
PHISHING WEBSITE DETECTION BY MACHINE LEARNING TECHNIQUES

MENTOR NAME:

DR. MAHESH MANCHANDA

MENTEE NAME:

PIYUSH CHAND

M.TECH C.S.E

2410013



INTRODUCTION

1. Phishing is the most commonly used social engineering and cyber attack.
2. Phishing is a type of online scam where fake websites are used to steal personal information.
3. Attackers trick people into thinking the website is real, so they enter things like passwords or bank details.
4. To stop phishing, we can:
 - Teach people how to spot fake websites,
 - Block known fake sites,

Objective

1. A phishing website is a common social engineering method that mimics trustful uniform resource locators (URLs) and web pages.
2. The objective of this project is to train machine learning models and deep neural nets on the dataset created to predict phishing websites.

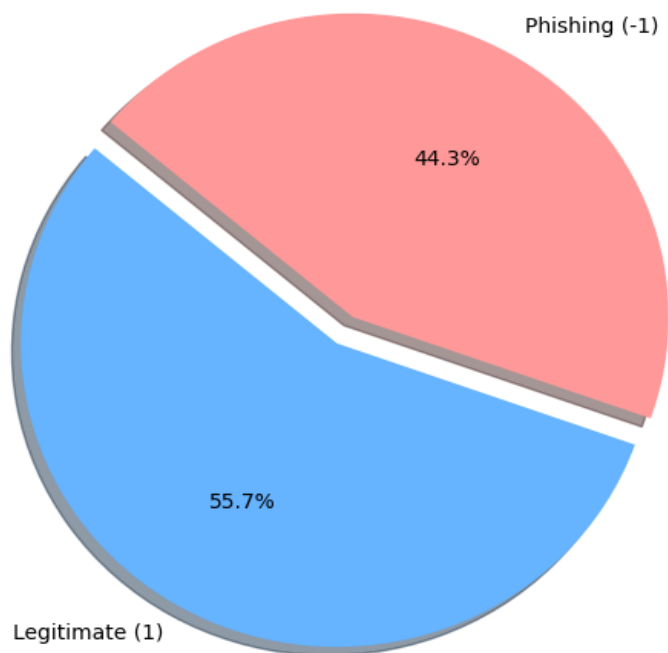
Data Collection

1. Source :
<https://www.kaggle.com/datasets/sai10py/phishing-websites-data?resource=download>
2. Number of record : 11055
3. Total Feature: 31 (30 Features + 1 Target)
4. Target Variable: Result (Indicates whether a website is phishing or legitimate)

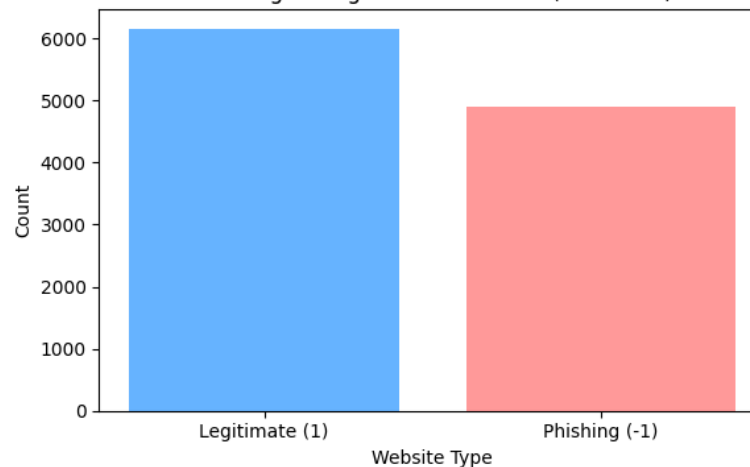
Phishing vs Legitimate Websites

Pie Chart & Bar Chart

Phishing vs Legitimate Websites (Pie Chart)

