

# Contaminant Removal for Android Malware Detection Systems

Lichao Sun<sup>\*</sup>, Xiaokai Wei<sup>†</sup>, Jiawei Zhang<sup>‡</sup>, Lifang He<sup>§</sup>, Philip S. Yu<sup>\*</sup> and Witawas Srisa-an<sup>¶</sup>

<sup>\*</sup>University of Illinois at Chicago, Chicago, IL <sup>†</sup>Facebook, Menlo Park, CA

<sup>‡</sup>IFM Lab, Florida State University, FL <sup>§</sup>Cornell University, New York City, NY

<sup>¶</sup>University of Nebraska - Lincoln, Lincoln, NE

Email: {lsun29, xwei2, psyu}@uic.edu, {jwzhanggy, lifanghescut}@gmail.com, witty@cse.unl.edu

**Abstract**—A recent report indicates that there is a new malicious app introduced every 4 seconds. This rapid malware distribution rate causes existing malware detection systems to fall far behind, allowing malicious apps to escape vetting efforts and be distributed by even legitimate app stores. When trusted downloading sites distribute malware, several negative consequences ensue. First, the popularity of these sites would allow such malicious apps to quickly and widely infect devices. Second, analysts and researchers who rely on machine learning based detection techniques may also download these apps and mistakenly label them as benign since they have not been disclosed as malware. These apps are then used as part of their benign dataset during model training and testing. The presence of contaminants in benign dataset can compromise the effectiveness and accuracy of their detection and classification techniques.

To address this issue, we introduce PUDROID (Positive and Unlabeled learning-based malware detection for Android) to automatically and effectively remove contaminants from training datasets, allowing machine learning based malware classifiers and detectors to be more effective and accurate. To further improve the performance of such detectors, we apply a feature selection strategy to select pertinent features from a variety of features. We then compare the detection rates and accuracy of detection systems using two datasets; one using PUDROID to remove contaminants and the other without removing contaminants. The results indicate that once we remove contaminants from the datasets, we can significantly improve both malware detection rate and detection accuracy.

**Keywords**-Mobile Security; Malware Detection; Noise Detection; Android Malware; PU Learning;

## I. INTRODUCTION

Android is currently the most used smart-mobile device platform in the world, occupying 87.6% of market share and over 1.4 billion Android devices in deployment [1]. Unfortunately, the popularity of Android also makes it a popular target for cyber-criminals to create malicious apps that can steal sensitive information and compromise systems [2]. During the first three months of 2016, Kaspersky Lab uncovered over 2 million malware samples including trojans, worms, exploits, and viruses. On average, a malicious app is introduced in every 3.79 seconds [3]. Some types of malicious apps have more than 50 variants, making detecting all of them very challenging [4].

There have been several approaches to detect these malicious Android apps. Most approaches focus on the attack be-

haviors, and use static or dynamic analysis to build detection

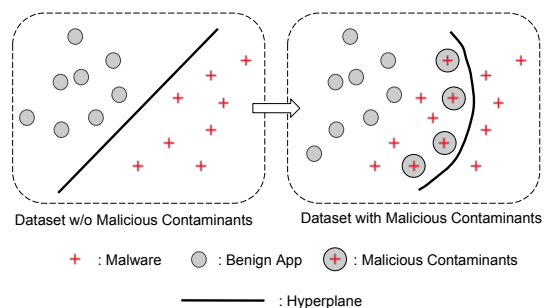


Figure 1. Left figure shows machine learning can classify malware and benign apps well without malicious contaminants. Right figure shows that the machine learning cannot work well for malware detection with malicious contaminants. tools that rely on approaches known to work well for desktop environments [5]. However, static analysis approaches in general can produce a large number of false positives while dynamic analysis approaches need adequate input suites to sufficiently exercise execution paths. Therefore, neither of them will work well for Android malicious app detection. Another emerging approach is to build detection techniques based on data mining and machine learning techniques [6], [7], [8].

For example, DREBIN [6] utilizes multi-view features by combining static analysis and supervised learning to accurately detect malware. SIGPID [7] improves upon DREBIN [6] by using many more features for training and detection. DROIDCLASSIFIER [8] uses traffic flow information and unsupervised learning to detect the malware and classify the family of each malicious app.

When machine learning techniques are used to help with malware detection, the detection effectiveness and accuracy are highly dependent on the quality of the training datasets. To create such dataset, researchers typically label a set of malicious apps and a set of benign apps. To build the malicious dataset, researchers manually label these malicious apps one by one based on known information from various malware analysis and collection sources (e.g., virusshare.com). To build the benign dataset, researchers download apps from trusted distribution sources such as Play Store and verify that those apps have not been recently disclosed as malware. However, as previously mentioned,

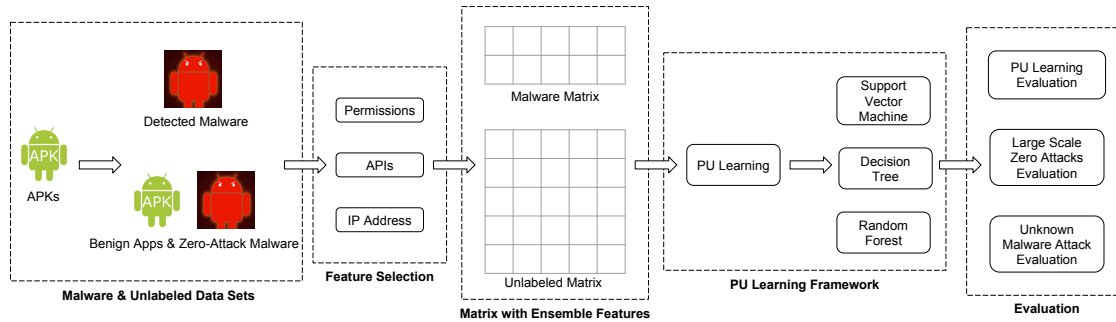


Figure 2. An Overview of PUDroid Approach

malicious apps are created every few seconds so malware detection sites often fall behind in disclosing new malware. As such, trusted distribution sites have been known to distributed repackaged malware [9].

