# Heuristic Based Approach For Fraud Detection using Machine Learning.

 <sup>1</sup>Nida Khan, <sup>2</sup>Prashantkumar rai, <sup>3</sup>Manliv Kaur, <sup>4</sup>Riddhi Panchal, <sup>5</sup>Mr. Nilesh Rathod
<sup>1,2,3,4</sup>Dept of I.T, Rajiv Gandhi Institute of Technology, Mumbai, India
<sup>5</sup> Professor, Dept of I.T, Rajiv Gandhi Institute of Technology, Mumbai, India
<sup>1</sup>Rajiv Gandhi Institute of Technology
Off Juhu Versova Link Road, Versova, Andheri (W), Mumbai, India

Abstract::Internet has become a useful part of our regular day to day life as we do almost all of our social and financial activities online. Today, everyone is heavily reliable on internet and online activities such as online shopping, online Banking, online booking, online recharge and many more. Phishing is a type of social engineering attack that targets a user sensitive information through a phony website that appears similar to a legitimate site, or by sending a phishing email .Heuristic based approach is to produce a solution in a reasonable time that is good enough for solving the problem .Heuristic approach defines that it may produce results by themselves, or they may be used in conjunction with optimization algorithms to improve their efficiency (e.g., they may be used to generate good seed values). With the limitation in the existing system we are introducing additional features through the heuristic approach which is simpler and effective than the earlier approaches. This is mainly used for real-world applications and one of this is used in fraud detection on an online platform. Internet environment and diversification of available web services, web attacks have increased in quantity and advanced in quality. Heuristics approach through machine learning underlie the whole field of Artificial Intelligence and the computer simulation of thinking, as they may be used in situations where there are no known algorithm. The heuristic-based detection technique analyses and extracts phishing site features and detects phishing sites using that information. It is imperative to detect and act on such threats in a timely manner. However, blacklists cannot be exhaustive, and lack the ability to detect newly generated malicious URLs. To improve the generality of malicious URL detectors, machine learning techniques have been explored with increasing attention in recent years.

Index terms: heuristic approach, phishing sites, blacklist detection, machine learning, detectors, DNS module.

### I. INTRODUCTION

The advent of new communication technologies has had tremendous impact in the growth and promotion of businesses spanning across many applications including online-banking, ecommerce, and in social networking. Unfortunately, the technological advancements come coupled with new sophisticated techniques to attack and scam users. The most common method to detect malicious URLs deployed by many antivirus groups is the black-list method. Blacklists are essentially a database of URLs that have been confirmed to be malicious in the past. This database is compiled over time. The limitations of traditional security management technologies are becoming more and more serious given this exponential growth of new security threats, rapid changes of new IT technologies, and significant shortage of security professionals. The first module is the URL and DNS matching module which contains a whitelist, which is used to increase the running time and decrease the false negative rate. Our white-list maintains two parameters, domain name and corresponding IP address. Whenever a user accesses a website, then the system matches the domain name of the current website with white-list. If the domain of the current website is matched with the white-list, then the system matches the IP address to take the decision. When the user access a website which is already present in the white-list, then our system matches the IP address of the corresponding domain to check the DNS poisoning attack. Our white-list starts with zero; it means that at the beginning, there is no domain in the list and the white-list starts increasing once a user accesses the new webpages. If the user is accessing the website for the first time, then the domain of the website will not be present in the white-list. In that case, our second module starts working. The second module is the phishing identification module, which checks whether a webpage is phishing. this then extract the hyperlinks from the webpage and apply our phishing detection algorithm. Our phishing detection algorithm examines the features from the hyperlinks to take the decision. After checking the legitimacy, if the website is phishing, then the system shows the warning to the user. Moreover, if the website is legitimate, then the system updates the domain in the white-list.