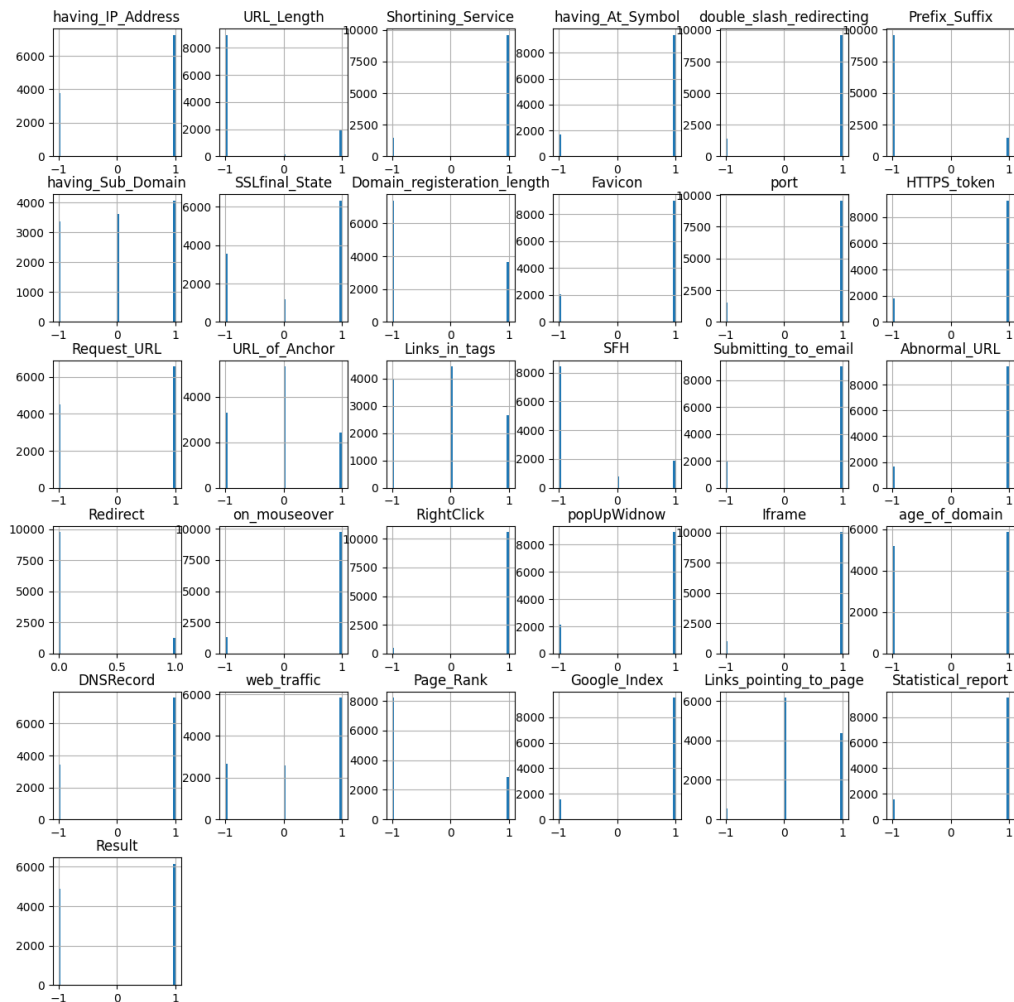


Phishing vs Legitimate Websites (Bar Chart)



