

**MINI PROJECT REPORT**  
**ON**  
**“DETECTION OF PHISHING WEBSITE”**  
**SUBMITTED IN PARTIAL FULLFILLMENT OF THE REQUIREMENTS OF**  
**DEGREE OF**  
**BACHELOR OF ENGINEERING**  
**BY**

**ROHAN SHINDE                      TE B 52**

**PRANJAL SONAWANE      TE B 58**

**PRANAY TATE                      TE B 61**

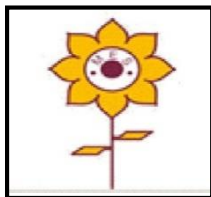
**GAURAV THAKUR              TE B 63**

**SUPERVISOR**

**Ms. SNEHAL SHINDE**



**DEPARTMENT OF COMPUTER ENGINEERING**  
**PILLAI HOC COLLEGE OF ENGINEERING AND TECHNOLOGY,**  
**PILLAI'S HOC EDUCATIONAL CAMPUS, HOCL COLONY,**  
**RASAYANI, TAL: KHALAPUR, DIST RAIGAD 410207**  
**UNIVERSITY OF MUMBAI**  
**[2022-23]**



**Mahatma Education Society's**  
**Pillai HOC College of Engineering and Technology,**  
**Rasayani-410207**  
**2022-23**

## **Certificate**

This is to certify that the Mini Project-2A entitled “**Detection Of Phishing Website**” is a bonafide work of **Rohan Shinde, Pranjal Sonawane, Pranay Tate & Gaurav Thakur** submitted to the University of Mumbai in partial fulfilment of the requirement for the award of the degree of “**Undergraduate**” in “**Computer Engineering**”.

---

**Ms. Snehal Shinde**  
(Supervisor)

---

**Ms. Snehal Chitale**  
(Project Coordinator)

---

**Ms. Rohini Bhosale**  
(Head of Department)

---

**Dr. J. W. Bakal**  
(Principal)

## **Mini Project-2A Report Approval**

This project report entitled “**Detection of Phishing Website**” submitted by “**Rohan Shinde, Pranjal Sonawane, Pranay Tate & Gaurav Thakur**” is approved for the degree of **Bachelor of Engineering in Computer Engineering**.

### **Examiners**

1. \_\_\_\_\_

2. \_\_\_\_\_

**Date**

**Place:**

# Declaration

We declare that this written submission represents our ideas in our own words and where others ideas or words have been included. We have adequately cited and referenced the original sources. We also declare that we have adhered to all principles of academic honesty and integrity and have not misrepresented or fabricated or falsified any idea/data/fact/source in our submission. We understand that any violation of the above will because for disciplinary action by the Institute and can also evoke penal action from the sources which have thus not been properly cited or from whom proper permission has not been taken when needed.

---

**Rohan Shinde**

---

**Pranjal Sonawane**

---

**Pranay Tate**

---

**Gaurav Thakur**

**Date:**

# Abstract

Phishing attacks are the simplest way to obtain sensitive information from innocent users. Aim of the phishers is to acquire critical information like username, password and bank account details. Cyber security persons are now looking for trustworthy and steady detection techniques for phishing websites detection. machine learning technology is used for detection of phishing URLs by extracting and analyzing various features of legitimate and phishing URLs. Decision Tree, K-Means Clustering, Random forest and Naïve Bayes algorithms are used to detect phishing websites. Aim is to detect phishing URLs as well as narrow down to the best machine learning algorithm by comparing accuracy rate, false positive and false negative rate of each algorithm.

**Keywords:** Machine learning, Phishing, Cyber Awareness, Phishing detection.

# **List of Abbreviations**

ML: Machine Learning

DB: Database

URL: Uniform Resource Locator

HTTP: Hypertext Markup Transfer Protocol

HTTPS: Hypertext Markup Transfer Protocol Secure

DNS: Domain Name Server

## List of Figures

<b>Figure no</b>	<b>Figure Name</b>	<b>Page No</b>
Fig 1.1	Existing System	6
Fig 1.2	System Architecture	11
Fig 1.3	Features Extraction	12
Fig 1.4	Home Page	15
Fig 1.5	Search Results for Phishing Website	15
Fig 1.6	Search Result for Non-Phishing Website	16
Fig 1.7	Accuracy Result	16

