

# “How Does It Detect A Malicious App?” Explaining the Predictions of AI-based Android Malware Detector

Zhi Lu<sup>a,\*</sup>, Vrizlynn L.L. Thing<sup>a</sup>

<sup>a</sup>*Cyber Security Strategic Technology Centre, ST Engineering, Singapore 567710*

---

## Abstract

AI methods have been proven to yield impressive performance on Android malware detection. However, most AI-based methods make predictions of suspicious samples in a black-box manner without transparency on models' inference. The expectation on models' explainability and transparency by cyber security and AI practitioners to assure the trustworthiness increases. In this article, we present a novel model-agnostic explanation method for AI models applied for Android malware detection. Our proposed method identifies and quantifies the data features relevance to the predictions by two steps: i) data perturbation that generates the synthetic data by manipulating features' values; and ii) optimization of features attribution values to seek significant changes of prediction scores on the perturbed data with minimal feature values changes. The proposed method is validated by three experiments. We firstly demonstrate that our proposed model explanation method can aid in discovering how AI models are evaded by adversarial samples quantitatively. In the following experiments, we compare the explainability and fidelity of our proposed method with state-of-the-arts, respectively.

*Keywords:* Explainable AI, Cyber security, Machine Learning

---

\*Corresponding author

*Email addresses:* [lu.zhi@stengg.com](mailto:lu.zhi@stengg.com) (Zhi Lu), [vrizlynn.thing@stengg.com](mailto:vrizlynn.thing@stengg.com) (Vrizlynn L.L. Thing)

## 1. Introduction

Artificial intelligence (AI) technologies, in particular shallow machine learning methods (e.g., logistic regression, SVM and random forest, etc.) and deep neural networks (e.g., CNN and LSTM, etc.), have been widely used to fight against cyber attacks today [1] [2]. However, explaining the reasoning behind the predictions made by AI models remains one of the most challenging problems in both AI and cyber security communities. The model explanation is to identify the relevance of data features to the model’s predicted class (Fig. 1). It is not only to facilitate the cyber security practitioners’ understanding on how the AI model identifies the cyber threats, where the confidence on the model’s stability is expected, but also helps the AI researchers to find out the models’ vulnerabilities from certain examples [3]. The main factors affecting the insights provided by a model explainer are: the challenges of data, such as unbalanced data and high dimensionalities for different applications (e.g., malware detection), and the model complexity that is about the number of parameters and various types of model structures (e.g, SVM, CNN and LSTM). These factors have motivated the interest on the *model-agnostic* explainers research [4] [5] that are capable of explaining any types of AI models.

The AI-based malware detection for Android operating system is a well studied problem [6] [7] [8] [9] [10]. The popularity and open nature of Android [11] and its official apps market, i.e., Google Play [12], have aroused large interests from malware developers, and accordingly make Android users more vulnerable to such attacks by thousands of malicious apps. In 2019, AV-TEST [13] reported 3.16 million Android malware developed. The efforts on Android malware detection research pay more attention on a secure AI model and the detection performance (e.g., detection rate), where the model explanation is rarely covered. Demontis et. al. [6] introduced a secure SVM on Android malware detection that can detect those malware apps evaded from standard SVM-based detectors, through the constraints on the model’s parameters selection. Zhang et. al. [14] combined the n-gram analysis and the online classifier techniques

