return 0
120
121 def https_token(url):
122    subDomain, domain, suffix = extract(url)
123    host =subDomain +'.' + domain + '.' + suffix
124    if(host.count('https')): #attacker can trick by putting https in domain
    part
125        return 1
126    else:
127        return -1
128
129 def request_url(url):
130    try:
131        subDomain, domain, suffix = extract(url)
132        websiteDomain = domain
133        opener = urllib.request.urlopen(url).read()
134        soup = BeautifulSoup(opener, 'lxml')
135        imgs = soup.findAll('img', src=True)
136        total = len(imgs)
137        linked_to_same = 0
138        avg =0
139        for image in imgs:
140            subDomain, domain, suffix = extract(image['src'])
```

```python
141              imageDomain = domain
142              if(websiteDomain==imageDomain or imageDomain==''):
143                  linked_to_same = linked_to_same + 1
144          vids = soup.findAll('video', src=True)
145          total = total + len(vids)
146          for video in vids:
147              subDomain, domain, suffix = extract(video['src'])
148              vidDomain = domain
149              if(websiteDomain==vidDomain or vidDomain==''):
150                  linked_to_same = linked_to_same + 1
151          linked_outside = total-linked_to_same
152          if(total!=0):
153              avg = linked_outside/total
154          if(avg<0.22):
155              return -1
156          elif(0.22<=avg<=0.61):
157              return 0
158          else:
159              return 1
160      except:
161          return 0
162
163
164 def url_of_anchor(url):
165      try:
166          subDomain, domain, suffix = extract(url)
167          websiteDomain = domain
168          opener = urllib.request.urlopen(url).read()
```

```python
169            soup = BeautifulSoup(opener, 'lxml')
170            anchors = soup.findAll('a', href=True)
171            total = len(anchors)
172            linked_to_same = 0
173            avg = 0
174            for anchor in anchors:
175                subDomain, domain, suffix = extract(anchor['href'])
176                anchorDomain = domain
177                if(websiteDomain==anchorDomain or anchorDomain==''):
178                    linked_to_same = linked_to_same + 1
179            linked_outside = total-linked_to_same
180            if(total!=0):
181                avg = linked_outside/total
182            if(avg<0.31):
183                return -1
184            elif(0.31<=avg<=0.67):
185                return 0
186            else:
187                return 1
188        except:
189            return 0
190 def Links_in_tags(url):
191        try:
192            opener = urllib.request.urlopen(url).read()
193            soup = BeautifulSoup(opener, 'lxml')
194            no_of_meta =0
195            no_of_link =0
196            no_of_script =0
```

```python
197        anchors=0
198        avg =0
199        for meta in soup.find_all('meta'):
200            no_of_meta = no_of_meta+1
201        for link in soup.find_all('link'):
202            no_of_link = no_of_link +1
203        for script in soup.find_all('script'):
204            no_of_script = no_of_script+1
205        for anchor in soup.find_all('a'):
206            anchors = anchors+1
207        total = no_of_meta + no_of_link + no_of_script+anchors
208        tags = no_of_meta + no_of_link + no_of_script
209        if(total!=0):
210            avg = tags/total
211
212        if(avg<0.25):
213            return -1
214        elif(0.25<=avg<=0.81):
215            return 0
216        else:
217            return 1
218    except:
219        return 0
220
221 def sfh(url):
222    #ongoing
223    return 0
224
```
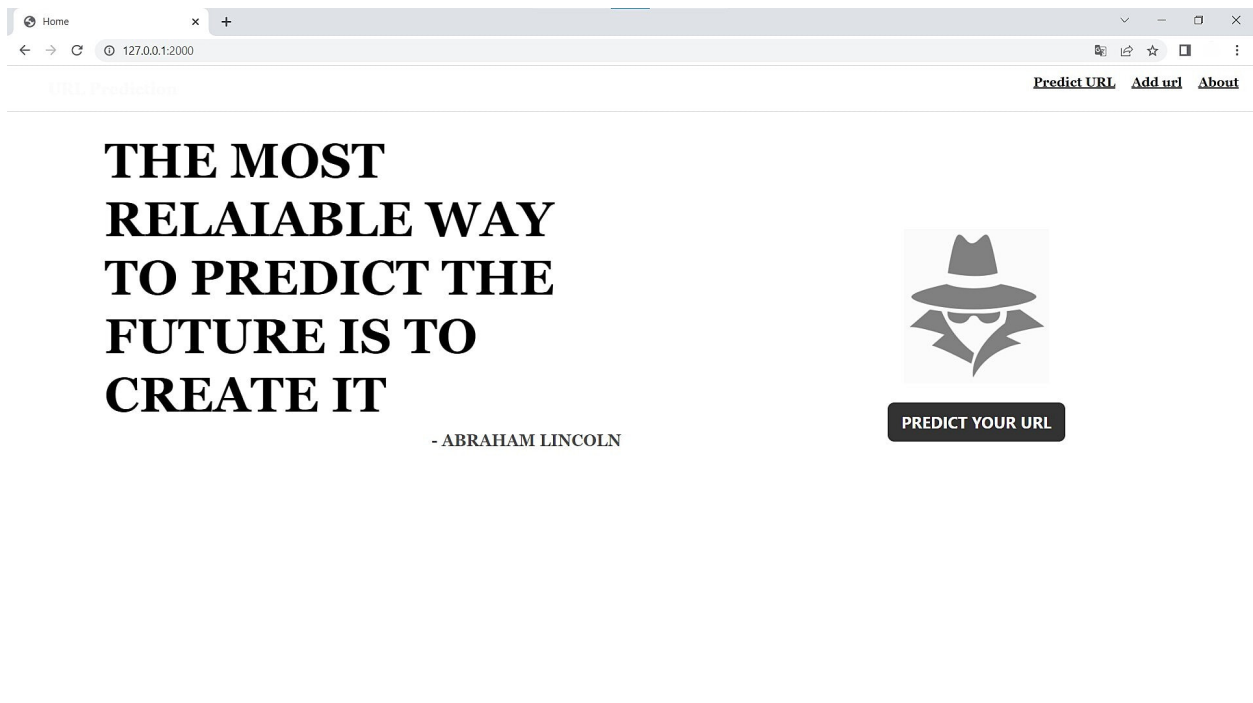
```python
225 def email_submit(url):
226     try:
227         opener = urllib.request.urlopen(url).read()
228         soup = BeautifulSoup(opener, 'lxml')
229         if(soup.find('mailto:')):
230             return 1
231         else:
232             return -1
233     except:
234         return 0
235
236 def abnormal_url(url):
237     #ongoing
238     return 0
239
240 def redirect(url):
241     #ongoing
242     return 0
243
244 def on_mouseover(url):
245     #ongoing
246     return 0
247
248 def rightClick(url):
249     #ongoing
250     return 0
251
252 def popup(url):
```