# TABLE OF CONTENTS

<b>Content</b>	<b>Page No.</b>
<b>Abstract</b>	<b>i</b>
<b>List of Abbrevations</b>	<b>ii</b>
<b>List of Figures</b>	<b>iii</b>
<b>1. Introduction</b>	<b>1</b>
1.1 Background	2
1.2 Motivation	3
<b>2. Literature Survey</b>	<b>4</b>
2.1 Basic Terminologies	5
2.2 Existing System	6
2.2 Problem Statement	7
<b>3. Requirement Gathering</b>	<b>8</b>
3.1 Software and Hardware Requirements	9
<b>4. Plan of Project</b>	<b>10</b>
4.1 Implemented system architecture	11
4.2 Methodology	12
<b>5. Result analysis</b>	<b>14</b>
5.1 Result and discussion	15
<b>6. Conclusion</b>	<b>17</b>
<b>References</b>	<b>19</b>



# **Chapter 1**

## **Introduction**

## 1.1 Background

Nowadays, Phishing has become a main area of concern for security researchers because it is not difficult to create fake websites which look so close to legitimate websites. Experts can identify fake websites but not all the users can identify the fake website and such users become the victim of phishing attacks. Main aim of the attacker is to steal bank account credentials. In United States businesses, there is a loss of US\$2billion per year because their client become victims of phishing. In Microsoft Computing Safer Index Report released in February 2014, it was estimated that the annual worldwide impact of phishing could be as high as \$5 billion. Phishing attacks are becoming successful because of lack of user awareness. Since phishing attacks exploit the weaknesses found in users, it is very difficult to mitigate them but it is very important to enhance phishing detection techniques. The general method to detect phishing websites by updating blacklisted URLs, Internet Protocol (IP) to the antivirus database which is also known as the "blacklist" method. To evade blacklists attackers use creative techniques to fool users by modifying the URL to appear legitimate via obfuscation and many other simple techniques including: fast-flux, in which proxies are automatically generated to host the web-page; algorithmic generation of new URLs; etc. Major drawback of this method is that it cannot detect zero-hour phishing attacks.

## 1.2 Motivation

Due to the rise of phishing attacks and phishing websites it is quite difficult to browse the web without any worries. This phishing attacks can be fatal as they try to steal sensitive information. Also it not so obvious to catch such websites as their UI/UX is exactly similar to original websites and manually detecting such websites can be time consuming. Hence, to avoid such phishing attacks we have developed phishing website detection system to help people detect such websites automatically.

# **Chapter 2**

## **Literature Survey**

## 2.1 Basic Terminologies

### **Phishing**

Phishing is a type of cybersecurity attack during which malicious actors send messages pretending to be a trusted person or entity. Phishing messages manipulate a user, causing them to perform actions like installing a malicious file, clicking a malicious link, or divulging sensitive information such as access credentials. Phishing is the most common type of social engineering, which is a general term describing attempts to manipulate or trick computer users. Social engineering is an increasingly common threat vector used in almost all security incidents. Social engineering attacks, like phishing, are often combined with other threats, such as malware, code injection, and network attacks.

### **Machine Learning**

Machine learning (ML) is a field of inquiry devoted to understanding and building methods that 'learn', that is, methods that leverage data to improve performance on some set of tasks. It is seen as a part of artificial intelligence. Machine learning algorithms build a model based on sample data, known as training data, in order to make predictions or decisions without being explicitly programmed to do so. Machine learning algorithms are used in a wide variety of applications, such as in medicine, email filtering, speech recognition, and computer vision, where it is difficult or unfeasible to develop conventional algorithms to perform the needed tasks.

## 2.2 Existing System

In earlier systems users need to detect and block phishing websites manually, which is quite difficult nowadays because it very hard to tell the difference in legitimate website and the phishing website as they look exactly similar and it's not so obvious to tell the difference between them based on their UI.

In recent times most of the phishing attacks are done via emails as it is one of the most widely used technology and have become integral part of our life. Attackers sends various links to thousands of users everyday which leads users to phishing websites. Various software and spam filters can be used to avoid such type of attacks but it can be difficult to setup such software and they can also be costly sometimes. Also installation of anti phishing and antivirus software can take excess of storage in system which will ultimately slow down the system.

