A Study on Phishing Awareness on Facebook Anonymous Applications

Tasmina Jahan Nishi (171-115-005),Sabbir Ahmed (163-115-019) tasminanishi@gmail.com, shawon.mu11@gmail.com
Metropolitan University Sylhet,Bangladesh

November 2019

Abstract

Social media has become a way of life for many people. Facebook is above all. Mass people share their information. Most of them are not aware of phishing website. We always see a warning when we leave facebook and visit another website. There are many unauthorized applications in facebook like when will you be in the next five years? Or when will you get married? Or who is your best friend? etc. When we enter in an application they get all data that we share on facebook. Even they get our email password. They can blackmail us. In this work we have conducted a study to detect phishing or facebook application, which are anonymous. Finally, we have focused on further scope at this study.

Keywords: Weka, Latex, Facebook, Privacy, Phishing.

1 Introduction

Social networking sites (SNSes), especially Facebook, have become an integral part of life for many people. Despite the complex use and plethora of information on these networks, a large fraction of users are lacking in security knowledge and awareness about how to navigate SNSes securely. On top of that, Some SNSes, Facebook in particular, have complicated systems of security and privacy settings due to their complex structure. Phishing attacks exploit human errors in online navigation. The attacker's goal is to either collect login credentials from the victim in order to gain access to their online accounts or have the victim visit a crafted malicious site with a drive-by download. Recently, there is also an increase in phishing attacks on SNSes using fake or compromised accounts. In SNSes, attackers can improve their chance of being clicked by creating targeted attack using information shared on the platform or using a link shortened (e.g. bitly.com) or specialized obfuscation services to disguise their malicious destinations. 2 although this requires more time and effort, it is generally more successful and harder to be detected by current defense systems and users. With over one billion active users, a successful attack in Facebook is worth the additional effort. To date, there has been little research into understanding the efficacy of attackers strategies in carrying out phishing attacks over SNS, which is important for understanding how to improve SNS defense mechanisms and user awareness. Our proposed study aims to fill this gap. In particular, we will investigate the importance of different aspects of a post (phishing or otherwise) that influence the user's decision of whether or not to click the link. This paper briefly describes the design of the study in relation to prior work.

2 Related work

We now discuss prior studies on phishing on Social Networking Sites and users' vulnerabilities. Phishing on Social Networking Sites: Dhamija et al. [1] showed a correlation between the success of the attacks and the low knowledge level of the users of the users as well as with the level of authenticity in the look and feel of the spoofed email and website. Chhabra et al. [2] discussed the rise of attacks using shortened URLs on Twitter. Shortened URL through third-party services such as bitly.com and owl.ly are widely used to reserve character space and provide memorable links for advertisement or personal use. However, attackers can use this service to misdirect their victims, fooling them by redirecting to a phishing website instead of the real one. According to the study, 89% of references on Twitter were reported to be inorganic. Vishwanath [3] pointed out the lack of statistics provided by Facebook regarding phishing attacks and fake account statistics, which makes it hard for researchers to conduct formal studies on the platform. They reported that approximately 1 in 10 Facebook accounts is a fake or a duplicate account. He stated that SNS has become a very attractive attack vector because of its continuing success. According to their study, attacks on Facebook have an approximately 40% success rate, compared to a success rate of just 1% for traditional email phishing. His findings indicate that attackers typically either post malicious links on a newsfeed, mimicking something of interest to the victims, or personally contact the victims through a private message. Alam et al. [4] noted that the success of targeted phishing is correlated with the amount of information the attacker has. Therefore, if an attacker is a friend with the victim or uses a compromised account of a friend of the victim, they will have little difficulty in fooling the victims without getting noticed. Since SNS users expose a lot of personal information through the site, particularly to their connections, the high success rates reported by Vishwanath may be considered unsurprising.

