



are often used to lure victims [1][2]. Phishing emails may contain links to websites that distribute malware. Detecting phishing websites often include lookup in a directory of malicious sites. Since most of the phishing websites are short lived, the directory cannot always keep track of all, including new phishing websites. So, the problem of detecting phishing websites can be solved in a better way by machine learning techniques. Based on a comparison of different ML techniques, the random forest classifier seems to perform better.

One way for an end user to benefit from this is to implement detection in a browser plugin. So that the user can be warned in real time as he browses a phishing site. However, browser extensions have restrictions such as they can be written only in JavaScript and they have limited access to page URLs and resources. Existing plugins send the URL to a server, so that the classification can be done in the server and the result is returned to the plugin. With this approach, user privacy is questioned, and the detection may be delayed due to network latency and the plugin may fail to warn the user in right time. As it is an important security problem and considering the privacy aspects, we decided to implement this on a chrome browser plugin which can do the classification without an external server [3][4][5].

2. SYSTEM ARCHITECTURE

Random Forest classifier on 17 features of a website is used to classify whether the site is phishing or legitimate. The dataset raff file is loaded using python raff library and 17 features are chosen from the existing 30 features [6]. Features are selected on basis that they can be extracted completely offline on the client side without being dependent on a web service or third party. The dataset with chosen features are then separated for training and testing. Then the Random Forest is trained on the training data and exported to the above mentioned JSON format. The JSON file is hosted on a URL. The client-side chrome plugin is made to execute a script on each page load, and it starts to extract and encode the above selected features. Once the features are encoded, the plugin then checks for the exported model JSON in cache and downloads it again in case it is not there in cache. The System Architecture is shown in Fig. 1

With the encoded feature vector and model JSON, the script can run the classification. Then a warning is displayed to the user, in case the website is classified as phishing. The entire system is designed lightweight so that the detection will be rapid.

2.1 User Interface (UI) Design

A simple and easy to use User Interface has been designed for the plugin using HTML and CSS. The UI contains a large circle indicating the percentage of the legitimacy of the website in active tab. The circle also changes its color with respect to the classification output (Green for legitimate website and Light Red for phishing). Below the circle, the analysis results containing the extracted features are displayed in the following color code.

Green - Legitimate
Yellow - Suspicious
Red - Phishing



2.3 Module Design

Preprocessing

The dataset is downloaded from UCI repository and loaded into a numpy array. The dataset consists of 30 features, which needs to be reduced so that they can be extracted on the browser. Each feature is experimented on the browser so that it will be feasible to extract it without using any external web service or third party. Based on the experiments, 17 features have been chosen out of 30 without much loss in the accuracy on the test data. More number of features increases the accuracy and on the other hand, reduces the ability to detect rapidly considering the feature extraction time.

Thus a subset of features is chosen in a way that the tradeoff is balanced.

1. IP address
2. Degree of subdomain
3. Anchor tag href of domains
4. URL length
5. HTTPS or not
6. Script & link tag domains
7. URL shortener
8. Favicon domain
9. Empty server form handler
10. '@' in URL
11. TCP Port
12. Use of mailto
13. Redirection with '//'
14. HTTPS in domain name
15. Use of iFrame
16. "-" in domain
17. Cross domain requests

The plugin also displays an alert warning in case of phishing to prevent the user from entering any sensitive information on the website. The test results such as precision, recall and accuracy are displayed in a separate screen. The UI is shown in Figure 3

