

Detection and Analysis of Drive-by Downloads and Malicious Websites

Saeed Ibrahim¹, Nawwaf Al Herami¹, Ebrahim Al Naqbi¹, and Monther Aldwairi¹

College of Technological Innovation
Zayed University, Abu Dhabi, UAE 144534
{m80007514,m80006805,m80006809,monther.aldwairi}@zu.ac.ae

Abstract. A drive-by download is a download that occurs without users action or knowledge. It usually triggers an exploit of vulnerability in a browser to download an unknown file. The malicious program in the downloaded file installs itself on the victims machine. Moreover, the downloaded file can be camouflaged as an installer that would further install malicious software. Drive-by downloads is a very good example of the exponential increase in malicious activity over the Internet and how it affects the daily use of the web. In this paper, we try to address the problem caused by drive-by downloads from different standpoints. We provide in-depth understanding of the difficulties in dealing with drive-by downloads and suggest appropriate solutions. We propose machine learning and feature selection solutions to remedy the drive-by download problem. Experimental results reported 98.2% precision, 98.2% F-Measure and 97.2% ROC area.

Keywords: drive by download, malware detection, web security

1 Introduction

Miscreants make use of malicious web content to perform attacks targeting web clients. Drive-by downloads (DBD) are unintentional downloads of malware or virus on to a mobile device or a computer. Due to the increased population of several web applications, DBD have become one of the most common malware spreading methods, thereby leading the security threats to cyber community. According to [29], query search results from Google contain more than 1.3% of the web pages that do DBD attacks. These downloads are located on normal-looking, but malicious websites [4]. They exploit vulnerabilities in out-of-date apps, browsers, plugins, or operating systems. Over the years, hackers have become much more sophisticated that just opening such web page could allow malicious code to be installed on the device without the knowledge and consent of the user. Downloaded malware takes complete control of the victims platform [13]. Once the attacker gets full control, he can download and execute any code and run malicious activities on the victim's platform such as joining botnets, sending spam emails, and participating in distributed denial of service attacks

[14]. Attackers may also record keystrokes, steal passwords, and can access sensitive information. Use of DBD to steal confidential data is also a major threat to the financial companies and banks.

A DBD attack occurs in four steps. First, the attacker compromises a genuine website and uploads malicious content to it. When a user visits that website, the malicious program is downloaded by browser, installed by itself, and the attacker gets full control [14]. APT programs and methods used by cybercriminal groups to attack businesses make them more dangerous.

In 2015, almost two million cases of malware infections to steal money were registered, while 34.2 % of computer users were exposed to at least one such attack through the year [12]. In order to ensure protection against such attack, there is a vital need for new methods and technologies that can safeguard the users from DBD attacks [31]. There are couple of existing techniques to detect and prevent such attacks. The detection of attacks can be performed by tracking web addresses with a history of malicious behavior [29]. According to Microsoft, Bing normally detects huge numbers of DBD pages every month. However, after getting blocked by Bing, the attackers switch servers and thus the same attacks are reborn but with different domain names [33]. Intrusion detection systems monitor traffic and system activities and may be used to detect attacks [7].

In order to counter the innovative tactics employed by the hackers, there is a vital need to develop efficient techniques that could potentially counter DBD attacks. In this paper, we proposed a novel design, which uses machine learning to detect and prevent DBD attacks. We selected nine attributes from a dataset of benign URLs from University of California Irvine (UCI) machine learning repository and malicious URLs from malware domain list [5]. Each attribute was chosen carefully to measure its effectiveness on different characteristics of malicious URLs. Furthermore, we employed several machine learning models for the training the system to detect malicious URLs. However, after empirical performance evaluation of these models, we selected Naive Bayes (NB), JRip, and J48 classifiers.

The rest of the paper is structured as follows: Section 2 contains the related work. Section 3 describes the methodology. Results are discussed in section 4 and section 5 concludes our work.

2 Related Work

Below we classify the most relevant work on detecting DBDs.

