# Phishing Attack Detection Based on Random Forest with Wrapper Feature Selection Method

Musbau Dogo Abdulrahaman, John K. Alhassan, Olawale Surajudeen Adebayo,
Joseph A. Ojeniyi and Morufu Olalere
[1,2,3,4,5] Department of Cyber Security Science, Federal University of Technology, Minna, Nigeria,
[1]Department of Information and Communication Science, University of Ilorin, Ilorin, Nigeria
rahaman.md@unilorin.edu.ng, jkalhassan@futminna.edu.ng, waleadebayo@futminna.edu.ng,
ojeniyija@futminna.edu.ng and lerejide@futminna.edu.ng

## Abstract

*Phishing website is a fake web page that mimics legitimate website using social engineering techniques to lure unsuspicious users to web page for the purpose of stealing their personal information such as credit card, username and password. Phishing attack is on the increase on a daily basis due to the inability of the existing systems to fully identify a phishing web page from legitimate. Machine learning technique is a trending and intelligent approach for detecting a phishing web page. but, identifying efficient algorithm with the ability to classify and identify web page as either legitimate or phishing in real-time continues to pose a challenge as the existing systems are characterized by misclassifications resulting into low detection rate, high false positive, high running time. The purpose of this study is to develop an efficient machine learning based model with the ability of detecting whether a web page is phishing or not. A performance analysis of some popular classification algorithms was performed and revealed Random Forest as the best classifier on the phishing dataset. A machine learning based model for the detection of phishing attack was built based on Random Forest with wrapper based on classifier attributes evaluator and ranker (CAER) feature selection method. The performance of the proposed model was evaluated using phishing dataset that comprises of static and dynamic features of websites. The experimental results show that the proposed Random Forest based model with feature selection outperformed some of the existing solutions including the best performance of the Random Forest when the full features were used with high accuracy of 97.3% in addition to better precision, sensitivity and lower false positive rate of 0.03 achieved.*

**Keywords:** *Phising_Attack_Detection, Phising_Website,Random_Forest, Wrapper_based_Feature_Selection, Machine_Learning_Techniques.*

## 1. Introduction

The recent development in the field of Information Technology (IT) has led to massive growth in the number of web based services such as e – commerce, online blogs, forums, banking, shopping, gaming, and file sharing. This has created many opportunities for individuals and organizations to interact, transact businesses, and handle payment services conveniently (Varshney, Misra, & Atrey, 2016). Due to the popularity of web applications and its openness in nature, there is an increase in the number of cyber-attacks through web

applications such as Structured Query Language (SQL) injection, Cross-Site Scripting (XSS), malware infections, and phishing attack (Liu, Pan, & Qu, 2016).

Phishing is one of the leading cyber-attacks which refers to a fake website that impersonate original or legitimate websites for the purpose of eliciting sensitive information from users such as username, password, credit card number, bank account number, mobile number among others (Li, Yang, Chen, Yuan, & Liu, 2019). Phishing attacker (phisher) uses social engineering technique to mimics legitimate website and lure users to web pages through different ways such as email, website, short message service (SMS), malware, and voice (Rao and Pais, 2019). Typically, an attacker sends an attack vector commonly inform of email, blog post, chat session which contains malicious link that directs unsuspicious user to a phishing website (Verma & Das, 2017). Phishing attack is a serious threat to web-based services including electronic commerce and continues to affect financial organizations and individuals on a daily basis (Li *et al.,* 2019; Sananse, 2015). According to RSA online fraud report for 2016, the phishing attacks have cost global organizations not less than $4.6 billion in loss in 2015 (RSA, 2017). Most of these attacks are as a result of vulnerabilities in web applications and users. In addition, according to the 2016 report of the internet security threat, about seventy – eight (78) percent of websites in 2015 had web related vulnerabilities that may be used by attackers to launch web attacks (Wang, Zhu, Tan, & Zhou, 2017). In order to curb the phishing attack menace, many security experts in industry and academia have devised solutions with different techniques.

Generally, phishing attack detection techniques can be categorized into two: the blacklist-based and heuristic-based. The blacklist-based technique compares the requested URL with those in the list of phishing sites and detects whether the web page is for phishing or benign activities. This technique relies on black-listed phishing web pages generated by security experts. Recently, several studies have established that this method is not effective with respect to the number of websites hosted daily (Liu *et al*., 2016; Verma & Das, 2017). In contrast, a heuristic based technique uses machine learning based algorithm to extract features from web pages and classifies every instance of the web page as either legitimate or phishing. This method is considered more effective, fast and reliable, because of its ability to detect a freshly created phishing website. (Liu *et al*., 2016; McCluskey, Thabtah, & Mohammad, 2014; Mohammad, Thabtah, & McCluskey, 2014b). However, the accuracy of the machine learning based detection systems depend on picking a set of appropriate features that can genuinely distinguish the websites. In addition, heuristic-based technique is categorized as either URL analysis-based or web content analysis-based. In the past few years, some of the early studies considered URL based extracted features as appropriate for phishing website classification, base on the fact that it is faster and does not require visiting the website for features or searching the internet for the purpose of retrieving and analyzing contents and network level features (Liu *et al*., 2016). Recently, studies have emphasized the need to include web contents analysis as it provides better and improved phishing detection capability (Wang *et al*., 2017). Meanwhile, due to the robustness and versatility of machine learning algorithms in determining the nature of an entity effectively, predicting benign or phishing status of a website can now be performed efficiently.