Currently, repackaging benign apps with existing malicious components is the leading approach employed by cybercriminals to create Android malware [1]. Thus, it is quite common for undisclosed malware to have previously known basic attack behaviors. However, the sheer volume of recently created malware also makes it possible for benign dataset to contain undisclosed malware as illustrated in the right side of Fig. 1. The presence of these malicious apps in a benign dataset (we refer to these malicious apps as contaminants) can negatively affect the accuracy and effectiveness of machine learning-based detection tools. As such, we need a mechanism to effectively remove such contaminants from benign dataset.

In this paper, we present PUDROID, an approach that leverages the *Positive and Unlabeled (PU)* learning to remove contaminants as part of building benign dataset for robust malware detection in Android. Specifically, we assume that the positive group contains malicious apps, and the unlabeled group contains benign apps and some contaminants consisting of unlabeled malicious apps. We make use of a feature selection strategy to build the desired input for PU learning. In a nutshell, PUDROID works as an additional validator to ensure that presumably benign apps downloaded from trusted distribution channels are indeed benign. That is, they do not contain any undisclosed but detectable malicious components due to repackaging. We then conduct empirical evaluation of PUDROID and find that it allows to remove nearly 100% of all contaminants as part of building benign dataset.

An overview of our PUDROID approach is illustrated in Fig. 2. The main contributions are summarized as follows:

(1) To the best of our knowledge, PUDROID is the first attempt to use PU learning to remove contaminants as part of building benign dataset for robust malware detection in Android.

(2) We introduce a feature selection strategy to analyze and select “explainable” features that are useful for effective malware detection. Our strategy reduces the number of embedding features by 93% when we compare to those used by existing approaches, while yielding nearly the same level of effectiveness in malware detection.

(3) We evaluate PUDROID using a large dataset containing 5,560 malware samples with 2,200 features. We also investigate a scenario in which PUDROID needs to detect different magnitudes and types of contaminants. The results show that PUDROID is very effective in such situations.

(4) We also evaluate PUDROID with different classification methods to fine tune the accuracy of PUDROID. The results indicate that the performance of a classification model is sensitive to the number of malware samples present in a dataset.

## II. BACKGROUND AND MOTIVATION

In this section we briefly introduce common machine learning approaches employed by researchers to build malware detection frameworks for Android. These techniques employ many features ranging from static information such as permissions to dynamic information such as network traffics.

### A. Permission-Based Learning

Android security system offers permission control mechanisms as one of most important components. Android app developers need to declare the permissions for each app to have access to resources such as text messaging and address book. The declaration can be found in the manifest file. When users try to install apps, they can choose to approve or decline the requested permissions.

In the case of installing apps from third party stores, users may need to root their devices, making their systems more vulnerable. Furthermore, we have also seen that developers tend to request more permissions than the apps actually need. These behaviors make Android security mechanism vulnerable to malicious attacks. To help a user determine

if requested permissions can make a system vulnerable to attacks, Android provides "protection Level" to help characterize the potential risks of these permissions [10]. The four protection levels are "normal", "dangerous", "signature", and "signature or system". Permissions such as "WAKE\_LOCK" which keep processor from sleeping or screen from dimming are considered as low risk. However, some permissions such as "WRITE\_SMS" are considered dangerous due to its ability to leak information and cause financial damages such as texting to premium services. Work by Sun et al. [7] finds that normal permissions, in addition to those classified as "dangerous" can also significantly contribute to malicious behaviors.

### B. API-Based Learning

Previously, we introduced the Android permission system and how it can be used to detect the potential malware. However, permission-based approaches may not always yield accurate results as many permissions are commonly used by both benign and malicious apps. To improve accuracies of these approaches, API (Application Programming Interface) information can also be used to add context that can help distinguishing between benign and malicious apps [6]. For example, by analyzing calling contexts leading to dangerous APIs such as WRITE\_SMS, one can compare benign calling contexts and malicious calling contexts.

### C. Network Information Based Learning

Dynamic information such as network traffic can be used as learning feature to detect the malware. To do so, input generation techniques such as Monkey are used to generate inputs to execute benign and malicious programs. The network traffic is then recorded and then analyzed as features for machine learning. For example, Shabtai et al. [11] present a Host-based Android malware detection system to target the repackaging attacks. They conclude that deviations of some benign behaviors can be regarded as malicious ones. Narudin et al. [12] introduce a TCP/HTTP based malware detection system. They extracted basic information, (e.g. IP address), content based, time based and connection based features to build the detection system. Their approach can determine if an app is malicious or not.

### D. Motivation

In this paper, we propose a PU learning-based approach (named as PUDROID) that aims to remove contaminants as part of building benign dataset for robust malware detection in Android. We show that PUDROID can protect the system when the benign dataset are infected by malware. To improve the performance of the PUDROID and to prevent various malicious behavior, we also use embedding features including permissions, APIs, and network information to efficiently detect the malware.

## III. PUDROID APPROACH

In this section, we introduce the main steps of our PUDROID approach including feature selection and PU learning. First, we present a feature selection strategy to help us find the informative features from a variety of Android features. Then, we use the resulted features to leverage PU learning process along with different classifiers to remove contaminants for robust malware detection.

### A. Feature Selection for Android Data Generation

Android apps can collect malicious and benign datasets, which usually contain various types of features such as permissions, APIs, IP address, activities, requested URLs, and services. However, not all of these features are effective for malware detection and many features such as activities, requested URLs, and services are difficult to explain due to large variability in these features. We need to perform feature selection to find the informative features from these available features. By using domain knowledge of security, our feature selection strategy only uses explainable and helpful features to detect malware. In particular, we only focus on permissions, APIs, IP address and URLs. We first convert the URLs into IP address, and then introduce our feature selection strategy on how to select features from permissions, APIs and IP address features.

IP address can provide similar information as URLs but with fewer variability. For example, multiple URLs can point to the same IP address if they are alias. Part of IP address, e.g., the three most significant bytes, can also provide company information (e.g., 216.59.192.xx tells us that the IP address in this range belong to Google). Additionally, note that both IP address and URLs can suffer from spoofing but IP address provides the same information with much fewer data points. Therefore, it is beneficial to convert the URLs to a more effective representation (IP address). To achieve this goal, we remove the invalid URLs as they tend not to have corresponding IP address, and use the socket API to acquire IP address.

