Phishing Websites Features

Rami M. Mohammad School of Computing and Engineering University of Huddersfield Huddersfield, UK. rami.mohammad@hud.ac.uk Fadi Thabtah E-Business Department Canadian University of Dubai Dubai, UAE. fadi@cud.ac.ae Lee McCluskey School of Computing and Engineering University of Huddersfield Huddersfield, UK. t.l.mccluskey@hud.ac.uk

1. Phishing Websites Features

One of the challenges faced by our research was the unavailability of reliable training datasets. In fact, this challenge faces any researcher in the field. However, although plenty of articles about predicting phishing websites using data mining techniques have been disseminated these days, no reliable training dataset has been published publically, maybe because there is no agreement in literature on the definitive features that characterize phishing websites, hence it is difficult to shape a dataset that covers all possible features.

In this article, we shed light on the important features that have proved to be sound and effective in predicting phishing websites. In addition, we proposed some new features, experimentally assign new rules to some well-known features and update some other features.

1.1. Address Bar based Features

1.1.1.Using the 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 as shown in the following link "http://0x58.0xCC.0xCA.0x62/2/paypal.ca/index.html".

<u>*Rule*</u>: IF { If The Domain Part has an IP Address \rightarrow Phishing Otherwise \rightarrow Legitimate

1.1.2.Long URL to Hide the Suspicious Part

Phishers can use long URL to hide the doubtful part in the address bar. For example:

 $http://federmacedoadv.com.br/3f/aze/ab51e2e319e51502f416dbe46b773a5e/?cmd=_home&dispatch=11004d58f5b74f8dc1e7c2e8dd4105e811004d58f5b74f8dc1e7c2e8dd4105e8@phishing.website. html$

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.

 $Rule: IF \begin{cases} URL \ length < 54 \rightarrow feature = Legitimate \\ else \ if \ URL \ length \ge 54 \ and \ \le 75 \rightarrow feature = Suspicious \\ otherwise \rightarrow feature = Phishing \end{cases}$

We have been able to update this feature rule by using a method based on frequency and thus improving upon its accuracy.

1.1.3.Using URL Shortening Services "TinyURL"

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".

 $\underline{\textit{Rule}}: IF \begin{cases} \text{TinyURL} \rightarrow \text{Phishing} \\ \text{Otherwise} \rightarrow \text{Legitimate} \end{cases}$

1.1.4.URL's having "@" Symbol

Using "@" symbol in the URL leads the browser to ignore everything preceding the "@" symbol and the real address often follows the "@" symbol.

$\label{eq:Rule:IF} \left\{ \begin{array}{l} \mbox{Url Having @ Symbol} \rightarrow \mbox{Phishing} \\ \mbox{Otherwise} \rightarrow \mbox{Legitimate} \end{array} \right.$

1.1.5.Redirecting using "//"

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.

$Rule: IF \begin{cases} The Position of the Last Occurrence of "//" in the URL > 7 \rightarrow Phishing \\ Otherwise \rightarrow Legitimate \end{cases}$

1.1.6. Adding Prefix or Suffix Separated by (-) to the Domain

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/.

$Rule: IF \left\{ \begin{matrix} \text{Domain Name Part Includes (-) Symbol} \rightarrow \text{Phishing} \\ \text{Otherwise} \rightarrow \text{Legitimate} \end{matrix} \right.$

1.1.7.Sub Domain and Multi Sub Domains

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) if it exists. Finally, we count the remaining dots. If the number of dots is greater than one, then the URL is classified as "Suspicious" since it has one sub domain. However, if the dots are greater than two, it is classified as "Phishing" since it will have multiple sub domains. Otherwise, if the URL has no sub domains, we will assign "Legitimate" to the feature.