However, the identification of efficient classification algorithm for building phishing attack detection system is a serious challenge as the existing models are characterized by classification errors resulting into low detection rate, high false positive and high running time. The aim of this paper is to perform a comparative analysis of the widely used classifiers for developing an efficient machine learning based model for detecting phishing attack. The performance analysis of the algorithms will help in the choice of the best performing classifier for the model. The model was evaluated using phishing dataset comprises of lexical, static and dynamic website features. The remaining sections of this paper are organized as follows: The literature review is provided in section 2, while section 3 provides methods adopted for building an efficient phishing attack detection model. Section 4 presents the results and discussion of the findings and the section 5 concludes the studies and provides recommendation for future research.

## 2. Literature Review

This section presents some related research works in the area of phishing attack detection and highlights major techniques used as well as their achievements.

### 2.1 Review of Related Concepts

Phishing is a fake web page with the aim of stealing personal information of users such as credit/debit cards, online banking password, other financial data (Aydin & Baykal, 2015). It uses social engineering technique (such as email spam, rogue software, blog post, chat session) to trick user to click a link that will direct victim to the fake website. Several attempts have been made by security experts to solve the problem of phishing attack using different techniques which can be categorized into blacklist-based and heuristic-based. The major problem of blacklist technique is its inability to cope with the number of sophisticated and dynamic websites currently hosted on a daily basis and not efficient for phishing attack detection. It makes website visitors vulnerable to phishing through phishing websites that have not been blacklisted. Heuristic based solution extracts features from website and analyze using different methods for the purpose of categorizing the web page as either phishing or benign (Mohammad, Thabtah, & McCluskey, 2014a). Heuristic based phishing website detection is considered to be reliable, fast and more efficient due to its ability to detect a freshly hosted phishing websites.

Heuristic based phishing website detection involves two common approaches, URL analysis and web page contents analysis. The URL analysis extracts features from web page link, analyze and determine whether it is a phishing or benign web page. On the other hand, web page content analysis based technique extracts either static features (such as lexical ) or dynamic contents of web page (such as JavaScript) to determine the behavior of a web page (Wang *et al*., 2017). However, some of the early studies in the field of phishing attack detection suggest that a URL based analysis is more efficient than the web content analysis, because it does not contribute to run-time latency and does not expose users to the browser based vulnerabilities (Liu *et al*., 2016). In a URL analysis, it is believed that if the characteristics of a phishing URL are known, it is possible to prevent users from visiting such web page. In addition, such technique does not increase the computational requirements to solve the problem of phishing attack.

A typical website has a URL that points to a specific web page, and every URL has a common syntax that includes protocol, hostname, and path. For example, "http://www.google.com/index.php" can be broken down as in figure 1
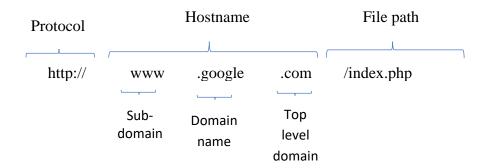


**Figure 1:** Typical URL syntax.

The protocol part represents which network protocol (such as Hypertext Transfer Protocol (http), HTTP with Transport Layer Security (https), and File Transfer Protocol (ftp)). The hostname part identifies the web server on the internet such as www.google.com, where "www" is the sub-domain, while "google" is the domain name, and ".com" shows that it is a Top Level Domain (TLD). The file path represents the path of any file on the computer and it may contain punctuation mark of different kinds like slashes, dots, dashes. In addition, the text after the first "?" symbol indicates the URL query part, the text after the first "@" symbol indicates the parameter part of the URL, while the text after the first "#" symbol indicates the fragment part of the URL.

## 2.2 Review of Related Works

In their quest to solving phishing attack problem, Feroz and Mengel (2015) describes an approach that automatically classifies URL as either phishing or benign in real-time based on their lexical and host based features. A system with the ability to perform URL clustering, categorizing, and ranking was proposed with the ability of adapting to evolving trends in URL features. To develop the system, a clustering was performed on the entire dataset using K-means algorithm to derive a cluster ID used for classification processes. An online URL reputation service was used for URL categorization, which in turn used as a source of information for URL ranking. The extracted and ranked URL features are classified using online algorithm with the help of Mahout as training and learning tool and Microsoft Reputation Services (MRS) as URL categorization. The online algorithm classifiers achieved a reasonable performance with accuracy ranging from 93% to 98% by detecting a large number of phishing websites with a modest false positive rate (FPR).

To determine the important features for classifying legitimate and phishing URLs, Jeeva and Rajsingh (2016) built a malicious URL detection system consisting of two phases, namely, URL searching and feature extraction. The essence of the URL search phase is to reduce the unnecessary computational bottleneck and improves the overall response time before the feature extraction stage. During the feature extraction phase, some heuristic rules are defined to extract features from URLs and produced 14 features that were subjected to association rule mining (apriori and predictive apriori) to determine whether a web page is legitimate or phishing. The dataset was collected from PhishTank database, while the features were extracted using PHP and classified in WEKA. The experiment shows some of the important features of phishing websites to include absence of transport layer security (https), unavailability of the top level domain in the URL, abnormal keyword within the path portion of the URL, dot in the host portion of the URL as well as length of the URL.

