WEB PHISHING DETECTION

IBM-Project-922-1658330945

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

Project Report

TEAM ID: PNT2022TMID27211

TEAM MEMBERS:

- 1. VIDHYALAKSHMI. R-310819205093
 - 2. YAZHINI. S-310819205097
 - 3. SNEHA. A-310819205083
- 4. SHARMILA. D-310819205079

BACHELOUR OF TECHNOLOGY IN INFORMATION TECHNOLOGY

JEPPIAAR ENGINEERING COLLEGE CHENNAI- 600 100

INDEX

1. INTRODUCTION

- 1. Project Overview
- 2. Purpose

2. LITERATURE SURVEY

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

3. IDEATION & PROPOSED SOLUTION

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

4. REQUIREMENT ANALYSIS

- 1. Functional requirement
- 2. Non-Functional requirements

5. PROJECT DESIGN

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

6. PROJECT PLANNING & SCHEDULING

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

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

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

8. TESTING

- 1. Test Cases
- 2. User Acceptance Testing

9. RESULTS

- 1. Performance Metrics
- 10. ADVANTAGES & DISADVANTAGES
- 11. CONCLUSION
- 12. FUTURE SCOPE
- 13. APPENDIX

Source Code GitHub & Project Demo Link

1. INTRODUCTION

1.1 Project Overview

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

1.2 Purpose

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

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

2. LITERATURE SURVEY

2.1 Existing problem

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

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

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

2.2 References

S.NO	PAPER TITLE	PAPER CONCEPT	ADVANTAGE	DISADVANTAGE	
1	LongfeiWu etal,	In this paper, author	Author propose MobiFish, a	Existing schemes	
	"Effective	did a comprehensive	novel automated lightweight	designed for web	
	Defense Schemes	study on the	anti- phishing scheme for	phishing attacks on PCs	
	for	Security	mobile platforms. MobiFish	c annot effectively	
	Phishing Attacks	vulnerabilities	verifies the validity of web	address the various	
	on Mobile	caused by mobile	pages, applications, and	phishing attacks on	
	Computing	phishing attacks,	persistent accounts by	mobile devices.	
	Platforms, " IEEE	including the web	comparing the actual		
	2016, pp.6678-	page phishing	Identity to the claimed		
	6691.	attacks.	identity		
2	Surbhi Gupta etal., To fool an online user		The paper discusses various	Every organization has	
	" A	into elicit personal	types of Phishing attacks	security issues that have	
	Literature	Information.	such as Tab-napping,	been of great concern to	
	Survey on Social	The prime objective	spoofing emails, Trojan	u sets, sited developers,	
	Engineering	of this review is to do	horse, hacking and how to	and specialists, in order	
	Attacks: Phishing	literature survey	prevent them.	to defend the	
	Attacks," in	on social engineering		confidential data from	
	International	attack:		this type of social	
	Conference on	Phishing attacks and		engineering attack.	
	Computing,	techniques to detect			
	Communication	attack.			
	and				
	Automation(ICCC				
	A2016),201				
	6, pp. 537-540.				

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

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

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

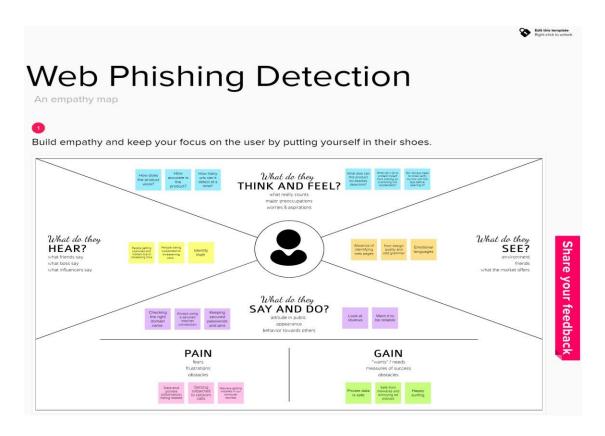
all approach.

