

Leveraging Google SafeBrowsing to Characterize Web-based Attacks

Peter Likarish, Eunjin Jung
Dept. of Computer Science
The University of Iowa
Iowa City, IA 52242
{plikaris, ejjung}@cs.uiowa.edu

Abstract

As the World Wide Web expands, it also hosts an increasing number of attacks designed to install malicious software on users' computers. These attacks are often transitory, as soon as they are detected the attack can be moved to a new location and repeated. As a result, it is difficult to study attacks in an online manner. In this paper, we present preliminary results that make use of tools provided by Google to determine connectivity among the attack sites, especially the sites that redirect from one to another.

1. Introduction

Malicious software used to spread primarily through email attachments, automated worms and P2P downloads. However, the ubiquitous nature of web browsers and the proliferation of 3rd party software provides malware distributors with a large number of attack vectors, including: drive-by download sites, fake codec installation requests, malicious advertisements and spam messages on blogs or social network sites. According to WebSense's 2008 report, the number of websites that host malware increased 46% between January 1st, 2008 and January 1st, 2009 and 77% of these sites were legitimate sites that have been compromised [6]. Furthermore, Symantec has found that malware authors target mainstream websites for better exposure, while obfuscation techniques reduce the chance that security auditors will detect their exploits [4].

Any web-based attack detection based on blacklists of attack-hosting domains provide only limited protection. When an attack is detected, it is easy for the attackers to select a new vulnerable website or to host their attack at a

newly registered domain. We followed the domains that are known to host malicious web-based attacks and monitored redirection relationships among the domains. Being able to understand the relationship between domains that host web-based attacks could aid security experts in identifying and detecting attacks and also possibly allow experts to predict which domains are at a greater risk of hosting an attack and to apply proactive measures.

1.1 Leveraging Google's Resources

While it is possible to crawl the internet in order to detect attacks as they are happening, this tends to be expensive in terms of the time and resources it requires. Additionally, the crawler must be fairly sophisticated, incorporating a honeyclient or other means of identifying attacks. Amongst others, Moshchuk et al [2] and Wang et al [5] have both devised systems to automatically identify web exploits.

Even then, attackers avoid basic automated collection methods by obfuscating their exploit code, utilizing multiple redirections and by requiring "human proofs" before installing their exploits. As an alternative, we attempt to offload the work of discovering malicious code to Google by utilizing their free SafeBrowsing Diagnostic tool in order to capture a snapshot of the malicious attacks Google has detected.

In this paper, we first describe the information available through Google's SafeBrowsing Diagnostic tool. We then describe the information we gained as a result of querying the tool for a total of 3,465 domains that had been observed hosting malicious software in the last 90 days (the maximum length of time for which Google's SafeBrowsing Diagnostic tool maintains records). One of our goals is to attempt to recreate *attack-domain graphs* based on information gathered from our queries. An attack-domain graph is where a node is a domain that either redirects to attack-hosting domain or hosts web-based attacks, and a directed edge is redirection from one domain to another. In general, attack-domain graphs are a valuable mechanism for identifying vulnerabilities. Shahriari et al provides one such example of the usefulness of applying attack-domain graphs at the client/server level [3] whereas this work is concerned with capturing redirections leading to malware infections.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Copyright 200X ACM X-XXXXX-XX-X/XX/XX ...\$5.00.

We conclude by discussing possible future directions for our work.

2. Google’s SafeBrowsing Diagnostics

The SafeBrowsing Diagnostic tool is available through a simple HTTP GET request at: <http://www.google.com/safebrowsing/diagnostic?site=xxx>. After supplying a domain as the site variable, Google returns a summary of the malicious activity (or lack thereof) that Google has observed occurring at that domain. This page is intended to provide webmasters and interested parties with pertinent information about any malicious activity observed at the domain. Some of this information includes:

1. Is the website currently listed as suspicious in Google’s search results?
2. How many pages on this domain Google has visited in the last 90 days and how many resulted in malicious software being downloaded/installed without consent.
3. The dates on which Google last visited the domain and the last time malicious content was found.
4. The network(s) on which the domain has been hosted.
5. If the website has hosted malware in the last 90 days.
6. The number of domains it has infected.
7. How many exploits have been found on this domain?
8. If the domain acted as an intermediary for malware distribution or if a website has redirected to this domain.

