



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 8, August 2015

Automatic Phishing Detection System

Pinky M S, Neethu Tom

M.Tech Student, Dept. of CS., Mangalam College, M G university, Kottayam, Kerala, India

Assistant Professor, Dept. of CS., Mangalam College, M G university, Kottayam, Kerala, India

ABSTRACT: Phishing is an e-mail fraud method in which the attacker sends out legitimate-looking email to the user to gather personal and financial information from recipients. Typically, the messages appear to come from well known and trustworthy web sites. Many methods are used to detect phishing attacks. The most recent method is to using PhishStorm, it is an automated phishing detection system. PhishStorm using URL word extraction, feature computation and finally prediction to detect whether the URL is Phishing URL or not. Content based approach is used with PhishStorm to detecting phishing websites. Content based approach is based on the LDA (Latent Dirichlet Allocation). Content based approach and PhishStorm can be implementing in to the web browser to improve accuracy. Implement these approaches in to Mozilla Firefox because it is an open source web browser.

KEYWORDS: phishing attack; PhishStorm; content based approach; rank calculation; LDA.

I. INTRODUCTION

Phishing is the attempt to acquire sensitive information such as usernames, passwords, and credit card details or phishing is an e-mail fraud method in which the perpetrator sends out legitimate-looking email in an attempt to gather personal and financial information from recipients. Typically the messages appear to come from well known and trustworthy websites.

Phishing is a continual threat that keeps growing to this day. The risk grows even larger in social media such as facebook, twitter, and google+. Hackers commonly take advantage these sites to attack people using them at their workplace, homes, or in public in order to take personal and security information that can affect the user or company. Phishing takes advantage of the trust that the user may have since the user may not be able to tell that the site being visited, or program being used, is not real; therefore, when this occurs, the hacker has the chance to gain the personal information of the targeted user, such as passwords, usernames, security codes, and credit card numbers, among other things.

The phishing is similar to phishing in a lake, but instead of trying to capture fish, phishers attempt to steal your personal information. They send out e-mails that appear to come from legitimate web sites such as eBay, PayPal, or other banking institutions. The emails state that your information needs to be updated or validated and ask that you enter your username and password, after clicking a link included in the e-mail. Some e-mails will ask that you enter even more information, such as your full name, address, phone number, social security number, and credit card number. However, even if you visit the false website and just enter your username and password, the phisher may be able to gain access to more information by just logging in to your account.

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 8, August 2015

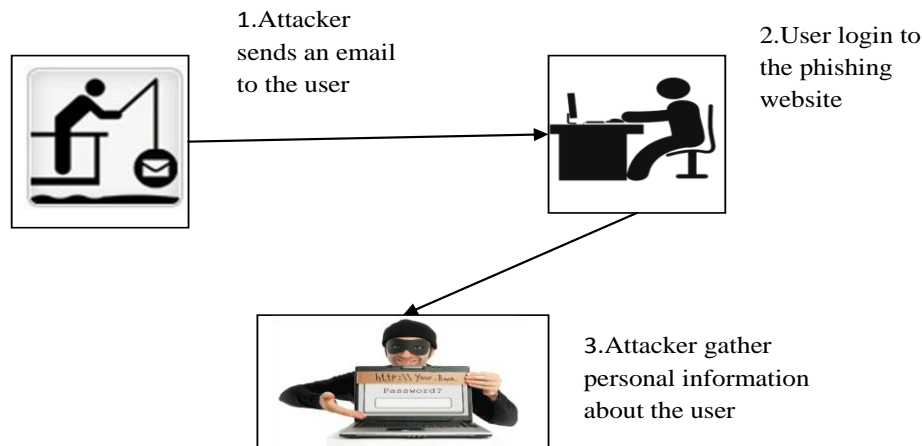


Fig 1. Phishing attack.

Fig 1 shows an example of phishing attack. Here the attacker created fake URL and sends an email to the user with this fake URL. If the user click on that URL then that user will go to the phishing website. Then the attacker can gather personal information about the user. The attacker aim is to gather sensitive information's like usernames, passwords and credit card details.

The main objective of this paper is to improve accuracy. In the proposed method content based approach and PhishStorm can be implementing in the web browser to improve accuracy. Here content based approach is based on the LDA. These methods using Page Rank, Alexa Ranking, @ representation in the URL, IPaddress using in the URL, age of the web sites, domain and sub domain to identify the phishing and legitimate web sites.

