

Improved Method Text Feature Extraction for Detection of Phishing Website

^[1]Ganesh S. Gangthade, ^[2]Shyam R. Gote, ^[3] Sachin P. Khengte

^[4] KMukund K. Kharade, ^[5] Prof. Rashmi Tundalwar,

^[1]gangthadeganeshs@gmail.com, ^[2] shyamgote30@gmail.com,

^[3]sachinpkhengte@gmail.com, ^[4]muks6063@gmail.com

Department Of Computer Engineering, Dhole Patil College of Engineering,Pune

Abstract :- — Phishing is a technique of gaining personal information of users from various web-sites. Sometimes it redirects the user to phish webpage to gain information of user like username, password, account and credit card details etc. Our main ambition here is to design system to provide safeguard to users against phishing attacks. Our work is mainly focuses on use of terms and URLs from web page to detect possible phishing patterns from web pages of phishing websites. Process initiates with pars-ing of web page to extract plain text terms and URLs. Further detected terms are fed to TF-IDF and URL weighting system to identify importance of each detected term. Later search engine lookup is carried out for most important terms which help to detect possible victim URLs for given input website. Finally WHOIS lookup is used to compare registration details of websites to correctly categorize website as phishing or legitimate one

Keywords: -Illegitimate, K-means, Naïve Bayes Algorithm, , Phishing detection

I. INTRODUCTION

Phishing is a relatively new internet crime in comparison with other forums, e.g., virus and hacking. Due to the requirement of internet users to facilitate them for 24/7 for banking, housekeeping activities and various many more needs, phishing attacks keep growing. More and more phishing web pages have been found in recent years in an accelerative way (Fu, et al., 2006). Its impact is the breach of information security through the compromise of confidential data and the victims may finally suffer losses of money or other kinds.

A phishing website is a broadly launched social engineering attack that attempts to defraud people of their personal information including credit card num-ber, bank account information, social security number, and their personal credentials in order to use these details fraudulently against them (James, 2006). Phishing web-sites use a number of different techniques to hide the fact that they are not authentic including overwriting or disguising the true URL shown in the browser, overlaying the genuine web site with a crafted pop-up window, drawing fake padlock images on top of the browser window to give the impression that SSL is enabled, and registering SSL certificates for domain names similar to the real organization etc.

II. LITERATURE SURVEY

In this section survey of various methodologies used for detecting phishing websites are provided. Even though different techniques are available (e.g. user's browser based dynamic security, predefined rules for web page creation by Website Company, visual and DOM tree similarity based approach and comparing URLs with blacklisted sites) our main focus is to incorporate use of text mining technique for classifying website as legitimate and fake. Thus following review illustrates application and observation of only text mining based website

No	Paper Details	Methodology	Observation
1.	CANTINA: A Content-Based Approach to Detecting Phishing Web Sites. (WWW 2007)	<ul style="list-style-type: none"> ❑ Traditional TF-IDF approach is used to identify more important words from testing web pages. ❑ Detected top 5 words with highest weight age used to retrieve search engine results. ❑ If domain name match found with top-n search engine results websites, website considered as legitimate otherwise phishing. 	<ul style="list-style-type: none"> • Simple approach to detect phishing websites. • To get better results from traditional TF-IDF approach it is required that web page contains large no. of words.

4	Phish Net: Predictive Blacklisting to Detect Phishing Attacks (IEEE 2010)	<ul style="list-style-type: none"> ❑ Technique for generating suspicious phishing URLs from already available URL blacklist and approximate URL matching suggested ❑ Incorporates mechanism for checking validity of generated URL's and matching content of suspicious URL websites with possible victim websites 	<ul style="list-style-type: none"> • Easy approach to predict URL's which can be generated by phishers
---	---	--	---

III. PROPOSED SYSTEM

We are proposing phishing website detection system which can categorize website as either phishing or legitimate. Preferred use of term weighting based phishing pattern detection reduces the false consideration of phishing website as legitimate one and vice versa (Refer Fig.3.1). Following are the main objectives of the system:

- Extract terms and URLs from web page using DOM parser.
- Identify important terms (brand name) using TF-IDF and URL weighting scheme.
- Search results for brand name using search engine API.
- Identify victim website for detected phishing website.
- In following section in detail explanation of proposed system architecture is given which helps to achieve mentioned objectives.

Module:

(A) Webpage parsing: Web page parsing phase is divided into three sub modules. Functionality of each module described in following

- 1) Plain text extraction
- 2) URL Extraction
- 3) Domain name Extraction

(B) Shroff's word Frequency: The shroff word frequency is used to assign the weight to the words that have been obtained from the web page parser. After the generation the weight is calculated and mostly important keywords are searched for that

TF-IDF [Term Frequency-Inverse Document Frequency]: The TF-IDF weight is a weight utilized as a part of data recovery and text mining. This weight is

a statistical measure used to evaluate how important a word is to a document in a collection or corpus.