3 Data

In this experiment, we are using 550 data. Through experiment we will come to know how to detect a phishing website. We are using 31 attributes here. Those are given below:

- 1. Having_IPhaving_IP_Address: If an IP address is used as an alternative of the domain name in the URL, such as "http://125.98.3.123/fake.html", users can be sure that someone is trying to steal their personal information. Sometimes, the IP address is even transformed into hexadecimal code.
- 2. URLURL_Length: Phishers can use long URL to hide the doubtful part in the address bar. For example:

http://federmacedoadv.com.br/3f/aze/ab51e2e319e51502f416dbe46b773a5e/?-cmd=_home&dispatch=11004d58f5b74f8dc1e7c2e8dd4105e811004d58f5b74f8-dc1e7c2e8dd4105e8@phishing.website.html To ensure accuracy of our study, we calculated the length of URLs in the dataset and produced an average URL length. The results showed that if the length of the URL is greater than or equal 54 characters then the URL classified as phishing. By reviewing our dataset we were able to find 1220 URLs lengths equals to 54 or more which constitute 48.8% of the total dataset size. Rule: IF We have been able to update this feature rule by using a method based on frequency and thus improving upon its accuracy.

- 3. Shortining_Service4: URL shortening is a method on the "World Wide Web" in which a URL may be made considerably smaller in length and still lead to the required webpage. This is accomplished by means of an "HTTP Redirect" on a domain name that is short, which links to the webpage that has a long URL. For example, the URL "http://portal.hud.ac.uk/" can be shortened to "bit.ly/19DXSk4".
- 4. Having_At_Symbol: Using "@" symbol in the URL leads the browser to ignore everything preceding the "@" symbol and the real address often follows the "@" symbol.
- 5. Double_slash_redirecting: The existence of "//" within the URL path means that the user will be redirected to another website. An example of such URL's is: "http://www.legitimate.com//http://www.phishing.com".

We examin the location where the "//" appears. We find that if the URL starts with "HTTP", that means the "//" should appear in the sixth position. However, if the URL employs "HTTPS" then the "//" should appear in seventh position.

- 6. Prefix_Suffix: The dash symbol is rarely used in legitimate URLs. Phishers tend to add prefixes or suffixes separated by (-) to the domain name so that users feel that they are dealing with a legitimate webpage. For example http://www.Confirme-paypal.com/.
- 7. Having_Sub_Domain: Let us assume we have the following link: http://www.hud.ac.uk/students/. A domain name might include the country-code top-level domains (ccTLD), which in our example is "uk". The "ac" part is shorthand for "academic", the combined "ac.uk" is called a second-level domain (SLD) and "hud" is the actual name of the domain. To produce a rule for extracting this feature, we firstly have to omit the (www.) from the URL which is in fact a sub domain in itself. Then, we have to remove the (ccTLD).
- 8. SSLfinal_State: The existence of HTTPS is very important in giving the impression of website legitimacy, but this is clearly not enough. The authors in (Mohammad, Tabatha and McCluskey 2012) (Mohammad, Tabatha and McCluskey 2013) suggest checking the certificate assigned with HTTPS including the extent of the trust certificate issuer, and the certificate age. Certificate Authorities that are consistently listed among the top trustworthy names include: "GeoTrust, Go Daddy, Network Solutions, Thawed, Comoro, Duster and VeriSign". Furthermore, by testing out our datasets, we find that the minimum age of a reputable certificate is two years.
- 9. Domain_registeration_length: Based on the fact that a phishing website lives for a short period of time, we believe that trustworthy domains are regularly paid for several years in advance. In our dataset, we find that the longest fraudulent domains have been used for one year only.
- 10. Favicon: A favicon is a graphic image (icon) associated with a specific webpage. Many existing user agents such as graphical browsers and newsreaders show favicon as a visual reminder of the website identity in the address bar. If the favicon is loaded from a domain other than that shown

in the address bar, then the webpage is likely to be considered a Phishing Attempt?

- 11. Port: This feature is useful in validating if a particular service (e.g. HTTP) is up or down on a specific server. In the aim of controlling intrusions, it is much better to merely open ports that you need. Several firewalls, Proxy and Network Address Translation (NAT) servers will, by default, block all or most of the ports and only open the ones selected. If all ports are open, phishes can run almost any service they want and as a result, user information is threatened.
- 12.HTTPS_token: he phishes may add the "HTTPS" token to the domain part of a URL in order to trick users. For example, http://https-www-paypal-it-webapps-mpp-home.soft-hair.com/.
- 13. Request_URL: Request URL examines whether the external objects contained within a webpage such as images, videos and sounds are loaded from another domain. In legitimate WebPages, the webpage address and most of objects embedded within the webpage are sharing the same domain.
- 14. URL_of_Anchor: An anchor is an element defined by the tag. This feature is treated exactly as "Request URL". However, for this feature we examine: 1.If the tags and the website have different domain names. This is similar to request URL feature. 2. If the anchor does not link to any webpage, e.g.: A. B. C. D.
- 15. Links_in_tags: Given that our investigation covers all angles likely to be used in the webpage source code; we find that it is common for legitimate websites to use tags to offer metadata about the HTML document.
- 16. SFH: SFHs that contain an empty string or "about: blank" are considered doubtful because an action should be taken upon the submitted information. In addition, if the domain name in SFHs is different from the domain name of the webpage, this reveals that the webpage is suspicious because the submitted information is rarely handled by external domains.
- 17. Submitting_to_email: Web form allows a user to submit his personal information that is directed to a server for processing. A phishes might redirect