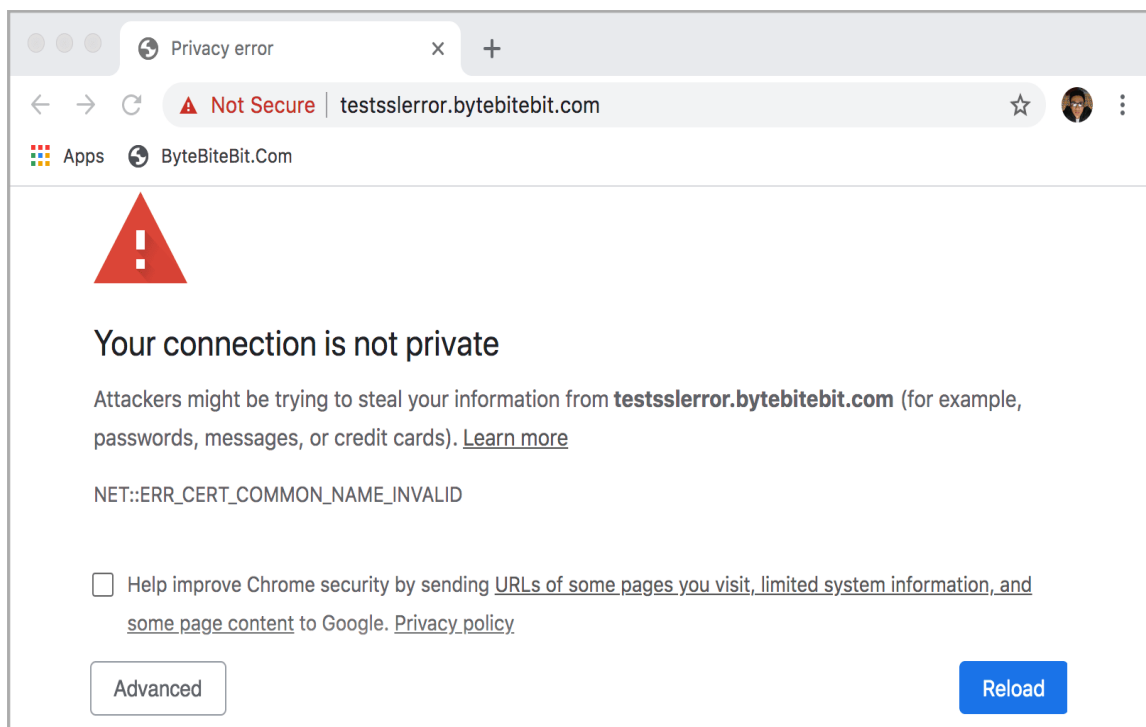


Fig : 1.1 Existing System

## 2.3 Problem Statement

Phishing attacks are the simplest way to obtain them. Phishing is one of the techniques which are used by the intruders to get access to the user credentials or to gain access to the sensitive data. This type of accessing is done by creating the replica of the websites which looks the same as the original websites which we use on our daily basis but when a user clicks on the link he will see the website and think its original and try to provide his credentials. To prevent such attack creating a system for detection of phishing websites.

# **Chapter 3**

## **Requirement Gathering**



## **3.1 Details of Hardware & Software**

### **1 Software Requirement**

This project is built using Visual Studio Code and Google Collab.

- Language: Python
- Framework: Flask
- Library: sklearn , pandas

### **2 Hardware Requirement**

We strongly recommend a computer fewer than 5 years old.

- Processor: Minimum 1 GHz; Recommended 2GHz or more
- Ethernet connection (LAN) OR a wireless adapter (Wi-Fi)
- Hard Drive: Minimum 32 GB; Recommended 64 GB or more
- Memory (RAM): Minimum 1 GB; Recommended 4 GB or above

# **Chapter 4**

## **Plan of Project**

## 4.1 Implemented System Architecture

To detect phishing website, we will collect the URL from user with help of web app and after the input of URL we will gather the required data of that website. After gathering of the data, we will proceed for the feature extraction step. In this step all 30 features of the websites are verified, Features which are packed in our program. Features like having IP Address, URL Lengthening At symbol, double slash redirecting, prefix suffix, Request URL, pop up window are few examples of those 30 features are verified, In next step various classification algorithms are used to train our module, algorithms such as Random forest, Naive Bayes etc. Among those algorithms which gives best accuracy is selected and used to train the module. After this process we will get an output of a Boolean value which will be 1 or 0. If output is 1 then website is phishing, if output is 0 then it is safe to visit the site without any worries. Based on this Boolean values , Output will be displayed to make user aware.