2.3 Problem Statement Definition

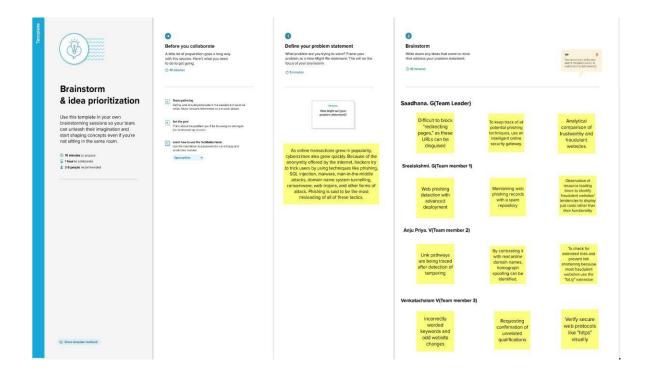
Human users' inability to recognize phishing sites allows phishing attacks to succeed. Past work in anti-phishing can be broadly divided into four categories: studies to understand why people fall for phishing attacks, strategies for teaching people not to fall for phishing attacks, user interfaces for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

3. IDEATION & PROPOSED SOLUTION

3.1 Empathy Map Canvas



3.2 Ideation & Brainstorming



3.3 Proposed Solution

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

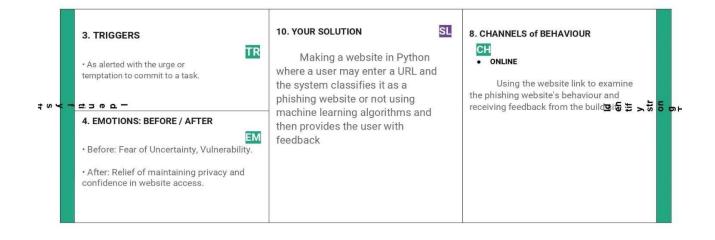
		variety of machine learningapproaches in order to produce a model with the highest level of accuracy and simplicity.
2.	Idea / Solution description	• The method includes the extraction of lexical features from collected webpages as well as host- and page-based feature extraction. The first stage is gathering phishing and legitimate websites. In the host-based technique, attribute extractions based on admiration andlexical bases are carried out to create a database of attribute value. This database contains knowledge that has been extracted using various machine learning methods. A selective classifier is chosen after comparing the methods, and it is put into practice in Python.
		The suggested approach gathered URLs of safe websites from sites likewww.alexa.com, www.dmoz.org, and browsing history. We gathered the phishing URLs from www.phishtak.com. 20000 benign URLs and 17000 phishing URLs make up the data collection.

3.	Novelty / Uniqueness	The dataset provided by UCI Machine Learning repository4 and compiled by Mohammad et al3 was used by the suggested system. The dataset contains6157 legal URLs and 4898 phishing URLsacross 11055 data points. Each data point had 30 features that weresorted into the three categories below: • Features extracted from the URL • Features based on the page's source code, such as URLs that are incorporated into the webpage and HTML and JavaScript-based features. • Features based on domains.
4.	Social Impact / Customer Satisfaction	The majority of the public (users) were assisted by the project in determining if a website was a phishing website or not. It assisted them in classifying the hazardous locations. Machine learning methods were employed in this research. The URL is entered, and it will recognize it and provideusers with precise results.

5.	Business Model (Revenue Model)	In the literature, a number of methods for phishing attack detection and filtering have been suggested. Researchers are still looking for a solution that can protect consumers from phishing attacks and produce better outcomes. It might be easierto spot phishing websites if we can recognize the specific traits and patterns they exhibit. The classification problem of identifying such traits can be resolved using machinelearning approaches.
6.	Scalability of the Solution	This project offers an effective method for phishing detection that pulls features from the URL and HTML source code of websites. In particular, we suggested a hybrid featureset that included features for the HTML source code's plaintext and noisy HTML data, different hyperlink information, and URL character sequence characteristics without the knowledge of experts. The suggested anti-phishing technique has demonstrated competitive performance onactual datasets in terms of several assessment statistics, according to extensive trials. The following criteria have been establishedfor our anti-phishing strategy. • Target independent • Real-time detection • High detection efficiency • Third-party independent

3.4 Problem Solution fit





4. REQUIREMENT ANALYSIS

4.1 Functional requirement

Following are the functional requirements of the proposed solution.