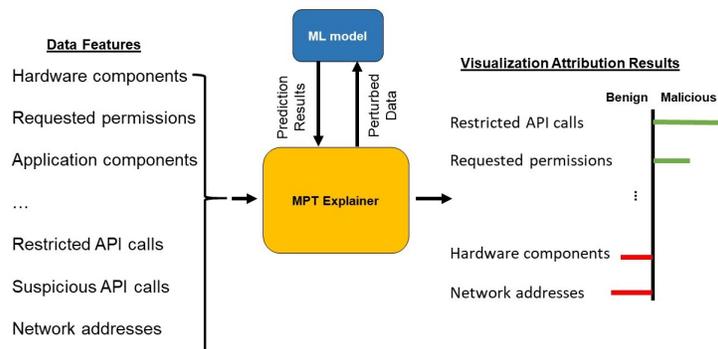


Figure 1: **Workflow of explaining an AI model for Android malware detection.** MPT explainer works together with the trained model (blue box) during the model explanation process. Firstly, MPT explainer perturb the input sample by features and observe the model’s prediction values. Second, MPT explainer optimizes the features attribution using modern portfolio theory. Finally, each feature’s attribution strength is visualized by the bar chart.

to detect the Android malware and attribute their malware families. The on-line classifier using incremental learning makes it scalable to the rapid evolved malware. The recently emerging detection methods based on deep neural networks do not only bring astonishing performance in terms of detection rate (i.e., true positive rate), but also remove the needs of features engineering manually, through an automatic joint end-to-end learning [7] [8] [9]. The large amount of parameters and hidden layers structure in deep neural networks, such as CNN and LSTM, make it more difficult for end users to understand the reasoning behind the predictions, even if the authors provided the insights on the effects of performance by the data and network structure manipulation [9]. However, only recently has the model explanation for malware detection been addressed [10].

Clearly, the primary aim of model explanation for Android malware detection is to identify the quantified “relevance” between the data features and the predictions. Existing explainable AI research [4] [5] [15] provides model-agnostic explainers for those models that may not be suitable for cyber security issues. LIME [4] explains a model through an optimization [16] of a linear model that is human-understandable over the perturbed data with the labels generated by the model’s inference. The high dimensionality and sparsity of Android malware

feature data may cease the explainability of the interpretable models in LIME that is too linear to underfitting. Integrated Gradients (IGs) [5] attributes the relevance of data features through the integration of predictions' gradients with respect to the data features varying from zero-values to the real values. However, the IGs method requires the availability of gradients computation for prediction scores with respect to the model input, which is not available in certain types of classifiers, such as random forests. The expectation, therefore, on accurate and efficient AI model explanation for cyber security issues, such as malware detection, is still remaining.

In this article, we propose a novel model explanation methodology capable of exploring the reasoning behind the predictions for AI-based Android malware detectors. The main goal of our methodology is to identify and quantify the relevance of data features with respect to the predictions for any AI-based malware detectors, as compared to the detectors with few human-understandable explanation for specific types of models [10]. Furthermore, this methodology explains an AI model by a simple data perturbation technique and efficient linear optimization that is based on the modern portfolio theory (MPT) [17]. The main advantage of the method based on optimization for features attribution lies in the potential to explain a model accepting a high-dimension data and no requirement on the gradients of prediction scores w.r.t. the inputs. The proposed methodology can be divided into two stages: an initial data perturbation technique followed by a linear optimization for feature attribution values (i.e., relevance of features). The first stage consists of the following steps: i) data features perturbation by multiplying a scaling factor in  $[0, 1]$  to have the small changes of feature values (i.e., perturbed data); and ii) associated prediction scores are obtained by the forward process of the model over these perturbed data. The second stage then minimizes the *variance* of the multiplication of these prediction scores and the attribution values, with the constraint that the sum of all features' attribution is one. The **main contribution** of this article is the formation of the optimization for features attribution, which assumes that the small changes of relevant features affect the prediction scores significantly,

and thus are assigned high attribution values.

In addition, results from the quantitative evaluation show that our proposed technique is sensitive to identify the relevant data features to the predictions made by the AI-based Android malware apps detectors, such as SVM [10] [18] [19] and BERT [20]. This is verified by the facts that 1) the feature with the largest attribution value is helpful for the clustering algorithm to group the apps from the same malware family into the same cluster, and more features do not increase the performance significantly; 2) a significant ratio of the camouflaged features activated for the adversarial samples, which evade the malware detection by the classifiers, can be identified.

