

Anti-Phishing Prediction

Sachin Satyanarayan Yadav

Theem College of Engineering

Abstract

Anti-phishing refers to efforts to block phishing attacks. Phishing is a kind of cybercrime where attackers pose as known or trusted entities and contact individuals through email, text or telephone and ask them to share sensitive information. Typically, in a phishing email attack, the message will suggest that there is a problem with an invoice, that there has been suspicious activity on an account, or that the user must login to verify an account or password. Users may also be prompted to enter credit card information or bank account details as well as other sensitive data. Once this information is collected, attackers may use it to access accounts, steal data and identities, and download malware onto the user's computer.

As technology is growing, phishing methods have started to progress briskly and this should be avoided by making use of anti-phishing techniques to detect phishing. Machine learning is a authoritative tool that can be used to aim against phishing assaults. There are several methods or approaches to identify phishing.

The machine learning approaches to detect phishing websites have been proposed earlier and have been implemented. The central aim of this project is to implement the system with high efficiency, accuracy and cost effectively. That is been achieved. The project is implemented using 2 machine learning supervised classification models. The two classification models are Logistic regression and random forest classifier. It was established that the RandomForest classifier provides best accuracy for the selected dataset and gives an accuracy score of 98 percent.

Objective

The project's objectives are as follows:

- To study various automatic phishing detection methods.
- To identify the appropriate machine learning techniques and define a solution using the selected method.
- To apply appropriate algorithms to achieve the solution to phishing attacks
- Make operations on dataset:

This dataset contains 48 features extracted from 5000 phishing webpages and 5000 legitimate webpages, which were downloaded from January to May 2015 and from May to June 2017. An improved feature extraction technique is employed by leveraging the browser

automation framework (i.e., Selenium WebDriver), which is more precise and robust compared to the parsing approach based on regular expressions.

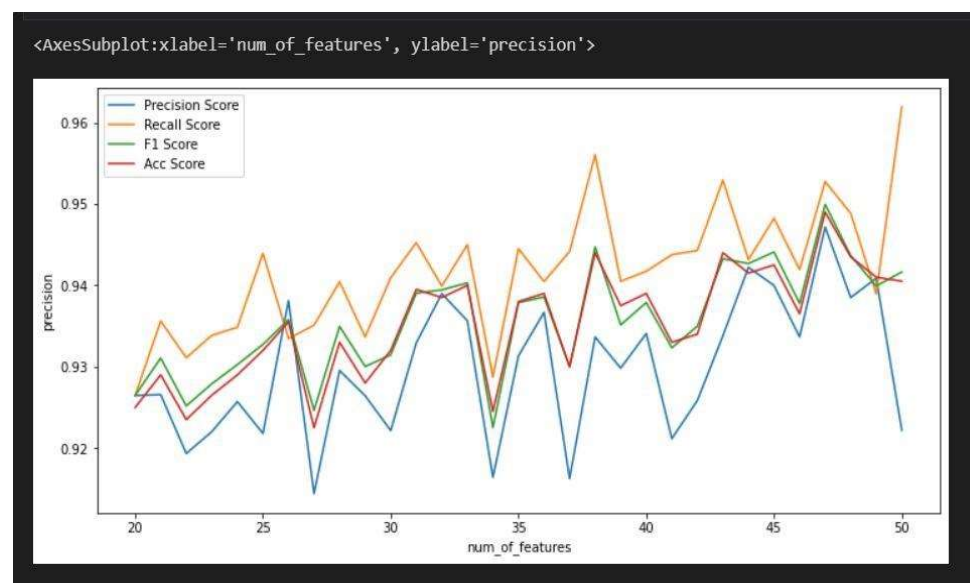
Models/Algorithm

We will first use logistic regression as for baseline, then try to beat the baseline using random forest classifier.

Our evaluation metrics will be accuracy, precision, recall and f1 score.

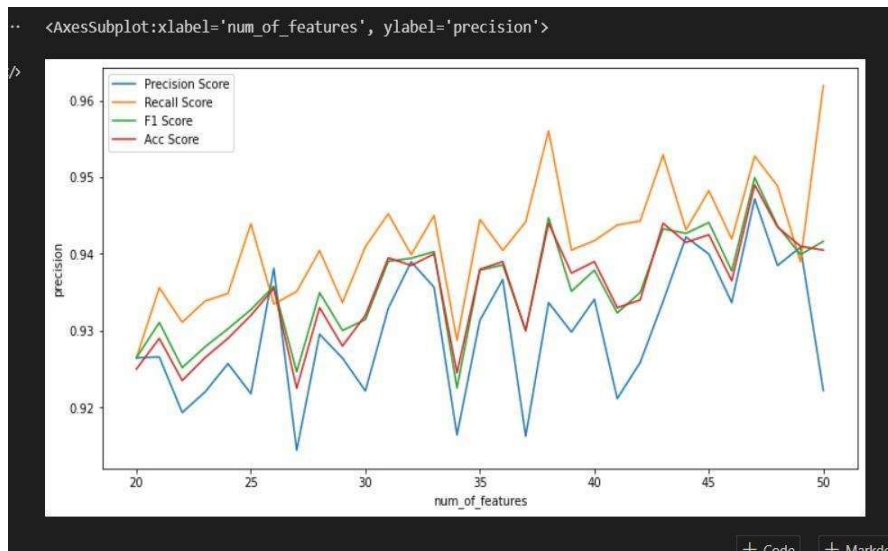
1. Logistic Regression Model

This method is to perform a repetitive training process using logistic regression model, the purpose for this is to find the optimal number of features that can be used to find the best fitted model without adjusting much of the hyper-parameters, hence the idea here is to go with Data-Centric training, basically the method takes number of top N features to be used for training the model and all the evaluation metrics are returned for evaluation purpose.



2. Random Forest Classifier

It is the same method as logistic regression, the only difference is that we are now using random forest classifier for training and trying to beat the logistic baseline. The Random forest classifier creates a set of decision trees from a randomly selected subset of the training set. It is basically a set of decision trees (DT) from a randomly selected subset of the training set and then It collects the votes from different decision trees to decide the final prediction.



Output/Accuracy

The model is now capable of predicting at up to 97% accuracy, this shows the model has high confidence in predicting phishing or non-phishing site.

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 0.97 | 0.97 | 0.97 | 986 |
| 1 | 0.97 | 0.97 | 0.97 | 1014 |
| accuracy | | | 0.97 | 2000 |
| macro avg | 0.97 | 0.97 | 0.97 | 2000 |
| weighted avg | 0.97 | 0.97 | 0.97 | 2000 |