the user's information to his personal email. To that end, a server-side script language might be used such as "mail ()" function in PHP. One more client-side function that might be used for this purpose is the "mailto:" function.

- 18. Abnormal_URL: This feature can be extracted from WHOIS database. For a legitimate website, identity is typically part of its URL.
- 19. Redirect: The fine line that distinguishes phishing websites from legitimate ones is how many times a website has been redirected. In our dataset, we find that legitimate websites have been redirected one time max. On the other hand, phishing websites containing this feature have been redirected at least 4 times.
- 20. On_mouseover: Phishes may use JavaScript to show a fake URL in the status bar to users. To extract this feature, we must dig-out the webpage source code, particularly the "onMouseOver" event, and check if it makes any changes on the status bar.
- 21. RightClick: Phishes use JavaScript to disable the right-click function, so that users cannot view and save the webpage source code. This feature is treated exactly as "Using onMouseOver to hide the Link". Nonetheless, for this feature, we will search for event "event. Button==2" in the webpage source code and check if the right click is disabled.
- 22. PopUpWidnow: It is unusual to find a legitimate website asking users to submit their personal information through a pop-up window. On the other hand, this feature has been used in some legitimate websites and its main goal is to warn users about fraudulent activities or broadcast a welcome announcement, though no personal information was asked to be filled in through these pop-up windows.
- 23. Iframe: IFrame is an HTML tag used to display an additional webpage into one that is currently shown. Phishes can make use of the "iframe" tag and make it invisible i.e. without frame borders. In this regard, phishes make use of the "frame Border" attribute which causes the browser to render a visual delineation.
- 24. Age_of_domain: This feature can be extracted from WHOIS database (Who is 2005). Most phishing websites live for a short period of time. By

reviewing our dataset, we find that the minimum age of the legitimate domain is 6 months.

- 25. DNSRecord: For phishing websites, either the claimed identity is not recognized by the WHOIS database (Who is 2005) or no records founded for the hostname (Pan and Ding 2006). If the DNS record is empty or not found then the website is classified as "Phishing", otherwise it is classified as "Legitimate".
- 26. Web traffic: This feature measures the popularity of the website by determining the number of visitors and the number of pages they visit. However, since phishing websites live for a short period of time, they may not be recognized by the Alexia database (Alexia the Web Information Company., 1996). By reviewing our dataset, we find that in worst scenarios, legitimate websites ranked among the top 100,000. Furthermore, if the domain has no traffic or is not recognized by the Alexia database, it is classified as "Phishing". Otherwise, it is classified as "Suspicious".
- 27. Page Rank: Page Rank is a value ranging from "0" to "1". Page Rank aims to measure how important a webpage is on the Internet. The greater the Page Rank value the more important the webpage. In our datasets, we find that about 95% of phishing WebPages have no Page Rank. Moreover, we find that the remaining 5% of phishing WebPages may reach a Page Rank value up to "0.2".
- 28. Google_Index: This feature examines whether a website is in Google's index or not. When a site is indexed by Google, it is displayed on search results (Webmaster resources, 2014). Usually, phishing WebPages are merely accessible for a short period and as a result, many phishing WebPages may not be found on the Google index.
- 29. Links_pointing_to_page: The number of links pointing to the webpage indicates its legitimacy level, even if some links are of the same domain (Dean, 2014). In our datasets and due to its short life span, we find that 98% of phishing dataset items have no links pointing to them. On the other hand, legitimate websites have at least 2 external links pointing to them.
- 30.Statistical_report: Several parties such as Phish Tank (Phish Tank Stats,

2010-2012), and StopBadware (StopBadware, 2010-2012) formulate numerous statistical reports on phishing websites at every given period of time; some are monthly and others are quarterly. In our research, we used 2 forms of the top ten statistics from Phish Tank: "Top 10 Domains" and "Top 10 IPs" according to statistical-reports published in the last three years, starting in January2010 to November 2012. Whereas for "StopBadware", we used "Top 50" IP addresses.