The rest of the paper is structured as follows: we present current work on AI-based Android malware apps detection and the model explanation methods in Section 2. The two learning-based classifiers, i.e., SVM and BERT, are presented in Section 3, followed by the introduction to the proposed model-agnostic MPT explainer in Section 5. We then validate the proposed methodology with the comprehensive experiments in Section 6. Finally, the conclusion of this methodology is discussed in Section 7.

## 2. Literature Review

In this section, we firstly review the research on Android malware apps detection using AI-based methods, where the early efforts in their work on these models explanation are also covered. Secondly, the general model explanation methods are introduced in terms of their strength and weakness.

### 2.1. Android Malware Detection

The vulnerabilities in Android system are still emerging [21], and the malware is in rapid evolution. This requires the development of AI-based models that can predict unknown malware samples in certain degrees. The AI-based methods described in the Android malware detection literature have focused on four basic features extraction approaches [22]. In this section, we reviewed

two main feature extraction techniques: *the dynamic analysis* and *the static analysis* that are relevant to our work. The first method for features extraction is *dynamic analysis* that obtain the features via the apps execution, such as system calls and file system usage, etc. For example, Xiao *et. al.* [9] propose a deep learning method that analyses the API calls obtained dynamically through two Long-term Short Memory (LSTM) networks trained for malware and non-malware apps, respectively.

The second and most commonly used method is the *static analysis technique* that analyses the application source code and resources without execution and requires small running time overhead. For example, Samra *et. al.* [23] use the clustering algorithm in the malware detection that the *permissions* data were extracted from the applications as the features statically. In [24], a mobile malware analysis tool, AndroDialysis, is developed and used to assess the malware detection performance with *intents* feature (explicit and implicit), which is a run-time binding message object for the inter-process communication in Android framework. The conclusion drawn from the experiment results shows that intents as the only features in the classification are more effective than that with permissions features. Zhang *et. al.* [25] convey the semantic information in malware classification via the graph construction of the dependencies among the *API calls* as the features. Similarly, Aafer *et. al.* [26] capture the semantic information about the apps' behavior at the API level. In addition to classification by single features, there are also a number of AI-based methods that classify the malware apps using a combination of features extracted by the static analysis technique. For instance, Zhu *et. al.* [27] detect malware apps by building up an ensemble classifier, Rotation Forests [28] with the features extracted from the APK archive file that includes permissions, sensitive APIs and monitoring system events, etc. The classifier was trained on a dataset of 2,130 samples of Android apps that malware and benign apps take 50% respectively, and the results outperforms a SVM classifier. Arp *et. al.* [10] learn a linear SVM classifier with the application features, including permissions, API calls and network addresses. They propose a lightweight features extraction method

to extract from *manifest.xml* file and the *dex bytecode*. In the above methods, only DREBIN [10] provided a limited explanation of the model’s prediction. However, their explanation is based on the weights of the SVM models, which cannot be generalised into other models.

Deep learning methods have been successfully applied on the cyber security issues. For example, Xiao et. al. [9] proposed a deep learning method that detects the Android malware through training two Long-term Short Memory (LSTM) networks for malware and non-malware apps features (i.e., system calls) respectively. Vinayakumar et. al. [29] tried different LSTM network topologies in Android malware detection. However, the hidden layers and cells structure in LSTM make the inference process in a black-box mechanism for users, and thus impossible to have the insights on the reasoning. This reduces the trust on the model’s prediction for an example, because it has the risk for the model to be misled to make a wrong decision by an carefully designed example, such as that in the stop sign recognition problem [3], and the user cannot know how the decision is made.

