

Malicious Link Detector on Website: Review

*Muhammad Saidu Aliero¹, Zulkifilu Muhammad¹, Abdurashid Allami¹, Basiru Umar Aliero², Salisu Adamu Aliero³,

¹Department of ICT Science, Faculty of Engineering, Kebbi State University of Science and Technology Aliero, Nigeria

²Department of Guidance and Counseling, Faculty of Education, Kebbi State University of Science and Technology Aliero, Nigeria.

³University Library, Kebbi State University of Science and Technology Aliero, Nigeria.
msaidua2000@gmail.com

Date of Submission: 15-10-2020

Date of Acceptance: 15-11-2020

ABSTRACT: Malicious link is one of the most common web-based application vulnerabilities that can be invaded to redirect users from trusted web application to malicious application result in unauthorized access and unauthorized data modification. Researchers have proposed many methods to tackle this problem, however these methods fail to address the whole problem of malware injection attack, because most of the approaches are vulnerable in nature, cannot resist sophisticated attack or limited to scope of subset particular malware type. In this paper we provide a detailed background of malware, and discuss most commonly used method by programmers to defend against malware using analytical evaluation approach lastly we conclude and recumbent future improvements .

Keywords: Malicious link, web application, web application security, web application vulnerably dynamic approach, analytical evaluation

I. INTRODUCTION

Malicious URLs have been widely used to mount various cyber-attacks including spamming, phishing and malware. Detection of malicious URLs and identification of threat types are critical to thwart these attacks. Automatic security assessment tools are used to automatically detect existence of defect, weakness or security flaws that can be exploited by potential attackers (MS Aliero et al, 2015).

These tools provide automatic way of security assessment either by examining the source code of applications or through penetration testing. Security assessment report published by Software Administration, Network and Security and Institute of Computer Security/FBI, show almost 500 computer security analyst concludes that 55% of web penetration tester use automatic approach for

testing and evaluating effectiveness of their applications (Antunes & Vieira, 2010).

This present a major concern that needs to be investigated as web application is today home to billions of users over the globe and Adversaries have used the Web as a vehicle to deliver malicious attacks such as phishing, spamming, and malware infection. For example, phishing typically involves using a fake website seemingly from a trustworthy source to trick people to click a link that takes you to a counterfeit webpage.

To address this challenge many researches in academia and industries have been working on malicious links detection typically through inspecting link to detect their malicious content (McAfee, 2011). Today, many malicious links are found on the web. And because of the high rate of joblessness in Nigeria, job seeking websites are becoming popular as they are convenient and economical. One can easily get a space with a publisher such as Jobbeman.com and place adverts for job vacancies. On the other hand, a lot of people rely heavily on those websites to have information about available job vacancies in other to apply. Unfortunately this bless can also turn to a curse: hackers and attackers have found web links to be low- cost and highly effective means to conduct malicious activities.

This research will focus on analytical analysis on proposed solution to detect malicious links this analysis would help administrator in selecting the best tool in the market or research who wants improve existing solutions.

II. EXISTING SOLUTION

Today's age, it is almost mandatory to have an online presence to run a successful venture. As a result, the importance of the World Wide Web has continuously been increasing especially in job seeking websites.

Unfortunately, the technological advancements come coupled with new sophisticated techniques to attack and scam users in job seeking websites. Such attacks include rogue websites that provide work counterfeit goods, financial fraud by tricking users into revealing sensitive information which eventually lead to theft of money or identity, or when installing malware in the user's system.

This section represent review of the state-of-the-art of malicious URL detection in literature. The research specifically focus on the contributions made for crawling activities used in detecting malicious links. The work also present categories of different various types of crawlers representation used for creating the training data for this task, and also categorize various learning algorithms used to learn a good prediction model. Most of common way by which user can be compromised is through spreading compromised URLs (or the spreading of such URLs forms a critical part of the attacking operation) (Hong, J. 2012).

Compromised URLs that are used for cyber-attacks are termed as malicious URLs. In fact, it was noted that close to one-third of all websites are potentially malicious in nature (Patil, D. R., & Patil, J. B. 2015) demonstrating rampant use of malicious URLs to perpetrate cyber-crimes.