Based on the above results, we define an unbalanced feature selection strategy for malicious and benign datasets. A large difference between the numbers of malware and benign samples can lead to unbalanced removal of less contributing features. To overcome this issue, our feature removal process biases the criteria based on the ratio between the number of benign samples and malware samples. For example, if there are twice as many benign samples as malware samples, we then remove a seldom occurring feature from malware dataset based on a threshold  $tm$ . For example, if feature occurs 5 times when we consider the entire malware dataset, we would then set  $tb$  for removal for the benign dataset to two (i.e., twice as many as  $tm$  or 10 times). This can be formulated as:

$$\frac{tm}{tb} = \eta \cdot \frac{\#benign\ samples}{\#malware\ samples} \quad (1)$$

Table I  
LIST OF BASIC SYMBOLS

| Symbol          | Definition and description  |
|-----------------|---|
| $\mathbf{P}$    | positive group/marked malware set   |
| $\mathbf{U}$    | unlabeled group/mixed malware and benign apps   |
| $\mathbf{1}$    | a vector of all ones  |
| $\mathbf{x}(s)$ | resulted feature vector of an app $s$   |
| $z(s)$          | 1 means the app $s$ is labeled, 0 otherwise   |
| $y(s)$          | 1 means the app $s$ is true malware, 0 otherwise  |
| $p(\cdot)$      | the probability of an app   |
| $f(\cdot)$      | malicious probability of an app without PU learning   |
| $g(\cdot)$      | malicious probability of an app with PU learning  |
| $M_d$           | classifier model without PU learning, if $f(\cdot) > 0.5$ , the app is malicious, otherwise is benign |
| $M_h$           | classifier model with PU learning, if $g(\cdot) > 0.5$ , the app is malicious, otherwise is benign    |
| $\sim$          | denotes the equivalent relation   |

where  $\eta$  is the coefficient that controls threshold selection and it is usually at least 2. Here we set  $\eta = 2$ , which means any feature must be used at least two malware and  $tm/tb$  benign apps will not be removed after feature selection.

After applying above feature selection strategy, we use  $\mathbf{x}_1$ ,  $\mathbf{x}_2$ , and  $\mathbf{x}_3$  to represent the selected feature subsets of permissions, APIs, and IP address, respectively. Our final resulted feature set is:

$$\mathbf{x} := \mathbf{x}_1 \cup \mathbf{x}_2 \cup \mathbf{x}_3$$

### B. PU Learning For Malware Detection

As mentioned earlier, unreliable negative examples (or contaminants) can be unknowingly included in training dataset due to the prolific rate of malware creation [9]. We leverage PU learning to detect and remove these contaminants from training datasets. Table I lists the basic symbols that will be used throughout this section.

PU learning is a semi-supervised technique for building a binary classifier on the basis of positive and unlabeled samples only. It is useful when we have not determined if an app is malicious or benign. In small collections of apps, labeling an app as malicious or benign can be done manually without too much effort. In large collections, however, manual identification may not be feasible. In this case, PU learning can be applicable. In addition, in a scenario that the identification process may produce inaccurate results (e.g., mistakenly identifying malicious apps as benign apps), PU learning can also help to identify and remove these unreliable negative samples.

In order to use PU learning for malware detection, we divide our dataset into the positive group ( $\mathbf{P}$ ) and the unlabeled group ( $\mathbf{U}$ ), where the positive group contains malicious apps, and the unlabeled group contains benign apps and some contaminants consisting of unlabeled malicious apps. To differentiate positive and unlabeled apps, we define “discovery state” ( $z$ ) to indicate whether an app is labeled or not in the dataset. For a given app  $s$  in the group  $\mathbf{P}$ , if  $s$  is marked as malicious, then  $z(s) = 1$ ; otherwise,  $z(s) = 0$ . As

a result, the “discovery states” of apps in groups  $\mathbf{P}$  and  $\mathbf{U}$  are:  $z(\mathbf{P}) = \mathbf{1}$  and  $z(\mathbf{U}) = \mathbf{0}$ . Besides, each app has another label called “hidden malware state” ( $y$ ), which can expose whether an app is actually malicious or benign. Here an app with 1 is a known malware app, and 0 is a benign app. For example, if an app  $s$  has been detected as malware, then  $y(s) = 1$ ; otherwise  $y(s) = 0$ . Based on this, we can check the “hidden malware state” of each app in group  $\mathbf{P}$  and group  $\mathbf{U}$ . Since every app in group  $\mathbf{P}$  is a known malware app,  $y(\mathbf{P}) = \mathbf{1}$ . However,  $y(\mathbf{U})$  can be either  $\mathbf{1}$  or  $\mathbf{0}$  as both malicious contaminants and benign apps can be in  $\mathbf{U}$ . As such, we have:

$$p(z(s) = 1 | \mathbf{x}(s), y(s) = 0) = 0. \quad (2)$$

where  $p(\cdot)$  is the probability, and  $\mathbf{x}(s)$  is the feature vector extracted for app  $s$  (i.e.,  $[0, 1, 0, 0, \dots, 1]$ , 1 denotes that the app requests the permission and 0 otherwise).

The goal of malware detection is to build a malware discovery model  $M_d \sim f(\mathbf{x}(s)) : \mathbb{R}^d \rightarrow \{1, 0\}$  from  $\mathbf{P}$  and  $\mathbf{U}$ , where  $d$  is the dimension of the  $x(s)$ , which represents the number of the features. 1 means  $s$  is identified as malware, 0 means  $s$  is not identified as malware. Given an app  $s$ , by applying  $M_d \sim f(\mathbf{x}(s))$ , the discovery probability of  $s$  as a malware is:  $p(z(s) = 1 | \mathbf{x}(s))$ . While our ultimate goal is to infer the true label of a given malicious app  $s$  (i.e.,  $y(s)$ ). Besides the *discovery probability*, we also need to build a hidden malware detection model  $M_h \sim g(\mathbf{x}(s))$  based on the detected malware and the unlabeled app sets. Formally, we define the probability that an app  $s$  is indeed malicious (i.e.,  $y(s) = 1$ ) as:  $p(y(s) = 1 | \mathbf{x}(s))$ . In the following, we discuss the details.