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

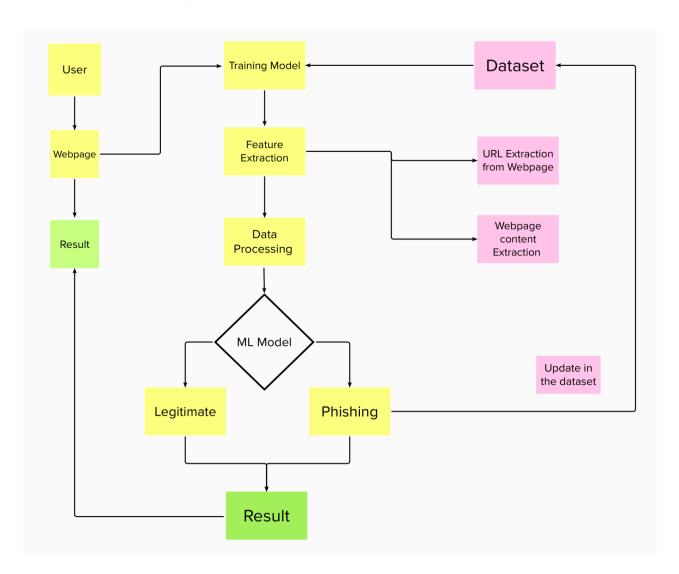
4.2 Non-Functional requirements

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

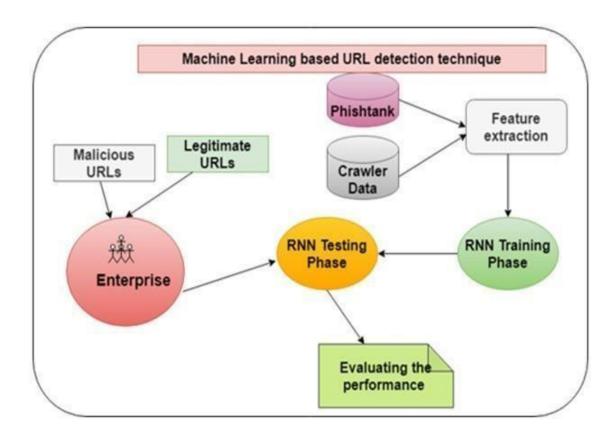
FR NO	Non-Functional Requirements	Description
NFR-1	Usability	System is easy to configure and is efficient in carrying out user tasks.
NFR-2	Availability	System is available to work asrequired when it is required.
NFR-3	Reliability	System will perform the tasks it was designed to do.
NFR-4	Performance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

5. PROJECT DESIGN

5.1 Data Flow Diagrams:



5.2 Solution & Technical Architecture



5.3 User Stories:

User Type	Functional Requireme nt(Epic)	User Story Number	User Story / Task	Acceptance criteria	Priority	Release
Custo mer (Mobi le 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 emailonce I have registered for the application	I can receive confirmation email & click confirm	High	Sprint-1

		USN-3	As a user, I can register for the application through Facebook	I can register & access the dashboard with Facebook Login	Low	Sprint-2
		USN-4	As a user, I can register for the application through Gmail		Medium	Sprint-1
	Login	USN-5	As a user, I can log into the application by entering email & password		High	Sprint-1
Custome r (Web user)	User Input	USN-1	As a user, I can enter the required URL in thebox while awaiting validation.	I can access the website without any problem	High	Sprint-1
Custom er Care Executi ve	Feature Extraction	USN-1	In the event that nothing is discovered during comparison, we can extract features using a heuristic and a visual similarity technique.	As a user I can have compari son between websites for security	High	Sprint-1
Administr ator	Prediction	USN-1	The model will use machine learning algorithms like a logistics regression and KNNto forecast the URLs of the websites.	I can accurately forecast the specific algorithms in this way.	High	Sprint-1
	Classifier	USN-2	To create the final product, I will now feed all of the model output to classifier.	I'll use this to identify the appropriate classifier for generating the outcome.	Medium	Sprint-2