2.1 Review on Malicious links detector

A Malicious URL or a malicious web site hosts a variety of unsolicited content in the form of spam, phishing, or drive-by-exploits in order to launch attacks. Unsuspecting users visit such web sites and become victims of various types of scams, including monetary loss, theft of private information (identity, credit-cards, etc.), and malware installation. Popular types of attacks using malicious URLs include: Drive-by Download, Phishing and Social Engineering, and Spam (Cova, M., Kruegel, C., & Vigna, G. 2010). Drive-by-download (Heartfield, R., & Loukas, G. 2016) refers to the unintentional download of malware upon just visiting a URL. Such attacks are usually carried out by exploiting vulnerabilities in plugins or inserting malicious code through JavaScript.

Phishing and Social Engineering attacks (OpenDNS, L. L. C. 2016) trick the users into revealing private or sensitive information by pretending to be genuine web pages. Spam is the usage of unsolicited messages for the purpose of advertising or phishing. These types of attacks occur in large numbers and have caused billions of dollars' worth of damage every year. Effective systems to detect such malicious URLs in

a timely manner can greatly help to counter large number of and a variety of cyber-security threats.

Consequently, researchers and practitioners have worked to design effective solutions for Malicious URL Detection. The most common method to detect malicious URLs deployed by many antivirus groups is the black-list method. Black-lists are essentially a database of URLs that have been confirmed to be malicious in the past. This database is compiled over time (often through crowd-sourcing solutions, e.g. PhishTank (Sinha, S., Bailey, M., & Jahanian, F. 2008) as and when it becomes known that a URL is malicious. Such a technique is extremely fast due to a simple query overhead, and hence is very easy to implement. Additionally, such a technique would (intuitively) have a very low false-positive rate (although, it was reported that often blacklisting suffered from non-trivial false-positive rates (Garera, S., Provos, N., Chew, M., & Rubin, A. D. 2007).

However, it is almost impossible to maintain an exhaustive list of malicious URLs, especially since new URLs are generated every day. Attackers use creative techniques to evade blacklists and fool users by modifying the URL to "appear" legitimate via obfuscation. (Garera, S., Provos, N., Chew, M., & Rubin, A. D. 2007) identified four types of obfuscation: Obfuscating the Host with an IP, Obfuscating the Host with another domain, obfuscating the host with large hostnames, and misspelling. All of these try to hide the malicious intentions of the website by masking the malicious URL. Recently, with the increasing popularity of URL shortening services, it has become a new and widespread obfuscation technique (hiding the malicious URL behind a short URL) (Garera, S., Provos, N., Chew, M., & Rubin, A. D. 2007), Alshboul, Y., Nepali, R., & Wang, Y. 2015).

Once the URLs appear legitimate, and users visit them, an attack can be launched. This is often done by malicious code embedded into the JavaScript. Often the attackers will also try to obfuscate the code so as to prevent signature based tools from detecting them. Attackers use many other simple techniques to evade blacklists including: fast-flux, in which proxies are automatically generated to host the web-page; algorithmic generation of new URLs; etc. Additionally, attackers can often simultaneously launch more than one attack, which alters the attack-signature, making it undetectable by tools that focus on specific signatures. Blacklisting methods, thus have severe limitations, and it appears almost trivial to bypass them, especially

due to the fact that blacklists are useless for making predictions on new URLs.

To overcome these issues, in the last decade, researchers have applied machine learning techniques for Malicious URL Detection (Patil, D. R., & Patil, J. B. 2015), Chhabra, S., Aggarwal, A., Benevenuto, F., & Kumaraguru, P. 2011), (Ma, J., Saul, L. K., Savage, S., & Voelker, G. M. 2009 and (Hoi, S. C., Wang, J., & Zhao, P. 2014. Machine Learning approaches, use a set of URLs as training data, and based on the statistical properties, learn a prediction function to classify a URL as malicious or benign. This gives them the ability to generalize to new URLs unlike blacklisting methods. The primary requirement for training a machine learning model is the presence of training data. In the context of malicious URL detection, this would correspond to a set of large number of URLs. Machine learning can broadly be classified into supervised, unsupervised, and semi-supervised, which correspond to having the labels for the training data, not having the labels, and having labels for limited fraction of training data. Labels correspond to the knowledge that a URL is malicious or benign.