**Assumption:** (Malware Discovered at Random): Assume malware samples are randomly detected by analysis, the probability of detection is not relevant with the features created, then it has:

$$p(z(s) = 1 | \mathbf{x}(s); y(s) = 1) = p(z(s) = 1 | y(s) = 1). \quad (3)$$

To build the proposed hidden malware detection model  $M_h$ , we should ideally know which apps are actually malware. For this purpose, we first prove Lemma 1 [13].

*Lemma 1:* Suppose the “Malware Discovered at Random” assumption holds, then

$$p(y(s) = 1 | \mathbf{x}(s)) = \frac{p(z(s) = 1 | \mathbf{x}(s))}{p(z(s) = 1 | y(s) = 1)}. \quad (4)$$

*Proof:* By holding the assumption, we have:

$$\begin{aligned} & p(z(s) = 1 | \mathbf{x}(s)) \\ &= p(z(s) = 1 | \mathbf{x}(s)) \cdot p(y(s) = 1 | \mathbf{x}(s), z(s) = 1) \\ &= p(y(s) = 1, z(s) = 1 | \mathbf{x}(s)) \\ &= p(y(s) = 1 | \mathbf{x}(s)) \cdot p(z(s) = 1 | y(s) = 1, \mathbf{x}(s)) \\ &= p(y(s) = 1 | \mathbf{x}(s)) \cdot p(z(s) = 1 | y(s) = 1). \end{aligned} \quad (5)$$

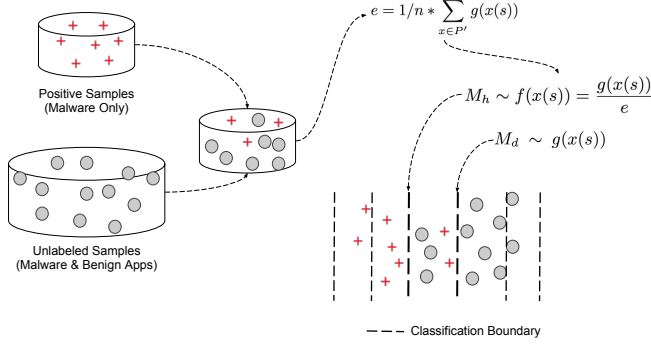


Figure 3. An Overview of Our PU Learning Process

Then, by dividing both sides of Eq. (5) by  $p(z(s) = 1|y(s) = 1)$ , we arrive at Lemma 1. ■

From Eq. (4), if we want to build  $M_h$  by calculating the *hidden malware probability*, we can use  $f(\mathbf{x}(s))/p(z(s) = 1|y(s) = 1)$ , where  $f(\mathbf{x}(s))$  is the malware probability of the an app  $s$  of  $M_d$ . Specifically, we have:

$$f(\mathbf{x}(s)) = p(z(s) = 1|\mathbf{x}(s)). \quad (6)$$

On the other hand,  $p(z(s) = 1|y(s) = 1)$  can be calculated by a validation set with applying  $M_d$ . First, we randomly select apps from  $\mathbf{P} \cup \mathbf{U}$  to be the validation set  $\mathbf{V}$ , and choose the subset  $\mathbf{P}'$  from  $\mathbf{V}$  with the positive label, i.e.  $z(\mathbf{P}') = 1$ . The estimator of  $p(z(s) = 1|y(s) = 1)$  is the average value of  $g(\mathbf{x}(s))$  for  $\mathbf{x}(s)$  in  $\mathbf{P}'$ , then we have:

$$p(z(s) = 1|y(s) = 1) \sim e = 1/n \cdot \sum_{\mathbf{x} \in \mathbf{P}'} f(\mathbf{x}(s)). \quad (7)$$

where  $n$  is the cardinality of  $\mathbf{P}'$  and the estimator  $e$  is the average value of  $f(\mathbf{x}(s))$  for  $\mathbf{x}$  in  $\mathbf{P}'$ . Since  $e$  is based on a certain number of data instances, it has a low variance and is preferable in practice [13]. Notice that to compute  $e$ , we need to specify the size of the set  $\mathbf{P}'$  and the size of the validation set  $\mathbf{V}$ . We explain how to set the size and evaluate our system in Section IV.

With  $e$  and the classifier model  $M_d \sim f(\mathbf{x}(s))$  on labels  $z(s)$ , we can adjust to a classifier  $M_h \sim g(\mathbf{x}(s))$  on relation labels  $y(s)$  as follows:

$$M_h \sim g(\mathbf{x}(s)) = p(y(s) = 1|\mathbf{x}(s)) = \frac{f(\mathbf{x}(s))}{e}. \quad (8)$$

Fig. 3 shows an overview of our PU learning process.

**Discussion:** Here we give an example to show how the Eq. (8) works for positive and unlabeled datasets. If all apps can be clearly separated into malicious and benign groups, (i.e., no malicious apps in the unlabeled group) the malicious probability ( $g(\mathbf{x}(m))$ ) of most random malware  $m$  in positive group should be close to 1. However, when the unlabeled group contains a large number of malicious contaminants, the malicious probability of most random

malware  $m$  in positive group will decrease. The worst case is that the  $g(\mathbf{x}(m))$  close or even less than 0.5. We know that if  $g(\mathbf{x}(m)) < 0.5$ , the malware  $m$  will be classified as a benign app by the detector.

To address this problem, when we first use a training set to build the detector, we select a subset  $\mathbf{P}_M$  from positive group where  $y(\mathbf{P}_M) = 1$ . Then we calculate the average malicious probability of  $\mathbf{P}_M$  by using  $g(\mathbf{x}(\mathbf{P}_M))$ . If we find the average malicious probability of  $\mathbf{P}_M$  is close to 0.5 rather than 1.0, we increase the malicious probability for each app in the training set. Then more malicious contaminants will be detected from the unlabeled group. If the system is faced with malicious contaminants, the malicious probability of every app will decrease including benign apps by the detector. So most benign apps will still be classified as benign after we increase the malicious probability for them. This is the framework of PU learning to prevent the malicious contaminants.

In order to apply PU learning for the dataset with group  $\mathbf{P}$  and group  $\mathbf{U}$ , we need to implement the learning algorithms to build the malware detection system. Various learning algorithms have been used in mobile security before [7], we choose three most frequently used learning algorithms: Support Vector Machine, Decision Tree and Random Forest.

