

ANALYSING ANDROID ENCRYPTED NETWORK TRAFFIC

¹Brij Mohan Goel*, ²Shefali Saini

¹Research Scholar, ²Supervisor (Professor)

¹⁻²Department of Computer Science and Engineering,

Baba MastNath University, Rohtak, Haryana, India

Email ID: shefalisaini9@gmail.com

Accepted: 21.01.2022

Published: 28.02.2022

Keywords: Android, Network Traffic.

Abstract

Signature-based malware detection algorithms are facing challenges to cope with the massive number of threats in the Android environment. In this paper, conversation-level network traffic features are extracted and used in a supervised based model. This model was used to enhance the process of Android malware detection, categorization, and family classification. The model employs the ensemble learning technique in order to select the most useful features among the extracted features. A real-world dataset called CICAndMal2017 was used in this paper. The results show that Extra-trees classifier had achieved the highest weighted accuracy percentage among the other classifiers by 87.75%, 79.97%, and 66.71% for malware detection, malware categorization, and malware family classification respectively. A comparison with another study that uses the same dataset was made. This study has achieved a significant enhancement in malware family classification and malware categorization. For malware family classification, the enhancement was 39.71% for precision and 41.09% for recall. The rate of enhancement for the Android malware

categorization was 30.2% and 31.14% for precision and recall, respectively.

Paper Identification



*Corresponding Author

I. INTRODUCTION

Nowadays, smartphones are not only for making phone calls as it was before. It is now a tool for holding personal information, health care, payment, and more eservices. As a result, the number of smartphone users in 2019 has increased by 5.9% more than in 2018 [1]. According to a report prepared by International Data Corporation (IDC) [2], the Android operating system is the most popular operating system for smartphones in 2019. It has an 86.7% market share more than any other smartphone's operating system. Android is a Linux based open-source operating system developed by Google [3]. It was invented in 2003, whereas the first Android smartphone was

invented in 2008 [4]. “Google Play” is an official market store offered by Google. Google Play offers more than two and a half million applications [5]. However, this store is not the only source of Android applications; many other unofficial third-party application developers exist. Along with such a massive number of applications, the number of potential security and privacy issues by malware is increased [6]. In order to reduce the risks of malware and other malicious applications, Google released a machine learning ecosystem under the name of “Play Protect” [7]. It is designed to detect malware before and after uploading applications to the market. In spite of this control process, more than 132 thousand malware was detected in the first quarter of 2018 [8], two million Android users were infected by “False-Guide” botnet in 2016 as well as, half-million users were infected by thirteen different malware from applications that were uploaded to Google play market [9]. Unfortunately, Android smartphones still a target for cybercriminals. The risks of malware are growing, and so are the efforts to mitigate their risks. In this respect, security researchers employ two methods to detect malware: the static-based method that aims to analyze the malware without running it and the dynamic-based method that monitors the malware behavior inside an isolated environment (i.e., monitoring the generated traffic of malware) [10]. Both methods can be used by machine learning to enhance malware detection. The contribution of this study is a) Distinguishing the most effective network traffic features. b) Determining the best machine learning algorithm (out of three classifiers) for detecting, categorizing and classifying of malware. c) Providing a comparison between this study and related studies that used the same dataset. d) Enhancing the used dataset “CICAndMal2017”.

CICANDMAL2017

Dataset After reviewing the most comprehensive and coherent set of related publications, we found that the Canadian Institute for Cybersecurity (CIC) [11] provides a competent real-world dataset called CICAndMal2017. CIC build their dataset in a way that reduces the drawbacks and shortcomings of the earlier dataset. As a starting point, the CIC collected more than four thousand malware applications from different resources, such as VirusTotal [12] and Contagiodumpst [13]. Moreover, more than six thousand benign applications published during 2015, 2016, and 2017 and uploaded to Google play market were collected. However, CIC has only managed to install 5 thousand of them (malware 429 and benign 5,065) on real android smartphones to conduct a real-world environment. Finally, CIC connected Android smartphones to a hotspot computer to capture network traffic in Capture Packet (PCAP) format using TCPDUMP software [14]. The CIC offered the dataset in two formats: CSV files generated by CIC Flow meter (2126 CSV file) and PCAP files (more than 20 gigabytes of captured network traffic). In this paper, the PCAP files are used. To deal with some advanced malware that uses the time delay technique to escape dynamic analysis, the CIC captured the network traffic in three different times: After malware installation directly, fifteen minutes before rebooting and fifteen minutes after rebooting. CICAndMal2017 has multiple levels of labeling. At the first level, the PCAP files are grouped into two categories: benign or malware. In the second level, malware types are categorized into four categories:

Adware: Adware automatically displays advertising materials and aims to collect the highest number of clicks or views on unwanted advertisement banners
 Ransomware: A malicious application aims to block access over computer resources. For example, it can encrypt users’ files to extort them to pay money for decrypting their files or unlocking their devices [12].

Scareware: This type of malware tries to scare users to let them purchase unnecessary and potentially dangerous software applications [14].

SMS malware: A malware that makes unauthorized calls or/and sends SMS messages without user consent. The malware owner can operate the infected devices as a premium channel for SMS services [11].

Table 1. Malware category and family types.

Category	Family Type		
Adware	Ewind	koodous	Kemoge
	Dowgin	Mobidash	Youmi
	Feiwo	Selfmite	Shuanet
	Gooligan		
Ransomware	Charger	Pletor	LockerPin
	Jisut	PornDroid	Svpeng
	Koler	RansomBO	WannaLocker
	Simplocker		
Scareware	AndroidDefender	FakeAV	FakeApp.AL
	AndroidSpy.277	FakeJobOffer	FakeAV
	AV for Android	FakeTaoBao	FakeApp
	AVpass	Penetho	VirusShield
SMS Malware	BeanBot	Jifake	FakeNotify
	Biige	Mazarbot	SMSsniffer
	FakeInst	Nandrobox	FakeMart
	Plankton		

II. RELATED WORK

In [11], a new method was introduced to detect android malware using edge computing and traffic clustering. First, the authors sent the android devices traffic to the edge server. Second, the edge server extracts mobile traffic content features (i.e., extracted plaintext from HTTP flows) and traffic behavior (i.e., packet intervals) and sent them to the cloud platform. Finally, they calculated the similarities between applications and clusters to detect the malware automatically. They used similarity methods: TF-IDF algorithm and cosine similarity. For evaluating their method, they used 400 android application. Note that the data set was not published online because of privacy concerns. The final average accuracy for their model was 96.9%. In [7], the author used the Long Short Term Memory (LSTM) based deep learning framework on detecting malware of type ransom ware.

Two categories of CICAndMal2017 dataset were selected in this study: benign and ransomware. Furthermore, they selected the top 19 flow-level features using 8 feature selection algorithms such as Chi-Square and information gain. The accuracy, recall, and F1-score results of their model were 97%. In [9] a part of the CICAndMal2017 dataset was used; the researchers choose one PCAP file for each malware family. Their chosen samples were taken randomly. Features were extracted from PCAP files using two steps. The first step: a Java program was developed to separate network flows using the flow level technique. Then, fifteen features were extracted, using a python program. One of the features was the minimum size of the sent packet within a flow. Three supervised machine-learning classifiers were used. The classifiers were K-Nearest Neighbours, Random forest and Decision Tree. They classify instances into two categories: malware and benign. Then, classifying malware instances into three categories: Adware, Ransom ware, and Scareware. The authors use three measures: recall, precision, and F-score.