2.1 Using web crawler to detect drive by downloads

Harley and Pierre-Marc work does not offer a solution to DBDs, but tries to provoke more research in the area by suggesting possible ideas [15]. It provokes researchers to pay more attention to attacks that are large scale in nature and which do not use codes that are self-propagating. This is because current attacks are sophisticated and, therefore, a long-lasting solution may be one that uses

the fault-tolerant and robust software in addition to ensuring the monitoring of web pages. A web crawler can be used to identify distribution points, however, due to the complexity of this detection, false positives risks can be lessened by either digital signing or obfuscating techniques can be avoided. Some of the characteristics of this web crawler would be; ability to analyze HTML pages as well as follows its links; ability to imitate web cookies; ability to imitate scripting languages in order to decode obfuscated code; and ability to use heuristics in the detection of possible exploits in web pages. It concludes that measures, which are semi-effective and multi-layered, and those that accept specific risks of both false positives and negatives offer much protection.

2.2 Antivirus software to detect drive-by downloads malware

Narvaez et al., studied how antivirus software can be useful in the detection of drive-by malware installation by studying the effectiveness of the current antivirus tools [27]. A sample of malware was collected by use of a honeypot. The sample of the malware was categorized into whether the malware used either delivered payload or downloader. An evaluation of the results was made by common antivirus software to determine their effectiveness in detecting exploits. After 30days, the sample of the malware was scanned again as it was expected that the antivirus would have made an update of signature databases. According to the initial results, Norton detected 66% of the collected malware, Kaspersky 91%, CA 61%, ClamWin 62% [9] and TrendMicro 69% [22]. The next scan, after 30days, showed an increase in the rate of detection with Norton having 90%, Kaspersky 98%, TrendMicro 70%, ClamWin 75% and CA 81%. However, even though there was an improvement in the second scan, signature-based antivirus may not perform well in reality. This is because just as they had an opportunity to perform an update on their signatures similarly would attackers update malware. The initial detection, which was low, shows that malware authors use polymorphic capabilities. In 84% of attacks, downloaders are used instead of payloads. Antivirus products struggle to keep their signature databases up to date with the continuously changing threat landscape [32].

2.3 BrowserGuard as a behavior-based solution

Hsu et al. [17] proposed a behavior-based BrowserGuard, which detects secret downloads and blocks the malware from being executed. BrowserGuard uses two phases to provide protection to its host. The first is the filtration phase, whereby BrowserGuard makes a distinction between malicious and benign files depending on the situations in which they are downloaded. The second is the prohibition phase, whereby a request for the execution of malicious files is denied. In order to test the technique in terms of false positives, BrowserGuard visited the 500 top-ranked websites from Alexa. As expected BrowserGuard did not issue any attack alert, therefore, BrowserGuard had zero false positives. To measure the false negatives, Metasploit framework was used to generate ten malicious web pages that are then hosted on a remote server. BrowserGuard blocked all

ten pages, therefore, the authors claimed zero false negatives. To assess the performance overhead of BrowserGuard, the time to download five web pages, from Alexa, was measured 2000 times. BrowserGuard introduced a fixed delay time and the worst performance overhead was 2.5%. Unfortunately, we believe the test samples are insufficient to support the conclusions and BrowserGuard only works for Windows Internet Explorer 7.0.

2.4 A framework for DBD attacks with users voluntary monitoring of the web

Matsunaka et al. [24] proposed participative monitoring framework that fights DBDs with voluntary monitoring of websites by users and expert analysts. The framework provided a security ecosystem whereby users allow monitoring of their web activities, while security analysts do an inspection of the information in order to detect threats, devise countermeasures and provide feedback to the users. The framework enables users to provide data via the sensors and security analysts to give feedback through analyzing the data available at the center. The sensors are located in web proxies, DNS servers, and web browsers. Additionally, a web crawler was used to inspect web pages that are suspicious. The real-time data enabled the framework to previously detect unknown malicious web pages. However, advertisement hosts can cause false positives and further work is needed to address that.

2.5 HTML and JavaScript feature for detecting the drive-by download

Priya et al. [28] provided a static approach to detect DBDs using JavaScript and HTML features [28]. A sample dataset was created with 311 malicious URLs, from www.malwaredomainlist.com, and 654 benign URLs from Alexa were used to test different classifiers. To view the source code of benign sites you just open the URL, however, opening a malicious web page is a problem because it will cause malware to be installed on the computer. Therefore, MATLAB parser was developed to extract the malicious source code without visiting and executing the code. The HTML code was parsed and JavaScript and HTML features were extracted. They used both WEKA and MATLAB to evaluate the classifiers performance with 92% best case detection accuracy.