#### IV. EVALUATION

We conducted experiments to evaluate the performance of PUDROID in identifying contaminants in Android datasets. We conducted our experiments using a large-scale real-world collection of apps. Our experiments were done to answer four research questions (RQs).

**RQ 1:** How effective is PUDROID in removing malicious contaminants from benign dataset? In this experiment, we create contaminated benign dataset and then compare the performances between a system with PUDROID and without PUDROID.

**RQ 2:** How do changes in the number of malicious contaminants in benign data sets affect the performance of PU-DROID? In this experiment, we systematically increased the number of malicious contaminants in our benign datasets.

**RQ 3:** How effective is PUDROID in removing unknown malicious contaminants from benign dataset? In this experiment, we include a family of malware into each benign dataset and evaluate if PUDROID can detect these contaminants.

**RQ 4:** How effective is PUDROID in removing benign contaminants from malicious dataset? Our first two focuses were on building accurate benign dataset. In this experiment, we reverse the contamination pattern by including benign apps in the malicious dataset and observe the performance of PUDROID to detect these benign contaminants.

In all experiments, we then compared our results with those from several state-of-the-art detection methods.

### A. Data Sets

We used the dataset based on prior work by Arp et al. [6]. The dataset includes 5,560 malware and 123,453 benign apps. We then extracted relevant features using a strategy that differs from theirs. In their work, they used many features including permission and API information. However, not all features that they used are effective for malware detection and many features such as activities, requested URLs, and services are difficult to explain due to large variability in these features (e.g., there are over 10,000 activities and services and 200,000 URLs in the dataset). Having too many features can result in very sparse matrix, leading to high machine learning overhead and over-fitting issues. Then we can apply feature selection to select informative features.

After feature selection, our system kept 2,200 features. Drebin, on the other hand, used over 300,000 features. Even with the reduced number of features, our dataset contains too many dimensions. To address this problem, we use PCA to project the whole dataset to 2-dimension representation of the dataset, Fig. 4 shows the distribution of our whole dataset. We can find most malware and benign apps can be separated even with low dimension representation with our features. That's the reason the machine learning can help to separate the malware and benign apps when dataset is without any contaminants.

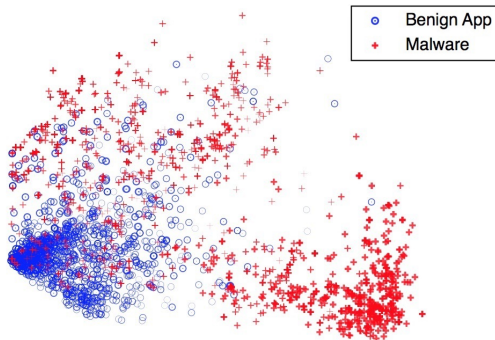


Figure 4. PCA Distribution of Whole Dataset

## V. RESULTS

In this section, we report the experimental results to answer the four proposed research questions. We used three measures, *accuracy*, *AUC* and *F-measure* to evaluate the performance of PUDROID.

### A. RQ 1: Effectiveness of PUDroid

In order to evaluate PUDROID, we first choose 1/3 of all malware samples and 1/3 of benign apps as the testing set. We then use the remaining malware and benign apps as the training set. To create an unlabeled group, we also include malware into the benign apps in multiple iterations using the following process. Initially, we have no malicious

contaminants, meaning there is no unlabeled samples in this case. Every sample is negative in the unlabeled group, which indicates benign sample. In the next iteration, we randomly select 100 malware samples from the training set as the malicious contaminants. We then add these malware samples to the benign portion of the training set and remove them from the malware portion of the training set. That is, we make these samples negative. Now, we have a positive group and unlabeled group which contains both malware and benign apps. We then repeat above step and randomly select more malware samples in the subsequent iterations. In summary, to build unlabeled dataset, we randomly include in the benign dataset,  $100 \cdot N$  malware samples as malicious contaminants in the  $N_{th}$  iteration.

We report the performance of PUDROID using our dataset in Fig. 5. Note that, due to limited space, we only report *AUC* values, but the conclusion is consistent across other measures. Our evaluation is based on the comparison of the performance of PUDROID and a classification method (denoted as PU in the figure) with that of using the same classification method alone (denoted in the figure as No PU or NPU). We used all three classification learning methods. Based on the results, we can make the following conclusions.

- As shown in Fig. 5, PUDROID can work well with all three classification methods.
- If there are no malicious contaminants in the testing dataset, PUDROID yields the same performance as using methods without PU learning.
- The performance of PUDROID with Support Vector Machine (SVM) is initially stable but degrades quickly when the training dataset contains a large number of malicious contaminants.
- The performance of Random forest is slightly better than that of Decision Tree.
- When the dataset contains only a small number of malicious contaminants, both of PUDROID system and methods without PU learning work well.

### B. RQ 2: Variable Levels of Contaminants

In this section, we evaluate the malware detection performance of PUDROID when the benign dataset contains a large portion of contaminants. To do so, we apply the following method to create the unlabeled dataset.

In  $N_{th}$  case, we set the number of malicious contaminants in the training dataset to be  $N$  times the number of malware in the training dataset. For example, in the first case, we set the ratio of the number of malicious contaminants and the number of malware in the training set to be 1:1. In the  $N_{th}$  case, the ratio becomes  $N : 1$ . The resulting unlabeled dataset should test our system's ability to deal with a large number of unlabeled malware samples.