 $\text{Rule: IF} \begin{cases} \text{Dots In Domain Part} = 1 \rightarrow \text{Legitimate} \\ \text{Dots In Domain Part} = 2 \rightarrow \text{Suspicious} \\ \text{Otherwise} \rightarrow \text{Phishing} \end{cases}$

1.1.8.HTTPS (Hyper Text Transfer Protocol with Secure Sockets Layer)

The existence of HTTPS is very important in giving the impression of website legitimacy, but this is clearly not enough. The authors in (Mohammad, Thabtah and McCluskey 2012) (Mohammad, Thabtah 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, GoDaddy, Network Solutions, Thawte, Comodo, Doster and VeriSign". Furthermore, by testing out our datasets, we find that the minimum age of a reputable certificate is two years.

1.1.9. Domain Registration 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.

1.1.10. 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.

$\label{eq:Rule:IF} \begin{array}{l} \mbox{Favicon Loaded From External Domain} \rightarrow \mbox{Phishing} \\ \mbox{Otherwise} \rightarrow \mbox{Legitimate} \end{array}$

1.1.11. Using Non-Standard 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, phishers can run almost any service they want and as a result, user information is threatened. The most important ports and their preferred status are shown in Table 2.

$Rule: IF \begin{cases} Port \ \# \ is \ of \ the \ Preffered \ Status \rightarrow \ Phishing \\ Otherwise \rightarrow \ Legitimate \end{cases}$

Table 1 Common ports to be checked

PORT	Service	Meaning	Preferred Status
21	FTP	Transfer files from one host to another	Close
22	SSH	Secure File Transfer Protocol	Close
23	Telnet	provide a bidirectional interactive text-oriented communication	Close
80	HTTP	Hyper test transfer protocol	Open
443	HTTPS	Hypertext transfer protocol secured	Open
445	SMB	Providing shared access to files, printers, serial ports	Close
1433	MSSQL	Store and retrieve data as requested by other software applications	Close
1521	ORACLE	Access oracle database from web.	Close
3306	MySQL	Access MySQL database from web.	Close
3389	Remote Desktop	allow remote access and remote collaboration	Close

1.1.12. The Existence of "HTTPS" Token in the Domain Part of the URL

The phishers 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/.

$Rule: IF \begin{cases} Using \text{ HTTP Token in Domain Part of The URL} \rightarrow \text{ Phishing} \\ \text{Otherwise} \rightarrow \text{ Legitimate} \end{cases}$

1.2. Abnormal Based Features

1.2.1. 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.

$\label{eq:Rule:IF} Rule: \mathrm{IF} \begin{cases} \mbox{ \% of Request URL } < 22\% \rightarrow \mbox{ Legitimate} \\ \mbox{ \% of Request URL } \geq 22\% \mbox{ and } 61\% \rightarrow \mbox{ Suspicious} \\ \mbox{ Otherwise} \rightarrow \mbox{ feature = Phishing} \end{cases}$

1.2.2. URL of Anchor

An anchor is an element defined by the $\langle a \rangle$ tag. This feature is treated exactly as "Request URL". However, for this feature we examine:

- 1. If the <a> 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.

```
\underbrace{\textit{Rule}}_{i}: \ \text{IF} \begin{cases} \text{\% of URL Of Anchor} < 31\% \rightarrow \textit{Legitimate} \\ \text{\% of URL Of Anchor} \geq 31\% \text{ And} \leq 67\% \rightarrow \text{Suspicious} \\ \text{Otherwise} \rightarrow \text{Phishing} \end{cases}
```

1.2.3. Links in <Meta>, <Script> and <Link> 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 <Meta> tags to offer metadata about the HTML document; <Script> tags to create a client side script; and <Link> tags to retrieve other web resources. It is expected that these tags are linked to the same domain of the webpage.

Rule:

```
% of Links in " < Meta > "," < Script > " and " < Link>" < 17% \rightarrow Legitimate IF % of Links in < Meta > "," < Script > " and " < Link>" \geq 17% And \leq 81% \rightarrow Suspicious Otherwise \rightarrow Phishing
```

1.2.4. Server Form Handler (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.

$\label{eq:Rule:IF} \text{Rule: IF} \left\{ \begin{array}{ll} \text{SFH is "about: blank" Or Is Empty} \rightarrow \text{Phishing} \\ \text{SFH Refers To A Different Domain} \rightarrow \text{Suspicious} \\ \text{Otherwise} \rightarrow \text{Legitimate} \end{array} \right.$

1.2.5. Submitting Information to Email

Web form allows a user to submit his personal information that is directed to a server for processing. A phisher 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.