2.6 Approach to detect drive-by download based on characters

Matsunaka et al. [23] proposed FCDBD that includes monitoring sensors on the client side and analysis center on the network. The sensors include web browsers, web sensors or DNS sensors. The browser sensors extracted the user's data while DNS and web sensors monitored DNS-/HTTP- related traffic [3]. The analysis center collects the logs and analyzes them, if malicious websites are detected, the information is reported to monitoring sensors so the users may not access

the websites. The approach was evaluated using D3M 2013 dataset. According to the results, false positives only occur when a transition of a sequence of web pages is terminated before the malware is downloaded. To compensate for that, advertisement or affiliates scripts are obfuscated and referrer field is empty.

2.7 Enhanced approach for malware downloading:

Adachi et al. [1] used two approaches to predict DBD through opcode and vulnerability evaluation. The first approach identified vulnerabilities CVE-IDs in the web pages to predict the of malware download. For analysis, Wepawet was employed to identify CVE-IDs in the web pages, and the National Vulnerability Database (NVD) provided information concerning the CVEs. To improve detection rates they are reduced unnecessary information by a grouping algorithm [26]. Features were then extracted and the prediction model computed malware-downloading probabilities. The second approach combined opcode with the first approach one because opcode by itself fails to detect attacks that do not use JavaScript. Pages from 2011-2014 D3M datasets and AlexaTop500 were used. The first approach had 83% prediction accuracy and low FPs rate, however it had high FNs rate. The second approach had a 92% prediction accuracy, 11% FNs and 6% FPs using Random Forest.

2.8 Analyze redirection code for mining URLs:

Takata et al. [30] MineSpider performed an analysis on JavaScripts that include browser fingerprinting and redirection code and extracted possible URLs through the execution of the redirection code. MineSpider applied program slicing to JavaScript in order to extract execution paths, the extracted code fragments are executed by an interpreter and URLs are extracted. The outcome is just URL extraction and no detection was done. However, the URLs extracted by this method can be analyzed for malice using other approaches. MineSpider could extract more than 30,000 URLs in seconds compared to other methods.

2.9 Visualize the flow of HTTP traffic

Kikuchi et al. [20] used decision trees to classify DBDs by using features such as object size and redirection methods. The first premise was that many code variations modify words that are user-defined without the structure of the script being affected. Second, the characteristics of the scripts do not protect from DBDs because of disguised transformations fabrication. Additionally, they used the prediction of latent behavior to detect large-scale DBDs by using the drive-by disclosure method, which bridges the gap in between static and dynamic approaches. The method captured models and learned latent behaviors as opposed to scanning web pages for content that is malicious. To evaluate the efficiency of the approach 50 malicious and 50 legitimate sessions were obtained from Alexa. It was found that the method had no false positives but had 0.06 chance of false

negatives. The results showed that drive-by exposure can filter out scripts that are benign in nature, detect malicious scripts, and detect a variety of obfuscated patterns of DBDs as well as sort-out scripts that are disguised. In comparison to other high-tech solutions, drive-by disclosure was doubling accurate when compared to Cujo and it outdid JSAND by 29%.

2.10 Drive-by download as a large scale web attacks

Jodavi et al. drive-by disclosure [2] used anomaly DbD hunter approach to train and detect using a collection of classifiers. In the training stage, inputs of benign web pages are run in a browser. Then, JavaScript byte codes are logged for the web pages and a feature vector generated for the sequence. The feature vectors are then used to construct the classifiers baseline. The detection stage involved logging JavaScript byte codes for web pages, after which a feature vector is generated and applied to all base classifiers. The detection performance of DbD hunter was evaluated and was found that it increased the rate of detection by 12.44%, while decreasing rates of false alarms by approximately 48.13%. It had an accuracy of 97%, a detection rate of 96.3% and false alarm rate of 1.8%. Anomaly detection approach [8] have been used to detect DBD. According to [19], attacks by DBDs make use of browser exploit packs (BEPs) that are deployed on compromised servers to spread malware. BEPs that are widely used include sweet orange, Black Hole, Angler, Nuclear, Sakura, Fiesta, Hunter, Magnitude and Styx. The study makes an analysis of features that are built-in, which allow successful attacks by DBDs. The study conclude that just as attacks by DBDs increase in sophistication, so should the solutions.