In previous section, we showed that PUDROID with SVM can work well when the number of malicious contaminants

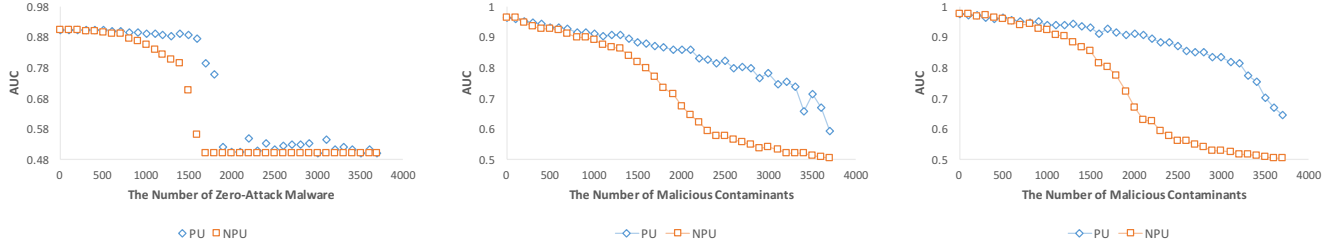


Figure 5. AUC of PUDroid with SVM (left), Decision Tree (center), Random Forest (right)

Table II  
EVALUATION OF HIGHLY CONTAMINATED DATASETS

| Scale Rate | Random Forest/PU | Random Forest/NPU | Difference | Decision Tree/PU | Decision Tree/NPU | Difference |
|------------|------------------|-------------------|------------|------------------|-------------------|------------|
| 100%       | 83.32%           | 44.25%            | 39.07%     | 73.18%           | 47.06%            | 26.12%     |
| 200%       | 76.39%           | 14.30%            | 62.10%     | 62.81%           | 16.46%            | 46.36%     |
| 300%       | 71.51%           | 8.69%             | 62.82%     | 57.15%           | 11.83%            | 45.32%     |
| 400%       | 64.91%           | 6.19%             | 58.72%     | 51.07%           | 7.97%             | 43.10%     |
| 500%       | 64.74%           | 3.32%             | 61.42%     | 49.70%           | 5.86%             | 43.85%     |
| 600%       | 59.70%           | 3.47%             | 56.24%     | 47.23%           | 4.79%             | 42.44%     |
| 700%       | 54.22%           | 2.50%             | 51.71%     | 44.32%           | 4.66%             | 39.65%     |
| 800%       | 52.14%           | 2.14%             | 50.00%     | 43.54%           | 3.82%             | 39.72%     |

is small to moderate. However, it does not work well when the number of malware is large. As such, we only evaluate Random Forest and Decision Tree in this scenario.

From Table II, we evaluate the ratios from 1:1 to 8:1, and report the malware detection rate (also referred to as True Positive Rate or Recall). Based on the results, we have the following observations:

- Random Forest or Decision Tree alone cannot work well when they face a large number of malicious contaminants. For example, when the ratio is 3:1, the detection rate of Random Forest and Decision Tree are 8.69% and 11.83%, respectively. However, when we applied PUDROID, the detection rates increases to 71.51% and 57.15%, respectively.
- PUDROID with Random Forest yields higher detection rates than those produced by PUDROID with Decision Tree. However, if the PU learning process is not used, Decision Tree performs better than Random Forest. For example, Decision Tree yields the detection rate of 11.83% while Random Forest only yields 8.69% when the ratio is 3:1.
- In extreme situations, Random Forest yields a detection rate of 52.14% when the ratio is 8:1. To put this into perspective, this result is better than using Random Forest alone when the ratio is 1:1.

In summary, PUDROID is very effective in identifying contaminants in highly contaminated datasets.

### C. RQ 3: Unknown Contaminants

In the dataset, we have 178 families of malware, but many families only have one malware sample. We only

found that only 5% of the total families contains more than 100 malware and about 50% of families contain 3 samples or fewer. Such a distribution pattern indicates that some families of malware are less likely to be repackaged and redeployed as new variations while others are popularly repackaged and redistributed as new variations.

While the focus of this particular investigation is on PUDROID’s performance when the benign dataset is contaminated with unknown malware, we also want the contaminants to be representative of the current practice of creating variations of malicious behaviors to promote rapid infections of devices. As such, our contaminants are from families that have more than 300 samples in our malware dataset. In this case, the families include *FakeInstaller*, *DroidKungFu*, *Plankton*, *Opfake*, *GinMaster* and *BaseBridge*. Next, we briefly describe the malicious behaviors of these families.

- *FakeInstaller* is the malware family with the largest number of variations in our dataset (925 malware samples). Malware authors simply repackaged commonly distributed apps (e.g., *Facebook*) with malicious functionalities. These malicious apps send SMS messages to premium rate numbers, without the users consent, passing itself off as the installer for a legitimate application. As previously reported, over 60% of Android malware samples processed by McAfee belong to this family [14].
- *DroidKungFu* family exploits several known but unpatched vulnerabilities in earlier Android versions to gain root access to a device and steal sensitive information. Stolen information is sent to remote command

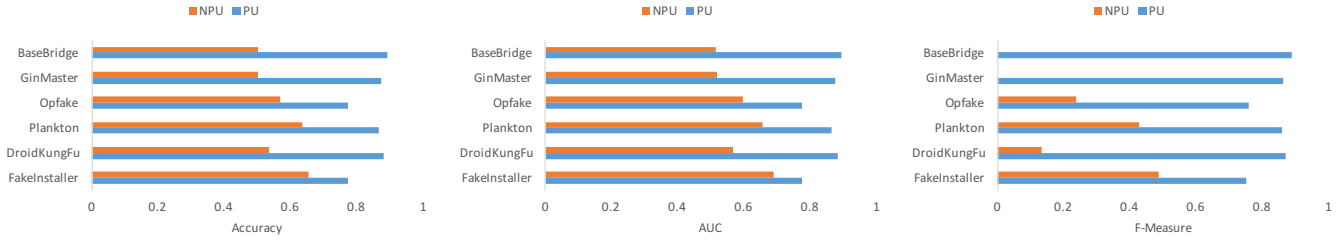


Figure 6. Performance of PUDroid when faced with new malware contaminants

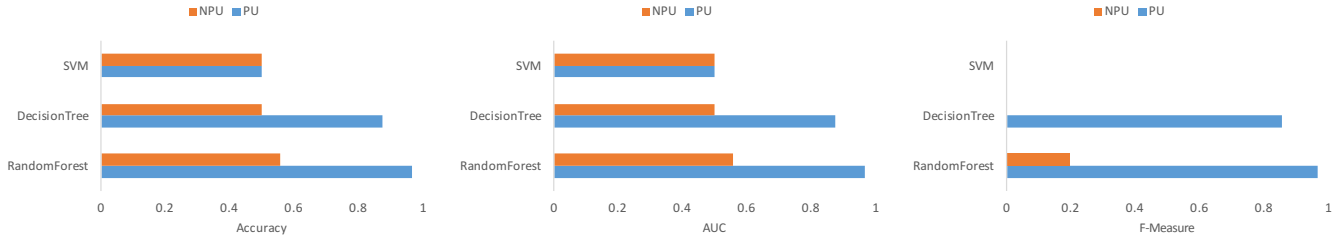


Figure 7. Performance of PUDroid when faced with contaminated data in malware dataset

and control (C&C) servers.