After the training data is collected, the next step is to extract informative features such that they sufficiently describe the URL and at the same time, they can be interpreted mathematically by machine learning models. For example, simply using the URL string directly may not allow us to learn a good prediction model (which in some extreme cases may reduce the prediction model to a blacklist method).

Thus, one would need to extract suitable features based on some principles or heuristics to obtain a good feature representation of the URL. This may include lexical features (statistical properties of the URL string, bag of words, n-gram, etc.), host based features (WHOIS info, geo-location properties of the host, etc.), etc. These features after being extracted have to be processed into a suitable format (e.g. a numerical vector), such that they can be plugged into an off-the-shelf machine learning method for model training.

The ability of these features to provide relevant information is critical to subsequent machine learning, as the underlying assumption of machine learning (classification) models is that the feature representations of the malicious and benign URLs have different distributions. Therefore, the quality of feature representation of the URLs is critical to the quality of the resulting malicious URL predictive model learned by machine learning. Finally, using the training data with the

appropriate feature representation, the next step in building the prediction model is the actual training of the model. (Sahoo, D., Liu, C., & Hoi, S. C. 2017)

There are plenty of classification algorithms that can be directly used over the training data (Naive Bayes, Support Vector Machine, Logistic Regression, etc.). However, there are certain properties of the URL data that may make the training difficult (both in terms of scalability and learning the appropriate concept). For example, the number of URLs available for training can be the order of millions (or even billions).

As a result, the training time for traditional models may be too high to be practical. Consequently, Online Learning (Canali, D, et al 2011), March) a family of scalable learning techniques have been heavily applied for this task. Similarly, for this task, URLs are represented using the bag-of-words (BoW) features. These features basically indicate whether a particular word (or string) appears in a URL or not - as a result every possible type of word that may appear in any URL becomes a feature. This representation may result in millions of features which would be very sparse (most features are absent most of the time, as a URL will usually have very few of the millions of possible words present in it). Accordingly, a learning method should exploit this sparse property to improve learning efficiency and efficacy.

Despite the promising generalizing ability of machine learning approaches, one potential shortcoming of these approaches for malicious URL detection may be their resource intensive nature (especially while extracting features that are non-trivial and expensive to compute), reducing their practical value when requiring real-time security assurance compared to blacklisting methods.

2.2 Review on Web Application Crawling

In this section, the study provides a brief introduction to crawling and reviews the previous works in the literature. Furthermore, the section also emphasizes on web crawling activities which is one of the major phases in detection of malicious links in web applications together with the relationships among document groups; more specifically, classes, to improve the performance of crawling activities. In particular, the study considers factors that are linked to low coverage of crawling activities by most of the previously proposed crawlers.

With the very fast growth of WWW, the quest for locating the most relevant answers to

users' information needs becomes more challenging. In addition to general purpose Web directories and search engines, several domain specific Web portals and search engines also exist, which essentially aim to cover a specific domain/topic (e.g., education), product/material (e.g., product search for shopping), region (e.g., transportation, hotels etc. at a particular country Chakrabarti, (S., Van den Berg, M., & Dom, B., 1999) or media/file type (e.g., mp3 files or personal homepages (Pant, G., & Srinivasan, P. 2006).

Such specialized search tools may be constructed manually-by also benefiting from possible assistance of the domain experts or automatically. Some examples of the automatic approaches simply rely on intelligent combination and ranking of results obtained from traditional search tools (just like meta search engines), whereas some others first attempt to gather the domain specific portion of the Web using web crawling techniques and then apply other operations (e.g., information extraction, integration, etc.).

(Chakrabarti, S., Van den Berg, M., & Dom, B., 1999) were the first to propose a soft-focus crawler, which obtains a given page's relevance score (i.e., relevance to the target topic) from a classifier and assigns this score to every URL extracted from that page. We refer to this soft-focus crawler as the baseline focused crawler. In a more recent work, they have proposed using a secondary classifier to refine the URL scores and increase the accuracy of this initial soft-focused crawler (Chakrabarti, S., Punera, K., & Subramanyam, M. 2002). An essential weakness of the baseline focused crawler is its inability to model tunneling; that is, it cannot tunnel toward the on-topic pages by following a path of off-topic pages (Barfouroush, A. et al, 2002)..