3 Methodology

We develop a novel mechanism to counter DBD attacks that employs machine learning techniques. The proposed mechanism is able to classify the URLs into benign and malicious categories accurately. The benign category refers to websites that are safe, whereas the malicious category relates to the websites created by attackers to gain access or retrieve sensitive information. We used Waikato Environment for Knowledge Analysis (WEKA) [16] to classify the URLs based on different attributes using machine learning based models. WEKA is a popular machine learning suite developed at the University of Waikato, New Zealand and is licensed under the GNU General Public License (GPL). It contains machine learning algorithms for data mining related tasks. Integration feature helps to integrate these algorithms with the application code. It also supports data pre-processing, classification, regression, clustering, association rules, and visualization. The following subsection summarize the methodology used to classify DBDs and evaluate the performance.

3.1 dataset

We collected benign URLs from open source UCI Machine Learning Repository [21] and we used a list of 63 updated malware and spyware URLs from Malware Domain List [25].

3.2 Feature selection

Feature selection, also known as variable selection or attribute selection, is a process to select relevant features from predictive models. Each instance of the dataset used by machine learning algorithms is represented by the same set of features. These features can be continuous, categorical, or binary. We selected multiple effective features to build our proposed model. Given a single URL, its features were extracted and categorized into eight attributes (plus class) that were used by WEKA as itemized below.

- HostRank: the URLs global Amazon Alexa ranking [10].
- CountryRank: the URLs Amazon Alexa website rank by country [11].
- ASNNumber: The autonomous system number (ASN), which is assigned to the URLs domain, and used in BGP routing. [18].
- DotsInURL: number of dots in URLs [6].
- Lengthofurl: length of the URL.
- IPaddress: is the host name using ip address rather than name address.
- Lengthofhostname: length of host name.
- Safe Browsing: rating of Google safe browsing.

Two attribute evaluators: Correlation Attributes Evaluation (CAE) and Information Gain Attributes (IG) have been used on the dataset. Correlation Attributes Evaluation is used to choose best attributes for model training. It measures the correlation between attribute and the class and evaluates its worth. Information Gain picks attributes by measuring IG with respect to the class. For this work, eight features were selected to be used with WEKA. Referring to the figure 1, most of the attributes have scored a high ranking except IPaddress and ASNnumber for which, IG was 0.0521 and 0.1691, respectively. On the CAE, the IPaddress and ASNnumber scored 0.247 and 0.148, which are the lowest scores in the precision test. Thus, these two attributes were eliminated from the attribute set. We finalized six features that include Host Rank, Country Rank, Dots in URL, Length of the URL, length of the host name, in addition to the class: malicious or benign.

3.3 Classification

Many classifiers were chosen to train on the selected dataset, however, NB, JRip, and J48 outperformed all others. Therefore, we experimentally determined that those three are the best classifiers based on their performance on a given dataset. To evaluate the trained model, we employed 10 folds cross validation. Cross-validation is a technique to evaluate predictive models by splitting the original

Rank	Attribute	Rank	Attribute ranking
0.921	1 hostrank	0.899	8 SafeBrowsing
0.8823	5 lengthoftheurl	0.699	5 lengthoftheurl
0.7094	8 SafeBrowsing	0.574	4 dotsinURL
0.6565	2 countryrank	0.532	7 lengthofhostname
0.6115	7 lengthofhostname	0.345	1 hostrank
0.5162	4 dotsinURL	0.247	6 IPaddress
0.1691	3 ASNnumber	0.148	2 countryrank
0.0521	6 IPaddress	0.148	3 ASNnumber

(a) Information Gain Ranking Filter (b) Correlation Ranking Filter

Fig. 1. Information Gain and Correlation Ranking Attributes Ranking

dataset sample into a training and test sets to train and evaluate the model respectively. The process is repeated k times, with each of the k sub-samples used exactly once as the validation data. For this problem, data was split into 10 sets of size $n/10$, training with 9 subsets and testing on the remaining one subset. This process was repeated ten times while using a different subset for the test each time. The final results were then calculated by taking the mean accuracy of ten tests.