Liu *et al*. (2016) proposed a learning based website classification technique that categorizes web pages as either benign, phishing, or malware. A total number of 13 lexical features were extracted based on the analysis of 7,017 URLs from phishing website, 20,976 URLs from benign websites, and 9,285 URLs from malware websites. Similarly, 3 site popularity based features were extracted and 5 host based features were also extracted. The dataset was collected from PhishTank, DMOZ Directory Project and DNS-BH project. For the experimentation, a 5-fold cross validation method was used for training and testing, while Chi-Square method and virtualization tool in WEKA was used to select the most informative features. Three learning algorithms comprising of J48 decision tree, logistic regression and support vector machine were used to train the dataset. The TPR values for Benign, Malware, and Phishing respectively are; DT (98.3%, 99.9%, and 90.7%). SVM (98.9%, 84.8%, and 44%), LR (97.5%, 100%, and 83.2%), and FPR values for Benign, Malware, and Phishing are; DT (3.6%, 0.2%, and 1.1%). SVM (30.6%, 0.7%, and 0.3%), LR (6.4%,0%, and 1.8%). The experiment clearly shows that the DT algorithm achieves the best classification accuracy of 97.53% with the least run time latency due to the feature selection method used.

Similarly, Aydin and Baykal (2015) proposed a safer framework for phishing website detection with high accuracy and lesser time. Phishing and legitimate URLs were collected from PhishTank and Google search engine respectively. Important features for the classification of websites were extracted and categorized into feature matrix compose of five different analyses that include character analysis (alpha_numeric),

keyword analysis, security analysis, domain identity analysis and rank based analysis. In order to identify the most prominent features, CFS method and consistency subset based features selection method were used which produced two different feature metrics 17 features and 25 features respectively and then evaluated on two algorithms Naïve Bayes and Sequential Minimal Optimization (SMO). The performance of the two algorithms was evaluated in WEKA using default settings and standard 10-fold cross validation to divide the training and testing data. The performance metrics include Accuracy, True Positive (TP) Rate, False Positive (FP) Rate and Precision with Naïve Bayes algorithm showing its high performance of 88.17% accuracy, compares to the SMO with 95.39% accuracy.

According to Dedakia & Mistry (2015), recent studies have established that URL based analysis is not enough to detect a phishing website. The author modified an existing multi-label classifier that was based on associative classification algorithm to include additional content and page style features. Some of the proposed features depends on spelling error, copying website, using forms with submit button, disabling right-click, and using pop-ups windows. The proposed phishing detection system was developed using data mining technique with 21 features that include 16 from existing works and additional 5 extracted from web contents. The phishing and normal URL dataset was collected from PhishTank and Id column database respectively and was used for 4 different set of experiments. The proposed system with 21 features performed better in term of accuracy (ranging from 92.85 to 94.29) when compared with the existing work with less features that did not include web content based features.

Similarly, Wang *et al*. (2017) proposed a new method for phishing web page detection based on hybridizing static and dynamic website analyses. Static features of web pages were extracted, trained and tested on classification algorithm at the first stage. Based on the pre-test threshold set, the algorithm then classified the web pages to either benign or malicious. Any web page that falls below the set threshold is regarded as "unknown" and will be put to undergo dynamic analysis stage, where the source code of such web page will be ran in an emulated environment. The emulated environment is equipped to detect any Shell code embedded in JavaScript. If the Shell code is detected, the web page will be classified as malicious otherwise benign. Java was used for the extraction of static features of malicious websites and then categorized into URL, Hypertext Markup Language (HTML) document, and JavaScript in the source code. Feature selection was done using correlation-based feature selection (CFS) and then classified by Decision Tree. For the dynamic analysis, a tool HTML Unit was used to emulate a full browser environment in order to interpret JavaScript code, Document Object Code as well as Shell code analysis, after which the output of HTML Unit is fed into another open source tool "scdbg' for dynamic analysis. The results of the analysis highlight the importance of feature selection in malicious web page analysis. The performance of the hybrid method was compared to the ordinary static and dynamic analyses and revealed the new proposed method outperforms others in term of precision, recall and F-sore as 0.952, 0.882, and 0.916 respectively and follow by static analysis. In addition, when compared with two commercial based anti viruses namely, Avira and 360 Total Security, the proposed method outperformed them in term of detection accuracy and time requirement for web page analysis.

The approach in Mohammad, Thabtah, & McCluskey (2014a) describes the important features for distinguishing phishing websites from a benign ones by applying rule – based classification data mining technique. PHP and JavaScript programs were built to extract features from about 2,500 phishing websites dataset from PhishTank and Millersmile archive, and about 450 legitimate URLs were collected from yahoo directory and categorized into abnormal based features, address based features, HTML based features, JavaScript based features, and domain based features. To build the classifier for phishing website detection, different rule based algorithms such as C45, RIPPER, PRISM, and Apriori were compared and evaluated on a randomly selected 450 legitimate URLs against 450 phishing URLs in WEKA. The results show that C4.5 algorithm outperformed other ones in term of predicting the class of the URL as either normal or abnormal with 5.76% average error rate. In order to reduce classification error, irrelevant features were removed using Chi-Square method and yields features like Request URL, Age of Domain, multi-sub domain, HTTPS not so important in malicious URL classification.

The feature selection shows an improvement in the prediction performance of most of the algorithms.