II. RELATED WORK

Many methods have been developed to prevent connection to the phishing websites. In [2] Phish Net is used as a phishing detection system and it contains two components. In the first component (URL prediction component), it proposes five heuristics to enumerate simple combinations of known phishing sites to discover new phishing URLs. The second component (approximate URL matching component) consists of an approximate matching algorithm that dissects a URL in to multiple components that are matched individually against entries in the blacklist. The following steps are used to identify phishing URLs [3]. First, we carefully select lexical features of the URLs that are resistant to obfuscation techniques used by attackers. Second, we evaluate the classification accuracy when using only lexical features, both automatically and hand-selected, vs. when using additional features. We show that lexical features are sufficient for all practical purposes. Third we thoroughly compare several classification algorithms, and we propose to use an online method (AROW) that is able to overcome noisy training data. Based on the insights gained from our analysis, we propose PhishDef, a phishing detection system that uses only URL names and combines the above three elements. PhishDef is a highly accurate method, lightweight, proactive and resilient to training data inaccuracies. Introduce PhishScore [4], an automated real-time phishing detection system. We observed that phishing URLs usually have few relationships between the part of the URL that must be registered and the remaining part of the URL. Hence, we define this concept as intra-URL relatedness and evaluate it using features extracted from words that compose a URL based on query data from Google and Yahoo search engines. These features are then used in machine learning based classification to detect phishing URLs from a real dataset. Introduce a novel parameter tuning framework based on a neuron-fuzzy [6] with six feature-sets, and identified different numbers of membership functions, different number of epochs, different sizes of feature-sets on a single platform. Parameter tuning based on neuron-fuzzy system with comprehensive features can enhance system performance in realtime. The outcome will provide guidance to the researchers who are using similar techniques in the field. It will decrease difficulties and increase confidence in the process of tuning parameters on a given problem. Propose a new technique [9] that leverages semantic and natural

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 8, August 2015

language processing tools in order to analyze large volumes of DNS data. The main research novelty consists in detecting malicious and dangerous domain names by evaluating the semantic similarity with already known names. This process can provide valuable information for reconstructing network and user activities. The efficiency of the method on experimental real datasets gathered from a national passive DNS system. Fast flux [8] is an evasion technique that cyber-criminals and internet miscreants use to evade identification and to frustrate law enforcement and anticrime efforts aimed at locating and shutting down web sites used for illegal purposes. Fast flux hosting supports a wide variety of cyber-crime activities and is considered one of the most serious threats to online activities today. One variant of fast flux hosting is double flux, exploits the domain name registration and name resolution services. Develop a methodology [7] to detect “domain fluxes” in DNS traffic by looking for patterns inherent to domain names that are generated algorithmically, in contrast to those generated by humans. Introduce PhishScore [5], an automated real-time phishing detection system. The phishing URLs usually have few relationships between the part of the URL that must be registered and the remaining part of the URL, Hence define this concept as intra-URL relatedness and evaluate it using features extracted from words that compose a URL based on query data from Google and Yahoo search engines. These features are then used in machine learning based classification to detect phishing URLs from a real dataset. Introduce PhishStorm [1], an automated phishing detection system that can analyze in real time any URL in order to identify potential phishing sites.

III. PROPOSED APPROACH

The attacker can able to gather user’s personal information’s like username, password, etc. PhishStorm is used to prevent this type of attacking. PhishStorm is an automated phishing detection system. In figure 2. PhishStorm and content based approach can be implement in to the web browser to improve accuracy. The content based approach is based on LDA (Latent Dirichlet Allocation).