For malware-benign classification, the results show that the Random Forest classifier has obtained the highest results by 92% for F-score and 95% for precision as well as recall. The other classifiers gained more than 85% of all used measures. Formal ware classification, the selected classifiers achieved more than 80% for the chosen measures. Similar to malware benign classification, the Random Forest classifier gained the highest results by 84% for recall, precision and F-score. The researchers did not show the results if the full dataset was used. In [11], CICAndMal2017 dataset was used. CIC researchers extracted network traffic flow-level features. Two algorithms for feature selection were used: Information Gain (IG) and Correlation-based Feature Selection (CFS). The two algorithms select nine features. Three machine learning classifiers were used to evaluate their model, namely: Decision Tree, Random Forests, and K-

Nearest (KNN). The classifiers categorized malware in three scenarios: malware binary detection, malware category classification and malware families' characterization. The results show that network traffic flow-level features are useful for binary detection, but not for the other scenarios. For clarification, the three classifiers gained 85% precision on average and 88% for recall measure for binary detection. On the other hand, malware category classification achieved less than 50% for precision and recall, and less than 20% for precision and recall for the family classification. In [38], the Decision Tree (J48) algorithm was used to detect malware traffic. They used 700 samples; 200 samples are malware from Drebin [6] and Contagiodumpst datasets, and 500 benign samples from Google play market. Network traffic of these samples was captured on a real smartphone using "tcpdump" [8]. The authors extracted seven features from the captured traffic. Finally, they calculated the accuracy for Drebin and Contagiodumpst datasets, and the results were 98.4% and 97.6%, respectively. In [10], a new model was proposed to detect and categorize Android malware based on network traffic features. The authors collected the generated network traffic of 1500 benign applications as well as 400 general malware and adware. Next, they used feature selection algorithms such as IG and CFS to select the most useful features. Finally, supervised machine learning classifiers were used to detect and categorize malware. The proposed model achieved more than 90% average accuracy and precision. In [4], network traffic features of Android malware are prioritized based on IG and Chi-Square tests. Next, network traffic features were minimized using a proposed algorithm to enhance the detection accuracy and reduce the time for training and testing phases. Statistical analysis techniques were used to rank features. The proposed algorithm finds that 9 out of 22 features are adequate for higher detection accuracy. Likewise, the study results show that it can

reduce the time for training and testing phases 50% and 30%, respectively.

III. EXPERIMENTATION

All experiments have been conducted on the Microsoft Windows 10 Professional (64-bit) version with a second-generation 2.20 GHz Intel Core i7 processor and 16 GB of memory. Python 3.7.0 was chosen for data pre-processing, feature selection, and model building because of its productive and useful libraries for such tasks. One tool that uses the conversation-level technique is the PeerShark tool [12]. The available version of this tool extracts six features only. Since it is an opensource tool, it can be enhanced to adopt new features. Therefore, a new fourteen conversation-level features were developed (E1-E14). Table 2 list these features. After executing the cleaning phase, the number of removed instances was 798 instances for the first scenario and 456 instances for the second and the third scenario. Only two identification features were removed: source IP and destination IP. In the feature selection phase and after executing the ensemble learning technique, nine features were selected.

Table 2. Peershark basic and extended features.

#	Feature name	Description
1	SourceIP	The source IP of the conversation
2	DestinationIP	The destination IP Source
3	NoOfPackets	Number of packets during the conversation
4	NoOfBytes	Number of bytes during the conversation
5	InterArrivaltime	Median of inter-arrival time of packets
6	DurationInSeconds	Duration time of a connection in seconds
E1	NoOfPacketSizeFWD	Number of byte per forward packet
E2	NoOfPacketSizeBWD	Number of byte per backward packet
E3	PacketPerSecFWD	The forward packet per second
E4	PacketPerSecBWD	The backward packet per second
E5	PKTFwdnum	Total number of forward packets
E6	PKTBwdnum	Total number of backward packets
E7	ByteFwdnum	Total number of forward bytes
E8	ByteBwdnum	Total number of backward bytes
E9	BytePerFlow	Number of byte per flow within the conversation
E10	PacketPerFlow	Number of packets per flow within the conversation
E11	FWDBytePerFlow	Number of forward bytes per flow within the conversation
E12	BWDBytePerFlow	Number of backward bytes per flow within the conversation
E13	FWDPacketPerFlow	Number of forward packets per flow within the conversation
E14	BWDPacketPerFlow	Number of backward packets per flow within the conversation

