

# URL-Based Phishing Detection Using Machine Learning

## INTRODUCTION

### 1.1 Overview

There are a number of users who purchase products online and make payments through e-banking. There are e-banking websites that ask users to provide sensitive data such as username, password & credit card details, etc., often for malicious reasons. This type of e-banking website is known as a phishing website. Web service is one of the key communications software services for the Internet. Web phishing is one of many security threats to web services on the Internet.

Web phishing aims to steal private information, such as usernames, passwords, and credit card details, by way of impersonating a legitimate entity. It will lead to information disclosure and property damage. Large organizations may get trapped in different kinds of scams.

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

In order to detect and predict e-banking phishing websites, we proposed an intelligent, flexible and effective system that is based on using classification algorithms. We implemented classification algorithms and techniques to extract the phishing datasets criteria to classify their legitimacy. The e-banking phishing website can be detected based on some important characteristics like URL and domain identity, and security and encryption criteria in the final phishing detection rate. Once a user makes a transaction online when he makes payment through an e-banking website our system will use a data mining algorithm to detect whether the e-banking website is a phishing website or not.

### 1.2 Purpose

The main purpose of the project is to detect the fake or phishing websites who are trying to get access to the sensitive data or by creating the fake websites and trying to get access of the user personal credentials. We are using machine learning algorithms to safeguard the sensitive data and to detect the phishing websites who are trying to gain access on sensitive data.

It can be used to preserve the confidentiality. To protect the user from phishing websites. To develop a user-friendly environment. To prevent or mitigate harm or destruction of computer networks, applications, devices, and data.

## LITERATURE SURVEY

### 2.1 Existing problem

The existing system handles the only one kind of phishing attacks. If that was a phishing site then the existing system only warns the user. The active warning gives the user options to close the window or displaying the website. The passive warning displays the popup dialog box.

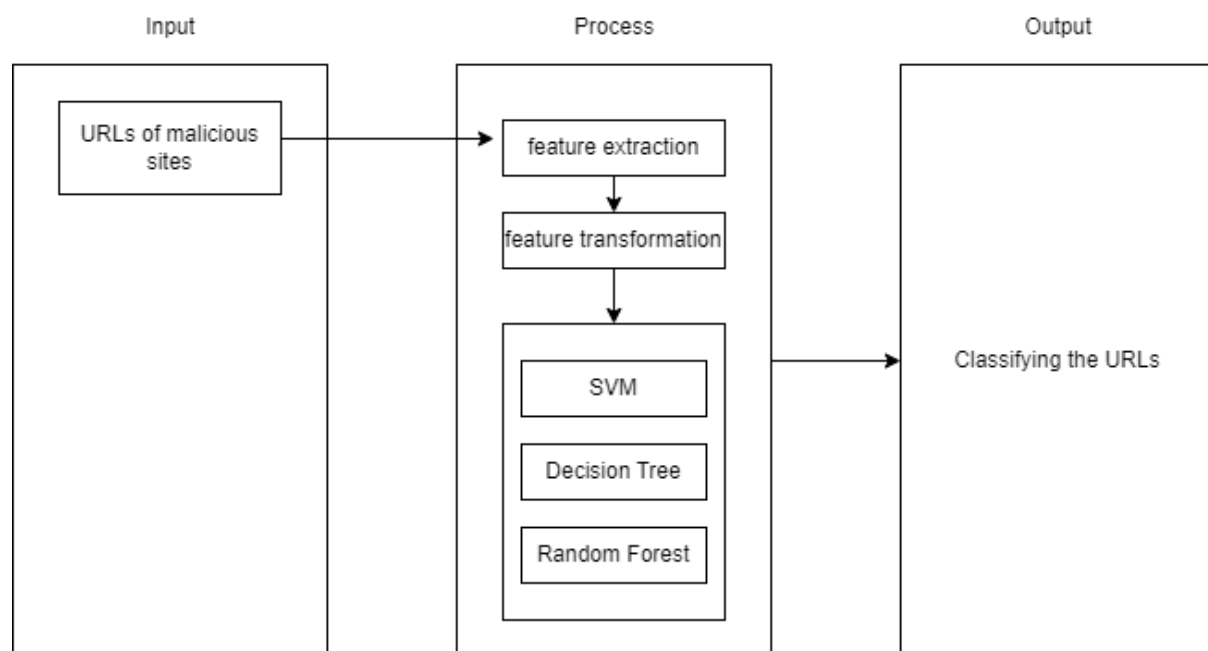
### 2.2 Proposed solution

The dataset of phishing and legitimate URL's is given to the system which is then pre-processed so that the data is in the useable format for analysis. The features have around 30 characteristics of phishing websites which is used to differentiate it from legitimate ones. Each category has its own characteristics of phishing attributes and values are defined. The specified characteristics are extracted for each URL and valid ranges of inputs are identified. These values are then assigned to each phishing website risk. For each input the values range from 0 to 10 , while for output range is from 0 to 100. The phishing attributes values are represented with binary no 0 and 1 which indicates the attribute is present or not.