In addition, Adewole, Akintola, Salihu, Faruk and Jimoh (2019) proposed a hybrid rule induction based algorithm that combines the inductive and intelligent capability of JRIP and Projective Adaptive Resonance Theory (PART) for phishing website detection, using two publicly available phishing datasets from UCI repository database. The approach was evaluated based on the performance metrics, namely, Accuracy, Kappa Statistic, Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE). The evaluation results show that the proposed hybridized method outperformed that of the JRIP and PART in term of accuracy of 0.9453 and 0.9908 respectively on the two datasets. However, the proposed model was not compared with the existing approaches in the literature such as those with adaptive and supervised machine learning methods for phishing website detection.

Despite the efforts of the existing web page filtering techniques, the need for an efficient method with the ability to cope with the dynamic nature and evolvement of the phishing techniques still persists. Therefore, with the insights from the reviewed literature, this study built on the extracted and labeled phishing detection dataset extracted by (Mohammad et al., 2014a) to build an efficient machine learning based model for the detection of phishing websites with high accuracy and low false positive rate (FPR).

## 3. Research Methodology

This section describes the method used in building phishing website detection model including the dataset, algorithms and evaluation metrics.

### 3.1 Phishing Dataset

This article utilizes the phishing website dataset made extracted by (Mohammad et al., 2014a) and made available on the UCI repository database ( https://archive.ics.uci.edu/ml/machine-learning-databases/00327/) for the purpose of building an efficient phishing website detection model. The dataset was built from the phishing URLs collected from PHishTank, Millersmiles archive, and legitimate URLs from Yahoo directory and starting point directory. In generating the dataset, a JavaScript program and PHP script were used to extract important features of phishing and legitimate websites. The extracted features then categorized into four including address based features (12 features), abnormal based features (6 features), HTML and JavaScript based features (5 features), as well as domain based features (7 features). The table 1 shows the features of the UCI phishing websites dataset and its description which consists of 30 features with additional one column as its class.

### 3.2. Phishing Attack Detection Development

In developing the proposed phishing attack detection model using machine learning technique, the following process was followed:
1) **Training and Validation Data Selection:** For the purpose of developing new phishing attack detection model, a phishing training dataset was collected from popular UCI repository database comprises a total number of 11,055 instances and used for training and testing of the model. A popular data mining and machine learning tool (WEKA) version 3.8 was adopted for the experimentation.
2) **Classification Algorithm Selection:** This involves performance analysis of some popular classification algorithms including AdaBoostM1, Bagging, K- Nearest Neighbor (KNN), Sequential Minimal Optimization (SMO), Multilayer Perceptron (MLP), Naïve Bayes (NB), Random Forest (RF), RepTree, JRip, and Decision Table was performed.

**Table 1:** UCI Phishing Attack Dataset

| S/N | Phishing website features | Feature Description | Category |
|---|---|---|---|
| 1 | having_IP_Address | Using of IP address instead of the domain name in the URL. | |
| 2 | URL_Length | Using of long URL to hide the doubtful part in the address bar | |
| 3 | Shortining_Service | Using URL shortening service which redirect user to long URL | **Address bar based Features** |
| 4 | having_At_Symbol | Using "@" symbol in the URL leads the browser to ignore everything preceding the symbol | |
| 5 | double_slash_redirecting | Using "//" symbol which redirect user to another website | |
| 6 | Prefix_Suffix | URL with dash symbol is rarely used in legitimate URLs | |
| 7 | having_Sub_Domain | URL having (dots) sub-domain | |
| 8 | SSLfinal_State | URL with HTTPS | |
| 9 | Domain_registration_length | Period of time of existence of domain | |
| 10 | Favicon | Domain where Favicon is loaded | |
| 11 | Port | Port status | |
| 12 | HTTPS_token | If HTTPS token is real or not | |
| 13 | Request_URL | Determines the domain where external object (e,g image, video) is loading from | |
| 14 | URL_of_Anchor | Determines the nature of anchor tag | |
| 15 | Links_in_tags | Determines if tags are linked to the same domain of the web page. | |
| 16 | SFH | Determines if the domain name in Server Form Handler is different from that of the web page. | **Abnormal Based Features** |
| 17 | Submitting_to_email | Determines if the form is not redirecting requests to different email | |
| 18 | Abnormal_URL | If URL does not look normal | |
| 19 | Redirect | By determining how many times a website has been redirected | **HTML and JavaScript based Features** |
| 20 | on_mouseover | Checks on MouseOver event if it makes any changes on the status bar | |
| 21 | RightClick | Checks if the right client function is disabled or not | |
| 22 | popUpWidnow | How popup Window is used | |
| 23 | Iframe | If iFrame is used or not | |
| 24 | age_of_domain | The information about the age of domain | |
| 25 | DNSRecord | If the URL DNS Record is empty or not | |
| 26 | web_traffic | The popularity of the webpage | |
| 27 | Page_Rank | Based on how the page is ranked | **Domain based Features** |
| 28 | Google_Index | If website is Google index or not | |
| 29 | Links_pointing_to_page | The number of links pointing to the webpage | |
| 30 | Statistical_report | Statistical report on URL based on reports from reputable bodies. | |

3) **Iterative Training and Evaluation of Classification Model:** based on the performance of Random Forest on the phishing dataset, this study adopts the use of FR for the identification of phishing websites. In order to train and evaluate the model, a standard 10 – fold cross validation method was used.