31. Result: Here we see the website we get is either legitimate or phishing.

4 Experimental Result

To evaluate classifier quality, First we can use Logistic Regression. Consider this algorithm, for the dataset, we obtain 2x2 confusion matrix and correctly classified instances are 543 that is 90.65% and incorrectly classified instances are 56 that is 9.36% where cross-validation is 10 and the percentage split is 80%.

In the confusion matrix, the first column contains all the samples which our model think 'a'-239 in total. The second column contains all the samples which our model thinks 'b'-360 in total. The first row contains all the samples which really are 'a'-243 of them and second row contains all the samples which really are 'b'-356 of them. In the top-left, 213 are really 'a' and bottom-left 26 our model think 'a' but are really 'b' j- one kind error. In the top right, 30 are things that our model thinks are 'b' but which are really 'a' j- so no error and bottom-right 330 are things that our model thinks 'b' are really 'b'.

```
a b <-- classified as
213 30 | a = phising
26 330 | b = legitimate
```

Ai Paper"Photo 1.pdf"

--- Detailed Accuracy By Class ---

	0.877	0.091		0.006	80C Area 0.970 0.970	0.949	Class phising legitimate
Ai Paper" Photo 2.pdf" Weighted Avg.			0.907		0.970		

Now we can use another algorithm that is Voted Perception regression. Consider this algorithm, for the dataset, we obtain 2x2 confusion matrixes and correctly classified instances are 550 that is 91.81% and incorrectly classified instances are 49 that is 8.18% where cross-validation is 10 and the percentage split is 80%.

In the confusion matrix, the first column contains all the samples which our models think 'a'-230 in total. The second column contains all the samples which our models think 'b'-369 in total. The first row contains all the samples which really are 'a'-243 of them and second row contains all the samples which really are 'b'-356 of them. In the top-left, 212 are really 'a' and bottom-left 18 our model think 'a' but are really 'b'; one kind error. In the top right, 31 are things that our model thinks are 'b' but which are really 'a'; so no error and bottom-right 338 are things that our model thinks 'b' are really 'b'.

=== Confusion Matrix ===

```
a b <-- classified as
212 31 | a = phising
18 338 | b = legitimate</pre>
```

Ai Paper"Photo 3.pdf"

--- Detailed Accuracy By Class ---

	TP Bate	FP Rate	Precision	Recall.	F-Measure	MCC	ROC Area	PBC Area	Class
	0.072	0.051	0.922	0.072	0.096	0.030	0.924	0.074	phising
	0.949	0.128	0.916	0.949	0.932	0.830	0.937	0.926	legitimate
Weighted Avg.	0.918	0.096	0.918	0.918	0.918	0.830	0.932	0.905	

Ai Paper" Photo 4.pdf" - Confusion Matrix -

Now we can use another algorithm that is lazy. IBK regression. Consider this algorithm, for the dataset, we obtain 2x2 confusion matrix and correctly classified instances are 526 that is 87.13% and incorrectly classified instances are 73 that is 12.18% where cross-validation is 10 and the percentage split is 80%.

In the confusion matrix, the first column contains all the samples which our models think 'a'-244 in total. The second column contains all the samples which our models think 'b'-355 in total. The first row contains all the samples

which really are 'a'-243 of them and second row contains all the samples which really are 'b'-356 of them. In the top-left, 207 are really 'a' and bottom-left 37 our model think 'a' but are really 'b'; one kind error. In the top right, 36 are things that our model thinks are 'b' but which are really 'a'; so no error and bottom-right 319 are things that our model thinks 'b' are really 'b'.

=== Confusion Matrix ===

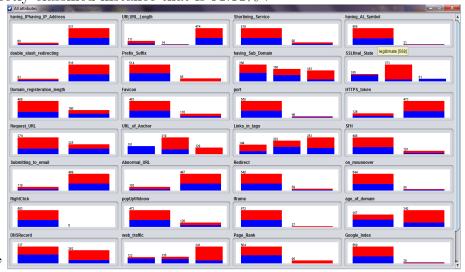
a b <-- classified as
207 36 | a = phising
37 319 | b = legitimate</pre>

Ai Paper"Photo 5.pdf"

Detailed Accuracy By Class ---Precision 0.848 0.852 0.850 0.918 phising 0.596 0.148 0.599 0.898 0.597 0.747 0.915 0.924 legitimate

Ai Paper"Photo 6.pdf"

we have used three classification algorithm here and Voted Perception has the best correctly classified instance that is 91.81%.