After this the data is trained we shall apply a relevant machine learning algorithm to the dataset. The machine learning algorithms are already explained in previous section. After this we use a hybrid classification in which we combine two of the classifier namely Naive Bayes and Random forest to predict the accuracy of the detection of the phishing URL, hence we get our desired result. This is also called a hybrid approach to test the data, in this method we propose to use the combination of two classifiers, as mentioned above. We shall then test the data and evaluate the prediction accuracy which shall be more than the existing system. We shall now see the different classifiers and discuss the hybrid combination used for our proposed system.

## THEORITICAL ANALYSIS

### 3.1 Block diagram



### 3.2 Hardware / Software designing

#### Hardware Requirements:

The following is the hardware requirements of the system for the proposed system:

- Processor : Any Processor above 500 MHz
- RAM :8 GB
- Hard Disk :1 TB
- Input device : Standard keyboard and mouse

#### Software Requirements:

The following is the software requirements of the system for the proposed system:

- OS : Windows 10
- Platform : Jupyter Notebook
- Language : Python
- IDE/tool : Anaconda 3-5.0.3

### EXPERIMENTAL INVESTIGATIONS

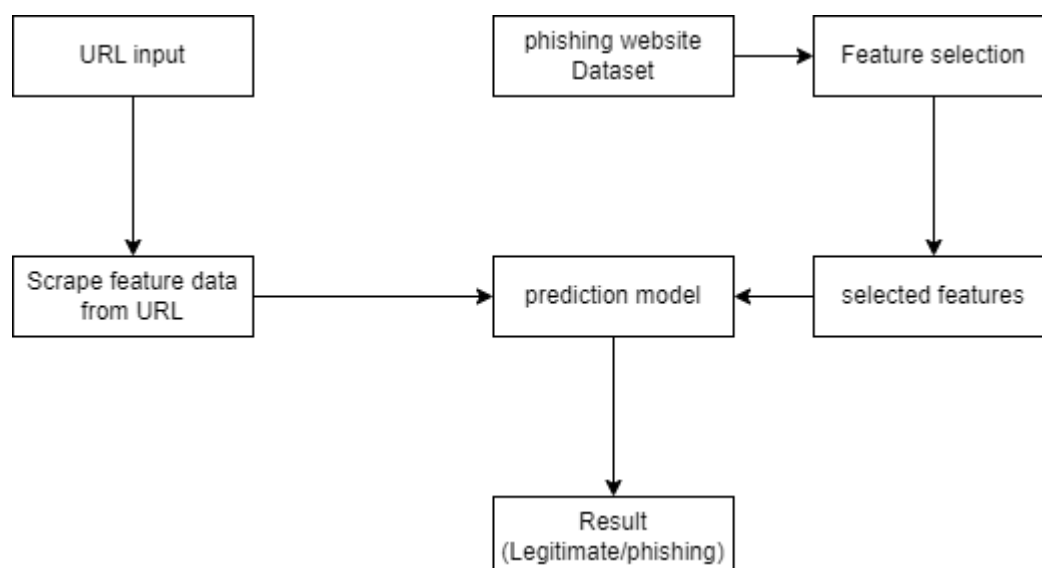
Scikit-learn tool has been used to import Machine Learning algorithms. Dataset is divided into training set and testing set in 70:30 ratio. Each classifier is trained using training data and testing data which is used to analyse the performance of classifier. Performance of classifier has been evaluated by calculating classifier's accuracy score.