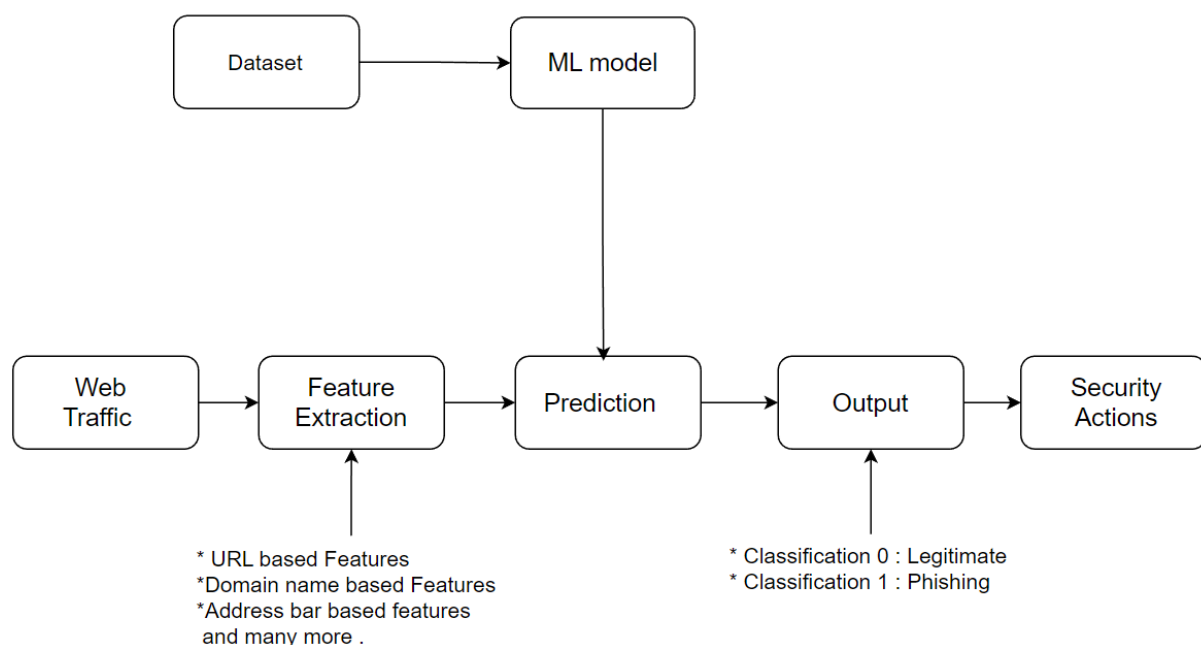


fig 1.2 System Architecture

## 4.2 Methodology.

Approach for building ML model for prediction system of phishing website start with feature extraction. In which 30 different features are consider for classification of phishing website these features play an important role in training the ML model.



Fig 1.3 Features Extraction

## **Algorithm used to train model**

- **K-Means Clustering Algorithm**

K-Means Clustering is an iterative algorithm that divides the unlabelled dataset into k different clusters in such a way that each dataset belongs only one group that has similar properties.

- **Naïve Bayes Classifier Algorithm**

Naïve Bayes algorithm is a supervised learning algorithm, which is based on Bayes theorem and used for solving classification problems. It is mainly used in text classification that includes a high-dimensional training dataset. Naïve Bayes Classifier is one of the simple and most effective Classification algorithms which helps in building the fast machine learning models that can make quick predictions. It is a probabilistic classifier, which means it predicts on the basis of the probability of an object.

- **Decision Tree Classification Algorithm**

Decision Tree is a Supervised learning technique that can be used for both classification and Regression problems, but mostly it is preferred for solving Classification problems. It is a tree-structured classifier, where internal nodes represent the features of a dataset, branches represent the decision rules and each leaf node represents the outcome. In a Decision tree, there are two nodes, which are the Decision Node and Leaf Node. Decision nodes are used to make any decision and have multiple branches, whereas Leaf nodes are the output of those decisions and do not contain any further branches.

