

Mining and Detection of Android Malware Based on Permissions

Abdirashid Ahmed Sahal
Dept. of Computer Engineering
Gebze Technical University
Kocaeli, Turkey
asahal@gtu.edu.tr

Shahid Alam
Dept. of Computer Engineering
Adana Science and Technology University
Adana, Turkey
salam@adanabtu.edu.tr

Ibrahim Soğukpınar
Dept. of Computer Engineering
Gebze Technical University
Kocaeli, Turkey
ispinar@gtu.edu.tr

Abstract—Due to the open app distribution and more than two billion active users, Android platform continues to serve as low-hanging fruit for malware developers. According to the McAfee threat report, the number of malware families found in the Google Play increased by 30% in 2017. Permission-based access control model is one of the most important mechanisms to protect Android apps against malware. In this paper, we propose a new permission-based model that enhances the efficiency and accuracy of Android malware analysis and detection, and has the capability of potentially detecting previously unknown malware. In this new model, we improve the feature selection by introducing a new weighting method, named TF-IDFCF, based on the class frequency (CF) of the feature. The results of our experiments show that our proposed method has a detection rate of greater than 95.3% with a low false positive rate, when tested with different classifiers.

Index Terms—Android, Permissions, Malware Analysis and Detection, TF-IDF, Machine Learning.

I. INTRODUCTION

Recently, Android has become the most selling operating system on mobile devices [1]. Android OS monthly has over two billion active users. Malware writers are actively and continuously developing malware programs to target Android platform. This continuous evolution and the diversity of malware pose a major threat to Android applications. According to the McAfee threat report, number of malware families found in the Google play increased by 30% in 2017 [2]. Different solutions have been proposed to protect mobile users from the increasing threats of Android malware. Permission-based access control model is the most important mechanism for Android protection against malware apps. In this paper, we use multiple machine learning algorithms with permission datasets to build and train models to classify Android malicious apps. We improved the Term Frequency-Inverse Document Frequency (TF-IDF) method, introduced new feature selection method that increases the detection accuracy. In this paper, we propose a new permission-based static analysis framework for the classification of Android applications into benign and malware. We improve upon

other permission-based approaches by introducing a feature selection method based on TF-IDF. This method improves the efficiency of malware analysis and detection, and obtains a high accuracy. The results of our experiments show that our proposed framework has a detection rate of more than 95.3% using most of the basic classifiers, such as SVM, J48, Naive Bayes and KNN. Our contribution of this paper is an improvement of TF-IDF weighting on vector space model. The TF-IDF method considers both TF and IDF [3]. If the TF is high and the term only appears in some part of the applications, then this term has a very good ability to differentiate the applications. A feature occurring frequently in the applications within same class represents more characteristics of the class. Therefore, we use total feature occurrence as a new parameter and enhance the TF-IDF to improve the efficiency of our classifier.

The remaining parts of the paper are organized as follows. In section 2 related works are presented. The proposed model is introduced in section 3. In section 4 experimental results are given. The last section is conclusions.

II. RELATED WORK

Many research has been performed on Android malware classification, using permission related features. We highlight some of related works of permission-based malware detection.

X. Liu and J. Liu proposed a framework that considers both requested and used permissions in the Android applications [4]. This framework is two layered malware detections and uses machine learning techniques to get high detection accuracy with the potential of detecting Android malware applications based on permissions. P. Rovelli and Ý. Vigfússon proposed a simple, client-server architecture malware detection system based on permissions which views requested permissions as behavioral markers [5]. The system has the server-side and client-side permission checker parts. The client-side part extracts the permissions from the Android apps and forwards the extracted permissions to the server-side part. The server-side part classifies the application as benign or malware.

write, and send SMS (such as contacts, banking information etc.) without notifying the user. **CALL_PHONE**: This permission allows an application to make calls on behalf of a user without confirming from the user, and hence is dangerous to the privacy of a user. **READ_HISTORY_BOOKMARKS**: This permission allows an application to read a user's browsing history and bookmarks, and hence is a risk to the privacy of a user. **ACCESS_COARSE_LOCATION** and **ACCESS_FINE_LOCATION**: These permissions allow an application to access the coarse (e.g., Cell-ID, Wi-Fi) and fine (e.g., GPS) locations citesarma2012android of a smartphone, and hence pose a risk to the privacy of a user.