## 2.2. AI Model Explanation

The research efforts in the model explanation literature are mainly to reveal the transparency of the reasoning behind the model’s predictions based on the data features, a.k.a., input variables. There are two major categories of the explanation approaches [30]: i) *transparency-based methods*; and ii) *post-doc interpretability methods*. The transparency-based methods [31] are the traditional model explanation strategy that conveys the interpretation of the predictions by interpretable AI models. That is, these models cannot be too complicated in terms of structure and parameters’ numbers such that humans cannot understand [32]. For instance, in [10], the authors tried to enhance the interpretation of the SVM classifier’s predictions through the presentation of the features with top ranked weights. Unfortunately, as the complexity of the AI models grows, it is inevitably more difficult for these methods to identify meaningful explanations. In addition, it is trade-off for a model between the transparency for

interpretability and the prediction performance [30].

The post-doc methods explored in the literature have raised more attention from both academia and industry. They extract the causality relationships between the prediction and the data features from the trained (learned) model. The early research in this field focused on the *model-specific* explanation that seek the rationale of predictions for specific models. For example, Zeiler *et. al.* [33] proposed to visualize and observe the neurons in a convolutional neural networks (CNNs) that shows how each neuron responses to different data instances. Xu *et. al.* [34] developed an attention-based caption generator for the images, where the attention mechanism shows the highlighted part of the image that is relevant to the generation of a particular word in the caption. Obviously, in the above methods, it is hard to extend the explanation to other models as the learning-based AI models are rapidly evolving today. Therefore, there are also several research that aim to develop the model explanation methods that can disclose the inference of predictions for any models, a.k.a., *model-agnostic* explanation. For instance, Samek *et. al.* [35] present two methods for deep learning models: (1) sensitivity analysis (SA) that measures how sensitive the predictions are w.r.t. each input variables, and (2) layer-wise relevance propagation (LRP) that propagates the prediction backward in the deep neural network, constrained by a set of propagation rules. Sundararajan *et. al.* [5] propose a feature attribution method, integrated gradients (IGs), that was inspired by a cooperative gaming theory [36]. It can explain the models that the input data are in different types [37], such as images, text and tabular data, only if the gradients of model’s prediction scores w.r.t. data features is provided. Besides, several efforts focusing on the complex model explanation that use model approximation. For example, LIME [4], a model-agnostic method for local explanation, explains any classifiers by learning an interpretable model, such as a linear model, with the generated data around the predictions. However, it has the risk that the linear model for interpretation through approximating the original model make the explanation is not capable of those models and datasets with complicated structure. Another approach presented by Wu *et. al.* in [38]

explains the deep models through the original model approximation with a decision tree with few nodes. The model explanation research efforts also focus on explaining the AI models in different tasks, such as common sense question answering (CQA) [39]. Recently, Guo et.al. [40] proposed LEMNA as the solution to LIME’s limitations that are caused by the linearity of its surrogate model. LEMNA is designed specific for the AI models that are used in cyber security tasks, which usually need to deal with sequential data, such as binary code analysis using RNN [41]. The linear approximation of the original model  $f(\mathbf{x})$  cannot take the dependencies among features into account that leads to limited capacity of explanation for sequential data. LEMNA solves this issue using the *fused lasso* [42], which forces the explanation considers the dependencies among features accordingly. In addition, LEMNA enhanced the fidelity of the local approximation on the non-linearity of the neighborhoods (sampled) around the input data sample  $\mathbf{x}$  by *mixture regression model* [43].

There are also several development tools for explainable AI (xAI) that are public available online. The What-If Tool [44] provides an interactive visual interface for users to inspect the models’ behaviors. It works with Tensorflow models [45] as a dashboard that allows users to compare models and find out the importance of features by manipulating the data feature values. However, this is a tool that requires models running on Google Cloud, and the user-interface is not friendly for those other than data scientists. Captum [46] is a Python library developed by Facebook recently that implements several model explanation methods, such as DeepLift [47], Input  $\times$  Gradient [48] and Integrated Gradients [5]. The functionalities of Captum are under expansion, and this model explanation library currently is only able to explain the deep learning models written in PyTorch [49]. Similarly, SecML [50] provides a Python library that focuses on the secure and explainable AI algorithms. There are limited number of model explanation methods provided in SecML library, because it is an integrated library that combines several machine learning, adversarial machine learning and explainable machine learning algorithms. In addition, SecML only accepts the data instances in the type (i.e., *CArray*) that is defined

in the library, which is not a requirement in our proposed method.