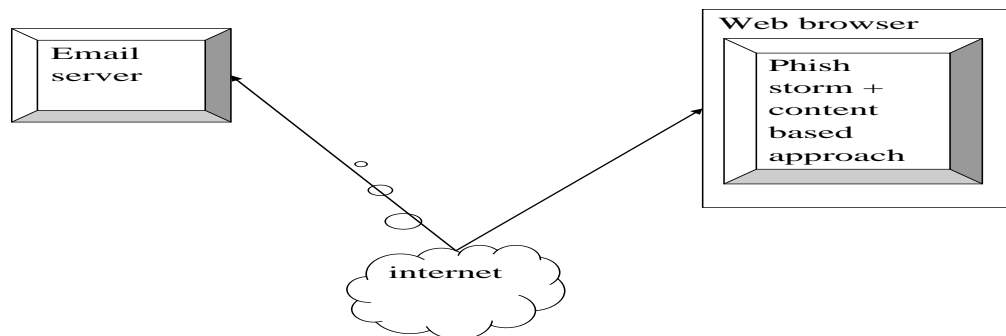


Figure 2. Web browser containing phishstorm and content based approach.

A. LDA:

LDA is a topic modelling technique and it is the most common topic modelling tool currently in use. It can discover the hidden topics in collections of documents using the words that appear in the documents.

B. Pattern Extraction:

Identify lexical patterns which represent semantic relationships between Phishing URL and legitimate URL.

C. Bloom Filter:

A Bloom filter is a space-efficient probabilistic data structure, that is used to test whether an element is a member of a set. Bloom proposed the technique for applications where the amount of source data would require an impracticably large hash area in memory if conventional error free hashing techniques were applied.

D. Rank Calculation:

The final rank can be calculated based on the PageRank, AlexaRank, IPaddress, @ symbol used in the URL, URL length, age of the web sites, and the domain name. The PageRank means it is a link analysis algorithm used by google

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 8, August 2015

and it is used to determine the relative importance of a website. Every web site will give Google PageRank score between 0 and 10 on an exponential scale. The AlexaRank is a ranking system by alexa.com and it is based on the amount of traffic recorded from the users. This traffic is based on such parameters as reach and page views. The reach refers to the number of users who visit a particular site in one day. Page view means it is the number of times a particular page (URL) is viewed by users. The first step of the ranking process is calculating the reach and number of page views for all the sites on the web on a daily basis. The Alexa ranking is obtained by performing the geometric mean of reach and page views, averaged over a predefined period of time. If the URL contains IPaddress and @ symbol, then it would be considered as phishing web site. The above methods are used to determine the phishing and legitimate web sites.

IV. PERFORMANCE EVALUATION

The attacker can gather personal information about the user. Here PhishStorm can be used to avoid this kind of attacking. PhishStorm and content matching can implemented in to the web browser to improve accuracy. If the user enter any URL in to the search space it will find whether this URL is phishing URL or not. If the URL is phishing URL then it will redirect in to the original web sites. This methods using ranking rate to find out whether this URL is Phishing URL or not. So this methods will find out the accurate results. The final ranking rate is calculate based on the PageRank, AlexaRank, age, @ presentation used in the web site, IPaddress used in the web site, domain and sub domain.

If the user enter the phishing web site like <http://www.paypal.shopping.uk/>, then the result shown in figure 3. Here PageRank and AlexaRank is 0, domainsim is 0.286, sub domainsim and age become 0, IPaddress is false that means IPaddress is not used in this URL, @ present is false because @ is not presented in this URL and finally large URL is true. So identified this URL is phishing URL.