In addition, Google provides up to three additional domains that have been infected by this domain, three intermediary domains that redirected to this domain and three domains to which this website has been observed redirecting to. Since there is a limit on the number of domains returned per query, we need to know a sufficient number of domains in a particular attack in order to observe the entirety of the attack.

3. Querying with Malicious Domains

In order to determine the structure of web attacks that have been captured by Google’s tool, we needed a set of seed sites that had been observed hosting malicious content. Malware Domains is a domain-based blacklisting project that attempts to curtail the spread of malware by providing a free list of domains that have been observed engaging in malicious activity [1]. We opted to use Malware Domains because it compiles its blacklist from a variety of reputable sources. From the Malware Domains blacklist, we select domains added during the past 90 days, the time period for which Google’s SafeBrowsing Diagnostic tool provides data. This yielded 3,465 domain names.

We wrote a python script to query Google’s SafeBrowsing Diagnostic tool, at set intervals, until we had retrieved

pages for these domains. In addition, every time we came across a new domain name on one of the retrieved pages, we retrieved the results for that page as well. When finished, we had retrieved information for 4,209 domains.

4. Characteristics of Observed Attacks Online

In Table 1, we present a summary of attack characteristics we were able to extract using the data gathered from Google’s SafeBrowsing Diagnostic tool and which are discussed in greater detail in this section.

Description	
Domains from Malware Domains blacklist	3,465
Total domains queried on <i>GSBDt</i>	4,209
Percentage of domains on <i>GSBDt</i> with no observed malicious activity	1,529 (44%)
Domains with malicious activity on <i>GSBDt</i> not on Malware Domains blacklist	620% (32%)
Number of disjoint attacks detected	1063
Average number of domains per attack	5.37
Number of singleton attack domains	345
Average number of domains per attack without singletons	7.47

Table 1. Basic characteristics of attacks gathered from Google’s SafeBrowsing Diagnostic tool (*GSBDt*).

4.1 Blacklist Coverage

We first investigate the *coverage* each blacklist achieves by evaluating what percentage of domains were solely present one of the two lists. This conception of coverage is relative between the two lists, as opposed to an estimation of the total number of web-based attacks the lists capture which would be more informative. Of the 4,209 domains queried on Google’s SafeBrowsing Diagnostic tool only 1,936 have been observed conducting some malicious activity (hosting malware, redirecting users to malware domains, or infecting other domains). This means that of the 3,465 domains blacklisted by Malware Domains in the past 90 days, in the case of 1,529 domains (44%) Google had not observed any kind of malicious activity.

Scanning the full Malware Domains blacklist, we discovered that 620 (32%) of the 1,936 domains Google had observed conducting malicious activity were never blacklisted on Malware Domain’s blacklist.

These results provide strong evidence to support the claim that no organization has a monopoly on knowledge of malicious activity online, reinforcing our earlier claim that observing such attacks is a difficult problem in and of itself. This also suggests that blacklisting may need to be supplemented with proactive detection methods in order to protect people before attacks are blacklisted.