Rule: IF $\{ Using "mail()" \text{ or "mailto:" Function to Submit User Information } \rightarrow Phishing Otherwise \rightarrow Legitimate \}$

1.2.6. Abnormal URL

This feature can be extracted from WHOIS database. For a legitimate website, identity is typically part of its URL.

$Rule: IF \left\{ \begin{matrix} \text{The Host Name Is Not Included In URL} \rightarrow \text{Phishing} \\ \text{Otherwise} \rightarrow \text{Legitimate} \end{matrix} \right.$

1.3. HTML and JavaScript based Features

1.3.1. Website Forwarding

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.

1.3.2. Status Bar Customization

Phishers 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.

$\label{eq:Rule:IF} \begin{array}{l} \mbox{(onMouseOver Changes Status Bar} \rightarrow \mbox{Phishing} \\ \mbox{It Does't Change Status Bar} \rightarrow \mbox{Legitimate} \end{array}$

1.3.3. Disabling Right Click

Phishers 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.

 $Rule: IF \begin{cases} \mbox{Right Click Disabled} & \rightarrow \mbox{Phishing} \\ \mbox{Otherwise} & \rightarrow \mbox{Legitimate} \end{cases}$

1.3.4. Using Pop-up Window

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.

$Rule: IF \left\{ \begin{matrix} \text{Popoup Window Contains Text Fields} \rightarrow \textbf{Phishing} \\ \textbf{Otherwise} \rightarrow \textbf{Legitimate} \end{matrix} \right.$

1.3.5. IFrame Redirection

IFrame is an HTML tag used to display an additional webpage into one that is currently shown. Phishers can make use of the "iframe" tag and make it invisible i.e. without frame borders. In this regard, phishers make use of the "frameBorder" attribute which causes the browser to render a visual delineation.

$Rule: IF \begin{cases} \textbf{Using iframe} \rightarrow \textbf{Phishing} \\ \textbf{Otherwise} \rightarrow \textbf{Legitimate} \end{cases}$

1.4. Domain based Features

1.4.1. Age of Domain

This feature can be extracted from WHOIS database (Whois 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.

$\label{eq:Rule:IF} \left\{ \begin{array}{ll} \mbox{Age Of Domain} \geq \mbox{6 months} \ \rightarrow \ \mbox{Legitimate} \\ \mbox{Otherwise} \ \rightarrow \ \mbox{Phishing} \end{array} \right.$

1.4.2. DNS Record

For phishing websites, either the claimed identity is not recognized by the WHOIS database (Whois 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".

Rule: IF{ no DNS Record For The Domain → Phishing Otherwise → Legitimate

1.4.3. Website 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 Alexa database (Alexa 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 Alexa database, it is classified as "Phishing". Otherwise, it is classified as "Suspicious".

1.4.4. PageRank

PageRank is a value ranging from "0" to "1". PageRank aims to measure how important a webpage is on the Internet. The greater the PageRank value the more important the webpage. In our datasets, we find that about 95% of phishing webpages have no PageRank. Moreover, we find that the remaining 5% of phishing webpages may reach a PageRank value up to "0.2".

1.4.5. 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.

1.4.6. Number of 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.

$Rule: IF \begin{cases} Of Link Pointing to The Webpage = 0 \rightarrow Phishing \\ Of Link Pointing to The Webpage > 0 and \le 2 \rightarrow Suspicious \\ Otherwise \rightarrow Legitimate \end{cases}$

1.4.7. Statistical-Reports Based Feature

Several parties such as PhishTank (PhishTank 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 PhishTank: "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.