Decision Tree : 0.9461782

Random Forest : 0.9104480

SVM : 0.541384

#### FLOWCHART



## RESULT

Result shows that Decision Tree algorithm gives better detection accuracy which is 94.61. it also shows that detection accuracy of phishing websites increases as more dataset used as training dataset.

```
In [32]: models = pd.DataFrame({
        'Model': ['SVM-rbf', 'Decision Tree', 'Random Forest'],
        'Test Score': [svm_rbf, dec_tree, rf]})
models.sort_values(by='Test Score', ascending=False)

Out[32]:
```

|   | Model         | Test Score |
|---|---------------|------------|
| 1 | Decision Tree | 0.946178   |
| 2 | Random Forest | 0.910448   |
| 0 | SVM-rbf       | 0.541384   |

```
In [23]: dtc.predict([[5,1,0,-1,1,1,-1,1,1,-1,1,1,1,1,0,0,-1,1,1,0,-1,1,-1,1,-1,-1,0,-1,1,1,1]])

C:\Users\Hp\anaconda3\lib\site-packages\sklearn\base.py:450: UserWarning: X does not have valid feature names, but DecisionTreeClassifier was fitted with feature names
  warnings.warn(

Out[23]: array([1], dtype=int64)
```

## ADVANTAGES & DISADVANTAGES

### Advantages

- This system can be used by many E-commerce or other websites in order to have good customer relationship.
- User can make online payment securely.
- Data mining algorithm used in this system 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.

### Disadvantages

- If Internet connection fails, this system won't work.
- All websites related data will be stored in one place.

## APPLICATIONS

Phishing is one of the most severe cyber-attacks where researchers are interested to find a solution. In phishing, attackers lure end-users and steal their personal in-formation. To minimize the damage caused by phishing must be detected as early as possible. There are various phishing attacks like spear phishing, whaling, vishing, smishing, pharming and so on. There are various phishing detection techniques based on white-list, black-list, content-based, URL-based, visual-similarity and machine-learning.

The importance to safeguard online users from becoming victims of online fraud, divulging confidential information to an attacker among other effective uses of phishing as an attacker's tool, phishing detection tools play a vital role in ensuring a secure online experience for users.

## CONCLUSION

The demonstration of phishing is turning into an advanced danger to this quickly developing universe of innovation. Today, every nation is focusing on cashless exchanges, business online, tickets that are paperless and so on to update with the growing world. Yet phishing is turning into an impediment to this advancement. Individuals are not feeling web is dependable now. It is conceivable to utilize AI to get information and assemble extraordinary information items. A lay person, completely unconscious of how to recognize a security danger shall never invite the danger of making money related exchanges on the web. Phishers are focusing on installment industry and cloud benefits the most. The project means to investigate this region by indicating an utilization instance of recognizing phishing sites utilizing ML. It aimed to build a phishing detection mechanism using machine learning tools and techniques which is efficient, accurate and cost effective. The project was carried out in Anaconda IDE and was written in Python. The proposed method used four machine learning classifiers to achieve this and a comparative study of the four algorithms was made. A good accuracy score was also achieved. All the three classifiers gave promising results with the best being Decision Tree Classifier with an accuracy score of 94.16%. The accuracy score might vary while using other datasets and other algorithms might provide better accuracy than decision tree classifier. This model can be deployed in real time to detect the URLs as phishing or legitimate

## FUTURE SCOPE

The proposed idea can be improvised by detecting the clickable pictures, malicious QR code, etc. The limitation is that all features are discrete. The other limitation is that the URL is to be copied and we have to search in the application then it will predict whether it is legitimate or not rather than redirecting the URL link to the application. If the URL is not there in the training and testing data set then it is difficult to predict that the URL is legitimate or not.

Further work can be done to enhance the model by using ensembling models to get greater accuracy score. Ensemble methods is a ML technique that combines many base models to generate an optimal predictive model. Further reaching future work would be combining multiple classifiers, trained on different aspects of the same training set, into a single classifier that may provide a more robust prediction than any of the single classifiers on their own. The project can also include other variants of phishing like smishing, vishing, etc. to complete the system. Looking even further out, the methodology needs to be evaluated on how it might handle collection growth. The collections will ideally grow incrementally over time so there will need to be a way to apply a classifier incrementally to the new data, but also potentially have this classifier receive feedback that might modify it over time.

## BIBLIOGRAPHY

- [1] Github. Github. <https://github.com/>
- [2] Rachna Dhamija, J. D. Tygar, and Marti Hearst. Why phishing works. In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, CHI '06, pages 581–590, New York, NY, USA, 2006. ACM.
- [3] npm. url-parse. <https://www.npmjs.com/package/url-parse>.
- [4] University of Dayton. Phishing, scams and spam. <https://udayton.edu/udit/safecomputing/spam.php>