4 Results

Figure 2 shows the comparison of each classifier for malicious, benign, and average instance by using precision metric. We observed that NB scored 97% Malicious, 99% Benign, and 98% Average whereas JRip scored 97% Malicious, 99% Benign and 98% Average. Finally, J48 scored 95% Malicious, 97% Benign and 96% Average. Among all the three classifiers, the J48 scored the lowest with the average score of 96%. Naive Bayes and JRip have scored the highest in the tests, with similar results of average being 98%. Therefore, NB and JRip classifiers are used in the following analysis.

4.1 Metrics

Confusion matrix The confusion matrix summarizes the performance of classification model. True Positive (TP), False Negative (FN), False Positive (FP), and True Negative (TN) are elements of confusion matrix as shown in Figure 3. Columns represent the predicted class while rows represent the actual class. Higher values in the main diagonal reflect better accuracy in the classification.

True positive rate A true positive rate is the proportion of positives that are correctly identified by classifier. The TP rate is defined as follows.

$$TPRate = \frac{TP}{TP + FN} \quad (1)$$

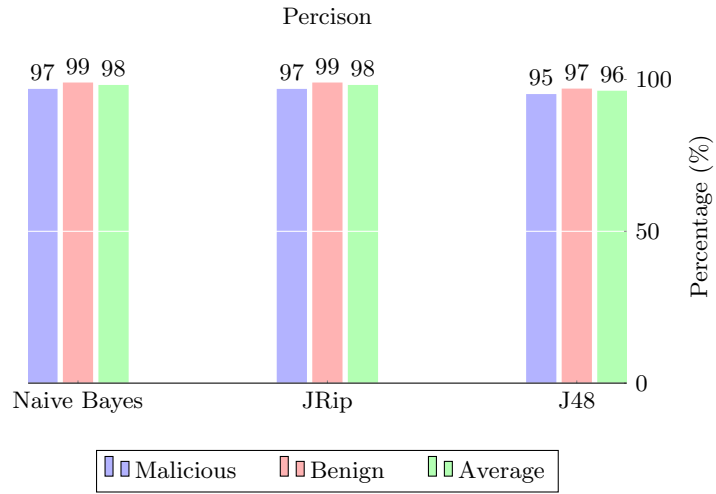


Fig. 2. Precision of different classifiers

		Predicted class	
		P	N
Actual Class	P	True Positives (TP)	False Negatives (FN)
	N	False Positives (FP)	True Negatives (TN)

Fig. 3. Confusion matrix

False positive rate A False Positive rate is the proportion of the outcome that is incorrectly predicted as yes (or positive) when it is actually no (negative). The FP rate is defined as follows.

$$FPRate = \frac{FP}{FP + TN} \quad (2)$$

Precision precision is the fraction of relevant instances among the retrieved instances.

$$Precision = \frac{TP}{TP + FP} \quad (3)$$

Recall Recall is the fraction of relevant instances among the retrieved instances.

$$Recall = \frac{TP}{TP + FN} \quad (4)$$

F-measure The F-measure is defined as a harmonic mean of precision and recall.

$$Precision = \frac{2 \times Precision \times Recall}{Recall + Precision} \quad (5)$$

Matthews correlation coefficient (MCC) MCC ranges from -1.0 (worst) to 1.0 (best) and is defined as follows.

$$MCC = \frac{(TP \times TN) - (FP \times FN)}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (6)$$

4.2 Naive Bayes

Table 1. Naive Bayes classifier results

	Malware	Benign	Average
TPR	98.40%	98%	98.20%
FPR	2%	1.60%	1.70%
Precision	96.90%	99%	98.20%
Recall	98.40%	98%	98.20%
F-Measure	97.60%	98.50%	98.20%
MCC	96.20%	96.20%	96.20%
ROC Area	98.70%	99.50%	99.20%
PRC Area	96%	99.70%	98.30%

The results of the NB classifier are shown in Table 1. The average score of TP is 98.20%, which indicates that the attributes have been correctly identified. The FP averaged 1.70%, which indicates that the result is scoring low on the error scale of the attributes. Therefore, the results can be identified as viable and true in this test. Table 2 shows the confusion matrix containing the details of the predicted and actual classes done by the NB classifier. Using these numbers we can calculate the TP and FP rates.