- *Plankton* family focuses on to collect every information on the smartphone, such as International Mobile Equipment Identity (IMEI) number, user ID, and the browser history. It then posts information via URL. It can also use remote control to modify browser’s bookmarks and install downloaded files on devices.
- *BaseBridge* family sends fake update notifications to users in hope of fooling them to install malicious software components to allow cyber-criminals to remotely control infected devices.

We included five out of six malware families in the malicious dataset for training. To create a contaminated dataset, we mixed samples from the remaining malware family and benign apps using the ratio of 1:1. Again, 1/3 of our dataset was used for testing and 2/3 was used for training. We repeated this process so that each of the six malware families was used as contaminants. We also adopted SVM classifier in this experiment because our earlier study indicates that SVM can work very well when the number of contaminants is moderate. We report the results in Fig. 6.

As shown, we can clearly see that PUDROID shows much better performance than simply applying SVM without PU learning when the system face any new malware contaminants especially in BaseBridge, GinMaster and DroidKungFu. In BaseBridge and GinMaster families, PUDROID achieves more than 45% accuracy. Moreover, without PUDROID achieves 0% in F-measure in families with smallest samples among the six families (Basebridge and Ginmaster).

In the case of FakeInstallers, SVM alone is effective since this is a very popular malware with the most common behavior. With PUDROID, we can achieve 10% higher accuracy. In general, PUDROID is effective in detecting repackaged contaminants.

#### D. RQ 4: Detecting Benign Contaminants in Malicious Datasets

It is quite possible that benign apps could be also mistakenly labeled as malicious and thus, are included in the malicious datasets. This is because to keep up with new malware samples, the datasets need to be updated periodically and this process can result in mislabeling by researchers (e.g., the number of benign apps is typically much larger than that of malicious apps). Furthermore, it is also possible for a malware repository to be intentionally compromised by someone (e.g., a cyber-attack to weaken training data). Regardless of the causes, the presence of benign apps in malicious datasets can make machine learning less effective.

To evaluate this scenario, we again used 1/3 of our dataset as the testing dataset. We set the ratio of benign and malicious app as 8:1 to simulate the case of highly contaminated datasets. In Fig. 7, we report the performances of PUDROID and using classifier alone. PUDROID with Random Forest achieves 96.36% accuracy when a malicious dataset is heavily contaminated with benign apps. PUDROID with Decision Tree is second. However, the results show that Random Forest without PU learning can still achieve 55.5% accuracy, but Decision Tree without PU learning achieves only about 50%, which is the same as making a random



guess. The results also confirm prior observation that SVM does not work well when the dataset is highly contaminated.

## VI. RELATED WORK

In this work, we combine several concepts and techniques to construct PUDROID, a framework to detect the malicious contaminants to increase accuracy of any machine-learning based detector. The framework performs feature selection and feature embedding techniques to increase robustness and efficiency of our proposed system.

Machine learning techniques have been used to build several robust and effective malware detection systems [6], [7]. For example, SIGPID [7] applies 67 machine learning algorithms to find which algorithm is better at classification based on permission features to detect the malware. Huang et al. [15] explore the use of machine learning of permission to detect malicious applications. Their detector is based on four common machine learning techniques. As these prior efforts have shown, different machine learning algorithms perform differently on different dataset, however, in general, machine learning has been effective in performing malware detection especially when datasets are properly labeled.

Typically, more features can help to improve the accuracy and detection rate, but also incur more time and space. Many advanced techniques, i.e. tensor, factorized machine [16], [17], [18], improve the performance by leveraging multi-view datasets. These techniques are frequently applied in many areas, such as nature language [19], recommendation [20], bio-medical [21], images [22], influence networks [23], behavioral detection [24].

Others have also used feature selection in different areas [7], [25], [26], [27]. By finding important or significant patterns, feature selection can improve the overall performance of a system. For example, SIGPID [7] proposes a multi-level mining technique to do feature selection, which ends up using only 22 of 135 permissions to detect malware.

In this work, we apply PU learning to help detect contaminants. PU learning is one kind of semi-supervised learning have been used to perform link prediction in the social network [28], and images [29]. MLI [28] uses the PU learning to help them deal with unlabeled link prediction in social networks. To use these unlabeled link information, MLI [28] use PU learning to identify reliable negative instances from the unlabeled set with the spy technique [30]. Our work, on the other hand, proposes to model the *hidden malware probability* based on the identified positive and unlabeled sets. We propose a non-parametric learning framework to detect malicious contaminants.

## VII. CONCLUSION

In this paper, we introduce PUDROID, a framework to detect and remove contaminants from training datasets used in machine learning based malware detection systems. We experimented with using a corpus of over 5,500 malware

and mixed them with benign apps to create various sizes of unlabeled datasets. We then evaluate the performance of PUDROID under four realistic settings. First, we evaluate the effectiveness of PUDROID using three commonly used classification techniques (SVM, Decision Trees, and Random Forest). Second, we evaluate PUDROID using highly contaminated benign datasets. Third, it is used to detect benign datasets contaminated with unknown malware. Fourth, we also consider a case in which benign apps are intentionally or accidentally included in the malware datasets.

We then compare the detection performances of a system with PUDROID and one without when contaminated datasets are used for training. The results indicate that PUDROID is effective at detecting malicious contaminants. PUDROID can improve malware detection rate by 62.82% over detectors that use only classifiers and no PU learning. We also observe that the detection accuracy is improved by 45% when a dataset is contaminated. To improve performance of PUDROID, we also apply feature selection and embedding features. For future work, we plan to explore other techniques to further improve the feature selection and embedding process.

## ACKNOWLEDGEMENTS

This work is supported in part by NSF through grants IIS-1526499 and CNS-1626432, NSFC through grants 61503253, 61672357 and 61672313, NIH through grant R01-MH080636, and the Science Foundation of Shenzhen through grant JCYJ20160422144110140.

