

Phishing Website Detection

Mr. Tahir Naquash H B¹, Sarang Rijul Prakash², Mohammed Arshad Usman², Mohammed Fahad², Mansoor Khan²

¹Assistant Professor, Department of Computer Science Engineering, HKBK College of Engineering, Bengaluru, Karnataka, India

²Department of Computer Science Engineering, HKBK College of Engineering, Bengaluru, Karnataka, India

ABSTRACT

Article Info

Volume 9, Issue 3

Page Number : 790-795

Publication Issue

May-June-2022

Article History

Accepted : 10 June 2022

Published : 30 June 2022

It is a crime to practice phishing by employing technical tricks and social engineering to exploit the innocence of unaware users. This methodology usually covers up a trustworthy entity so as to influence a consumer to execute an action if asked by the imitated entity. Most of the times, phishing attacks are being noticed by the practiced users but security is a main motive for the basic users as they are not aware of such circumstances. However, some methodologies are limited to look after the phishing attacks only and the delay in detection is mandatory. In this paper we emphasize the various techniques used for the detection of phishing attacks. We have also discovered various techniques for detection and prevention of phishing. Apart from that, we have introduced a new model for detection and prevention of phishing attacks.

Keywords : Detection, Phishing email, Filtering, Classifiers, Machine learning, Authentication

I. INTRODUCTION

The purpose or goal behind phishing is data, money or personal information stealing through the fake website. The best strategy for avoiding the contact with the phishing web site is to detect real time malicious URL. Phishing websites can be determined on the basis of their domains. They usually are related to URL which needs to be registered (low-level domain and upper-level domain, path, query). Recently acquired status of intra-URL relationship is used to evaluate it using distinctive properties extracted from words that compose a URL based on

query data from various search engines such as Google and Yahoo.

These properties are further led to the machine-learning based classification for the identification of phishing URLs from a real dataset. This paper focus on real time URL phishing against phishing content by using phish-STORM. For this a few relationship between the register domain rest of the URL are consider also intra URL relentless is consider which help to dusting wish between phishing or non phishing URL. For detecting a phishing website certain typical blacklisted URLs are used, but this

technique is unproductive as the duration of phishing websites is very short.

Phishing is the name of avenue. It can be defined as the manner of deception of an organization's customer to communicate with their confidential information in an unacceptable behavior. It can also be defined as intentionally using harsh weapons such as Spasm to automatically target the victims and targeting their private information. As many of the failures being occurred in the SMTP are exploiting vectors for the phishing websites, there is a greater availability of communication for malicious message deliveries. Along with the various criminal enterprises, if there is enough amount of money generated through the mode of phishing, hunting of various other systems of message delivery can be done, even though the errors are closed eventually in SMTP [5]. Along with the ever increasing dishonesty through phishing scams, organizations are getting more attention from their customers regarding the security of their personal information.

DIFFERENT KINDS OF PHISHING ATTACKS

- **Malware-Based Phishing:** - It refers to the execution of wicked software on the user's PC. Malwares are intruded along with an attachment in the email, as the downloadable files can trace the inputs from keyboard.
- **Deceptive Phishing:** - Actual meaning of phishing is secretarial stealing using direct communication but nowadays the most commonly used method is deceptive messaging. The text sent to the victim concerns about the need of verification of account details, system failure makes it mandatory to re-enter the details of users, fake charges, unfavorable changes in account, unexpected free provisions leading to fast actions, and a lot of more are being broadcasted to maximum number of recipients hoping that the innocents may fall in their trap.
- **System Reconfiguration:-** Attacks may apply unwanted changes in the user's machine for

wicked purposes. Illustration: Websites which are mentioned in mostly used files can be changed in such a way that same website is visited repeatedly.

- **Hosts File Poisoning:** - A URL is converted into an IP address before it is broadcasted over the Internet.
- **Data Shoplifting:** - PCs without security may consist of susceptible information being stored on protected servers.

II. LITERATURE SURVEY

A. Protecting user against phishing using Anti-phishing: -

Anti-Phish is used to avoid users from using fraudulent web sites which in turn may lead to phishing attack. Here, Anti Phish traces the sensitive information to be filled by the user and alerts the user whenever he/she is attempting to share his/her information to a untrusted web site. The much effective elucidation for this is cultivating the users to approach only for trusted websites. However, this approach is unrealistic. Anyhow, the user may get tricked. Hence, it becomes mandatory for the associates to present such explanations to overcome the problem of phishing. Widely accepted alternatives are based on the creepy websites for the identification of "clones" and maintenance of records of phishing websites which are in hit list.