6. PROJECT PLANNING & SCHEDULING

6.1 Sprint Planning & Estimation

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

6.2 Sprint Delivery Schedule

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

```
Velocity:

AV=Velocity/Duration = 35/7 =5

AV=Velocity/Duration = 15/8 =1.875A

V=Velocity/Duration = 25/8 =3.12
```

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

7.1 Feature 1

LOGIN

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

SIGNUP

```
@app.route('/signup/',methods=['POST'])
def signup():
    if request.method=="POST":
        userInfo={
        "fullName":request.form.get('fullName'),
        "email":request.form.get('email'),
        "phoneNumber":request.form.get('phoneNumber'),
        "password":request.form.get('password'),
        userInfo['password']=pbkdf2 sha256.encrypt(userInfo['password'])
        if(account.find one({"email":userInfo['email']})):
            flash("Sorry,user with this email already exist", "signupError")
            return redirect(url_for('index', signupError=True))
        if(account.insert one(userInfo)):
            return start session(userInfo)
    flash("Signup failed", "signupError")
    return redirect(url for('index', signupError=True))
```

ABOUT US

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

7.2 Feature 2

HISTORY PAGE

CONTACT US PAGE

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

FAQ

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

: ×

7.3 Database Schema (if Applicable)

Tables		New table +	×
☐ Name ▼	Schema	Properties	
ACCOUNT	YSX70667		
DETECTIONHISTORY	YSX70667		

Table definition : × ACCOUNT No statistics available. Length Name Data type Nullable Scale **FULLNAME** VARCHAR N 100 0 0 100 0 0 **EMAIL VARCHAR** N **PHONENUMBER** LONG VARCHAR N 32700 0 0 **PASSWORD** N 100 0 0 VARCHAR

Table definition

DETECTIONHISTORY

No statistics available

			NO Stati	stics availabl
Data type	Nullable	Length	Scale	
VARCHAR	N	100	0	0
VARCHAR	N	100	0	0
VARCHAR	N	100	0	0
VARCHAR	N	100	0	0
	VARCHAR VARCHAR	VARCHAR N VARCHAR N VARCHAR N	VARCHAR N 100 VARCHAR N 100 VARCHAR N 100	Data typeNullableLengthScaleVARCHARN1000VARCHARN1000VARCHARN1000

8. TESTING

8.1 Test Cases

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

			nothing is entered in the textbox		2.Enter nothing and click submit 3.An alert is displayed to provide proper input					
Prediction Page_TC_ OO1	Functional	Prediction form page	Verify user is able to see the result when URL is entered in the textbox	Run the flask app in local host	1.Enter URL and click go 2. Enter any URL and click submit 3. The result of the classification is displayed in a new page.	https://go ogle.com/	Result of classification will be displayed with a corresponding emoticon	Working as expected	Pass	

8.2 User Acceptance Testing

1. Defect Analysis:

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

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

2. Test Case Analysis:

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

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

9. RESULTS

9.1 Performance Metrics

S.No.	Parameter	Values	Screenshot
1.	Model Summary	Decision Tree ModelAccuracy – 97%	#2 Descrition Tree from silarm.tree import DecisionTreeClassifier from silarm.tree import DecisionTreeClassifier from the tree import DecisionTreeClassifier from the tree import Decision Tree import Decision train_core = del.core(t.grain y_train) train_core = del.core(t.grain y_train) train_core = test_core (i.e. d.,delellassides2eee) y_pred = dt.predict(t.test) y_pred = dt.predict(t.test) score = delellassides2eeeee score = delellassides2eeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeee
2.	Accuracy	Training Accuracy -Test	d) profess for from silann-tree diport DecisionTreeClassifier dt = DecisionTreeClassifier() dt.dt.(X_t, Y_t, X_t, Y_t, Y_t, X_t, Y_t, Y_t, Y_t, Y_t, Y_t, Y_t, Y_t, Y

