# Using machine learning, a Python-based Phishing URL detection

Meghashri B
Dept of Computer Applications
JNN College of Engineering,
Shimoga, Karnataka , India
meghashrib12298@gmail.com

Sandhya R
Dept of Computer Applications
JNN College of Engineering,
Shimoga, Karnataka , India
sandhya_r@jnnce.ac.in

**Abstract:** With the rapid growth of mobile image display, there has been a growing trend to shift nearly all real-world functions to the cyberworld. Even though this makes our daily lives easier, because of privacy, it also brings various security breaches. The Internet's structured Antivirus and firewall software are used. The majority of attacks can be avoided by using systems and experienced individuals.Attackers try to exploit computer users' vulnerable sites. To phishing them some of the pages on site are modular. Sites to steal include major banking, social networking sites, e-commerce, and other services. User-ids, passwords, and bank account information are examples of sensitive information. To get accountinformation, credit card numbers, and so on. Phishing detection is a difficult task.There are many various solutions presented to this difficult challenge such as list, rule-based detection, animateddetection on the market detection, and so on.

**Keywords**

*Website classification, Protection, Scam, Machine Learning, Phishing, Cyber Security, Spam.*

## 1. INTRODUCTION

In various ways, having a computer and connection to the web makes our work and personal lives easier. It enables us to performtransactions and operations in fields including trade, health, education, communication, finance, aviation, research, engineering, entertainment, and public services. In a timely of mobile and wireless technologies, users that require access to a local network can now connect to the Internet from anywhere and at any time. Cybercriminals, pirates, non-malicious (capped) attacks, and hacktivists are all capable of carrying out attacks .The object information in the computer attacks data it contain, are carried out on the various ways of attacks (Morris Worm) began in 1,988 and has been carried till now.

Fraud, forgery, coercion, shakedown, and hacking are only a few examples. Malware apps and illicit digital content are all examples of illegal digital content's major flaws. As a result, people in cyberspace must take precautions against potential cyber-attacks. as well as social engineering Attackers hope to obtain a large amount of information or money by contacting a wide number oftarget users. According to Kaspersky's data, the average cost of an assault in 2019is between $ 108K and $ 1.4 billion. Furthermore, roughly $ 124 billion is spenton worldwide security products and services.

"Phishing attacks" are the most common and dangerous types of attacks. Cybercriminals typically employ email or other social networking communication channels in this type of assault. Attackers

trick consumers into thinking the message came from a reputable source, such as a bank, an e-commerce site, or something similar. As a result, attempt to gain access to sensitive information.

## 2. REQUIREMENT ANALYSIS

### 2.1 Motivation

The COVID-19 epidemic has increased use of technologies in every sector, causing activities such as planning official meetings, taking courses,shopping, phishers so on to move from the physical to the virtual world. The phishers will have more opportunities to carry out attacks opportunities the victim financially, personally, and professionally. In 2013, about 5.9 billion USD was lost asa result of 450 thousand phishing attacks1.

77 percent of IT professionals believe their security teams are inadequate for today's cybersecurity issue, according to the Checkpoints Research Security Report 2018. According to the same survey, 64 percent of businesses have been the victim of a phishing assault in the last year.

### 2.2Objective

A phishing website is a common social engineering method  that mimics trustful uniform resource locators (URLs) and webpages. The objective of this project is to train machine learning models and deep neurons on the dataset created to predict phishing websites. Both the correct website URL and wrong URL are gathered to form a dataset and the required  URL and website  content-based  features  are extracted. The performance level of each model ismeasured and compared.

## 3. Design measured Data Collection

5000 random phishing Addresses are collected from this dataset to train the machine learning models.

### 3.1 Data Cleaning

Cleaning data is an essential part  of any  machine-learning  development.  In tabular data, you can investigate data using a set of analytical analysis and data visualization  approaches  to  find visualization activities.

Before moving on to more advanced methodologies, every machine-learning project should start with some basic data cleaning activities. These are so basic that even skilled machine learning practitioners sometimes forget them, but are so important that if are skipped, models may break or produce unduly optimistic performance results.

### 3.2 Data Wrangling

#### 3.2.1 Modeling

The dataset is partitioned into 80-20, or 8000 training samples and 2000 testing samples before the ML model is trained. It is evident from the dataset this is a monitored machine-learning problem. Classification and regression are the two main categories of supervisedclassification issues.

#### 3.2.2 Decision tree

Evaluating the label split in the training and testing into Both classesshould get an equal share of the pie.
Legitimate – 0 Phishing – 1

### 3.3 Data Visualization

The Word Embedding approach is used to display text data, with the size of each word indicating its frequency or relevance. A word cloud is being used in frequently analyzed textual documents. Data from fraud word cloud required using word cloud is obtained using matplotlib, pandas, and word cloud is required awarded word cloud in Python.

### 3.4 Data Manipulation

PIL (Python Imaging Library) is a fair and open extension python package system that example of adding several standards for accessing and saving a variety of image file types.

