

SEdroid: A Robust Android Malware Detector using Selective Ensemble Learning

Ji Wang
andrewwang@pku.edu.cn
Peking University
Beijing, China

Qi Jing*
jingqi@pku.edu.cn
Peking University
Beijing, China

Jianbo Gao
gaojianbo@pku.edu.cn
Peking University
Beijing, China

ABSTRACT

For the dramatic increase of Android malware and low efficiency of manual check process, deep learning methods started to be an auxiliary means for Android malware detection these years. However, these models are highly dependent on the quality of datasets, and perform unsatisfactory results when the quality of training data is not good enough. In the real world, the quality of datasets without manually check cannot be guaranteed, even Google Play may contain malicious applications, which will cause the trained model failure. To address the challenge, we propose a *robust* Android malware detection approach based on selective ensemble learning, trying to provide an effective solution not that limited to the quality of datasets. The proposed model utilizes genetic algorithm to help find the best combination of the component learners and improve robustness of the model. Our results show that the proposed approach achieves a more robust performance than other approaches in the same area.

CCS CONCEPTS

• **Security and privacy** → **Software and application security**;
• **Computing methodologies** → *Neural networks*.

KEYWORDS

Android malware detection, Machine learning, Selective ensemble, Genetic algorithm

ACM Reference Format:

Ji Wang, Qi Jing, and Jianbo Gao. 2019. SEdroid: A Robust Android Malware Detector using Selective Ensemble Learning. In *CCS '19: ACM Conference on Computer and Communications Security, November 11–15, 2019, London, UK*. ACM, New York, NY, USA, 3 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

In order to address the challenges of the rapid increase of Android malwares, many deep learning approaches have been presented in recent years. Some of them are proved to be efficient and of high accuracy. Sahs et al.[1] construct a SVM model using features

consisting of Android permission and control flow graphs, which achieves a 85% precision rate and a 95% recall rate. Zhu et al. [2] introduce DBN (Deep Belief network) to this problem and achieve great results with a 95.05% F1-score.

However, there still occurred some problems. On one hand, some researchers may find it difficult to reproduce a proven effective model in practice. On the other hand, the lack of dataset and low-quality training data may cause a model's failure. This is because of the high dependency of high-quality dataset for deep learning model. If the dataset is not good enough, the trained model may be not that good in training process and become hard to implement in application. This problem may become worse for large and high-quality dataset are commonly inaccessible for the reluctancy of holders to share the datasets for security or benefit reasons. So the robustness of the model becomes fairly important to this challenge.

In this paper, we try to exploit the architecture of selective ensemble to offer a solution with stronger generalization ability and robustness to this problem. For ensemble learning can increase memory use and computational cost, we selectively pick some of "good" ones and discard the "bad" ones. We consider DBN as component learners in our approach but we also compare it with SVM in our experiments. Since single deep neural network is easy to be over-fitting, we try to construct multiple neural networks and combine their predictions together to make a precise prediction. The main contribution can be summarized as follows:

- We propose a more robust Android malware detection approach based on selective ensemble learning and genetic algorithm, which is called SEdroid. Through comparative experiment, SEdroid showed a more robust and preminent capability in Android malware detection with a 98.3% precision rate and a 98.1% malware recall rate.
- The first time that selective ensemble learning is applied in Android malware detection. The algorithm we designed based on genetic algorithm considers both accuracy and diversity of the ensemble, which makes it easier to find the optimal ensemble combination and gives the model stronger generalization ability and robustness.

2 APPROACH

The general architecture is shown in Figure 1. First, we extract 3 types of features from Android application packages (.apk file) and vectorize them, which are android-permission, intent-action and API calls. Second, vectorized features are sampled to construct component neural networks using bootstrap sampling, which means that we randomly sample the data with replacement. Third, we designed a kind of genetic algorithm that consider both accuracy and diversity of the ensemble to find the best ensemble combination.

*Corresponding author.

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. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

CCS '19, November 11–15, 2019, London, UK

© 2019 Association for Computing Machinery.

ACM ISBN 978-1-4503-5693-0/19/11...\$15.00

<https://doi.org/10.1145/1122445.1122456>

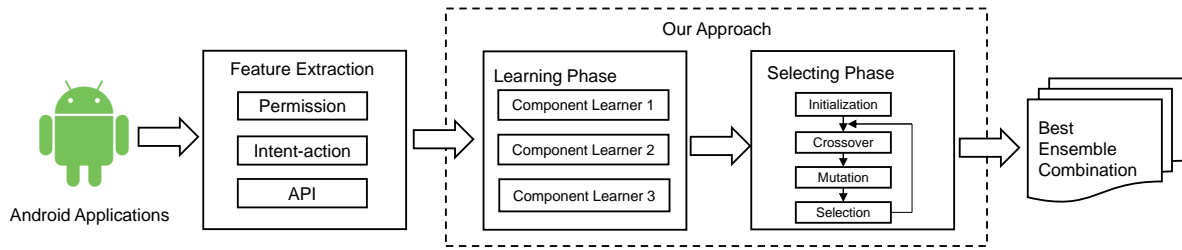


Figure 1: Overview of SEDroid

Finally, the selected component neural networks are combined to make up an ensemble via majority voting, which is used to detect Android malware.