### 3. Learning-based Android Malware Detector

The malware detection is cast into a binary classification problem, i.e., malware and non-malware (benign). DREBIN [10], an Android malware dataset that is used to train and validate the classifiers, is firstly introduced. Second, the feature representation based on TF-IDF weighting scheme [51] that extends the representative of the binary feature vector in [10] is illustrated. Finally, we briefly introduce two classifiers: (1) SVM and (2) BERT, that are used in the Android malware detection.

#### 3.1. DREBIN Dataset

DREBIN [10] was developed as a lightweight Android malware detector that utilises the *static analysis technique* to extract the features of suspicious Android apps. An Android malware dataset with the same name, DREBIN, was also released in public. Each Android application is represented by means of feature sets that were extracted through a linear sweep over the application’s manifest file, *AndroidManifest.xml*, and disassembled *dex code* from the byte-code for the Dalvik virtual machine of Android platform. These feature sets fall into 8 categories, which are numbered from  $S_1$  to  $S_8$ , such as  $S_2$ : *requested permissions*. DREBIN datasets collected 5,560 Android malware apps and 123,453 non-malware apps in total. An Android app is labeled as *malicious* only if it is detected by at least two scanners in VirusTotal [52], and *non-malicious* otherwise.

#### 3.2. Embedding Features Vector by TF-IDF

The features for an application are in the form of text data, where each feature can be treated as a single “word” or token. It is, thus, necessary to convert the set of features to the vector space that reveals the dependencies among features and the machine learning models can identify the patterns of malicious or non-malicious Apps. In DREBIN [10], the features of an App

$\mathbf{X}$  were mapped into a vector space with Boolean values only. The mapping function for binary vector values is  $\phi : \mathbf{X} \mapsto \{0, 1\} \in R^m$  that any dimension is marked as 1 refers to the associated feature contained by the App  $\mathbf{X}$ , and 0 otherwise, where  $m$  is the total number of features. However, the vector representation with Boolean values cannot reveal all the dependencies among features in an App, because of its limited capacity of representation.

We, therefore, extend the representation of a vector with a new mapping function:

$$\phi(x_i) = \begin{cases} \frac{f_{x_i, \mathbf{X}}}{|\mathbf{X}|} \cdot \log \frac{|\mathcal{D}|}{|\{\mathbf{X} \in \mathcal{D} : x_i \in \mathbf{X}\}|}, & \text{if feature } x_i \in \mathbf{X} \\ 0 & , \text{ otherwise} \end{cases} \quad (1)$$

for any feature  $x_i$ 's numeric representation (i.e., word embedding) in the vector space, where  $|\mathcal{D}|$  is the number of Apps in dataset  $\mathcal{D}$ ,  $|\mathbf{X}|$  is the amount of features in app  $\mathbf{X}$ ,  $f_{x_i, \mathbf{X}}$  is the frequency of  $x_i$  in app  $\mathbf{X}$ , and  $\frac{f_{x_i, \mathbf{X}}}{|\mathbf{X}|} \cdot \log \frac{|\mathcal{D}|}{|\{\mathbf{X} \in \mathcal{D} : x_i \in \mathbf{X}\}|}$  computes the term weight in TF-IDF weighting scheme. This weighting scheme guarantees the assignment of relatively low weights on those features appearing frequently among both malicious and non-malicious apps, and relatively high weights on the features that are used in a specific class of apps. For example, the feature `android.permission.SEND_SMS` are more commonly used in malicious apps [10].