Model Performance Comparison:

 $\underline{https://dataplatform.cloud.ibm.com/analytics/notebooks/v2/8188710d-09bc-\underline{4dd6-824d-}$

<u>ec5519a26fea/view?access_token=c66af8d9cede342710725423df04821dc21cc5</u> <u>3af0d4eae2f7496d9e1d3f5a7f</u>

10. ADVANTAGES & DISADVANTAGES:

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

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

11. CONCLUSION

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

12. FUTURE SCOPE:

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

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

13. APPENDIX

Source Code

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

```
. . .
   def login_required(f):
       @wraps(f)
       def wrap(*args, **kwargs):
           if('logged in' in session):
                return f(*args, **kwargs)
           else:
               return redirect('/')
       return wrap
11 def start session(userInfo):
       del userInfo['password']
       session['logged_in']=True
       session['user']=userInfo
       session['predicted']=False
       return redirect(url_for('index'))
19 @app.route('/login/',methods=['POST'])
20 def login():
       if request.method=="POST":
           email=request.form.get("email")
           password=request.form.get("password")
           verify_account = "SELECT * FROM account WHERE email =?"
           stmt = ibm_db.prepare(conn, verify_account)
           ibm db.bind param(stmt,1,email)
           ibm_db.execute(stmt)
           fetch_account = ibm_db.fetch_assoc(stmt)
           if(fetch_account):
                if(pbkdf2_sha256.verify(password,fetch_account['PASSWORD'])):
                   userInfo={
                        "fullName":fetch account['FULLNAME'],
                        "email":fetch account['EMAIL'],
                        "phoneNumber":fetch_account['PHONENUMBER'],
                        "password":fetch_account['PASSWORD'],
                   return start_session(userInfo)
               else:
                   flash("Password is incorrect", "loginError")
                    return redirect(url_for('index',loginError=True))
           flash("Sorry, user with this email id does not exist", "loginError")
           return redirect(url_for('index',loginError=True))
```

```
• • •
   @app.route('/signup/',methods=['POST'])
   def signup():
       if request.method=="POST":
            userInfo={
            "fullName":request.form.get('fullName'),
            "email":request.form.get('email'),
            "phoneNumber":request.form.get('phoneNumber'),
            "password":request.form.get('password'),
            userInfo['password']=pbkdf2_sha256.encrypt(userInfo['
   password'])
            sql = "SELECT * FROM account WHERE email =?"
            stmt = ibm_db.prepare(conn, sql)
           ibm_db.bind_param(stmt,1,userInfo['email'])
           ibm db.execute(stmt)
           account = ibm_db.fetch_assoc(stmt)
           if account:
                flash("Sorry, user with this email already exist","
   signupError")
               return redirect(url_for('index', signupError=True))
           else:
                insert sql = "
   INSERT INTO account(fullName, email, phoneNumber, password) VALUES
   (?, ?, ?, ?)
               prep_stmt = ibm_db.prepare(conn, insert_sql)
               ibm_db.bind_param(prep_stmt, 1, userInfo['fullName'])
               ibm_db.bind_param(prep_stmt, 2, userInfo['email'])
               ibm_db.bind_param(prep_stmt, 3, userInfo['phoneNumber'
   ])
                ibm_db.bind_param(prep_stmt, 4, userInfo['password'])
               ibm_db.execute(prep_stmt)
                return start_session(userInfo)
        flash("Signup failed", "signupError")
        return redirect(url_for('index', signupError=True))
32 @app.route('/logout/',methods=["GET"])
33 def logout():
       if request.method=="GET":
            session.clear()
       return redirect(url for('index'))
```

```
@app.route('/')
       if(session and '_flashes' in dict(session)):
           loginError=request.args.get('loginError')
           signupError=request.args.get('signupError')
           if(loginError):
               return render_template('./index.html',loginError=loginError,
   carousel_content=carouselData['carousel_content'],currentYear=datetime.date.today().
   year)
           if(signupError):
