

Detection of Hidden Fraudulent URLs within Trusted Sites using Lexical Features

Enrico Sorio Alberto Bartoli Eric Medvet

DIA - University of Trieste



September 5th, 2013

<http://machinelearning.inginf.units.it>

Table of Contents

- 1 Scenario and motivation
- 2 Our proposal
 - Requirements
 - Variants
- 3 Experimental evaluation
 - Dataset
 - Results



Security threats to web sites content

Means of attack:

- Fraudulent changes to existing pages
- Fraudulent creation of new *hidden* pages

Attacker's goals:

- defacement
- search/web spam
- malware spreading



Security threats to web sites content

Means of attack:

- Fraudulent changes to existing pages
- **Fraudulent creation of new *hidden* pages**

Attacker's goals:

- defacement
- search/web spam
- malware spreading



Fraudulent hidden pages

How long to actually detect them?

- 50% of defacements (non-hidden fraudulent pages) are detected later than 1 week after the incident (study says¹)

¹Bartoli, A.; Davanzo, G.; Medvet, E., *The Reaction Time to Web Site Defacements*, Internet Computing, IEEE , vol.13, no.4, pp.52,58, July-Aug. 2009

Fraudulent hidden pages

How long to actually detect them?

- 50% of defacements (non-hidden fraudulent pages) are detected later than 1 week after the incident (study says¹)
- no studies for hidden pages, but likely longer

¹Bartoli, A.; Davanzo, G.; Medvet, E., *The Reaction Time to Web Site Defacements*, Internet Computing, IEEE , vol.13, no.4, pp.52,58, July-Aug. 2009

Fraudulent hidden pages

How many of them?

- > 9% Italian PA domains contain fraudulent hidden pages (study says²)

²Sorio, E.; Bartoli, A.; Medvet, E., *A look at hidden web pages in Italian public administrations*, Computational Aspects of Social Networks (CASoN), 2012 Fourth International Conference on , vol., no., pp.291,296, 21-23 Nov. 2012



Fraudulent hidden pages

How many of them?

- > 9% Italian PA domains contain fraudulent hidden pages (study says²)
- a significant fraction of Zone-H entries are fraudulent hidden pages

²Sorio, E.; Bartoli, A.; Medvet, E., *A look at hidden web pages in Italian public administrations*, Computational Aspects of Social Networks (CASoN), 2012 Fourth International Conference on , vol., no., pp.291,296, 21-23 Nov. 2012



Fraudulent hidden pages

Are they risky?

- actually used for illicit activities, e.g., illicit drug trade



Fraudulent hidden pages

Are they risky?

- actually used for illicit activities, e.g., illicit drug trade
- HTTPS does not help: the whole server content is authenticated (by definition) → the fraudulent hidden page may live within a **trusted** site (phishing!)



Fraudulent hidden pages

Are they risky?

- actually used for illicit activities, e.g., illicit drug trade
- HTTPS does not help: the whole server content is authenticated (by definition) → the fraudulent hidden page may live within a **trusted** site (phishing!)

The screenshot shows a browser window with the title "Error - Login - PayPal". The address bar contains the URL <https://www.polisiohor.gov.my/templates/bee>. The page content includes the PayPal logo, a navigation menu with "Home", "Personal", and "Business" tabs, and a login form with fields for "Email address" and "PayPal password". A yellow warning box with a triangle icon contains the text: "Please make sure you enter your email address and password correctly. If you still can't Troubleshooting Tips below." Below the login form is a "Troubleshooting Tips" section with links for "Forgot your password?" and "Forgot your email address?".

Malaysian police

HTTPS!

PayPal phishing



Table of Contents

- 1 Scenario and motivation
- 2 Our proposal
 - Requirements
 - Variants
- 3 Experimental evaluation
 - Dataset
 - Results



Detection: requirements

Detecting if a page at URL u is **fraudulent and hidden**:

- without fetching the page
- without using the *domain* part of u

(~~http://foo:bar@www.here.com/path/to/there.htm?p=1#title~~)



Detection: requirements

Without using the *domain* part of u

- most existing detectors use the domain as key feature \Rightarrow URLs at domain d are either deemed *all* legitimate or *all* illegitimate
- we exclude the domain and evaluate approach on both legitimate and illegitimate on the *same* domain



Method variants

Three variants:

- lexical features
- lexical features and headers
- lexical features, headers and age



Method variants

Three variants:

- | | |
|-------------------------------------|--------------------------------|
| • lexical features | few features, low complexity |
| • lexical features and headers | ↓ |
| • lexical features, headers and age | many features, high complexity |



Variants complexity

- Lexical → no HTTP requests
- Lexical+headers → one HEAD request
- Lexical+headers+age → two HEAD requests



Variants complexity

- Lexical → no HTTP requests **can work offline!**
- Lexical+headers → one HEAD request
- Lexical+headers+age → two HEAD requests



Method phases

Two phases:

- 1 calibration, using a *training set* U of URLs and corresponding labels L
- 2 actual classification



Lexical

Calibration:

- 1 remove all up to domain from URLs
(~~http://foo:bar@www.here.com/path/to/there.htm?p=1#title~~)
- 2 compute unigrams ($u_i \in U \rightarrow f_i \in \mathbb{R}^n$)
- 3 train a SVM on $\langle f_i, l_i \rangle$



Lexical

Calibration:

- 1 remove all up to domain from URLs
(~~http://foo:bar@www.here.com~~/path/to/there.htm?p=1#title)
- 2 compute unigrams ($u_i \in U \rightarrow f_i \in \mathbb{R}^n$)
- 3 train a SVM on $\langle f_i, l_i \rangle$

URL	Lexical				/
	a	b	c	...	
/p/attc.txt?ac	2	0	2	...	+
/bb/user/car/38	1	2	1	...	-
/xXx.htm	0	0	0	...	+



Lexical

Classification of u :

- 1 remove all up to domain from u
- 2 compute unigrams ($u \rightarrow f$)
- 3 apply SVM to f



Lexical+headers

As above, but with more features:

- 1 obtain by HTTP HEAD header values of
 - Server (only major+minor version)
 - X-Powered-By (only major+minor version)
 - Content-Type
 - Content-Length
- 2 map (all but Content-Length) as binary features



Lexical+headers

As above, but with more features:

- 1 obtain by HTTP HEAD header values of
 - Server (only major+minor version)
 - X-Powered-By (only major+minor version)
 - Content-Type
 - Content-Length
- 2 map (all but Content-Length) as binary features

URL	Lexical				Server				Content-Type			...	/	
	a	b	c	...	Apache/2.2	Apache/2.4	IIS/7.2	...	text/plain	text/html	text/css	...		
c/p/attc.txt?ac	2	0	2	...	0	1	0	...	1	0	0	+
/bb/user/car/38	1	2	1	...	1	0	0	...	0	1	0	-
/xXx.htm	0	0	0	...	1	0	0	...	0	1	0	+



Lexical+headers+age

As above, but with one more feature (the *relative age*):

- 1 obtain the Last-Modified header value of u
- 2 obtain the Last-Modified header value of *home page* of u
(<http://www.carexperts.it/bb/user/car/38> → <http://www.carexperts.it>)
- 3 use the difference as relative age



Table of Contents

- 1 Scenario and motivation
- 2 Our proposal
 - Requirements
 - Variants
- 3 Experimental evaluation
 - Dataset
 - Results



Dataset

Evaluating a detector of hidden fraudulent pages

- no public dataset available
- we built one (and made publicly available³)

³<http://machinelearning.inginf.units.it/data-and-tools/hidden-fraudulent-urls-dataset>

Positives and negatives

Positive examples (+):

- hidden fraudulent pages

Negative examples (−):

- “normal” pages on “normal” sites
- “normal” pages on attacked sites



Positive examples (+)

Hidden fraudulent pages:

- validated phishing attacks from Phishtank⁴ ($U_+^P \rightarrow 6564$)
- validated defacement attacks from Zone-H⁵ ($U_+^D \rightarrow 2144$)

⁴<http://www.phishtank.com>

⁵<http://www.zone-h.org>

Positive examples (+)

Hidden fraudulent pages:

- validated phishing attacks from Phishtank⁴ ($U_+^P \rightarrow 6564$)
- validated defacement attacks from Zone-H⁵ ($U_+^D \rightarrow 2144$)

Hidden \rightarrow not reachable while crawling the site within 3rd level

⁴<http://www.phishtank.com>

⁵<http://www.zone-h.org>

Negative examples (–)

“Normal” pages on “normal” sites:

- 20 among 500 Alexa top sites, crawled up to 10th level ($U_- \rightarrow 3713$)

“Normal” pages on attacked sites:

- sites of positives phishing attacks, crawled up to 3rd level ($U_-^P \rightarrow 78388$)
- sites of positives defacements attacks, crawled up to 3rd level ($U_-^D \rightarrow 94370$)



Experiments procedure

Two suites (phishing, defacements)

- balanced training set and testing set
- 90% training, 10% testing
- 5 repetitions



Results

Phishing:

Method	Accuracy (%)			FNR on U_+^P (%)			FPR on U_-^P (%)			FPR on U_- (%)		
	Avg.	Dev.	Std.	Avg.	Dev.	Std.	Avg.	Dev.	Std.	Avg.	Dev.	Std.
Lexical	92.50	0.62		9.80	0.85		5.20	0.84		4.93	1.21	
Lexical+headers	95.34	0.48		4.93	1.05		4.38	0.62		0.79	0.70	
Lexical+headers+age	95.57	0.37		4.87	1.06		3.98	0.52		0.73	0.73	

Defacements:

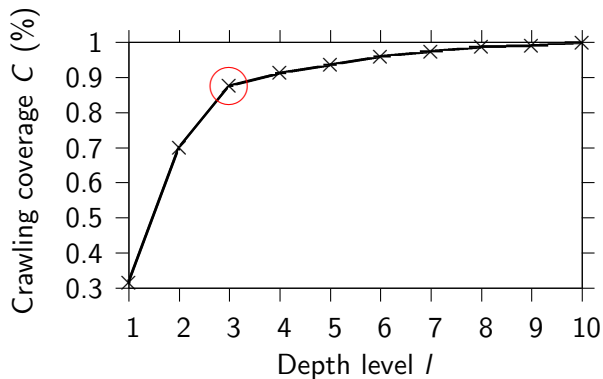
Method	Accuracy (%)			FNR on U_+^D (%)			FPR on U_-^D (%)			FPR on U_- (%)		
	Avg.	Dev.	Std.	Avg.	Dev.	Std.	Avg.	Dev.	Std.	Avg.	Dev.	Std.
Lexical	98.37	0.68		0.93	0.74		2.32	0.87		2.41	1.92	
Lexical+headers	99.35	0.3		0.37	0.61		0.93	0.57		0.93	0.65	
Lexical+headers+age	99.26	0.3		0.47	0.66		1.02	0.61		0.93	0.65	

Thanks!



Hidden flag: validation of the assumption

“Hidden → page not reachable within 3rd level of crawling”



- 20 among 500 Alexa top sites
- 88% of non-hidden pages are reachable within 3rd level



Positives and negatives

Positive examples (+):

- hidden fraudulent pages

Negative examples (−):

- non-hidden fraudulent pages
- non-hidden non-fraudulent pages



Positives and negatives

Positive examples (+):

- hidden fraudulent pages

Negative examples (−):

- non-hidden fraudulent pages
- non-hidden non-fraudulent pages

		Hidden flag	
		T	F
Fraudulent flag	T	+	−
	F		−

Fraudulent flag

T → fraudulent pages

- validated phishing attacks from Phishtank⁶
- validated defacement attacks from Zone-H⁷

F → legitimate pages

- 20 among 500 Alexa top sites, crawled up to 10th level
- Phishtank phished sites, crawled up to 3rd level
- Zone-H defaced attacks sites, crawled up to 3rd level

⁶<http://www.phishtank.com>

⁷<http://www.zone-h.org>

Hidden flag

T → hidden pages

- not reachable while crawling the site

F → non-hidden pages

- reachable while crawling the site



Hidden flag

T → hidden pages

- not reachable while crawling the site within 3rd level

F → non-hidden pages

- reachable while crawling the site within 3rd level



Summary

U_+^P , hidden phishing \rightarrow 6564

U_-^P , non-hidden phishing \rightarrow 78388

U_+^D , hidden defacements \rightarrow 2144

U_-^D , non-hidden defacements \rightarrow 94370

U_- , Alexa \rightarrow 3713

Summary

U_+^P , hidden phishing \rightarrow 6564
 U_-^P , non-hidden phishing \rightarrow 78388
 U_+^D , hidden defacements \rightarrow 2144
 U_-^D , non-hidden defacements \rightarrow 94370
 U_- , Alexa \rightarrow 3713

		Hidden flag	
		T	F
Fraudulent flag	T	U_+^P, U_+^D	U_-^P, U_-^D
	F		U_-