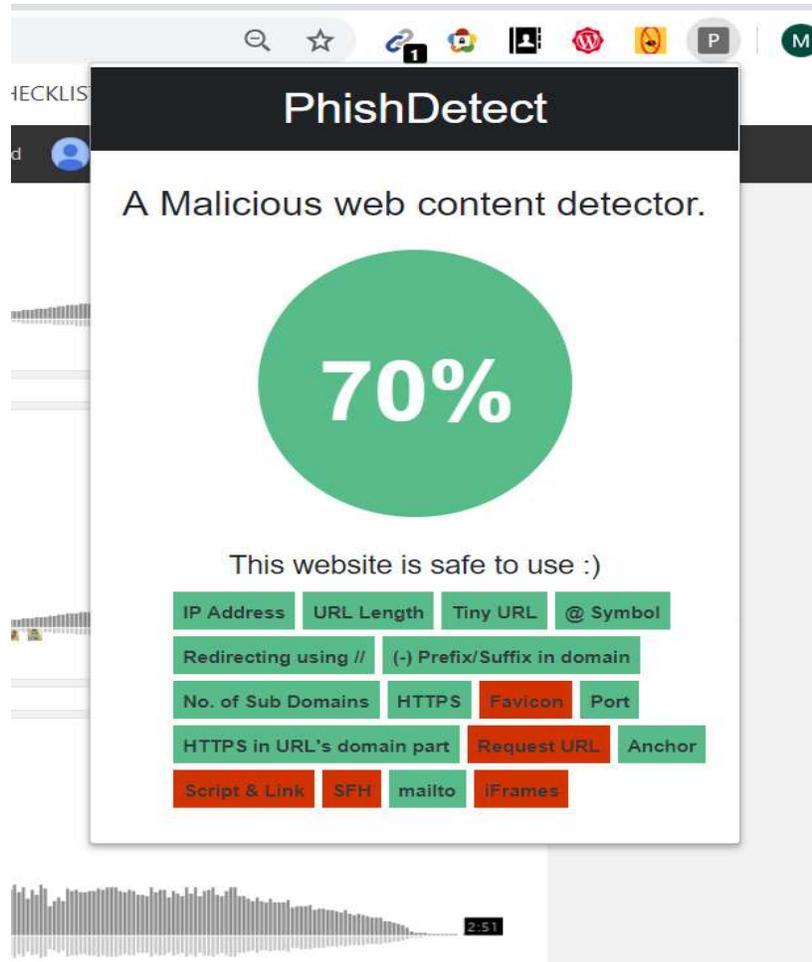


Figure 3: UI Design

Then the dataset is split into training and testing set with 30% for testing. Both the training and testing data are saved to disk.

3. TRAINING

The data from the preprocessing module is loaded from the disk. A random forest classifier is trained on the data using scikit-learn library. Random Forest is an ensemble learning technique and thus 19 an ensemble of 10 decision tree estimators is used. Each decision tree follows CART algorithm and tries to reduce the Gini impurity.

$$Gini(E) = 1 - \sum_{j=1}^c p_j^2 \quad (1)$$

The cross-validation score is also calculated on the training data. The F1 score is calculated on the testing data. Then the trained model is exported to JSON using the next module.



3.1 Exporting Model

Every machine learning algorithm learns its parameter values during the training phase. In Random Forest, each decision tree is an independent learner and each decision tree learns node threshold values and the leaf nodes learn class probabilities. Thus, a format needs to be devised to represent the Random Forest in JSON. The overall JSON structure consists of keys such as number of estimators, number of classes etc. Further it contains an array in Figure 4.4 Random Forest JSON structure 20 which each value is an estimator represented in JSON. Each decision tree is encoded as a JSON tree with nested objects containing threshold for that node and left and right node objects recursively.

3.2 Plugin Feature Extraction

The above mentioned 17 features needs to be extracted and encoded for each webpage in realtime while the page is being loaded. A content script is used so that it can access the DOM of the webpage. The content script is automatically injected into each page while it loads. The content script is responsible to collect the features and then send them to the plugin. The main objective of this work is not to use any external web service and the features needs to be independent of network latency and the extraction should be rapid. All these are made sure while developing techniques for extraction of features. Once a feature is extracted it is encoded into values {-1, 0, 1} based on the following notation.

-1 - Legitimate 0 - Suspicious 1 - Phishing

The feature vector containing 17 encoded values is passed on to the plugin from the content script.

3.3 Classification

The feature vector obtained from the content script is ran through the Random Forest for classification. The Random Forest parameters JSON is downloaded and cached in disk. The script tries to load the JSON from disk and incase of cache miss, the JSON is downloaded again. A JavaScript library has been developed to mimic the Random Forest behavior using the JSON by comparing feature vector against the threshold of the nodes. The output of the binary classification is based on the leaf node values and the user is warned if the webpage is classified as phishing.

4. COMPLEXITY ANALYSIS

4.1 Complexity of the Design Architecture

The complexity of the project lies in balancing the tradeoff between accuracy and rapid detection. Choosing a subset of features that will make the detection fast and at the same time without much drop-in accuracy.

- ❖ Porting of scikit-learn python object to JavaScript compatible format. For example, JSON.
- ❖ Reproducing the Random Forest behavior in JavaScript reduced the accuracy by a small margin.
- ❖ Many features are not feasible to extract without using an external web service. Use of an external web service will again affect the detection time.
- ❖ Maintaining rapid detection is important as the system should detect the phishing before the user submit any sensitive information.