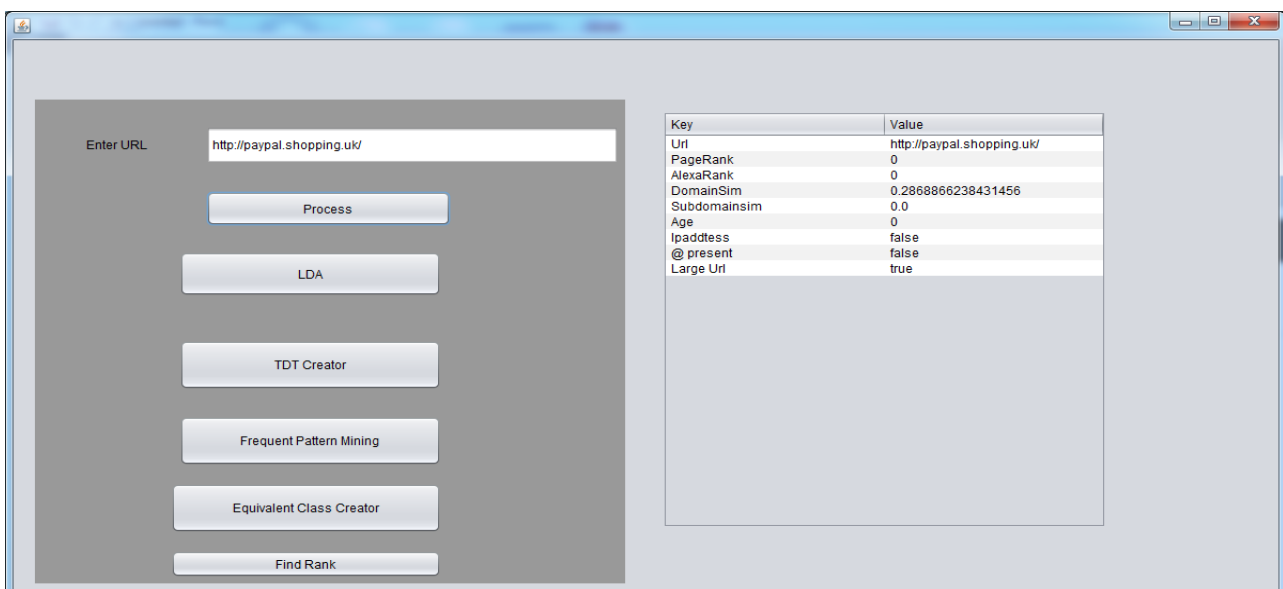


Figure 3. Details of the URL.

V. CONCLUSION

In this paper a content based approach and PhishStorm are introduced. These two approaches are implemented in to Mozilla Firefox; it is an open source web browser to improve the accuracy. PhishStorm is an automatic phishing detection system, this method using URL word extraction, feature computation and finally prediction to detect the phishing web sites. Content based approach using LDA, is a topic modeling technique and it is the most common topic modeling tool currently in use. It can discover the hidden topics in collections of documents using the words that appear



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 8, August 2015

in the documents. If the user enters any URL in to the search space it will accurately find weather this URL is phishing URL or not. If the URL is phishing URL then it will redirect in to the original web sites.

REFERENCES

1. Samuel Marchal, Jerome Francois, Radu State and Thomas Engel, "PhishStorm: Detecting phishing with streaming analytics",IEEE Trans.on network and service management,vol.11, no.4,december 2014.
2. P.Prakash, M. Kumar, R. Kompella, and M. Gupta, "PhishNet: Predictive blacklisting to detect phishing attacks," in Proc. IEEE INFOCOM, pp. 1–5, 2010.
3. G. A. Miller, "WordNet: A lexical database for English," Commun ACM, vol. 38, no. 11, pp. 39–41, Nov. 1995.
4. A. Le, A. Markopoulou, and M. Faloutsos, "PhishDef: URL names say it all," in Proc. IEEE INFOCOM, pp. 191–195, 2011.
5. S. Marchal, J. François, R. State, and T. Engel, "PhishScore: Hacking phishers' minds," in Proc. 10th Int. CNSM, pp. 46–54, 2014.
6. P. Barraclough, G. Sexton, M. Hossain, and N. Aslam, "Intelligent phishing detection parameter framework for e-banking transactions based on neuro-fuzzy," in Proc. SAI, pp. 545–555, 2014.
7. S. Yaday, A. K. K. Reddy, A. N. Reddy, and S. Ranjan, "Detecting algorithmically generated domain-flux attacks with DNS traffic analysis," IEEE/ACM Trans. Netw., vol. 20, no. 5, pp. 1663–1677, Oct. 2012.
8. ICANN Security and Stability Advisory Committee, Angeles, CA, USA, Tech. Rep. SAC 025, "SSAC advisory on fast flux hosting and DNS," 2008.
9. S. Marchal, J. Francois, R. State, and T. Engel, "Semantic based DNS forensics," in Proc. IEEE Int. WIFS, pp. 91–96, 2012.
10. Anti-Phishing Working Group, Lexington, MA, USA, Tech. Rep. 1H2014, "Global phishing survey: Trends and domain name use," 2014.

BIOGRAPHY

Pinky M S is an M.Tech student in the Computer Science Department, Mangalam College of Engineering, M G University. She received B.Tech degree in 2011 from Anna University, Trichy, Tamilnadu, India. Her interests are in Computer Networks (wireless Networks).