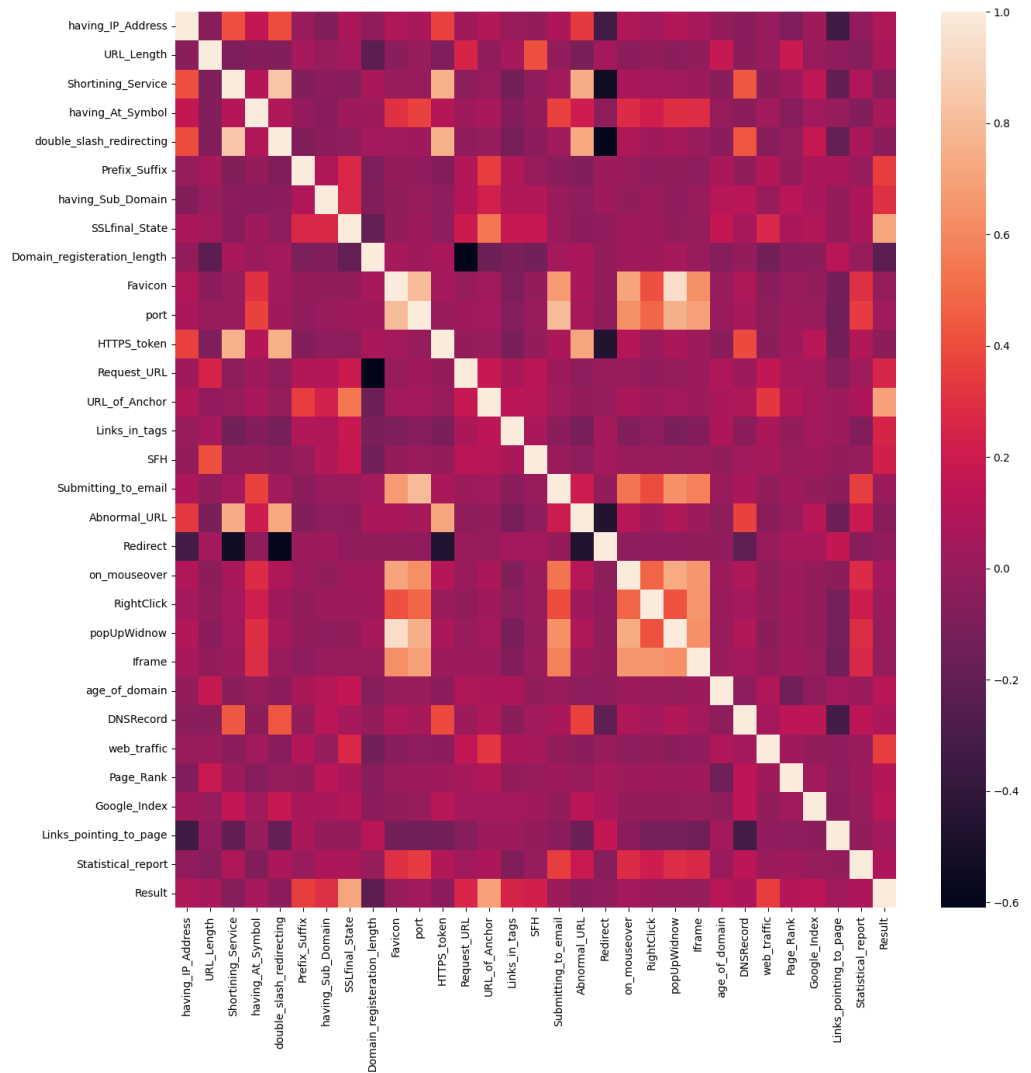
Feature Distribution in Phishing Dataset:

This plot shows the distribution of feature values (-1, 0, 1) across all attributes in the phishing dataset.



Feature Correlation Heatmap :

The Heatmap shows correlation values between all features in the phishing website dataset.



Machine Learning Models & Training

1. An algorithm that learns patterns from data to make predictions.
2. Based on the dataset above, we can identify this as a supervised machine learning task. Supervised learning problems are generally categorized into two main types: classification and regression."
3. This notebook explores various supervised learning classification models to train the given dataset:-
 - Decision Tree
 - Random Forest
 - Multilayer Perceptron
 - Support Vector Machines
 - XGBoost

Decision Tree

1. Decision Tree is a popular machine learning model used for classification and regression.
2. The goal is to learn the best sequence of questions to predict the correct outcome quickly and accurately.
3. Result :
 - Accuracy on training Data: 0.927
 - Accuracy on test Data: 0.926

Random Forest

1. Random Forest is an ensemble of decision trees used for both classification and regression.
2. The final prediction is made by averaging (for regression) or voting (for classification) across all trees.
3. Result :
 - Accuracy on training Data: 0.931
 - Accuracy on test Data: 0.934

Multilayer Perceptron (MLP):

1. MLPs are basic types of neural networks, also called feed-forward networks.
2. MLP consists of layers of neurons that process inputs in multiple stages.
3. Result:
 - Accuracy on training Data: 0.987
 - Accuracy on test Data: 0.972

Support Vector Machines

1. It works by finding the best boundary (Hyperplane) that separates classes in the data.
2. Commonly used for binary classification (two categories), but can be extended to multi-class.
3. Result:-
 - Accuracy on training Data: 0.929
 - Accuracy on test Data: 0.925

XGBoost Classifier

1. XGBoost stands for eXtreme Gradient Boosting.
2. Designed for high performance, with features like regularization and parallel processing.
3. Result:-
 - Accuracy on training Data: 0.989
 - Accuracy on test Data: 0.974

Comparison of Models

	ML Model	Train Accuracy	Test Accuracy
4	XGBoost	0.989	0.974
2	Multilayer Perceptrons	0.987	0.972
1	Random Forest	0.931	0.934
0	Decision Tree	0.927	0.926
3	SVM	0.929	0.925

➤ For the above comparison, XGBoost is the best performing model on this dataset in terms of both training and test accuracy.

So, saving the model for future use.

Future Scope

1. This project can be added to web browsers and mobile apps to block fake websites.
2. The model can be improved to learn from new phishing attacks automatically.
3. It can also be connected with global phishing databases for better accuracy.

THANK

YOU