Two other remarkable projects, the context-graph-based crawler (Diligenti, M. et al, 2000) and Cora's focused crawler (A. McCallum, 1999) achieve tunneling. The context-graph based crawler [56] also employs a best-search heuristic, but the classifiers used in this approach learn the layers which represent a set of pages that are at some distance to the pages in the target class (layer).

More specifically, given a set of seeds, for each page in the seed set, pages that directly refer to this seed page (i.e., parents of the page) constitute layer-1 train set, pages that are referring to these layer-1 pages constitute the layer-2 train set, and soon; up to some predefined depth limit. The overall structure is called the context graph, and the

classifiers are trained so that they assign a given page to one of these layers with a likelihood score.

The crawler simply makes use of these classifier results and inserts URLs extracted from a layer-i page to the layer-i queue, i.e., it keeps a dedicated queue for each layer. URLs in each queue are also sorted according to the classifier's score. While deciding the next page to visit, the crawler prefers the pages nearest to the target class that is, the URLs popped from the queue that correspond to the first nonempty layer with the smallest layer label. This approach clearly solves the problem of tunneling, but it requires constructing the context graph, which in turn requires finding pages with links to a particular page (back links). In contrast, our rule-based crawler uses forward links while generating the rules and transitively combines these rules to effectively imitate tunneling behavior.

In particular, reinforcement learning is used in CORA's focused crawler. CORA's crawler basically searches for the expected future reward by pursuing a path starting from a particular URL. The training stage of classifier(s) involves learning the paths that may lead to on-topic pages in some number of steps.

In contrast, our rule-based crawler does not need to see a path of links during training, but constructs the paths using the transitive combination and chaining of simple rules of length 1.

The focused crawler of Web Topic Management System (WTMS) fetches only pages that are close (i.e., parent, child, and sibling) to on-topic pages (Mukherjee, S. 2000). In WTMS, the relevancy of a page is determined by only using IR-based methods. In another work, Aggarwal et al. attempt to learn the Web's linkage structure to determine a page's likelihood of pointing to an on-topic page. However, they do not consider interclass relationships in the way we do in this study. Bingo! is a focused-crawling system for overcoming the limitations of initial training by periodically retraining the classifier with high quality pages (Diligenti, M., 2000).

Recently, Menczer et al. present an evaluation framework for focused crawlers and introduce an evolutionary crawler (Sizov, S., Graupmann, J., & Theobald, M. 2003). In another work, Pant and Srinivasan provide a systematic comparison of classifiers employed for focused crawling task (Menczer, F., Pant, G., & Srinivasan, P. 2004).

Two recent methods that exploit link context information are explored in (Menczer, F., Pant, G., & Srinivasan, P. 2004). In the first approach, so called text-window, only a number of

words around each hyperlink is used for determining the priority of that link. The second one, tagtree heuristic, uses the words that are in the document object model (DOM) tree immediately in the node that a link appears, or its parents, until a threshold is satisfied.

In (Craven, et al. 1998) we proposed a similar but slightly different technique, so-called page segmentation method, which fragments a Web page according to the use of HTML tags.

Focused crawling paradigm is employed in a number of prototype systems for gathering topic/domain specific Web pages. For instance, in (Pirkola, A. 2012), focused crawling is used for obtaining high quality pages on a mental health topic (depression). In (Pant, G., & Srinivasan, P. 2006), a prototype system is constructed that achieves focused crawling and multilingual information extraction on the laptop and job offers domains.

Black Box Testing Scanner (BBTS) was developed in (Chen and Wu, 2010) to find SQL injection vulnerabilities using dynamic approach. The BBTS uses state aware crawler to identify forms and links with injection parameters. In this approach the authors' use state aware scanner that will be able to recognize webpage that contain injection parameter.

This improvement enables scanner to save time by not downloading pages that do not contain injection parameter and also avoid false positives that result from attacking pages that do not have injection parameter. Author experimental evaluation shows the scanner achieved 100% accuracy on tested application in short period of time. However author did not evaluate his approach with available existing technique to evaluate both accuracy and efficiency of the proposed scanner.