B. Learning to Detect Phishing Emails: -

An alternative for detecting these attacks is a relevant process of reliability of machine on a trait intended for the reflection of the besieged deception of user by means of electronic communication. This approach can be used in the detection of phishing websites, or the text messages sent through emails that are used for trapping the victims[2]. Approximately, 800 phishing mails and 7,000 non phishing mails are traced till date and are detected accurately over 95% of them along with the categorization on the basis of

0.09% of the genuine emails. We can just wrap up with the methods for identifying the deception, along with the progressing nature of attacks

C. Phishing detection system for e-banking using fuzzy data mining: -

Phishing websites, mainly used for e-banking services, are very complex and dynamic to be identified and classified. Due to the involvement of various ambiguities in the detection, certain crucial data mining techniques may prove an effective means in keeping the e-commerce websites safe since it deals with considering various quality factors rather than exact values. In this paper, an effective approach to overcome the “fuzziness” in the e-banking phishing website assessment is used an intelligent resilient and effective model for detecting

e-banking phishing websites is put forth. The applied model is based on fuzzy logics along with data mining algorithms to consider various effective factors of the e-banking phishing website.

D. Collaborative Detection of Fast Flux Phishing Domains:-

Here, two approaches are defined to find correlation of evidences from multiple servers of DNS and multiple suspects of FF domain[3]. Real life examples can be used to prove that our correlation approaches expedite the detection of the FF domain, which are based on an analytical model which can quantify various DNS queries that are required to verify a FF domain. It also shows implementation of correlation schemes on a huge level by using a distributed model, that is more scalable as compared to a centralized one, is publish N subscribe correlation model known as LARSID. In deduction, it is quite difficult to detect the FF domains in a accurate and timely manner, as the screen of proxies is used to shield the FF Mother ship. A theoretical approach is used to analyze the problem of FF detection by calculating the number of DNS queries required to get back a certain amount of unique IP addresses.

E. A Prior-based Transfer Learning Method for the Phishing Detection: -

A logistic regression is the root of a priority based transferrable learning method, which is presented here for our classifier of statistical machine learning. It is used for the detection of the phishing websites depending on our selected characteristics of the URLs. Due to the divergence in the allocation of the features in the distinct phishing areas, multiple models are proposed for different regions. It is almost impractical to gather enough data from a new area to restore the detection model and use the transfer learning algorithm for adjusting the existing model. An appropriate way for phishing detection is to use our URL based method. To cope with all the prerequisites of failure of detecting characteristics, we have to adopt the transferring method to generate a more effective model [1]. Comparative study of the classifiers’ model-based features is shown in the table.

ALGORITHM

A. Decision Tree Algorithm:-

One of the most widely used algorithm in machine learning technology. Decision tree algorithm is easy to understand and also easy to implement. Decision tree begins its work by choosing best splitter from the available attributes for classification which is considered as a root of the tree. Algorithm continues to build tree until it finds the leaf node. Decision tree creates training model which is used to predict target value or class in tree representation each internal node of the tree belongs to attribute and each leaf node of the tree belongs to class label. In decision tree algorithm, gini index and information gain methods are used to calculate these nodes.

B. Random Forest Algorithm:-

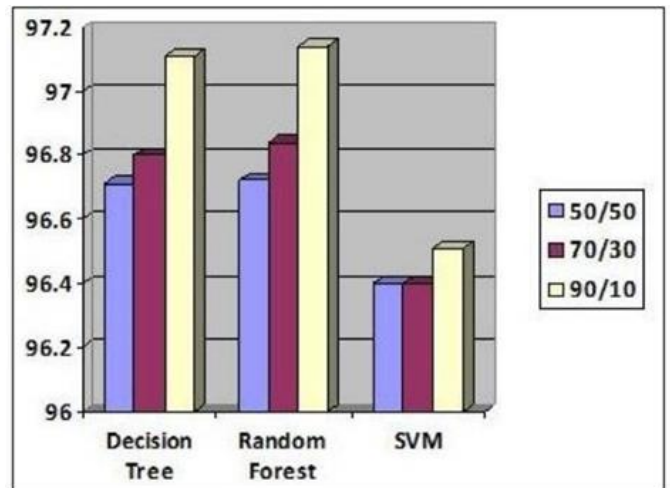
Random forest algorithm is one of the most powerful algorithms in machine learning technology and it is based on concept of decision tree algorithm. Random forest algorithm creates the forest with number of decision trees. High number of tree gives high

detection accuracy. Creation of trees are based on bootstrap method. In bootstrap method features and samples of dataset are randomly selected with replacement to construct single tree.