# **Chapter 5**

## **Result Analysis**

## 5.1 Results and Discussion

Below is the home page of the phishing detection website. As given below user needs to enter the URL of the website in order to check whether that website is safe to use or not. After entering the URL user can simply click on check button.

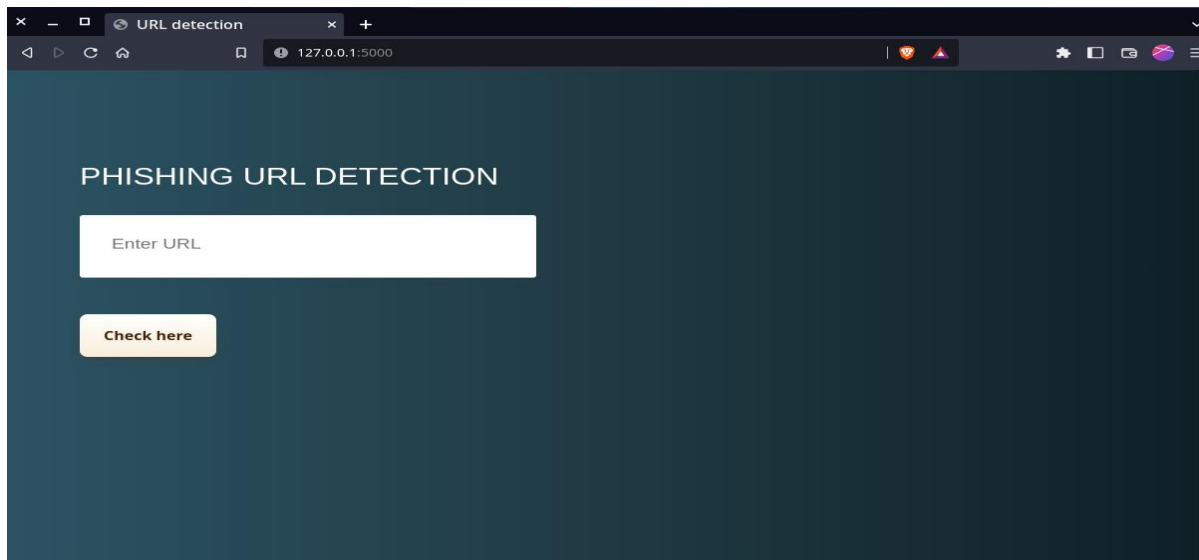


Fig 1.4 Home Page

After clicking the check button, the ML model will classify whether the website is phishing or not based on the 30 features which were used to train the model.

As we can see below the website is classified as unsafe to use as it is using http and not https which makes the website vulnerable.

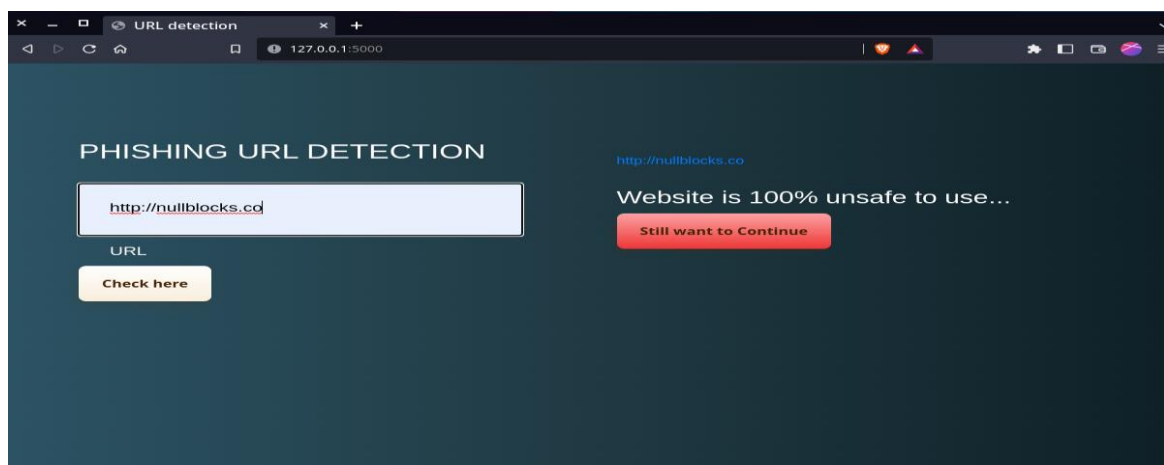


Fig 1.5 Search Results for Phishing Website

In this case the website is safe to use. Hence user can continue to use such websites. This result is also given on the basis of 30 features which were used to train the ML model.