Viper is black box scanner proposed in (Ciampa et al., 2010) that detects SQLi vulnerability through dynamic testing. Scanners that implement dynamic approach are required to have different number of attacks type as well as way of analyzing server response. Viper uses different number of predefined attack pattern to trigger hidden vulnerability in application but it is capable of performing single analysis on server response to attacks. Experiment shows that Viper was able to carry successful attacks that trigger blind SQL injection vulnerability but Viper was not able to report it, this is because analysis was not address to analyzed blind SQL injection vulnerability.

Static Scanner was proposed in (Shar and Tan, 2012) that detect vulnerable point by

characterizing input function into pattern of code attributes. Static code attributes are collected from backward static program slices of sensitive program points, so that to mine both input sanitization code patterns and input validation code patterns, from such static code attributes. This scanner uses vulnerabilities prediction model to enable scanner to predict vulnerable code for SQLi and cross site scripting (XSS) vulnerabilities. Authors evaluated their scanner with different PHP web-based application source code, which shows the effectiveness of their approach. However their approach can only be effective in PHP web-based applications.

4SQLi is the scanner that combines both static and dynamic approach to detect SQLi vulnerabilities (Appelt et al., 2014). 4SQLi uses a single or multiple mutation operators of different types that can be used as a single input parameter to generate desired inputs which will use latter for detecting subtle vulnerabilities that can only be triggered with an input generated by combining multiple mutation operators. For example, consider an application that alters inputs by searching for known attack patterns that can be generated using one of the behavior-changing operators. To form a successful attack, it is necessary to apply a behavior-changing operator and then apply one or more obfuscation operators.

III. DISCUSSION

This study presented explored existing solutions on malicious link detector and result of our analysis reveals that the existing solution faces two common challenges. One is bypassing login authentication; bypassing login application using SQL injection attack is highly dependent on how developer designed authentication query and type of backend database used in target application. Because it is possible single attack pattern that bypass login application in mysql version 5 may failed to bypass mysql version 5.1. A query that is performing check on numbers of return is more easier to bypass than query that is performing checking on single rows returned. Therefore, our malicious link detector perform brute force attack on attempt to bypass login authentication since it has no knowledge of knowing the type and version database is using on backend application.

Another challenge faced by our malicious link detector is inability to recognize malicious links underneath valid link. This is one of the reasons why our malicious link detector failed to identify any malicious link underneath valid link.

RECOMMENDATION

In short, the our recommendation on how to improve the current state of art is as follows:

- i. Enhancing crawler intelligence to perform full indexing of links in target application even if the page requires partial page refreshment.
- ii. Enhancing crawler intelligence to recognize the page as potential to malicious link attack as long as link is directly commutating with database and end users.
- iii. Shaping of pattern used in tautology attaks to limit the number of return record into one to bypass the safety check applied by the developer so as to bypass login application even in the presence of the safety check.
- iv. Updating database of SQL qeury related error messages return by different database server including Oracle, MySQL server and Sybase so as to avoid false positive on platforms other than PHP applications.