```python
253        #ongoing
254        return 0
255
256 def iframe(url):
257        #ongoing
258        return 0
259
260 def age_of_domain(url):
261        try:
262            w = whois.whois(url)
263            start_date = w.creation_date
264            current_date = datetime.datetime.now()
265            age =(current_date-start_date[0]).days
266            if(age>=180):
267                return -1
268            else:
269                return 1
270        except Exception as e:
271            print(e)
272            return 0
273 def dns(url):
274        #ongoing
275        return 0
276
277 def web_traffic(url):
278        #ongoing
279        return 0
280
```

```python
281 def page_rank(url):
282     #ongoing
283     return 0
284
285 def google_index(url):
286     #ongoing
287     return 0
288
289
290 def links_pointing(url):
291     #ongoing
292     return 0
293
294 def statistical(url):
295     #ongoing
296     return 0
297
298 def main(url):
299
300
301     check = [[url_having_ip(url),url_length(url),url_short(url),having_at_symbol(url),
302     doubleSlash(url),prefix_suffix(url),sub_domain(url),SSLfinal_State(url),
303     domain_registration(url),favicon(url),port(url),https_token(url),request_url(url),
304     url_of_anchor(url),Links_in_tags(url),sfh(url),email_submit(url),abnormal_url(url),
```

```
305    redirect(url),on_mouseover(url),rightClick(url),popup(url),iframe(url),

306    age_of_domain(url),dns(url),web_traffic(url),page_rank(url),google_index(url),

307              links_pointing(url),statistical(url)]]

308    return check

309
```
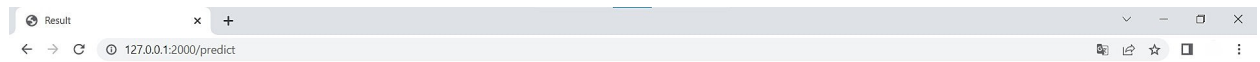
## SCREENSHOTS

127.0.0.1:2000/predicturl

URL Prediction                      Home   Add url   About



**Enter URL to predict**

Ex : https://abcde.com/

**Predict URL**   **Clear**

---
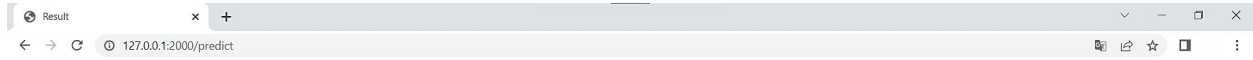
127.0.0.1:2000/predict



**TRUSTED SITE**

**Entered Site or URL is Not a phishing. So don't worry about this site.**

**Predict Another URL**

## GIT HUB & PROJECT DEMO LINK

➤ https://github.com/IBM-EPBL/IBM-Project-1483-1658390493