Ai Paper"scs.pdf"

5 Discussion

We have used three algorithms here. These are logistic, Voted Perception and lazy IBK. Where cross-validation is 10 and the percentage split is 80% In logistic correctly classified instance is 90.65% and incorrectly classified instance is 9.36% In Voted Perception correctly classified instance is 91.81% and incorrectly classified instance is 8.18% In lazy IBK correctly classified is 87.13% and incorrectly classified instance is 12.18% so we can say Voted Perception is the best one among these.

6 Future Scope

Nobody wants to fall prey to a phishing scam. There's a good reason that such scams will continue, though: They are successful enough for cybercriminals to make massive profits. Phishing scams have been around practically since the inception of the Internet, and they will not go away any time soon. Now we come to know how to detect a phishing website and we can avoid these type of websites in future.

7 References

[1] R. Dhamija, J. D. Tygar, and M. Hearst, "Why phishing works," in Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, ser. CHI '06. New York, NY, USA: ACM, 2006, pp. 581–590. [Online].

Available: http://doi.acm.org/10.1145/1124772.1124861

[2] S. Chhabra, A. Aggarwal, F. Benevenuto, and P. Kumaraguru, "Phi.sh/\$ocial: The phishing landscape through short urls," in Proceedings of the 8th Annual Collaboration, Electronic Messaging, Anti-Abuse and Spam Conference, ser. CEAS '11. New York, NY, USA: ACM, 2011, pp. 92–101. [Online]. Available: http://doi.acm.org/10.1145/2030376.2030387

[3] A. Vishwanath, "Habitual facebook use and its impact on getting deceived on social media," Journal of Computer-Mediated Communication, vol. 20, no. 1, pp. 83–98, 2015. [Online].

Available: http://dx.doi.org/10.1111/jcc4.12100

[4] S. Alam and K. El-Khatib, "Phishing susceptibility detection through social media analytics," in Proceedings of the 9th International Conference on Security of Information and Networks, ser. SIN '16. New York, NY, USA: ACM, 2016, pp. 61–64. [Online].

Available: http://doi.acm.org/10.1145/2947626.2947637

- [5] Tsikerdekis and S. Zeadally, "Online deception in social media," Commun. ACM, vol. 57, no. 9, pp. 72–80, Sep. 2014. [Online]. Available: http://doi.acm.org/10.1145/2629612
- [6] T. N. Jagatic, N. A. Johnson, M. Jakobsson, and F. Menczer, "Social phishing," Commun. ACM, vol. 50, no. 10, pp. 94–100, Oct. 2007. [Online].

Available: http://doi.acm.org/10.1145/1290958.1290968

- [7] Facebook, "Company info facebook newsroom," Dec 2016. [Online]. Available: http://newsroom.fb.com/company-info/
- [8] A. N. Joinson, "Looking at, looking up or keeping up with people?: Motives and use of facebook," in Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, ser. CHI '08. New York, NY, USA: ACM, 2008, pp. 1027–1036. [Online].

Available: http://doi.acm.org/10.1145/1357054.1357213

- [9] C. A. Lampe, N. Ellison, and C. Steinfield, "A familiar face(book): Profile elements as signals in an online social network," in Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, ser. CHI '07. New York, NY, USA: ACM, 2007, pp. 435–444. [Online]. Available: http://doi.acm.org/10.1145/1240624.1240695
- [10] S. Patil, "Will you be my friend?: Responses to friendship requests from strangers," in Proceedings of the 2012 iConference, ser. iConference '12. New York, NY, USA: ACM, 2012, pp. 634–635. [Online].

Available: http://doi.acm.org.ezproxy.rit.edu/10.1145/2132176.2132318

[11] Y. Boshmaf, I. Muslukhov, K. Beznosov, and M. Ripeanu, "The Social-bot network: When bots socialize for fame and money," in Proceedings of

the 27th annual computer security applications conference. ACM, 2011, pp. 93–102.