4) **Feature selection:** For the purpose of improving the performance of Random Forest, a feature selection based on wrapper method was used to identify the most relevant features and remove the redundant features in the dataset. A classifier attributes evaluator with ZeroR classifier and Ranker (CAER) as search method was used to rank and select the most relevant features for the phishing dataset classification. The wrapper based feature selection evaluates the worth of each attribute, while ranker generates the attribute ranking.

## 3.3. Design Flowchart of the Phishing Attack Detection Model

This section describes the design flowchart for building the proposed machine learning – based phishing attack detection model.

As represented in figure 2, the algorithm starts by loading a labeled phishing dataset comprises of both phishing and legitimate (benign) instances of web wages into the system. In order to select optimal subset features and ensure that only the relevant features are included in the phishing website classification process, all the 30 input features and the target (class) feature were evaluated with classifier attribute evaluator (CAE) using ZeroR as based classifier in conjunction with ranker (R) to know the worth and ranking of each attribute of phishing attack dataset. Series of experiments were performed by dropping each of the least ranked features until 28 most important features were discovered. For the classification stage, Random Forest algorithm was introduced to identify the class of each of the instances of web page in the dataset. The classifier then makes a decision by classifying each of the instances as either benign or phishing attack. Once a decision is made, a phishing attack is detected which makes the algorithm to stop.

## 3.4. Performance Evaluation Metrics

In order to evaluate the performance of the developed malicious website detection model, the confusion matrix is shown in table 2 and follows by the other popular machine learning performance metrics:

**Table 2:** Confusion Matrix

|  |  |  | Predicted Class | |
|---|---|---|---|---|
|  |  |  | Benign | Phishing |
| **Actual Class** | Benign page | Web | TN (4703) | FP(195) |
|  | Phishing Page | Web | FN(104) | TP(6053) |

Where:
a) **True positive (TP):** The total number of malicious webpage instances "correctly" labeled by the classifier

b) **True Negative (TN):** The total number of normal webpage instances "correctly" labeled by the classifier

c) **False positive (FP):** The total number of normal webpage instances "incorrectly" labeled by the classifier as phishing

d) **False Negative (FN):** The total number of phishing webpage instances "incorrectly" labeled by the classifier as normal

1) **Accuracy:** it measures how accurate a model can detect whether a webpage is legitimate or phishing. It can be expressed as follows:

$$\text{Accuracy (ACC)} = TP+TN/ (TP+FP+FN+TN) \tag{1}$$
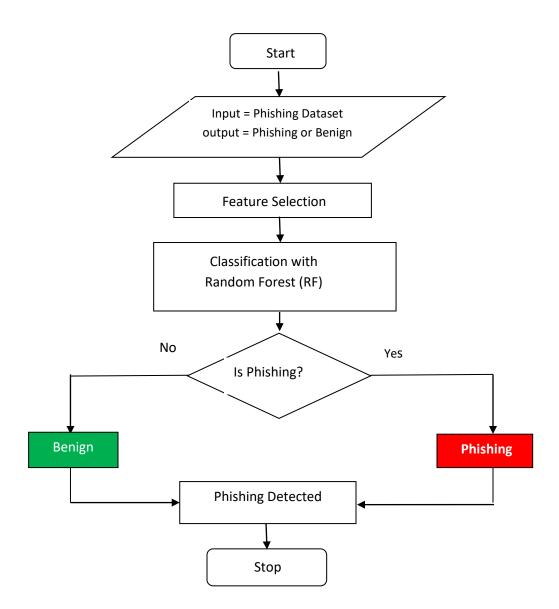


**Figure 2:** Design flow chart for the proposed Random Forest based Phishing Attack Detection

2) **Sensitivity (Recall):** This is also known as detection rate. It is the ratio between total numbers of phishing web pages detected by the model to the actual total number of phishing web pages present in the dataset. It can be expressed as follows:

$$\text{Sensitivity or (DR)} = TP/TP+FN \tag{2}$$

3) **Specificity:** this can be described as the ratio of total number of correctly classified as benign web pages to the actual number of the normal web pages. It is expressed thus:

$$Specificity = TN/TN+FP \qquad (3)$$

4) **Precision:** This can be expressed as the ratio of the total number of the correctly labeled instances of phishing web pages (TP) to the addition of total number of correctly classified phishing web pages (TP) and total number of misclassification as phishing (FP). This is expressed as follows:

$$Precision = TP/ TP + FP \qquad (4)$$

## 4. Results and Discussion

This section presents the results of the experiments conducted in order to develop an effective machine learning based model for phishing attack detection. It also discusses the findings of the experiments in detail.

### 4.1 Comparative Analysis of Classifiers

Phishing attack detection is a classification problem, and there exist several classifiers with different performance capability on different dataset. As established in literature, no single machine learning algorithm works for all problems. In order to be rightly guided in the choice of classifier for the proposed model by identifying the best classifier for phishing attack detection, a comparative analysis of ten (10) widely adopted classification algorithm including AdaBoostM1, Bagging, K- Nearest Neighbor (KNN), Sequential Minimal Optimization (SMO), Multilayer Perceptron (MLP), Naïve Bayes (NB), Random Forest (RF), RepTree, JRip, and Decision Table was performed and investigated. The result as presented in table 3 shows that Random Forest performed better than other classifiers in terms of true positive rate (0.973%), false positive rate (0.03%), Accuracy (97.2592%), Precision (0.973%), Recall (0.973%), F-measure (0.973%), Area under ROC curve (0.996%), and running time (3.14 sec). This was closely followed by K – nearest neighbor that achieved 0.972 as true positive rate, 0.030 as false positive rate, Accuracy of 97.1777%, Precision (0.972%), Recall (0.972%), F-measure (0.972%), Area under ROC curve (0.936%), and running time of 0 sec. Naïve Bayes produced the least performance with an accuracy of 92.9806% within 0.01 seconds among other parameters. For the experiment a 10 fold cross-validation method was used for training and testing processes. Though, as can be seen from table 3 the KNN has a very low running time compare to RF, it performance in term of other metrics is lower than the Random Forest.