REFERENCES

- [1]. Antunes, Nuno, & Vieira, Marco. (2011). Enhancing penetration testing with attack signatures and interface monitoring for the detection of injection vulnerabilities in web services. Paper presented at the Services Computing (SCC), 2011 IEEE International Conference on.
- [2]. Antunes, Nuno, & Vieira, Marco. (2012). Evaluating and improving penetration testing in web services. Paper presented at the Software Reliability Engineering (ISSRE), 2012 IEEE 23rd International Symposium on.
- [3]. Antunes, Nuno, & Vieira, Marco. (2015). Assessing and Comparing Vulnerability Detection Tools for Web Services: Benchmarking Approach and Examples. *Services Computing, IEEE Transactions on*, 8(2), 269-283.
- [4]. Appelt, Dennis, Nguyen, Cu Duy, Briand, Lionel C, & Alshahwan, Nadia. (2014). Automated testing for SQL injection vulnerabilities: an input mutation approach. Paper presented at the Proceedings of the 2014 International Symposium on Software Testing and Analysis.
- [5]. Bandhakavi, Sruthi, Bisht, Prithvi, Madhusudan, P, & Venkatakrishnan, VN. (2007). CANDID: preventing sql injection attacks using dynamic candidate evaluations. Paper presented at the Proceedings of the 14th ACM conference on Computer and communications security.
- [6]. Bau, Jason, Bursztein, Elie, Gupta, Divij, & Mitchell, John. (2010). State of the art: Automated black-box web application vulnerability testing. Paper presented at the Security and Privacy (SP), 2010 IEEE Symposium on.
- [7]. Boyd, Stephen W, & Keromytis, Angelos D. (2004). SQLrand: Preventing SQL injection attacks. Paper presented at the Applied Cryptography and Network Security.
- [8]. Buehrer, Gregory, Weide, Bruce W, & Sivilotti, Paolo AG. (2005). Using parse tree validation to prevent SQL injection attacks. Paper presented at the Proceedings of the 5th international workshop on Software engineering and middleware.
- [9]. Cenzic's. (2014). Application Vulnerability Trends Report: 2014. Retrieved 29/09/2015, from <https://www.info-point-security.com/sites/default/files/cenzic-vulnerability-report-2014.pdf>
- [10]. Chen, Jan-Min, & Wu, Chia-Lun. (2010). An automated vulnerability scanner for injection attack based on injection point. Paper presented at the Computer Symposium (ICS), 2010 International.
- [11]. Cheon, Eun Hong, Huang, Zhongyue, & Lee, Yon Sik. (2013). Preventing SQL Injection Attack Based on Machine Learning. *International Journal of Advancements in Computing Technology*, 5(9).
- [12]. Cho, Ying-Chiang, & Pan, Jen-Yi. (2015). Design and Implementation of Website Information Disclosure Assessment System. *PloS one*, 10(3), e0117180.
- [13]. Ciampa, Angelo, Visaggio, Corrado Aaron, & Di Penta, Massimiliano. (2010). A heuristic-based approach for detecting SQL-injection vulnerabilities in Web applications. Paper presented at the Proceedings of the 2010 ICSE Workshop on Software Engineering for Secure Systems.
- [14]. Cook, William R, & Rai, Siddhartha. (2005). Safe query objects: statically typed objects as remotely executable queries. Paper presented at the Software Engineering, 2005. ICSE 2005. Proceedings. 27th International Conference on.
- [15]. Djuric, Zoran. (2013). A black-box testing tool for detecting SQL injection vulnerabilities. Paper presented at the Informatics and Applications (ICIA), 2013 Second International Conference on.
- [16]. Halfond, William GJ, & Orso, Alessandro. (2007). Detection and prevention of sql

- injection attacks Malware Detection (pp. 85-109): Springer.
- [17]. Huang, Shih-Kun, Lu, Han-Lin, Leong, Wai-Meng, & Liu, Huan. (2013). Craxweb: Automatic web application testing and attack generation. Paper presented at 7th International Conference on the Software Security and Reliability(SERE), IEEE.
- [18]. Kumar, Praveen. (2013). The multi-tier architecture for developing secure website with detection and prevention of sql-injection attacks. *International Journal of Computer Applications*, 62(9), 30-35.
- [19]. Lawal, MA, Sultan, Abu Bakar Md, & Shakiru, Ayanloye O. (2016). Systematic Literature Review on SQL Injection Attack. *International Journal of Soft Computing*, 11(1), 26-35.
- [20]. Liban, Abdilahi, & Hilles, Shadi. (2014). Enhancing Mysql Injector vulnerability checker tool (Mysql Injector) using inference binary search algorithm for blind timing-based attack. Paper presented at the Control and System Graduate Research Colloquium (ICSGRC), 2014 IEEE 5th.
- [21]. Liu, Anyi, Yuan, Yi, Wijesekera, Duminda, & Stavrou, Angelos. (2009). SQLProb: a proxy-based architecture towards preventing SQL injection attacks. Paper presented at the Proceedings of the 2009 ACM symposium on Applied Computing.
- [22]. Livshits, V Benjamin, & Lam, Monica S. (2005). Finding Security Vulnerabilities in Java Applications with Static Analysis. Paper presented at the Usenix Security.
- [23]. McClure, Russell A, & Kruger, Ingolf H. (2005). SQL DOM: compile time checking of dynamic SQL statements. Paper presented at the Software Engineering, 2005. ICSE 2005. Proceedings. 27th International Conference on.
- [24]. Medhane, Munqath H Alattar SP. R-WASP: Real Time-Web Application SQL Injection Detector and Preventer.