               return render_template('./index.html',signupError=signupError,
   carousel_content=carouselData['carousel_content'], currentYear=datetime.date.today().
       if(session and '_flashes' not in dict(session)):
           if(session['logged_in']==True):
               return render_template('./index.html',userInfo=session['user'],
   carousel_content=carouselData['carousel_content'], currentYear=datetime.date.today().
               return render_template('./index.html',carousel_content=carouselData['
   carousel_content'],currentYear=datetime.date.today().year)
           return render_template('./index.html',carousel_content=carouselData['
   carousel_content'],currentYear=datetime.date.today().year)
   @app.route('/detect/', methods=['GET','POST'])
   @login_required
       if request.method == 'POST':
           title=request.form['title']
           url = request.form['url']
           checkprediction = inputScript.main(url)
           output=prediction[0]
           session['predicted']=True
           if(output==1):
               pred = "Wohoo! You are good to go."
session['status']='safe'
               session['pred'] = pred
               pred = "Oh no! This is a Malicious URL"
               session['status']='unsafe'
               session['pred'] = pred
           session['title']=title
           session['url']=url
           insert_detection_info_stmt="
   INSERT INTO DETECTIONHISTORY(email,title,url,status) VALUES(?,?,?,?)"
           insert_detection_info = ibm_db.prepare(conn, insert_detection_info_stmt)
           ibm_db.bind_param(insert_detection_info,1,session['user']['email'])
           ibm_db.bind_param(insert_detection_info,2,session['title'])
           ibm_db.bind_param(insert_detection_info,3,session['url'])
           ibm_db.bind_param(insert_detection_info,4,session['status'])
           ibm_db.execute(insert_detection_info)
           if(session and session['logged_in']):
               if(session['logged_in']==True):
                   return redirect(url_for('predictionResult'))
           return render template('./templates/predict-form.html',userInfo=session['user'
```

```
@app.route('/detection-result/')
   @login_required
   def predictionResult():
       if(session['predicted']==True):
           urlInfo={
           'message' :session['pred'] ,
           'title':session['title'],
            'url':session['url'],
           'status':session['status']
           return render_template("./templates/prediction-result.html", urlInfo
   =urlInfo,userInfo=session['user'])
       else:
           return redirect(url_for('predict'))
16 @app.route('/detection-history/')
   @login_required
18 def detectionHistory():
       if(session and session['logged_in']):
           if(session['logged_in']==True):
               get_detection_history_stmt = "
   SELECT title,url,status FROM detectionHistory where email=?"
               get_detection_history = ibm_db.prepare(conn,
    get_detection_history_stmt)
               ibm_db.bind_param(get_detection_history,1,session['user']['email
    '])
               ibm_db.execute(get_detection_history)
               fetch_detection_history = ibm_db.fetch_assoc(
   get_detection_history)
               detection_history = []
               ind = 0
               while fetch_detection_history != False:
                   detection_history.append(fetch_detection_history)
                   fetch_detection_history = ibm_db.fetch_assoc(
   get_detection_history)
               detection_history= detection_history[::-1]
               return render_template('./templates/detection-history.html',
   userInfo=session['user'],detectionHistory=detection_history)
   @app.route('/about/')
   def about():
       if(session and session['logged in']):
           if(session['logged_in']==True):
               return render_template('./templates/about.html',userInfo=session
   ['user'], aboutContents=aboutData['aboutContents'])
           else:
               return render template('./templates/about.html',aboutContents=
   aboutData['aboutContents'])
       else:
           return render_template('./templates/about.html',aboutContents=
   aboutData['aboutContents'])
```

GitHub and Demo link:

GITHUB: IBM-Project-922-1658330945

DEMO LINK:

https://youtube.com/watch?v=imX9n9UTXtA&feature=share