#### SIMULATIONS AND EXPERIMENTAL RESULTS

Platform/Parameter	Time(seconds)	Accuracy (%)
	10.00	( <b>7</b> ) <b>7</b>
Database handling unit	10.02	65.05
Detection Unit	20-30	75.02
Result generation unit	5-10	86.20

## II. PROPOSED ALGORITHM

*ID3 ALGORITHM:* Iterative Dichotomiser that is ID3 is a decision tree learning algorithm which was invented by Ross Quinlan which is used for generation of decision tree from datasets. ID3 is the precursor to the C4.5 algorithm and is typically used in fields like machine learning and natural language processing domains.

The ID3 algorithm consists of original set S as the root node. On each iteration of algorithm, it iterates through every unused attribute of the set S and calculates the entropy H(S) or the information gain IG(S) of that attribute. The attribute with the smallest entropy value or largest information gain value is selected. The set split S is then split by the selected attribute to produce subsets of the data. For example age is less than 50, age is between 50 and 100, age is greater than 100. The algorithm continues recursion on each subset, considering only attributes that are never slected before.

Recursion on a subset may stop in one of these cases:

- every element in the subset belongs to the same class (+ or -), then the node is turned into a leaf and labelled with the class of the examples
- there are no more attributes to be selected, but the examples still do not belong to the same class (some are + and some are -), then the node is turned into a leaf and labelled with the most common class of the examples in the subset
- $\circ$  there are no examples in the subset, this happens when no example in the parent set was found to be matching a specific value of the selected attribute, for example if there was no example with age >= 100. Then a leaf is created, and labelled with the most common class of the examples in the parent set.

Throughout the algorithm, the decision tree is constructed with each non-terminal node representing the selected attribute on which the data was split, and terminal nodes representing the class label of the final subset of this branch.

Summary of Algorithm and PSEUDOCODE is as follows:

- 1. Using data set S calculate the entropy of every attribute.
- 2. Split the set S into subsets using the attribute for which the resulting entropy (after splitting) is minimum (or, equivalently, information gain is maximum)
- 3. Make a decision tree node containing that attribute
- 4. Recurse on subsets using remaining attributes.

819

ID3 (Examples, Target Attribute, Attributes) Create a root node for the tree If all examples are positive, Return the single-node tree Root, with label = +. If all examples are negative, Return the single-node tree Root, with label = -. If number of predicting attributes is empty, then Return the single node tree Root, with label = most common value of the target attribute in the examples. Otherwise Regin A + The Attribute that best classifies examples. Decision Tree attribute for Root = A. For each possible value,  $v_i$ , of A, Add a new tree branch below Root, corresponding to the test A =  $v_{\rm f}$ . Let  $Examples(v_i)$  be the subset of examples that have the value  $v_i$  for A If Examples( $v_i$ ) is empty Then below this new branch add a leaf node with label = most common target value in the examples Else below this new branch add the subtree ID3 (Examples( $\nu_i$ ), Target Attribute, Attributes - {A}) End Return Root

# III. CONCLUSION AND FUTURE WORK

In this project, we proposed a URL based phishing attack technique that employs URL-based features .We have added new features by analyzing the websites which are phising websites along with URL based features that were used in the previous studies. We have generated classifiers using machine learning algorithms and found that ID3/DECISION TREE Algorithm are good classifiers.The technique which we have proposed in our project can help naïve users to detect the phising sites based on the features and also help in providing them with security for personal information and reduce damage caused by phising sites and phising attacks.It can detect new and temporary phishing sites that evade existing phishing detection techniques, such as the blacklist-based technique.