- [5] PhishTank. Phishtank. <https://www.phishtank.com/>
- [6] Whois. Whois domain lookup. <https://www.whois.com/whois/>
- [7] Panos Louridas and Christof Ebert. Machine learning. *IEEE Software*, 33:110–115, 09 2016.
- [8] Michael Jordan and T.M. Mitchell. Machine learning: Trends, perspectives, and prospects. *Science* (New York, N.Y.), 349:255–60, 07 2015.
- [9] Steven Aftergood. Cybersecurity: The cold war online. *Nature*, 547:30+, Jul 2017. 7661.
- [10] Aleksandar Milenkoski, Marco Vieira, Samuel Kounev, Alberto Avritzer, and Bryan Payne. Evaluating computer intrusion detection systems: A survey of common practices. *ACM Computing Surveys*, 48:12:1–, 09 2015.
- [11] Chirag N. Modi and Kamatchi Acha. Virtualization layer security challenges and intrusion detection/prevention systems in cloud computing: a comprehensive review. *The Journal of Supercomputing*, 73(3):1192–1234, Mar 2017.
- [12] Eduardo Viegas, Altair Santin, Andre Fanca, Ricardo Jasinski, Volnei Pedroni, and Luiz Soares de Oliveira. Towards an energy-efficient anomaly-based intrusion detection engine for embedded systems. *IEEE Transactions on Computers*, 66:1–1, Jan 2016. 53
- [13] Y. Xin, L. Kong, Z. Liu, Y. Chen, Y. Li, H. Zhu, M. Gao, H. Hou, and C. Wang. Machine learning and deep learning methods for cybersecurity. *IEEE Access*, 6:35365– 35381, 2018.
- [14] Neha R. Israni and Anil N. Jaiswal. A survey on various phishing and anti-phishing measures. *International journal of engineering research and technology*, 4, 2015.
- [15] Pingchuan Liu and Teng-Sheng Moh. Content based spam e-mail filtering. pages 218–224, 10 2016.
- [16] N. Agrawal and S. Singh. Origin (dynamic blacklisting) based spammer detection and spam mail filtering approach. In *2016 Third International Conference on Digital Information Processing, Data Mining, and Wireless Communications (DIPDMWC)*, pages 99–104, 2016.
- [17] Vikas Sahare, Sheetalkumar Jain, and Manish Giri. Survey:anti-phishing framework using visual cryptography on cloud. *JAFRC*, 2, 01 2015.
- [18] S. Patil and S. Dhage. A methodical overview on phishing detection along with an organized way to construct an anti-phishing framework. In *2019 5th International Conference on Advanced Computing Communication Systems (ICACCS)*, pages 588– 593, 2019.

## APPENDIX

### A. Source Code

```
import requests
```

```
# NOTE: you must manually set API_KEY below using information retrieved from your IBM Cloud account.
```

```
API_KEY = " "
```

```
token_response = requests.post('https://iam.cloud.ibm.com/identity/token', data={"apikey":
```

```
API_KEY, "grant_type": 'urn:ibm:params:oauth:grant-type:apikey'})

mltoken = token_response.json()["access_token"]

header = {'Content-Type': 'application/json', 'Authorization': 'Bearer ' + mltoken}

# NOTE: manually define and pass the array(s) of values to be scored in the next line

payload_scoring = {
    "input_data": [
        {
            "fields": [
                "index",
                "having_IPhaving_IP_Address",
                "URLURL_Length",
                "Shortining_Service",
                "having_At_Symbol",
                "double_slash_redirecting",
                "Prefix_Suffix",
                "having_Sub_Domain",
                "SSLfinal_State",
                "Domain_registration_length",
                "Favicon",
                "port",
                "HTTPS_token",
                "Request_URL",
                "URL_of_Anchor",
                "Links_in_tags",
                "SFH",
                "Submitting_to_email",
                "Abnormal_URL",
                "Redirect",
                "on_mouseover",
                "RightClick",
```

```

        "popUpWidnow",
        "Iframe",
        "age_of_domain",
        "DNSRecord",
        "web_traffic",
        "Page_Rank",
        "Google_Index",
        "Links_pointing_to_page",
        "Statistical_report" ] , "values": [[1,-1,1,1,1,-1,-1,-1,-1,-1,1,1,-1,1,-1,1,-
1,-1,-1,0,1,1,1,1,-1,-1,-1,-1,1,1,-1]]] } }

response_scoring = requests.post('https://us-
south.ml.cloud.ibm.com/ml/v4/deployments/d35291ad-94b0-4257-afd0-
86dc7f82c6eb/predictions?version=2023-01-30', json=payload_scoring,

headers={'Authorization': 'Bearer ' + mltoken})

print("Scoring response")

print(response_scoring.json())

pred=response_scoring.json()

output=pred['predictions'][0]['values'][0][0]

print("result is: ",output)

```