## 4. Result Analysis

### a. Sklearn

Scikit-learn (Sklearn): Scikit-learn (Sklearn) is the most usable and robust machine-learning package in Python. It uses a Python coherence interface to give a set of fast methods for machine learning model classification, clusters, and processing NumPy, SciPy, and Jupiter are the foundations of this library.

- Naive Bayes Strategy : It is a collection of unsupervised segregation algorithms which have independent assumptions between each pair of factors

- Sklearn model selection: A Python package that can be used for sorting, sorting, and model selection, among other things. Selection is a way of designing an analysis template and using it to measure new data.

- Sklearn.pipeline: Truly a representative of the Pipeline Object () {[native code]} function; does not require or allow ratings to be named. Instead, their identities will be automatically charged the capitals of the charge dress.

### b. Pickle

The pickle module supports binary serialization and calculation methods for Python object formats. Pickling converts the Python object component into a byte stream, while randomly converting the byte stream (from a binary file or something like bytes) back to the object system.

### c. Seaborn:

Based on the dataset is a data visualization package based on visualization linked with Python's panda's structures. The core panda of Seaborn is visualization which aids in data visualization and comprehensionUnivariate and bivariate distributions are visualized

### d. Beautiful soup :

Beautiful Soup is a Python package for reading HTML and XML files and extracting data from them. It integrates with h preferred parser to offer fluent navigation, search, and modification of theparse tree. It is normal for programmers to save hours or eve
n days of effort.

### e. Nltk:

It is a Python framework, for creating programs that interact with human language data to toss in quantitative natural language processing (NLP). It includes tokenization, parsing, categorization, stemming, tagging, and semantic reasoning libraries.

### f. Matplotlib :

It's a graphing library for the Python computer language and NumPy, the Python numerical mathematics extension. Matplotlib is an incredible Python visualization page for 2D array charts. Matplotlib is a multiplatform data via visualization package based on Numpyand designed to operate with the deformitystack as a whole.

### g.Flask :

Extension for object-relational all mappers, from valid framework-relatedly, different open authentication Pinterest and numerous com framework- relate lack-based are available for Flask interest and LinkedIn are two examples of Flask-based applications.

## 5. Conclusion :

To conclude, phishing is a significant danger to the web's security andsafety, and phishing identification is a significant issue domain. We looked at some of the conventional phishing detection methods, such as target list and inference evaluation methods, as well as their drawbacks. We evaluated the performance of two computer

algorithms on the Hacking Websites Dataset. The algorithms are chosen, formulates and a Chrome plugin for identifying phish websites are detected.

The plugin makes it simple for end users to apply our phishing detection technology.

We propose to construct the phishing detection system as a scalable web service learning , so that new phishing assault patterns can be learned quickly in the future.

## 6. References:

[1]. J. Shad and S. Sharma, A Novel Machine Learning Approach to Detect Phishing Websites Jaypee Institute of Information Technology, pp. 425430, 2018.

[2]. Y websites, T. Tuncer, H. Gajkal, and E. Avci, Phishing web sites feature classification base Forensics Secure learning machine, 6th Int. Symp. Digital Forensics Securer. ISDFS 2018 –Proceeding, vol. 2018 Jnua, pp. 15,2018.

[3]T. Peng, I. Harris, and Y. Saw, Detecting Phishing Attacks Using Natural Language Processing and Machine Learning, Proc. – 12th IEEE Int. Conf. Semant. Comput. ICSC 2018, vol. 2018Janua, pp. 300301, 2018.

[4]M.Karabakh and T. Mustafa, Performance comparison of classifierson reducing Forensics Secure website dataset, 6thInt. Symp. Digital Forensics Securer. ISDFS 2018 –Proceeding, vol. 2018Janua, pp. 15,2018.

[5] S. Parekh, D. Parikh, S.Kotak, and P. S. Sankhe, A New Method for Detection of PhishingWebsites: URL Detection, in 2018 Second International Conference on InventiveCommunication and Computational Technologies (ICICCT), 2018, vol. 0, no. Icicct, pp. 949952.

[6] K. Shima et al., Classification of URL bit streams using bagi of bytes, in 2018 21st Conference on Innovation in Clouds, Internet and Networks and Workshops (ICIN), 2018, vol. 91, pp. 15.

[7]WFadhilel, M. Abusharkh, and I. Abdel-Qader, On Feature Selection for the Prediction of Phishing Websites, 2017 IEEE 1Securentl Conf Dependable, AutonSecureur. Compute. 15th Intl Conf Pervasive Intell. Compute. 3rd IntlConf Big Data Intell. Compute. Cyber Sci. Technol. Congr., pp. 871876, 2017.

[8] X. Zhang, Y. Zeng, X. Jin, Z. Yan, and G. Geng, Boosting the Phishing Detection Performance by Semantic Analysis, 2017.

[9] L. Machado and J. Gadge, Phishing Sites Detection Based on C4.5 Decision Tree Algorithm, in 2017 International Conference on Computing, Communication, Control, and Automation, ICCUBE