Table 2. Confusion Matrix Naive Bayes

	a=Malicious	b=Benign
a=Malicious	61	2
b=Benign	1	100

Applying formula 1 and 2 to the confusion matrix of NB, we get the following results.

$$TPRate = \frac{61}{61 + 2} = 0.968 \quad (7)$$

$$FPRate = \frac{1}{1 + 100} = 0.009 \quad (8)$$

4.3 JRIP

Table 3 shows that average TP of 98.20%, which indicates that JRip is able to correctly classify the URLs. The FP score is 1.70%, which indicates the classification had a low number of errors.

Table 3. JRip classifier results

	Malware	Benign	Average
TPR	98.40%	98%	98.20%
FPR	2%	1.60%	1.70%
Precision	90.60%	99%	98.20%
Recall	98.40%	98%	98.20%
F-Measure	97.60%	98.50%	98.20%
MCC	96.20%	96.20%	96.20%
ROC Area	97.20%	97.20%	97.20%
PRC Area	92.80%	98%	96%

In Table 4 the confusion matrix is presented, which contains the details about the predicted and actual classification done by the JRip classifier. The count of

Table 4. Confusion Matrix JRip

	a=Malicious	b=Benign
a=Malicious	62	1
b=Benign	1	100

TP is 62, FN is 1, and FP is 1 whereas TN is equal to 100. Using these numbers we can calculate the TP rate and FP rate.

$$TPRate = \frac{62}{62 + 1} = 0.984 \quad (9)$$

$$FPRate = \frac{1}{1 + 100} = 0.009 \quad (10)$$

4.4 J48

From Table 5, we can deduce that the average TP is 96.30%, which indicates that most of the URLs are correctly classified. The FP score is 4.10%, which indicates that the classification had a low number of errors.

Table 5. J48 classifier results

	Malware	Benign	Average
TPR	95.20%	97%	96.30%
FPR	3%	4.80%	4.10%
Precision	95.20%	97%	96.30%
Recall	95.20%	97%	96.30%
F-Measure	95.20%	97%	96.30%
MCC	92.30%	92.30%	92.30%
ROC Area	95.60%	95.60%	95.60%
PRC Area	91.60%	96.10%	94.40%

Table 6 shows the confusion matrix, which contains the details about the predicted and actual classification done by the J48 classifier. Using these numbers we can calculate the TP rate and FP rate.

Table 6. Confusion Matrix J48

	a=Malicious	b=Benign
a=Malicious	60	3
b=Benign	3	98

$$TPRate = \frac{60}{60 + 3} = 0.952 \quad (11)$$

$$FPRate = \frac{3}{3 + 98} = 0.029 \quad (12)$$

5 Conclusions

In this paper, we proposed an approach to filter benign and malicious websites. The URL based analysis is performed that helped by removing the runtime latency and delay of loading the websites. Furthermore, the proposed design protects the users from attacks induced by browser vulnerabilities. The proposed approach can be applied via a blacklisting content and system-based evaluation of site content and behavior of the site. By selecting the right features and algorithms, our system has achieved 98% accuracy in detecting and classifying the malicious URLs. The limitation of the work include the small dataset, number of classifiers used and actual real time testing. Future work would include creating a browser plugin and testing the system with real data, using a much larger dataset and investigating deep learning methods.

Acknowledgment

This research was supported, in part, by Zayed University Research Office, Research Incentives Grant # R18054.

References

1. Adachi, T., Omote, K.: An approach to predict drive-by-download attacks by vulnerability evaluation and opcode. In: Information Security (AsiaJCIS), 2015 10th Asia Joint Conference on. pp. 145–151. IEEE (2015), <https://doi.org/10.1109/AsiaJCIS.2015.17>
2. Al-Taharwa, I.A., Lee, H.M., Jeng, A.B., Ho, C.S., Wu, K.P., Chen, S.M.: Drive-by disclosure: A large-scale detector of drive-by downloads based on latent behavior prediction. In: Trustcom/BigDataSE/ISPA, 2015 IEEE. vol. 1, pp. 334–343. IEEE (2015), <https://doi.org/10.1109/Trustcom.2015.392>
3. Aldwairi, M., Guled, M., Cassada, M., Pratt, M., Stevenson, D., Franzon, P.: Switch architecture for optical burst switching networks. In: Proceedings of the first workshop on Optical Burst Switching, OPTICOMM'03 (Oct 2003)
4. Aldwairi, M., Alsalman, R.: Malurls: Malicious urls classification system. In: Annual International Conference on Information Theory and Applications (2011), https://doi.org/10.5176/978-981-08-8113-9_ITA2011-29
5. Aldwairi, M., Alsalman, R.: Malurls: a lightweight malicious website classification based on url features. Journal of Emerging Technologies in Web Intelligence 4(2), 128–133 (2012), <https://doi.org/10.4304/jetwi.4.2.128-133>