(C) URL Based Weight: In this mainly URL's are weighted according to the occurrences in the document of the source code

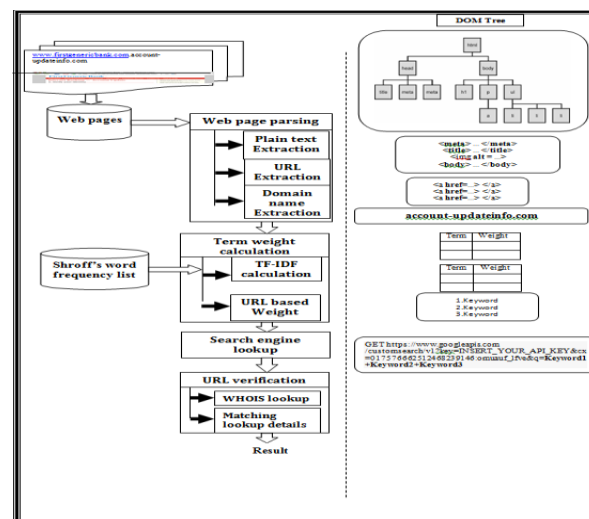
Search Engine lookup: In this lookup search engines used are "Google and yahoo" to search keyword related information. In this phase the three keywords are selected from the shroff's word frequency given to the search engine and the result is generated in the form of top 30 results related to the keyword. After that extract domain name of that related URLs is extracted and is given to next phase of the URL verification.

(D) URL Verification

In This phase we use the "WHOIS lookup Engine" for verifying matching the contents of the web pages with the DOM parser

WHOIS lookup: WHOIS is a program that will tell you the owner of any second-level domain name who has registered it with VeriSign (or with Network Solutions, which was acquired by VeriSign) [7]. Network Solutions was originally the only Internet registrar of the com, net, and org domain names) and many domain names are still registered with VeriSign. In this, the extracted domain name of webpage and extracted domain name of related URL's are given to WHOIS lookup. And also in this phase the comparison to justify us whether the URL is legitimate or phished. This is a step of comparing legitimate domain name is D1 and query of domain name is Dq.

IV. ARCHITECTURE



Architecture dia

V. APPLICATION

Current implementation of proposed system (e.g. as standalone application) is helpful for data mining researchers to identify various data patterns used for phishing websites.

Implemented system can be useful for end users if implemented as “Plug-in” to the browsers where browser can take care of doing mentioned process for user recommended websites e.g. banking, social networking & so on.

VI. FUTURE SCOPE

The future development of community detection system will be concentrated on improving phishing website detection system results by incorporating more than single search engine as it would reduce the chances of getting biased search engine results.

VII. ACKNOWLEDGEMENT

We would like to acknowledge our heartfelt gratitude to our guide prof rashmi tundalwar. Of dhole patil college of engineering, wagholi, pune. For her guidance and motivation

VIII. CONCLUSION

Implemented phishing website detection system considers term and URLs weights for identification of possible victims for phished web pages. Previous approaches rely on traditional term weighting approaches e.g. TF-IDF technique which does not provides accurate results due to lack of sufficient textual content and need of processing no. of web pages. Proposed system adopts URL and shroff's word list weights in weighting scheme which eliminates need of processing multiple web pages. Following parameters are going to be considered while evaluating system:

- (i) Identification of true phishing websites.
- (ii) Identification of true victims of phishing website.

REFERENCES

- [1] Choon Lin Tan, Kang Leng Chiewy, San Nah Sze, “Phishing Website Detection Using URL-Assisted Brand Name Weighting System”, *In IEEE International Symposium on Intelligent Signal Processing and Communication Systems (ISPACS)*, Pages 54-59, 2014
- [2] (2014, June) Phishing guide part 1. PayPal Inc. [Online]. Available: <https://www.paypal.com/au/webapps/mpp/security/generalunderstandphishing>
- [3] (2014, June) Phishing activity trends report, 2nd half / 2010.Anti-Phishing Working Group. [Online]. Available: http://docs.apwg.org/reports/apwg_report_h2_2010.pdf
- [4] Y. Zhang, J. I. Hong, and L. F. Cranor, “Cantina: A content-based approach to detecting phishing web sites,” in *Proceedings of the 16th International Conference on World Wide Web, ser. WWW '07*. New York, NY, US: ACM, 2007, pp. 639–648. [Online]. Available: <http://doi.acm.org/10.1145/1242572.1242659>