C. Classification Models

In this section, we evaluated using Weka tool the classification capability of our framework for Android APK files as benign or malware on two datasets, a training set of 700 and a testing set of 300 samples. This framework gets the original feature set by extracting permissions from Android APK files using reverse engineering tool and forms the feature set by using the newly introduced method named TF-IDF-CF. We used the techniques of machine learning to classify the Android applications. We used multiple learning algorithms such as J48, Naive Bayes, SVM and KNN to test our model. The metrics we used are as follows:

TP (True Positive): Correctly detected number of malware apps. **FP** (False Positive): Number of benign applications wrongly identified as malware. **TN** (True Negative): Correctly detected number of benign apps. **FN** (False Negative): Number of malware applications wrongly identified as benign.

We used the below metrics to evaluate the performance of our proposed permission-based detection framework.

True Positive Rate (TPR): Number of samples correctly classified as malware out of the total malware dataset ($TP / TP+FN$). **False Positive Rate (FPR)**: Number of samples wrongly classified as malware out of the total benign dataset ($FP / TN+FP$). **Overall Accuracy (ACC)**: Proportion of Android apps, that are correctly identified as either malicious or benign app. ($TP+TN / TP+TN+FP+FN$).

D. Evaluation

To check the accuracy of our permission-based detection method we also used dataset containing 1000 APK samples, including 500 malwares and 500 benign samples, with 10-fold cross-validation method. This method divides the dataset into ten parts and, takes one part the 10% as testing set and the rest 90% as training test.

To further evaluate and make a comparison with our proposed technique, we selected two other well-known information mining techniques, Principal Component Analysis (PCA) [17] and Information Gain [3]. Both these techniques reduce the number of attributes to achieve a

TABLE II
EXPERIMENTAL RESULTS

Dataset	Method	TPR	FPR	ACC
Training	J48	0.991 %	0.009 %	99.1429 %
Training	Naive Bayes	0.98 %	0.02 %	98 %
Training	SVM	1 %	0 %	100 %
Training	KNN	1 %	0 %	100 %
Test	J48	0.97 %	0.03 %	97 %
Test	Naive Bayes	0.97 %	0.03 %	97 %
Test	SVM	0.983 %	0.017 %	98.3333 %
Test	KNN	0.953 %	0.047 %	95.3333 %

high detection rate. PCA is a statistical method used to emphasize variation and bring out strong patterns in a dataset. Information Gain depends on the entropy of an attribute and selects a feature that provides the foremost information gain. Using the same dataset, we compared our method to Information Gain and PCA. Results are shown in table III. As seen in table III the use of our proposed method for assigning weights to each feature extracted from the APK has higher accuracy rate than the other two techniques.

TABLE III
COMPARED RESULTS

Information Gain Method				
Cross Validation	Classifier	TPR	FPR	ACC
10-fold	J48	0.86 %	0.14 %	86 %
10-fold	Naive Bayes	0.819 %	0.181 %	81.9 %
10-fold	SVM	0.864 %	0.136 %	86.4 %
10-fold	KNN	0.856 %	0.144 %	85.6 %
Principal Component Analysis Method				
Cross Validation	Classifier	TPR	FPR	ACC
10-fold	J48	0.864 %	0.136 %	86.4 %
10-fold	Naive Bayes	0.578 %	0.422 %	57.8 %
10-fold	SVM	0.815 %	0.185 %	81.5 %
10-fold	KNN	0.843 %	0.157 %	84.3 %
Our Proposed (TF-IDF-CF) Method				
Cross Validation	Classifier	TPR	FPR	ACC
10-fold	J48	0.975 %	0.025 %	97.5 %
10-fold	Naive Bayes	0.973 %	0.027 %	97.3 %
10-fold	SVM	0.986 %	0.014 %	98.6 %
10-fold	KNN	0.968 %	0.032 %	96.8 %

E. Limitations

In this work, our weight assigning method relies on the labeled (into classes) data. Before training, the data needs to be divided into classes, which enables us to calculate the class frequency (c_i) and allows us to use class-based instead of corpus based TF-IDF. Currently the only classes data is divided into, before training, are malware and benign.