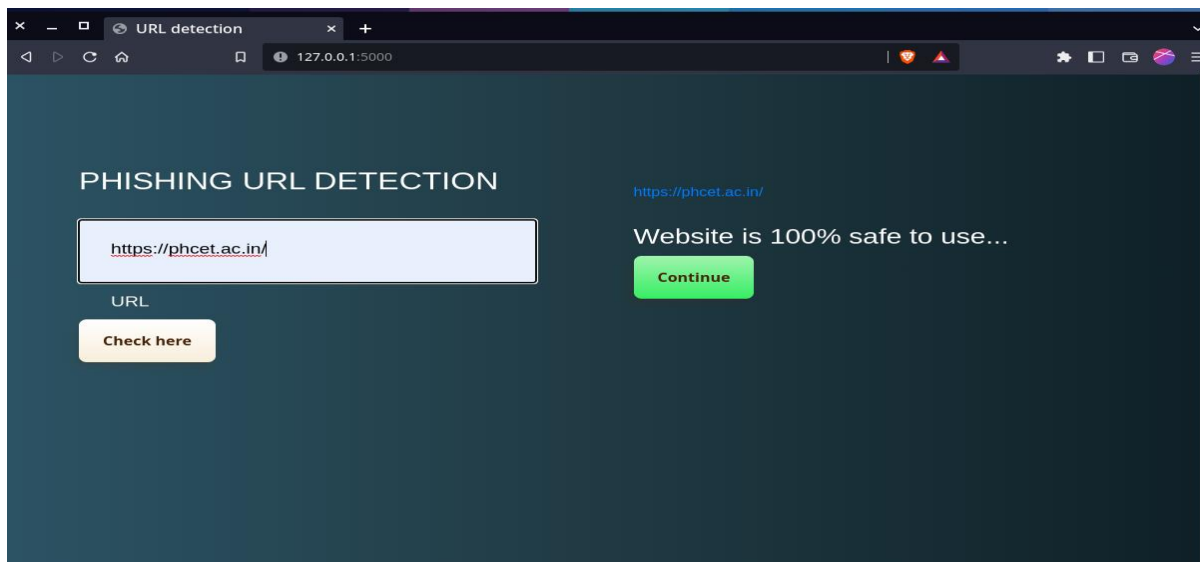


Fig 1.6 Search Result for Non Phishing Website

Below is the accuracy result of the different ML models. The model which is giving the highest accuracy is selected to detect the phishing websites.

```
#Sorting the dataframe on accuracy
sorted_result=result.sort_values(by=['Accuracy', ],ascending=False).reset_index(drop=True)
sorted_result
```

Model	Accuracy
Random Forest	0.967
Decision Tree	0.961
K-Nearest Neighbors	0.956
Naive Bayes Classifier	0.605

Fig 1.7 Accuracy Result



## **Chapter 6**

## **Conclusion**

## **Conclusion**

In this Project we explore the possibility of using machine learning to detect phishing URLs without performing additional scraping of the URL itself. Web app is used for detecting & blocking the phishing attacks and also malicious links on web pages & keeping users safe on the internet .

## References

- [1] Ankit Kumar Jain and B. B. Gupta . “Analysis of Visual Similarity Based Approaches” *IEEE Security and Communication Networks*, vol.2017, Issue : 10 Jan 2017 .
- [2] Wei KingTiong . “Utilisation of website logo for phishing detection” *IEEE Computers & Security*, vol 54 , Issue : October 2015 .
- [3] Weili Han. “Anti-phishing based on automated individual white-list” *IEEE Digital identity management*, vol 12, Issue : October 2008
- [4] Arathi Krishna , “Phishing Detection using Machine Learning based URL Analysis” *International journal of engineering research & technology (ijert)* , vol9 , Issue 02-08-2021