2.1 Feature Extraction

Three sensitive features of Android applications are considered in our approach, including permission, intent-action and API. Permission feature is our first selected feature because all the sensitive permissions that a malware need to operate on sensitive data must be declared explicitly in the AndroidManifest.xml file. Intent-action is mainly used to assist interaction and messaging between different components, which has been researched that it is highly related to sensitive operation[3, 4]. API calls are the classes and methods that have been decompiled to smali grammar, which include low-level system calls. Hence, it's reasonable to take these three static features into consideration. After extracting features, we vectorize the derived features and concatenate them to the final feature vector.

2.2 Selective Ensemble Method

Selective ensemble method was first proposed by Zhou et al [5]. An ensemble is a combination of predictions by multiple component learners like neural networks, which are trained for the same task. It is researched that an ensemble of neural networks can enhance the model's generalization ability[6, 7] and an ensemble of some of the component learners may be better than all[5]. Genetic algorithm is introduced to pick out the best combination of component neural networks for ensemble. Here, we suppose that we have N component learners which is denoted as $f_i (i = 1, 2, \dots, N)$. Then $f_i(x)$ denotes the i th component learner's prediction on sample x . If the prediction is malicious, we denote it as 1. If not, we denote it as -1 . Bootstrap sampling method is adopted to generate the ensemble, which means we randomly sample the data with replacement for N times. To denote the ensemble effect mathematically, we suppose a weight vector $\omega(\omega_1, \omega_2, \dots, \omega_N)$ that is N -dimensional with each dimension having the value of either 0 or 1. Then our selective ensemble combination predicted by component learners based on majority voting can be given as follows.

$$H(x) = \text{sgn}\left(\sum_{i=1}^N \omega_i f_i(x)\right) \quad (1)$$

When ω_i equals 0, the related component learner f_i cannot make any effect on the final ensemble combination $H(x)$. And the component learners whose ω_i is 1 make the final ensemble prediction

about whether the application is benign or malicious. To find the appropriate weight vector ω , we introduce genetic algorithm to help us search the global optimal solution and obtain the best ensemble combination.

2.3 Genetic Algorithm

Genetic algorithm is adopted to obtain the best ensemble combination. The pseudocode of the genetic training algorithm of our approach is presented in Algorithm 1.

Algorithm 1 Pseudocode of genetic algorithm

- 1: $population \leftarrow$ generate pop_size chromosomes randomly
 - 2: **for** iteration = 1 to max_iter **do**
 - 3: crossover($population$)
 - 4: mutation($population$)
 - 5: evaluate the fitness function of population
 - 6: $new_population \leftarrow$ select_newpop($population$)
 - 7: $population \leftarrow new_population$
 - 8: **end for**
-

Population Initialization. The weight vectors that are mentioned in *Selective Ensemble Method* section are naturally the chromosomes owing to their inherent binary attributes. Hence, during this procedure we initialize the population by randomly generating multiple 0-1 vectors as chromosomes.

Fitness Function Evaluation. We consider two factors in our fitness function. One is ensemble accuracy, the other is the diversity of the component learners. It is researched that component learners with good accuracy and high diversity can achieve excellent performance than any of the component [8]. We define our fitness function $F(x)$ as follows:

$$F(x) = Accuracy \times D \quad (2)$$

In equation (2), *Accuracy* denotes the accuracy of an ensemble on the training data. *D* is the diversity factor of an ensemble. Suppose i and j are separate component learners of ensemble E , x_k is the k th sample of the data and $p_t(x_k)$ denotes the prediction that component learner t makes on sample x_k . Then we consider defining diversity factor D using Euclidean distance:

$$D(x) = \frac{1}{N} \sum_{i < j} \sqrt{\sum_{k=1}^M (p_i(x_k) - p_j(x_k))^2} \quad (3)$$

In equation(3), N denotes the number of the component learners of an ensemble. It can be seen that both the accuracy factor and diversity factor constrain the fitness function. If the predictions of component learners in an ensemble are prone to be the same, the diversity factor D will decrease to 0, which make the fitness value so small to survive to the next generation. But if the predictions of the component learners are totally different, it will cause the accuracy factor decrease to 0, making the ensemble hard to survive. Only the ensembles with prediction difference on different part of the data and high accuracy are inclined to survive in the genetic algorithm.

3 PRELIMINARY RESULTS

In this section, we will elaborate the process of our experiment and prove the effectiveness of our approach in real world dataset.