IV. CONCLUSIONS

This research introduces an enhanced model for malware detection, categorization, and family classification in the android environment. The model extracts conversation-level network traffic features from a recent and real-world dataset named "CICAndMal2017. For the process of the feature extraction phase, conversation-level features were extracted using the PeerShark tool. Multiple stages of data pre-processing have been conducted to the dataset. The most useful features were selected using the ensemble learning technique by three feature selection algorithms: Random Forest, RFE, and LightGBM classifiers. Moreover, the developed model was trained and tested using three classifiers: Decision Tree, Random Forest, and Extra-trees. Finally, this study compared the provided model results with another model that used the same dataset. According to the final results, conversation-based features can enhance the detection, categorization, and family classification of Android malware. Furthermore, among the selected classifiers, the Extratrees algorithm achieved the maximum accuracy results. In comparison with a study from CIC, this model obtains better results in binary detection, and significant enhancement in malware categorization by 30.3% for precision and 31.14% for recall. Furthermore, the accuracy of malware family classification is improved by 39.71% for precision and 41.09% for recall.

RÉFÉRENCES

1. Hamandi K., Chehab A., Elhadj I., and Kayssi A., "Android SMS malware: Vulnerability and Mitigation," in Proceedings of the 27th International Conference on Advanced Information Networking and Applications Workshops, Barcelona, pp. 1004-1009, 2013.
2. IDC, "Smartphone Market Share," 2019. [Online]. Available: <https://www.idc.com/promo/smartphone-marketshare/os>, Last Visited, 2020.
3. Kashefi I., Kassiri M., and Saleh M., "Preventing Collusion Attack in Android," The International Arab Journal of Information Technology, vol. 12, no. 6A, pp. 719-727, 2015.
4. Lashkari A., Kadir A., Taheri L., and Ghorbani A., "Toward Developing a Systematic Approach to Generate Benchmark Android Malware Datasets and Classification," in Proceedings International Carnahan Conference on Security Technology, Montreal, pp. 1-7, 2018.
5. Lashkari A., Akadir A., Gonzalez H., Mbah K., and Ghorbani A., "Towards A Network-Based Framework for Android Malware Detection and Characterization," in Proceedings of 15th Annual Conference on Privacy, Security and Trust, PST, Calgary, pp. 233-242, 2018.
6. Liao Q., "Ransomware: a Growing Threat to SMEs," in Proceedings of Southwest Decision Sciences Institute's Annual Conference, Houston, pp. 360-366, 2008.
7. Narang P., Hota C., and Venkatakrishnan V., "PeerShark: Flow-Clustering and ConversationGeneration for Malicious Peer-To-Peer Traffic Identification," EURASIP Journal on Information Security, vol. 2014, no. 1, pp. 1-12, 2014.
8. Nauman M. and Khan S., "Design and Implementation of A Fine-Grained Resource Usage Model for the Android Platform," The International Arab Journal of Information Technology, vol. 8, no. 4, pp. 440-448, 2011.
9. Parkour M., "Contagio malware database," contagiodump. 2013. [Online]. Available: <http://contagiodump.blogspot.com/2011/03/ta>

- kesample-leave-sample-mobile-malware.html,
Last Visited, 2020.
10. Chen R., Li Y., and Fang W., “Android Malware Identification Based on Traffic Analysis,” in Proceedings of International Conference on Artificial Intelligence and Security, New York, pp. 293-303, 2019.
 11. Draper-Gil G., Lashkari A., Mamun M., and Ghorbani A., “Characterization Of Encrypted And VPN Traffic Using Time-Related Features,” in Proceedings of the 2nd International Conference on Information Systems Security and Privacy, Italy, pp. 407-414, 2016.
 12. F-Secure, “Android/Kmin.” 2012. [Online]. Available: https://www.f-secure.com/vdescs/trojan_android_kmin.shtml, Last Visited, 2020.
 13. Google, “Google Play Protect.” [Online]. Available: <https://www.android.com/playprotect/>Last Visited, 2020.
 14. Google, “Google Play Store.” [Online]. Available: <https://play.google.com/store/>, Last Visited, 2020.