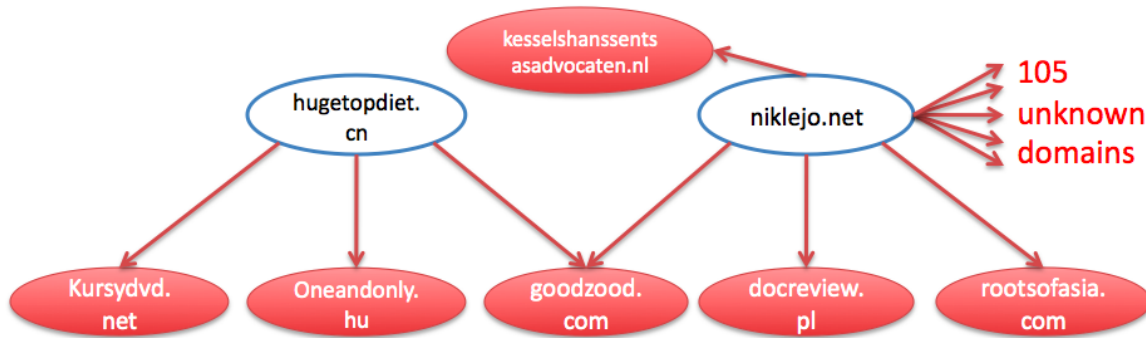


Figure 1. Example attack-domain graph

4.2 Domains and Attack Size

In order to provide a rough estimate of attack size, with regard to the number of domains involved in an attack, we kept track of which domains were associated with one another according to Google’s tool. According to this measure, we collected 1,063 attacks with an average size of 5 domains per attack. The number of attacks seemed suspiciously high to us and the number of domains per attack low. We discovered that we had 345 singleton domains which were not associated with any others. This is likely due to the fact that we did not provide a broad enough set of initial seeds to connect a large number of attacks and suspect we are overestimating the number of attacks and underestimating their average size.

When we removed the 345 singleton domains, leaving us with 718 attacks, and recalculated the average size we found an average of 7.47 domains per attack. This suggests that as future work we need a way be sure we are capturing entire attacks and to omit attacks we have capture partially when querying Google’s SafeBrowsing Diagnostic tool. We expect that when we have successfully done so, we will see a much larger average size and a much more connected attack structure than the current disjoint graphs (forrest).

5. Conclusion and Future Work

To the best of our knowledge, we are the first researchers to propose using Google’s SafeBrowsing Diagnostic tool to aid in reconstructing web-based attack graphs. As future work we intend to aim for better coverage of captured attacks by doing more extensive querying of the Google’s SafeBrowsing Diagnostic tool in order to refine the numbers presented here. We also intend to spend time examining the structure of links between domains hosting malicious content and to compare the attack structure to the underlying links of the legitimate websites when they are not under attack. We present the attack-domain graph for one attack in Figure 1 as an example of the types of graphs we hope to generate on a large scale. This is generated from partial knowledge of the attack topology. The clear circles represent domains redirecting users to filled circles, domains hosting malware. While we’ve capture all links from the leftmost blue circle, the other redirects to an as yet unde-

termined number of domains. We hope to devise a means of identifying structures that are likely the result of an attack.

We are also interested in how this work might interface with high-interaction honeyclients, for instance, by altering their crawling behavior.

6. References

- [1] DNS-BH. Malware domains. <http://www.malwaredomains.com/>.
- [2] A Moshchuk, T Bragin, S Gribble, and H Levy. A crawler-based study of spyware on the web. *Proceedings of the 2006 Network and Distributed System Security Symposium*, Jan 2006.
- [3] H Shahriari and R Jalili. Vulnerability take grant (vtg): an efficient approach to analyze network vulnerabilities. *Computers & Security*, Jan 2007.
- [4] Symantec. Web based attacks. http://www.symantec.com/content/en/us/enterprise/media/security_response/whitepapers/web_based_attacks_02-2009.pdf, 2009.
- [5] Y Wang, D Beck, X Jiang, and R Roussev. Automated web patrol with strider honeymonkeys: Finding web sites that exploit browser *Proceedings of the 13th Annual Network and . . .*, Jan 2006.
- [6] WebSense Security Labs. State of internet security, q3 – q4, 2008. http://securitylabs.websense.com/content/Assets/WSL_ReportQ3Q4FNL.PDF, 2008.