3.1 Datasets

We collected a dataset that is composed of 16000 real world Android applications in total, 8000 for malwares and 8000 for benign ones. Malicious Apps are obtained from Information Security Lab of Peking University and each of them has been detected with multiple malicious actions using Android malware detection platform like *VirusTotal*. Benign Apps are collected from Google Play, which have been manually checked. So we basically can ensure the correctness of our dataset. To simulate the circumstance in real world, we swap some proportion of malicious Apps and benign Apps, which is 10%. To distinguish this dataset from the original one, we will refer this dataset as *modified dataset*, and the other is referred as *original dataset*. SEdroid will be evaluated on both datasets. In our experiment, we split our dataset into three parts, 60% for training set, 20% for validation set and 20% for test set. For evaluation metrics, accuracy, precision rate, recall rate and F1-score are adopted to evaluate the model.

3.2 Evaluation on Original Dataset

Based on previous work, we consider SVM and DBN as component learners. As is shown in Table 1, when the dataset is correct or almostly correct, all the machine learning methods listed below can achieve good results on our original dataset. SVM is really useful, which achieves a 94.7% precision rate and 94.1% F1-score. DroidDeep[9] behaves even better, which has a 96.7% precision rate and 95.2% F1-score. For DBN and SVM, the selective ensemble method slightly enhance the performance at about 2%, and the ensemble of DBN (SEdroid) performs even better than the ensemble of SVM, which is probably caused by the stronger representation ability of DBN. In general, all machine learning methods above performs good on the original dataset, and our approach has a slightly superiority than DroidDeep.

3.3 Evaluation on Modified Dataset

We did 30 repeated experiments to evaluate SEdroid compared to DroidDeep on the modified dataset and the results are shown in Table 2. The performance of DroidDeep on modified dataset is not as good as on original dataset, which is caused by the data permutation in the modified dataset. The worst performance of DroidDeep is 89.7% of precision rate and 90.3% of recall rate. On

Table 1: Evaluation Results on Original Dataset

	Precision(M)	Recall(M)	F1-score(M)
SVM	0.947	0.936	0.941
DroidDeep (DBN)	0.967	0.937	0.952
Ensemble of SVM	0.960	0.972	0.966
SEdroid (DBN)	0.979	0.985	0.982

the other hand, the best performance of DroidDeep can reach to 96.7% of precision rate and 96.0% of recall rate. This suggests that the performance of DroidDeep has a relatively large variance on the modified dataset, which can be seen in Fig.6 intuitively. In contrast, the performance of SEdroid has superior performance on modified dataset than DroidDeep, which has a 98.3% precision rate and 98.1% recall rate on average, which proves superiority and robustness of SEdroid.

Table 2: Evaluation Results on Modified Dataset

	Precision	Recall	F1-score
DroidDeep (worst)	0.897	0.903	0.900
DroidDeep (best)	0.967	0.960	0.963
DroidDeep (average)	0.938	0.943	0.940
SEdroid (average)	0.983	0.981	0.982

4 CONCLUSION

In this paper, we propose an Android malware detection approach called SEdroid, which is based on selective ensemble learning and genetic algorithm. SEdroid takes into account both accuracy and diversity of the ensemble of component learners, which increases the model's generalization ability and robustness. In the evaluation, SEdroid managed to perform more effective and robust compared to other approaches.

REFERENCES

- [1] Justin Sahs and Latifur Khan. A machine learning approach to android malware detection. In *Intelligence and Security Informatics Conference*, pages 141–147, 2012.
- [2] Dali Zhu, Hao Jin, Ying Yang, Di Wu, and Weiyi Chen. Deepflow: Deep learning-based malware detection by mining android application for abnormal usage of sensitive data. pages 438–443, 2017.
- [3] Xin Su, Dafang Zhang, Wenjia Li, and Kai Zhao. A deep learning approach to android malware feature learning and detection. In *Trustcom/bigdata/iaANspa*, pages 244–251, 2017.
- [4] El Mouatez Billah Karbab, Mourad Debbabi, Abdelouahid Derhab, and Djedjiga Mouheb. Android malware detection using deep learning on api method sequences. 2017.
- [5] Zhi Hua Zhou, Jianxin Wu, and Wei Tang. Ensembling neural networks: Many could be better than all. In *ARTIFICIAL INTELLIGENCE*, 2002.
- [6] J. G. Carney, P. Cunningham, and U. Bhagwan. Confidence and prediction intervals for neural network ensembles. In *International Joint Conference on Neural Networks*, pages 1215–1218 vol.2, 1999.
- [7] David W. Opitz and Jude W. Shavlik. Generating accurate and diverse members of a neural-network ensemble. *Advances in Neural Information Processing Systems*, 8:535–541, 1996.
- [8] Sunil Choenni. Design and implementation of a genetic-based algorithm for data mining. In *VLDB 2000, Proceedings of International Conference on Very Large Data Bases, September 10–14, 2000, Cairo, Egypt*, pages 33–42, 2000.
- [9] Zi Wang, Juecong Cai, Sihua Cheng, and Wenjia Li. Droiddeeplearner: Identifying android malware using deep learning. In *Sarnoff Symposium, 2016 IEEE*, pages 160–165, 2017.