## REFERENCES

- [1] I. IDC Research, "Smartphone os market share, 2016 q2," in *IDC Research Report*, 2016.
- [2] G. Kelly, "Report: 97% of mobile malware is on android. this is the easy way you stay safe," in *Forbes Tech*, 2014.
- [3] I. Kaspersky Lab, "Statistics of mobile threats," in *IT THREAT EVOLUTION IN Q1 2016*, 2016.
- [4] Symantec, "Latest intelligence for march 2016," in *Symantec Official Blog*, 2016.
- [5] M. Grace, Y. Zhou, Q. Zhang, S. Zou, and X. Jiang, "Riskranker: scalable and accurate zero-day android malware detection," in *Proceedings of the 10th international conference on Mobile systems, applications, and services*. ACM, 2012, pp. 281–294.
- [6] D. Arp, M. Spreitzenbarth, M. Hübner, H. Gascon, K. Rieck, and C. Siemens, "Drebin: Effective and explainable detection of android malware in your pocket," in *Proceedings of the Annual Symposium on Network and Distributed System Security (NDSS)*, 2014.
- [7] L. Sun, Z. Li, Q. Yan, W. Srisa-an, and Y. Pan, "Sigpid: significant permission identification for android malware detection," in *2016 11th International Conference on Malicious and Unwanted Software (MALWARE)*. IEEE, 2016, pp. 1–8.

- [8] Z. Li, L. Sun, Q. Yan, W. Srisa-an, and Z. Chen, "Droidclassifier: Efficient adaptive mining of application-layer header for classifying android malware," in *International Conference on Security and Privacy in Communication Systems*. Springer, 2016, pp. 597–616.
- [9] S. Acharya, "Google Removes 13 Android Apps from Play Store Infected with Brain Test Malware," <http://www.ibtimes.co.uk/google-removes-13-android-apps-play-store-infected-brain-test-malware-1537049>, January 2016.
- [10] Google, "App manifest," in *API Guides*, 2016.
- [11] A. Shabtai, U. Kanonov, Y. Elovici, C. Glezer, and Y. Weiss, "andromaly: a behavioral malware detection framework for android devices," *Journal of Intelligent Information Systems*, vol. 38, no. 1, pp. 161–190, 2012.
- [12] F. A. Narudin, A. Feizollah, N. B. Anuar, and A. Gani, "Evaluation of machine learning classifiers for mobile malware detection," *Soft Computing*, vol. 20, no. 1, pp. 343–357, 2016.
- [13] C. Elkan and K. Noto, "Learning classifiers from only positive and unlabeled data," in *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2008, pp. 213–220.
- [14] F. Ruiz, "'FakeInstaller' Leads the Attack on Android Phones," <https://securingtomorrow.mcafee.com/mcafee-labs/fakeinstaller-leads-the-attack-on-android-phones/>, October 2012.
- [15] C.-Y. Huang, Y.-T. Tsai, and C.-H. Hsu, "Performance evaluation on permission-based detection for android malware," in *Advances in Intelligent Systems and Applications-Volume 2*. Springer, 2013, pp. 111–120.
- [16] L. He, X. Kong, P. S. Yu, X. Yang, A. B. Ragin, and Z. Hao, "Dusk: A dual structure-preserving kernel for supervised tensor learning with applications to neuroimages," in *Proceedings of the 2014 SIAM International Conference on Data Mining*. SIAM, 2014, pp. 127–135.
- [17] L. He, C.-T. Lu, G. Ma, S. Wang, L. Shen, S. Y. Philip, and A. B. Ragin, "Kernelized support tensor machines," in *International Conference on Machine Learning*, 2017, pp. 1442–1451.
- [18] L. He, C.-T. Lu, H. Ding, S. Wang, L. Shen, P. S. Yu, and A. B. Ragin, "Multi-way multi-level kernel modeling for neuroimaging classification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 356–364.
- [19] C. Zhang, S. Xie, Y. Li, J. Gao, W. Fan, and P. S. Yu, "Multi-source hierarchical prediction consolidation," in *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*. ACM, 2016, pp. 2251–2256.
- [20] L. Zheng, V. Noroozi, and P. S. Yu, "Joint deep modeling of users and items using reviews for recommendation," in *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining*. ACM, 2017, pp. 425–434.
- [21] B. Cao, L. He, X. Wei, M. Xing, P. S. Yu, H. Klumpp, and A. D. Leow, "t-bne: Tensor-based brain network embedding," SIAM, 2017.
- [22] Z. Hao, L. He, B. Chen, and X. Yang, "A linear support higher-order tensor machine for classification," *IEEE Transactions on Image Processing*, vol. 22, no. 7, pp. 2911–2920, 2013.
- [23] W. Shao, L. He, and S. Y. Philip, "Clustering on multi-source incomplete data via tensor modeling and factorization," in *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, 2015, pp. 485–497.
- [24] L. Sun, Y. Wang, B. Cao, P. S. Yu, W. Srisa-an, and A. D. Leow, "Sequential keystroke behavioral biometrics for mobile user identification via multi-view deep learning," in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases (ECML/PKDD)*, 2017.
- [25] X. Wei, B. Cao, and P. S. Yu, "Unsupervised feature selection with heterogeneous side information," in *Proceedings of ACM International Conference on Information and Knowledge Management (CIKM)*, 2017.
- [26] B. Cao, L. He, X. Kong, S. Y. Philip, Z. Hao, and A. B. Ragin, "Tensor-based multi-view feature selection with applications to brain diseases," in *2014 IEEE International Conference on Data Mining*. IEEE, 2014, pp. 40–49.
- [27] X. Wei, B. Cao, and P. S. Yu, "Multi-view unsupervised feature selection by cross-diffused matrix alignment," in *Proceedings of International Joint Conference on Neural Networks (IJCNN)*, 2017.
- [28] J. Zhang, P. Yu, and Z. Zhou, "Meta-path based multi-network collective link prediction," in *KDD*, 2014, pp. 1286–1295.
- [29] L. Cui, J. Zhang, Z. Chen, Y. Shi, and P. S. Yu, "Inverse extreme learning machine for learning with label proportions," in *Proceedings of IEEE International Conference on Big Data*, 2017.
- [30] B. Liu, Y. Dai, X. Li, W. Lee, and P. Yu, "Building text classifiers using positive and unlabeled examples," in *ICDM*, 2003, pp. 179–186.