**Table 3:** Comparative Analysis of Classifiers for Phishing Detection

| Classifier | TPR | FPR | Accuracy | Precision | Recall | F-Measure | Roc Area | Run Time |
|---|---|---|---|---|---|---|---|---|
| **AdaBoostM1** | 0.926 | 0.076 | 92.5825 | 0.926 | 0.926 | 0.926 | 0.981 | 0.4 |
| **Bagging** | 0.962 | 0.041 | 96.2008 | 0.962 | 0.962 | 0.962 | 0.993 | 2.63 |
| **KNN** | 0.972 | 0.030 | 97.1777 | 0.972 | 0.972 | 0.972 | 0.989 | **0** |
| **SMO** | 0.938 | 0.066 | 93.8037 | 0.938 | 0.938 | 0.938 | 0.936 | 49.18 |
| **MLP** | 0.969 | 0.033 | 96.9064 | 0.969 | 0.969 | 0.969 | 0.995 | 195.86 |
| **Naïve Bayes** | 0.930 | 0.076 | 92.9806 | 0.930 | 0.930 | 0.930 | 0.981 | 0.01 |
| Random Forest | **0.973** | **0.030** | **97.2592** | **0.973** | **0.973** | **0.973** | **0.996** | **3.14** |
| **RepTree** | 0.953 | 0.050 | 95.3324 | 0.953 | 0.953 | 0.953 | 0.985 | 0.36 |
| **JRip** | 0.950 | 0.054 | 95.0158 | 0.950 | 0.950 | 0.950 | 0.961 | 9.96 |
| **Decision Table** | 0.932 | 0.75 | 93.2429 | 0.933 | 0.932 | 0.932 | 0.979 | 7.36 |

.

## 4.2 Feature Selection for Phishing Detection

For the purposed of developing an efficient phishing detection model with a less running time, a wrapper subset evaluator (WSE) with ranker method was used for selecting the most relevant features for phishing website classification after several feature selection methods have been investigated. The wrapper based subset evaluator evaluated and ranked the features of phishing dataset according to their worth. The result shows that the first twenty – eight (28) out of the thirty (30) predictors (attributes) are essential in order to detect a phishing attack with an enhanced performance based on Random Forest.

The result of the feature selection method based on feature ranking is presented in table 4 in the following order:

**Table 4:** Phishing Attack dataset Features Ranking Based on Wrapper Subset Evaluator with Ranker Method

| Rank | Feature | Rank | Feature | Rank | Feature |
|---|---|---|---|---|---|
| 1 | Statistical_report | 11 | URL_Length | 21 | Submitting_to_email |
| 2 | Port | 12 | having_At_Symbol | 22 | Google_Index |
| 3 | Favicon | 13 | Prefix_Suffix | 23 | age_of_domain |
| 4 | HTTPS_Token | 14 | double_slash_redirecting | 24 | Iframe |
| 5 | SSLfinal_State | 15 | URL_of_Anchor | 25 | popUpWidnow |
| 6 | Request_URL | 16 | Links_in_tags | 26 | RightClick |
| 7 | Domain_registration_length | 17 | SFH | 27 | Abnormal_URL |
| 8 | having_Sub_Domain | 18 | web_traffic | **28** | **Redirect** |
| 9 | Links_Pointing_to_page | 19 | DNSRecord | **29** | **on_mouseover** |
| 10 | Shortining_Service | 20 | Page_Rank | 30 | having_IP_Address |

In addition, based on several experiments performed to eliminate the redundant features by dropping the least ranked features, it was discovered that removing **Redirect, and on_mouseover** (HTML and JavaScript based features) from the phishing dataset will enhance the performance of the Random Forest (RF) algorithm and reduce the running time compared to when all the features are used. It was observed that any further reduction of the features will deteriorate the performance of the algorithm.

## 4.3 Evaluation of the Proposed Phishing Detection Model

In order to evaluate the performance of the developed Random Forest based phishing attack detection model with classifier attributes evaluator and ranker (CAER)feature selection method, some state-of-the-art performance metrics in machine learning was used as presented in table 5.

The result shows that the developed model with feature selection achieved higher classification accuracy within a minimum running time compared to when all the features were used. This shows the importance of feature selection for the purpose of enhancing the performance of the classifier. It also revealed that Redirect and on_mouseover features are not necessary for the identification of phishing attacks. Specifically, the developed model achieved the following performance: True Positive Rate (0.973%), False Positive Rate (0.03%), Accuracy (97.2953%), Precision (0.973%), Recall (0.973%), F-measure (0.973%), Area under ROC curve (0.996%), and running time of 2.5 seconds.