6. Aldwairi, M., Alwahedi, A.: Detecting fake news in social media networks. *Procedia Computer Science* 141, 215 – 222 (2018), <http://www.sciencedirect.com/science/article/pii/S1877050918318210>, the 9th International Conference on Emerging Ubiquitous Systems and Pervasive Networks (EUSPN-2018) / The 8th International Conference on Current and Future Trends of Information and Communication Technologies in Healthcare (ICTH-2018) / Affiliated Workshops
7. Aldwairi, M., Hasan, M., Balbahaith, Z.: Detection of drive-by download attacks using machine learning approach. *International Journal of Information Security and Privacy (IJISP)* 11(4), 16–28 (2017), <https://doi.org/10.4018/IJISP.2017100102>
8. Aldwairi, M., Mardini, W., Alhowaide, A.: Anomaly payload signature generation system based on efficient tokenization methodology. *International Journal on Communications Antenna and Propagation (IRECAP)* 8(5) (2018), <https://doi.org/10.15866/irecap.v8i5.12794>
9. Aldwairi, M., Mhaidat, K., Flaifel, Y.: Efficient pattern matching hardware for network intrusion detection systems. In: *International Conference on Electrical, Electronics, Computers, Communication, Mechanical and Computing (EECCMC)*. IEEE (2018), <https://arxiv.org/pdf/2003.00405.pdf>
10. amazon: Alexa, <https://www.alexa.com/>
11. amazon: Alexa: The top 500 sites on the web, <https://www.alexa.com/topsites/countries>
12. Anton, I., Makrushin, D., Van Der Wiel, J., Garnaeva, M., Namestnikov, Y.: Kaspersky security bulletin 2015. overall statistics for 2015. Securelist Information about Viruses Hackers and Spam. AO Kaspersky Lab (2015)
13. Cova, M., Kruegel, C., Vigna, G.: Detection and analysis of drive-by-download attacks and malicious javascript code. In: *Proceedings of the 19th international conference on World wide web*. pp. 281–290. ACM (2010), <https://doi.org/10.1145/1772690.1772720>
14. Egele, M., Wurzing, P., Kruegel, C., Kirda, E.: Defending browsers against drive-by downloads: Mitigating heap-spraying code injection attacks. In: Flegel, U., Bruschi, D. (eds.) *Detection of Intrusions and Malware, and Vulnerability Assessment*. pp. 88–106. Springer Berlin Heidelberg, Berlin, Heidelberg (2009), https://doi.org/10.1007/978-3-642-02918-9_6
15. Harley, D., Bureau, P.M.: Drive-by downloads from the trenches. In: *2008 3rd International Conference on Malicious and Unwanted Software (MALWARE)*. pp. 98–103. IEEE (2008), <https://doi.org/10.1109/MALWARE.2008.4690864>
16. Holmes, G., Donkin, A., Witten, I.H.: Weka: A machine learning workbench. In: *Intelligent Information Systems, 1994. Proceedings of the 1994 Second Australian and New Zealand Conference on*. pp. 357–361. IEEE (1994), <https://doi.org/10.1109/ANZIIS.1994.396988>
17. Hsu, F.H., Tso, C.K., Yeh, Y.C., Wang, W.J., Chen, L.H.: Browserguard: A behavior-based solution to drive-by-download attacks. *IEEE Journal on Selected areas in communications* 29(7), 1461–1468 (2011), <https://doi.org/10.1109/JSAC.2011.110811>
18. Huston, G.: Exploring autonomous system numbers. *The Internet Protocol Journal* 9(1), 2–23 (2006), https://www.kiv.zcu.cz/~ledvina/DHT/ipj_9-1.pdf
19. Jodavi, M., Abadi, M., Parhizkar, E.: Dbdhunter: An ensemble-based anomaly detection approach to detect drive-by download attacks. In: *Computer and Knowledge Engineering (ICCKE), 2015 5th International Conference on*. pp. 273–278. IEEE (2015), <https://doi.org/10.1109/ICCKE.2015.7365841>