### IV. **REFERENCES**

- R. K. Nepali and Y. Wang, "You look suspicious!!: Leveraging visible attributes to classify malicious short urls on twitter," in 2016 49th Hawaii International C]
- C. Seifert, I. Welch, and P. Komisarczuk, "Identification of malicious web pages with static heuristics," in Telecommunication Networks and Applications Conference, 2008. ATNAC 2008. Australasian. IEEE, 2008, pp. 91–960nference on System Sciences (HICSS). IEEE, 2016, pp. 2648–2655.
- Huang, Huajun, Junshan Tan, and Lingxi Liu. "Countermeasure techniques for deceptive phishing attack." New Trends in Information and Service Science, 2009. NISS'09. International Conference on. IEEE, 2009
- PhishTank, [Online] Available:http://www.phishtank.com
- DMOZ, [Online] Available: http://rdf.dmoz.org/rdf/
- Anti Phishing Working Group. (2015. March.) APWG PhishingActivity Trend Report 2nd Quarter 2014. [Online]. Available: http://docs.apwg.org/reports/apwg\_report\_q2\_2010.pdf
- Huang, Huajun, Junshan Tan, and Lingxi Liu. "Countermeasuretechniques for deceptive phishing attack." New Trends in Information and Service Science, 2009. NISS'09. International Conference on.IEEE, 2009.
- Ma, Justin, et al. "Beyond blacklists: learning to detect malicious websites from suspicious URLs." Proceedings of the 15th ACM SIGKDDinternational conference on Knowledge discovery and data mining.
- ACM, 2009

 Nguyen, Luong Anh Tuan, et al. "A novel approach for phishingdetection using URL-based heuristic." Computing, Management and Telecommunications (ComManTel), 2014 International Conference on. IEEE, 2014.

- Wikipedia. (2015. March) Uniform Resource Loactor. Available:http://en.wikipedia.org/wiki/Uniform\_resource\_locator
- Kausar, Firdous, et al. "Hybrid Client Side Phishing Websites
- Detection Approach." International Journal of Advanced ComputerScience and Applications (IJACSA) 5.7 (2014).

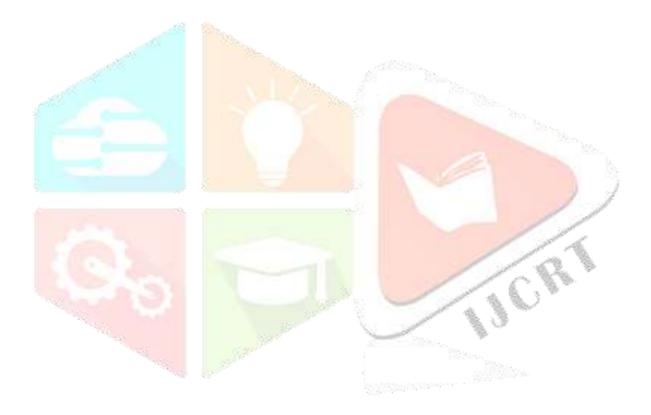
• Sunil, A. Naga Venkata, and Anjali Sardana. "A pagerank baseddetection technique for phishing web sites." Computers & Informatics (ISCI), 2012 IEEE Symposium on. IEEE, 2012.

• Mohammad, Rami M., Fadi Thabtah, and Lee McCluskey."Intelligent rule-based phishing websites classification." Information Security, IET 8.3 (2014): 153-160.

820

- Canali, Davide, et al. "Prophiler: a fast filter for the large-scaledetection of malicious web pages." Proceedings of the 20th international conference on World wide web. ACM, 2011.
- Xiang, Guang, et al. "Cantina+: A feature-rich machine learning
- framework for detecting phishing web sites." ACM Transactions onInformation and System Security (TISSEC) 14.2 (2011): 21.
- WANG, Wei-Hong, et al. "A Static Malicious Javascript Detection

- L. Ladha and T. Deepa, "Feature selection methods and algorithms,"International journal on computer science and engineering, vol 3, no5, 2011.
- Hou, Yung-Tsung, et al. "Malicious web content detection by
- machine learning." Expert Systems with Applications 37.1 (2010): 55-60.
- Cao, Ye, Weili Han, and Yueran Le. "Anti-phishing based on automated individual white-list." Proceedings of the 4th ACM workshop on Digital identity management. ACM, 2008.



Using SVM." strings. Vol. 40. 2013.