Among randomly selected features, random forest algorithm will choose best splitter for the classification and like decision tree algorithm; Random forest algorithm also uses gini index and information gain methods to find the best splitter. This process will get continue until random forest creates n number of trees. Each tree in forest predicts the target value and then algorithm will calculate the votes for each predicted target. Finally random forest algorithm considers high voted predicted target as a final prediction.

C. Support Vector Machine Algorithm :- Support vector machine is another powerful algorithm in machine learning technology. In support vector machine algorithm each data item is plotted as a point in n-dimensional space and support vector machine algorithm constructs separating line for classification of two classes, this separating line is well known as hyper plane. Support vector machine seeks for the closest points called as support vectors and once it finds the closest point it draws a line connecting to them.

Support vector machine then construct separating line which bisects and perpendicular to the connecting line. In order to classify data perfectly the margin should be maximum[4]. Here the margin is a distance between hyper plane and support vectors. In real scenario it is not possible to separate complex and non linear data, to solve this problem support vector machine uses kernel trick which transforms lower dimensional space to higher dimensional space.



Dataset Split ratio	Classifiers	Accuracy Score	False Negative Rate	False Positive Rate
50:50	Decision Tree	96.71	3.69	2.93
	Random Forest	96.72	3.69	2.91
	Support vector machine	96.40	5.26	2.08
70:30	Decision Tree	96.80	3.43	2.99
	Random Forest	96.84	3.35	2.98
	Support vector machine	96.40	5.13	2.17
90:10	Decision Tree	97.11	3.18	2.66
	Random Forest	97.14	3.14	2.61
	Support vector machine	96.51	4.73	2.34

III. CONCLUSION

This paper aims to enhance detection method to detect phishing websites using machine learning technology. We achieved 97.14% detection accuracy using random forest algorithm with lowest false positive rate.

Also result shows that classifiers give better performance when we used more data as training data. In future hybrid technology will be implemented to detect phishing websites more accurately, for which

random forest algorithm of machine learning technology and blacklist method will be used.

IV. REFERENCES

- [1]. Gunter Ollmann, "The Phishing Guide Understanding & Preventing Phishing Attacks", IBMInternet Security Systems, 2007.
- [2]. T. Churi, P. Sawardekar, A. Pardeshi, and P. Vartak, "A secured methodology for anti-phishing," Proc. 2017 Int. Conf. Innov. Information, Embed. Commun. Syst. ICIECS 2017, vol. 2018- Janua, pp. 1-4, 2018, doi: 10.1109/ICIECS.2017.8276081.
- [3]. Mahmoud Khonji, Youssef Iraqi, "Phishing Detection: A Literature Survey IEEE, and Andrew Jones, 2013
- [4]. Modeling and Preventing Phishing Attacks by Markus Jakobsson, Phishing detection system for e-banking using fuzzy data mining by Aburrous, M. Dept. of Comput., Univ. of Bradford, Bradford, UK ; Hossain, M.A. ; Dahal, K. ; Thabatah, F.
- [5]. Mohammad R., Thabtah F. McCluskey L., (2015) Phishing websites dataset. Available: <https://archive.ics.uci.edu/ml/datasets/Phishing+Websites> Accessed January 2016
- [6]. M. V. Kunju, E. Dainel, H. C. Anthony, and S. Bhelwa, "Evaluation of phishing techniques based on machine learning," 2019 Int. Conf. Intell. Comput. Control Syst. ICCS 2019, no. Icccs, pp. 963-968.

Cite this article as :

Mr. Tahir Naquash H B, Sarang Rijul Prakash, Mohammed Arshad Usman, Mohammed Fahad, Mansoor Khan , "Phishing Website Detection", International Journal of Scientific Research in Science and Technology (IJSRST), Online ISSN : 2395-602X, Print ISSN : 2395-6011, Volume 9 Issue 3, pp. 791-795, May-June 2022. Available at doi : <https://doi.org/10.32628/IJSRST2293116>
Journal URL : <https://ijsrst.com/IJSRST2293116>