20. Kikuchi, H., Matsumoto, H., Ishii, H.: Automated detection of drive-by download attack. In: Innovative Mobile and Internet Services in Ubiquitous Computing (IMIS), 2015 9th International Conference on. pp. 511–515. IEEE (2015), <https://doi.org/10.1109/IMIS.2015.71>
21. Lichman, M.: UCI machine learning repository (2013), <http://archive.ics.uci.edu/ml>
22. Masri, R., Aldwairi, M.: Automated malicious advertisement detection using virus-total, urlvoid, and trendmicro. In: 2017 8th International Conference on Information and Communication Systems (ICICS). pp. 336–341 (April 2017), <https://doi.org/10.1109/IACS.2017.7921994>
23. Matsunaka, T., Kubota, A., Kasama, T.: An approach to detect drive-by download by observing the web page transition behaviors. In: Information Security (ASIA JCIS), 2014 Ninth Asia Joint Conference on. pp. 19–25. IEEE (2014), <https://doi.org/10.1109/AsiaJCIS.2014.21>
24. Matsunaka, T., Urakawa, J., Kubota, A.: Detecting and preventing drive-by download attack via participative monitoring of the web. In: Information Security (Asia JCIS), 2013 Eighth Asia Joint Conference on. pp. 48–55. IEEE (2013), <https://doi.org/10.1109/ASIAJCIS.2013.15>
25. MDL: Malware Domain List featuring a list of malware-related sites plus a discussion forum on new threats (2009), <https://www.malwaredomainlist.com/mdl.php>
26. Mohammad, M.: A numerical solution of fredholm integral equations of the second kind based on tight framelets generated by the oblique extension principle. *Symmetry* 11(7), 854 (Jul 2019), <http://dx.doi.org/10.3390/sym11070854>
27. Narvaez, J., Endicott-Popovsky, B., Seifert, C., Aval, C., Frincke, D.A.: Drive-by-downloads. In: System Sciences (HICSS), 2010 43rd Hawaii International Conference on. pp. 1–10. IEEE (2010), <https://doi.org/10.1109/HICSS.2010.160>
28. Priya, M., Sandhya, L., Thomas, C.: A static approach to detect drive-by-download attacks on webpages. In: Control Communication and Computing (ICCC), 2013 International Conference on. pp. 298–303. IEEE (2013), <https://doi.org/10.1109/ICCC.2013.6731668>
29. Provos, N., Mavrommatis, P., Rajab, M.A., Monroe, F.: All your iframes point to us. In: Proceedings of the 17th Conference on Security Symposium. pp. 1–15. SS’08, USENIX Association, Berkeley, CA, USA (2008), <http://dl.acm.org/citation.cfm?id=1496711.1496712>
30. Takata, Y., Akiyama, M., Yagi, T., Hariu, T., Goto, S.: Minespider: Extracting urls from environment-dependent drive-by download attacks. In: Computer Software and Applications Conference (COMPSAC), 2015 IEEE 39th Annual. vol. 2, pp. 444–449. IEEE (2015), <https://doi.org/10.1109/COMPSAC.2015.76>
31. Vergelis, M., Shcherbakova, T., Demidova, N.: Kaspersky security bulletin. spam in 2014. *Secure List* p. 68 (2015)
32. Yaseen, Q., Jararweh, Y., Al-Ayyoub, M., Aldwairi, M.: Collusion attacks in internet of things: Detection and mitigation using a fog based model. In: 2017 IEEE Sensors Applications Symposium (SAS). pp. 1–5 (March 2017), <https://doi.org/10.1109/SAS.2017.7894031>
33. Zhang, J., Seifert, C., Stokes, J.W., Lee, W.: Arrow: Generating signatures to detect drive-by downloads. In: Proceedings of the 20th international conference on World wide web. pp. 187–196. ACM (2011), <https://dl.acm.org/doi/pdf/10.1145/1963405.1963435>