V. CONCLUSIONS

Permission is one of the most important features in Android security, and meaningful in malware detection. Our proposed permission-based framework uses machine learning algorithms to detect potentially malware apps. Also, to improve the efficiency of permission-based Android malware analysis and detection we introduced a new method based on TF-IDF to assign weight to each feature extracted from an Android APK file. We evaluated the new technique using different metrics and achieved a detection rate higher than 95.3% using different classifier algorithms. To further improve malware detection, in future we will add more features, such as, API calls etc., to our proposed framework. To improve the weight assigning method, in future, we are going to divide the data (especially malware samples) into more (than two) classes based on similarity measures.

ACKNOWLEDGEMENTS

This work was supported by the Scientific and Technological Research Council of Turkey (TÜBİTAK), Grant No: ARDEB-116E624.

REFERENCES

- [1] "The Statistics Portal: Statistics and studies from more than 22,500 sources," <https://www.statista.com/statistics/266210/number-of-available-applications-in-the-google-play-store/>, accessed: 2018-05-01.
- [2] "McAfee mobile threat report Q1," <https://www.mcafee.com/enterprise/en-us/assets/reports/rp-mobile-threat-report-2018.pdf>, accessed: 2018-05-29.
- [3] D. M. Christopher, R. Prabhakar, and S. Hinrich, *Introduction to information retrieval*. New York, NY, USA: Cambridge University Press, 2008.
- [4] X. Liu and J. Liu, "A two-layered permission-based android malware detection scheme," in *Mobile cloud computing, services, and engineering (mobilecloud), 2014 2nd ieee international conference on*. IEEE, 2014, pp. 142–148.
- [5] P. Rovelli and Ý. Vigfússon, "Pmds: Permission-based malware detection system," in *International Conference on Information Systems Security*. Springer, 2014, pp. 338–357.
- [6] Z. Aung and W. Zaw, "Permission-based android malware detection," *International Journal of Scientific & Technology Research*, vol. 2, no. 3, pp. 228–234, 2013.
- [7] Z. Xiaoyan, F. Juan, and W. Xiujuan, "Android malware detection based on permissions," *The Institute of Engineering and Technology*, 2014.
- [8] D. Arp, M. Spreitzenbarth, M. Hubner, H. Gascon, K. Rieck, and C. Siemens, "Drebin: Effective and explainable detection of android malware in your pocket." in *Ndss*, vol. 14, 2014, pp. 23–26.
- [9] Y. Aafer, W. Du, and H. Yin, "Droidapiminer: Mining api-level features for robust malware detection in android," in *International conference on security and privacy in communication systems*. Springer, 2013, pp. 86–103.
- [10] N. Bhargava, G. Sharma, R. Bhargava, and M. Mathuria, "Decision tree analysis on j48 algorithm for data mining," *Proceedings of International Journal of Advanced Research in Computer Science and Software Engineering*, vol. 3, no. 6, 2013.
- [11] T. R. Patil and S. Sherekar, "Performance analysis of naive bayes and j48 classification algorithm for data classification," *International Journal of Computer Science and Applications*, vol. 6, no. 2, pp. 256–261, 2013.
- [12] A. J. Smola and B. Schölkopf, "A tutorial on support vector regression," *Statistics and computing*, vol. 14, no. 3, pp. 199–222, 2004.
- [13] M.-L. Zhang and Z.-H. Zhou, "MI-knn: A lazy learning approach to multi-label learning," *Pattern recognition*, vol. 40, no. 7, pp. 2038–2048, 2007.
- [14] "Google Play Store google play store," <https://play.google.com/store/apps>, accessed: 2018-05-01.
- [15] "Contagio Mobile genome project," <http://contagiominiidump.blogspot.com.tr/>, accessed: 2018-05-01.
- [16] Y. Zhou and X. Jiang, "Android malware genome project," *Disponibile a http://www.malgenomeproject.org*, 2012.
- [17] I. T. Jolliffe, *Principal Component Analysis*. New York, NY, USA: Springer-Verlag, 2002.