By the new mapping function  $\phi(x_i)$ , the vector representation of a malicious app  $\mathbf{X}$  in the family `Plankton` may look like:

$$\begin{pmatrix} 0 \\ 0.13 \\ \dots \\ 0.08 \end{pmatrix} \begin{matrix} \text{m.facebook.com} & S_8 \\ \text{android.hardware.wifi} & S_1 \\ \dots & \dots \\ \text{android.hardware.touchscreen} & S_1 \end{matrix}$$

### 3.3. Learning-based Malware App Detection

We cast the malware detection problem as a binary classification task, where the *malware* is labeled as the positive class and the *benign* app (*non-malware*)

is the negative class. Specifically, two popular classifiers, i.e., SVM and BERT, are trained to separate the malicious apps from the non-malicious (benign).

**SVM** is a classical and popular classifier for binary classification before the deep learning era. It maps a sample into a high-dimension space that has a clear gap between two classes, and thus makes the prediction of the sample’s category more accurate. It has been proven effective on Android malware detection by many early research [10] [19] [18]. We train a SVM model with Radial basis function (RBF) kernel [53]. Each feature of the apps are taken as a token in text analysis, and vectorized by TF-IDF weighting scheme. The trained SVM produces the prediction scores for each class, which indicates its confidence on the prediction for a class. In addition, the state-of-the-art deep learning method for natural language processing (NLP), **BERT**, is employed as another classifier. In order to have the trade-off between the performance and the hardware resources, we train the base uncased BERT model implemented by Huggingface [54] with a maximum length of input text as 128. Both classifiers are trained on the balanced DREBIN dataset. The detection performance shown in Table 1 denote to the high detection rates (0.9684 for SVM and 0.9591 for BERT) with low false positive rates (0.0425 for SVM and 0.0420 for BERT).

#### 4. Evade Detection by Adversarial Samples

A machine learning model for malware detection can be vulnerable to attacks using well-crafted adversarial samples, because of the intrinsic assumption to build up the model that the training and test data are drawn from the same distribution [6] [55]. The adversarial sample violates such assumption by activating a certain number of features in the features vector to make it available to evade the AI-based malware detection without compromising of its malicious functionalities.

##### 4.1. Adversarial Samples Generation Algorithm

In this paper, we simulate an evasion attack using the adversarial samples to mislead SVM and BERT classifiers in malware detection. The attacker is

assumed to have full capacity of manipulating the data in the inference stage, including the knowledge of the feature space, and the model’s feedback, such as prediction scores. The adversarial samples are generated with the following settings: Firstly, the optimised set of features for the adversarial samples are selected/activated through a process of heuristic search that maximises the fitness function [56]. Then, the manipulated samples will be converted into the features space by calculating the TF-IDF values for the features in the adversarial samples, including the original and the activated features. Please note that we only activate the features under the “Permission” category, in order to preserve the malicious functionalities, as suggested as by [6]. We also extend the genetic algorithm [56] that was used in the adversarial samples generation, where the model’s prediction score of benign class is used as the fitness values. This guarantees that (1) the adversarial samples can evade the detection, because of the prediction scores  $> 0.5$  for benign class; (2) we have the quantified confidence that the classifier incorrectly classifies them into benign class, which will be used in later experiments (see Section 6.2). The pseudo code for the adversarial samples generation is illustrated in Algorithm 1. The MPT explainer is then used to explain why the classifiers are bypassed by these adversarial samples, where the activated features in the adversarial samples are expected to be attributed with high values.

Specifically, we run the genetic algorithm to generate the adversarial samples dataset by the following parameters: (1) there are 200 malware samples randomly selected from the test set that will be camouflaged to fool the classifiers in the experiments, which are called **base samples**; (2) the maximum iterations for the evolution of solutions in genetic algorithm are 500 loops, which is based on our observation in the early experiments that is enough to guarantee the adversarial samples can be generated for most base samples. In addition, the optimisation process of the genetic algorithm will be converged if any of the following conditions is fulfilled: (1) the maximum prediction score of the adversarial samples has been idle for at least 10 iterations that implies the no solutions (i.e., a set of activated features) with higher fitness values can be ob-

tained through the optimisation process in reasonable time; (2) the maximum fitness value  $> 0.99$ , i.e., the prediction score, which means the classifiers cannot distinguish the adversarial samples from the malware samples at all, and there is little benefits to increase the fitness value by the genetic algorithm; or (3) the solutions in the optimisation process have been evolved for 500 loops.

