

AN EFFICIENT ANDROID MALWARE PREDICTION USING ENSEMBLE MACHINE LEARNING ALGORITHMS

Neamat Al Sarah¹ Fahmida Yasmin Rifat¹ Md. Shohrab Hossain² Husnu S. Narman³

¹Department of Computer Science and Engineering, Military Institute of Science and Technology, Mirpur, Dhaka, Bangladesh.

²Department of Computer Science and Engineering, Bangladesh University of Engineering and Technology, Bangladesh.

³Weisberg Division of Computer Science, Marshall University, Huntington, WV, USA.



Bangladesh University of Engineering and Technology



Presentation Outline









Introduction



- ✓ Lack of trustworthiness review methods, developers can upload their Android apps including repackaged apps, ransomware, or trojans to the market easily in even Google's Android market
- ✓ Which posed serious threats to the smart phone users, such as stealing user credentials, auto-dialing premium numbers, and sending SMS messages without user's concern







Bangladesh University of Engineering and Technology







□ To automate the process of malware detection

□ To detect malware within short time

□ To achieve high accuracy with minimum number of false positives

□ To be able to detect malware from all families

□ To find out the best algorithm for predicting malware



Bangladesh University of Engineering and Technology



Brief Literature Review

Static Analysis

Relies on features extracted without executing code

Analysis for malware detection is done without running the app Static features such as manifest file components, API calls are used

Low resource consumption, fast detection and low real time requirements

Dataset



UNIVERSITY

Eight Features of DREBIN dataset

S1Hardware Components72S2Requested Permission3,812S3App contents218,952S4Filtered Intents6379S5Restricted API Calls733S6Used Permission70
S2Requested Permission3,812S3App contents218,952S4Filtered Intents6379S5Restricted API Calls733S6Used Permission70S7Suspicious API215
S3App contents218,952S4Filtered Intents6379S5Restricted API Calls733S6Used Permission70S7Suspicious API215
S4Filtered Intents6379S5Restricted API Calls733S6Used Permission70S7Suspicious API215
S5Restricted API Calls733S6Used Permission70S7Suspicious API215
S6Used Permission70S7Suspicious APL215
S7 Suspicious API 215
Calls
S8 Network Address

- ✓ We use DREBIN Dataset
- ✓ Contains 5560 applications from 179 different malware family
- ✓ Collected in the period of August 2010 to October 2012
- ✓ Available to us by the Mobile Sandbox project.



Methodology





Our Approach: Key Features





Ban Engir

Engineering and Technology



Feature Selection using RFE







Bangladesh University of Engineering and Technology

Significant Features



TO



Bangladesh University of Engineering and Technology

...



Why LightGBM?

Angladest Technology For advancement

- ✓ Gradient Boosting Ensemble Algorithm
- ✓ Fast, distributed, high-performance gradient boosting framework based on decision tree algorithm
- ✓ Grows tree leafwise while other algorithm grows level wise.
- \checkmark Chooses the leaf with max delta loss to grow.
- ✓ When growing the same leaf, Leaf-wise algorithm can reduce more loss than a level-wise algorithm



Parameter Tuning

- Exploring a range of possibilities
- Parameters are crucial
- Finding an optimal combination of parameters that minimizes a predefined loss function to give better results

For Faster Speed

- bagging by setting bagging_fraction and bagging_freq
- •feature sub-sampling by setting feature_fraction
- •small max_bin
- •save_binary to speed up data loading in future learning



Bangladesh University of Engineering and Technology





Parameter Tuning

For Better Accuracy

- ✓ large max_bin (may be slower)
- ✓ small learning_rate with large num_iterations
- ✓ large num_leaves (may cause over-fitting)
- ✓ bigger training data
- ✓ dart

Deal with Over-fitting

- ✓ small max_bin
- ✓ small num_leaves
- ✓ min_data_in_leaf and min_sum_hessian_in_leaf
- ✓ bagging by set bagging_fraction and bagging_freq
- ✓ feature sub-sampling by set feature_fraction
- ✓ bigger training data
- ✓ lambda_l1, lambda_l2 and min_gain_to_split for regularization
- ✓ max_depth to avoid growing deep tree













Evaluation Criteria & Results

1. Accuracy:







Thank You!



Bangladesh University of Engineering and Technology