**Table 5:** Performance of Random Forest with Feature Selection

| Classifier | TPR | FPR | Accuracy | Precision | Recall | F-Measure | Roc Area | Run Time |
|---|---|---|---|---|---|---|---|---|
| **Random Forest Without feature Selection** | 0.973 | 0.03 | 97.259 | 0.973 | 0.973 | 0.973 | 0.996 | 3.14 |
| Random Forest With Feature Selection | 0.973 | 0.03 | **97.295** | 0.97 | 0.973 | 0.973 | 0.996 | **2.5** |

As can be seen from table 5, the proposed model (Random Forest) with feature selection outperformed the Random forest without feature selection in term of Accuracy and Running time for building the model. It is very important for the phishing attack detection system to classify the web pages accurately within a very short period of time which the proposed model has demonstrated.

### 4.4 Proposed Model Comparison with existing Solutions

In order to further evaluate the proposed model, this section compares the performance the developed model with some existing phishing detection solutions in the literature. Table 6 shows that Random Forest with feature selection had superior performance over other benchmarked (existing) phishing attack detection models such as self-structuring neural network based model as in(Mohammad, Thabtah, & McCluskey, 2014b), Naïve Bayes and Sequential Minimal Optimization (SMO) in (Aydin & Baykal, 2015), stacking of Gradient Boosting Decision Tree (GBDT), XGBoost, and Light GBT in (Li *et al*., 2019), Ensemble of Gaussian naïve Bayes, SVM, KNN, Logistic Regression (LR), MLP, Gradient Boosting and Random Forest in (Ubing, Kamilia, Abdullah, Jhanjhi, & Supramaniam, 2019), Decision Tree in (Wang *et al*., 2017), and Multi-label classifier Based Associative Classification (MCAC) in (Dedakia & Mistry, 2015).

As depicted in table 6, the proposed model that consists of an enhanced Random Forest (RF) outperform the benchmarked machine learning based phishing detection models in terms of accuracy, precision, recall, true positive rate, and false positive rate, and F-measure. The results also show that most of the existing works mostly used accuracy and precision for evaluating the performance of their models, with less interest on other important performance metrics such as running time, Recall, F-measure, TPR, and FPR.

In term of accuracy, the proposed model achieved the higher performance (97.2953%)and closely followed by the stacking based model of (Li *et al*., 2019) with 96.45% accuracy. The least accuracy of 92.18% wais achieved by self-structuring neural network (NN) based model proposed by (Mohammad, Thabtah, & McCluskey, 2014b)**.**

**Table 6:** Comparison of the Proposed Model with existing Models

| Reference | Model | TPR | FPR | Accuracy | Precision | Recall | F-Measure | Run Time |
|---|---|---|---|---|---|---|---|---|
| Proposed Model | **RF + CAER** | **0.973** | **0.030** | **97.29** | **0.973** | **0.973** | **0.973** | **2.5** |
| **Ubing et al., (2019)** | Ensemble | - | - | 95.40 | 0.947 | 0.959 | 0.947 | - |
| **Aydin & Baykal(2015)** | SMO | 0.954 | 0.046 | 95.39 | 0.954 | - | - | - |
| **Li et al., (2019)** | Stacking | - | - | 96.45 | - | - | - | **-** |
| **Dedakia & Mistry (2015)** | MCAC | - | - | 94.29 | - | - | - | **-** |
| **Mohammad et al., (2014b)** | NN | - | - | **92.18** | - | - | - | **-** |
| **Wang et al., (2017)** | DT | - | - | - | 0.952 | 0.882 | 0.916 | **-** |

**Note:** (-) means Not Applicable or not reported in the reference.

In addition, the figure 3 and figure 4 further compare the proposed model pictographically with some existing phishing attack solutions in term of accuracy and precision.
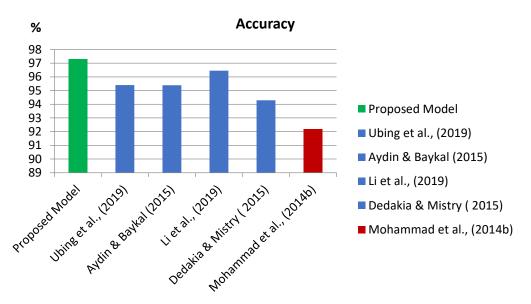


**Figure 3:** Models Accuracy Comparison

**Figure 3:** Models Precision Comparison

The remarkable performance achievement of the proposed model against other benchmarked phishing attack detections demonstrates the appropriateness of the model to enhancing phishing attack detection and shows a commendable contribution to knowledge and online security.

## 5. Conclusion and Recommendations

Phishing attack is one of the most sophisticated web attacks and it is considered as a serious threat to website users, electronic commerce and other web services. This paper proposed a Random Forest based model for the purpose of classifying phishing websites for detecting phishing attacks. The performance of the classification algorithm with feature selection based on classifier attributes evaluator and ranker (CAER) was evaluated using a phishing dataset consisting of the combination of URL, host and third-party based features. The result of the evaluation shows that the proposed model has high accuracy of 97.2953% and low error rates of 0.03% compare to when the whole features were used as well as other existing machine learning based models. The experimental results also revealed that the new proposed model performs well in term of high precision, sensitivity and running time. Finally, the higher classification accuracy result confirms the appropriateness of the hybridized dataset comprises of static and dynamic features.

For future work, it is hoped that more feature selection on the phishing dataset is performed in order to get most relevant features and further improves the performance of the phishing attack detection model. The authors also plan to use more phishing datasets to further investigate how best to combat phishing attack and make cyber-space safe for all.