By these settings, we generated 29,821 adversarial samples that can evade both SVM and BERT classifiers. The samples are distributed by their model’s prediction scores, where the interval is 0.1 and most of the adversarial samples evade the model’s prediction (SVM and BERT) with scores of  $\leq 0.9$ . We randomly pick up 100 samples for each prediction scores interval from the pool (i.e., totally, 500 samples for BERT and 499 for SVM that one sample is removed because it was false negative in the original prediction, respectively) for both SVM and BERT classifiers for a fair assessment on the model explanation experiment. The MPT explainer is then used to attribute the features for each adversarial sample such that the “relevant” features to the false negative predictions can be identified. More details about the explanation on the adversarial samples by MPT explainer will be discussed in Section 6.

## 5. Explaining Classifier’s Predictions

In this section, we firstly formulate the model explanation as a feature attribution problem that is to quantify the relevance of features to the classifier’s prediction. The data perturbation technique, then, for the generation of the candidate data in model explanation is introduced. Finally, we present the details of the proposed methodology, i.e., MPT Explainer, that is based on modern portfolio theory (MPT) [17].

### 5.1. Theory Basis

**Modern Portfolio Theory** (MPT) is an economic theory that maximize the expected returns through a portfolio of assets (e.g., an allocation of investment on a set of stocks), given a certain level of risk. Inspired by MPT, the

---

**Algorithm 1: Adversarial Samples Generation**

---

**Input** : A malware sample from DREBIN dataset.

SVM classifier

**Output:** A set of adversarial samples.

**Stage 0 - Candidate Features**

- 1 Collect features from training set under “Permission” category.
- 2 Exclude features in the input malware sample.

**Stage 1 - Construction of Adversarial Samples**

- 3 Initialize a “population” matrix that random features are activated in each solution.

$shape=(num\_solutions, num\_genes)$ .

One **solutions** refers to an adversarial sample candidate.

The **genes** means the candidate camouflaged features.

4 `iter_counter = 0`

5 `MAX_LOOP = 500`

6 **while** *True* **do**

7     Compute fitness values - **p**

8     Select matting pool

9     Crossover

10    Mutation

11    Update “population”

12    `iter += 1`

13    **if**  $iter\_counter \geq MAX\_LOOP$  **then**

14       | **break**

15    **if**  $max(p)$  idle for 10 iterations **then**

16       | **break**

17    **if**  $max(p) > 0.99$  **then**

18       | **break**

19 **end**

20 **return** Adversarial samples that fitness values  $> 0.5$

**and** prediction score of benign class by BERT  $> 0.5$

---

features in a data sample can be treated as the assets. The output attribution values for the features, which refer to the explanation, are similar with the allocation of the investment, given a certain level of the perturbation of the classifier’s prediction scores.

### 5.2. Problem Statement

**Problem statement** The classifier  $f(\mathbf{x}) : \rightarrow [0, 1]^{|C|} \in \mathcal{F}$  is a mapping function to map the data instance  $\mathbf{x}$  to a set of probabilities, which denote to the confidence of the classifier  $f$  to classify  $\mathbf{x}$  to each class. A data instance  $\mathbf{x} = (x_1, x_2, \dots, x_m) \in R^m$  is a vector with  $m$  features that a feature  $x_i$  can be the numerical representation of either a token in text data or a pixel (or a superpixel) in image data. The **explanation** of a classifier  $f$  is to seek the quantified relevance of each feature  $x_i$  in data instance  $x$  to the predicted class. Specifically, this is called *feature attribution* that is in the form of an attribution vector  $\mathbf{A} = (a_1, a_2, \dots, a_m) \in R^m$ , where  $a_i \in [-1, 1]$  represents a degree of relevance of feature  $x_i$  to the prediction, and the lower the less relevant.