# Reference

Adewole, K. S., Akintola, A. G., Salihu, S. A., Faruk, N., & Jimoh, R. G. (2019). *Hybrid Rule-Based Model for Phishing URLs Detection* (Vol. 1). https://doi.org/10.1007/978-3-030-23943-5_9

Amrutkar, C., Kim, Y. S., & Traynor, P. (2017). Detecting Mobile Malicious Webpages in Real Time.*IEEE Transactions on Mobile Computing*, *16*(8), 2184–2197. https://doi.org/10.1109/TMC.2016.2575828

Belgiu, M., & Dragut, L. (2016). Random forest in remote sensing : A review of applications and future directions. *Journal of Photogrammetry and Remote Sensing*, *114*, 24–26.

Breiman, L. E. O. (2001). Random Forests, *45*, 5–32.

Aydin, M., & Baykal, N. (2015). Feature extraction and classification phishing websites based on URL. In *2015 IEEE Conference on Communications and NetworkSecurity, CNS 2015* (pp. 769–770). https://doi.org/10.1109/CNS.2015.7346927

Dedakia, M. (2015). Phishing Detection using Content Based Associative Classification Data Mining, *4*(7), 209–214.

Fan, Y., Ye, Y., & Chen, L. (2016). Malicious sequential pattern mining for automatic malware detection. *Expert Systems with Applications*, *52*(January), 16–25. https://doi.org/10.1016/j.eswa.2016.01.002

Feroz, M. N., & Mengel, S. (2015). Phishing URL Detection Using URL Ranking. In *Proceedings - 2015 IEEE International Congress on Big Data, BigData Congress 2015* (pp. 635–638). https://doi.org/10.1109/BigDataCongress.2015.97

Jabbar, M. A., Aluvalu, R., Satyanarayana, S., & Reddy, S. (2017). RFAODE : A Novel Ensemble Intrusion Detection System. *Procedia Computer Science*, *115*, 226–234. https://doi.org/10.1016/j.procs.2017.09.129

Jeeva, S. C., & Rajsingh, E. B. (2016). Intelligent phishing url detection using association rule mining. *Human-Centric Computing and Information Sciences*, *6*(1). https://doi.org/10.1186/s13673-016-0064-3

Li, Y., Yang, Z., Chen, X., Yuan, H., & Liu, W. (2019). A stacking model using URL and HTML features for phishing webpage detection. *Future Generation Computer Systems*, *94*, 27–39. https://doi.org/10.1016/j.future.2018.11.004

Liu, H., Pan, X., & Qu, Z. (2016). Learning based Malicious Web Sites Detection using Suspicious URLs. In *34th Internationai Conference on Software Engineering.* (pp. 3–5). Retrieved from http://www.dtic.mil/cgi-

Machová, K., Barčák, F., & Bednár, P. (2006). A bagging method using decision trees in the role of base classifiers. *Acta Polytechnica Hungarica*, *3*(2), 121–132.

Mafarja, M., & Mirjalili, S. (2018). Whale optimization approaches for wrapper feature selection. *Applied Soft Computing*, *62*, 441–453.

MathWorks. (2018). Mastering Machine Learning A Step-by-Step Guide with MATLAB.

Mohammad, R. M., Thabtah, F., & McCluskey, L. (2014). Intelligent rule-based phishing websites classification. *IET Information Security*, *8*(3), 153–160. https://doi.org/10.1049/iet-ifs.2013.0202

Mohammad, R. M., Thabtah, F., & McCluskey, L. (2014). Predicting phishing websites based on self-structuring neural network. *Neural Computing and Applications*, *25*(2), 443–458. https://doi.org/10.1007/s00521-013-1490-z

Mazini, M., Shirazi, B., & Mahdavi, I. (2018). Anomaly network-based intrusion detection system using a reliable hybrid artificial bee colony and AdaBoost algorithms. *Journal of King Saud University - Computer and Information Sciences*. https://doi.org/10.1016/j.jksuci.2018.03.011

Rao, R. S., & Pais, A. R. (2019). Jail-Phish: An improved search engine based phishing detection system. *Computers and Security*, *83*, 246–267. https://doi.org/10.1016/j.cose.2019.02.011

RSA (2017). RSA Research and Thought leadership: Insights that Drive Change. Retrieved on March 15, 2017 from https://www.rsa.com/en-us/perspectives/industry/ online-fraud

Sananse, B. E. (2015). Phishing URL Detection : A Machine Learning and Web Mining-based Approach, *123*(13), 46–50.

Ubing, A. A., Kamilia, S., Abdullah, A., Jhanjhi, N., & Supramaniam, M. (2019). Phishing Website Detection: An Improved Accuracy through Feature Selection and Ensemble Learning. *International Journal of Advanced Computer Science and Applications*, *10*(1), 252–257. https://doi.org/10.14569/ijacsa.2019.0100133

Varshney, G., Misra, M., & Atrey, P. K. (2016). A survey and classification of web phishing detection schemes. *Security and Communication Networks*. https://doi.org/10.1002/sec.1674

Verma, R. (2017). What ' s in a URL : Fast Feature Extraction and Malicious URL Detection, 55–63.

Wang, R., Zhu, Y., Tan, J., & Zhou, B. (2017). Detection of malicious web pages based on hybrid analysis. *Journal of Information Security and Applications*, *35*, 68–74. https://doi.org/10.1016/j.jisa.2017.05.008