### 5.3. Data Perturbation

We look for the interpretation of the classifier’s prediction with respect to single data instances, which is also known as *local interpretation* [4]. The challenge of local interpretation is the limited number of data that can be used to attribute the features’ relevance, and the observation of the changes of prediction scores by different data becomes impossible. Therefore, the synthetic data generation through perturbing the values of each feature becomes necessary for further model explanation. Specifically, given a data instance  $\mathbf{x} = (x_1, x_2, \dots, x_m)$  from the dataset, the classifier outputs a prediction  $f(\mathbf{x})$  that is the probability of the predicted class (e.g, in binary classification, this is the larger probability). For each feature  $x_i$  ( $i = 1, 2, 3, \dots, m$ ), the  $K$  perturbed data over this feature are computed as:

$$x_i^{(k)} = \alpha_i^{(k)} \cdot x_i \tag{2}$$

where  $k \in \{1, 2, 3, \dots, K\}$  is the number of perturbed data, and  $\{\alpha_i\} \in [0, 1]$  are evenly spaced and the coefficients to linearly change the feature  $i$  values slightly. Based on our observation in the early experiments,  $K$  is set as 50 for the trade-off between capacity of data perturbation and computation efficiency.

By the data perturbation strategy, we simply obtain a set of  $K \times |\mathbf{x}|$  perturbed data, which is  $\mathbf{P}_\mathbf{x}$ . These data will be fed into the classifier  $f$  to have the prediction scores over the classes. Both the perturbed data and the associated prediction scores will be used in the optimization of the attribution values vector  $\mathbf{A}$  (See Section 5.4).

#### 5.4. Optimization for Features Attribution

The purpose of model explanation is to find the optimized attributions values for the features. The assumption behind is that the prediction scores are more sensitive to the changes of relevant features than those relatively less relevant. That is, a small changes of relevant features values will cause significant changes of prediction scores. With the perturbed dataset  $\mathbf{P}_\mathbf{x}$ , the changes of prediction scores between the  $\mathbf{x}^{(ik)}$  (i.e., the  $k$ -th perturbed data instance by perturbing the  $i$ -th feature value) and the original data instance  $\mathbf{x}$  is defined as the fraction of the prediction scores difference between  $f(\mathbf{x}^{(ik)})$  and  $f(\mathbf{x})$ :

$$r_i^{(k)} = \frac{f(\mathbf{x}^{(ik)}) - f(\mathbf{x})}{f(\mathbf{x})} \quad (3)$$

This forms a  $K$ -dimension vector of the prediction scores changes,  $\mathbf{r}_i = (r_i^{(1)}, r_i^{(2)}, \dots, r_i^{(K)})$ , when varying the values of feature  $x_i$  by equation (2).

We model the association between the attribution values of features, i.e.,  $\mathbf{A}$ , with the random prediction scores changes,  $\mathbf{r} = (r_1, r_2, \dots, r_m)$ , for all  $m$  features by their inner product:

$$\mathbf{A}^T \mathbf{r} = \sum_{i=1}^m r_i a_i \quad (4)$$

we compute the variance of this association as:

$$\begin{aligned}
\text{Var}[\sum_{i=1}^m r_i a_i] &= E[(\sum_{i=1}^m r_i a_i - E(\sum_{i=1}^m r_i a_i))^2] \\
&= E[(\sum_{i=1}^m r_i a_i - a_i \mu_i)(\sum_{j=1}^m r_j a_j - a_j \mu_j)] \\
&= \sum_{i=1}^m \sum_{j=1}^m a_i a_j \cdot E[(r_i - \mu_i)(r_j - \mu_j)] \\
&= \sum_{i=1}^m \sum_{j=1}^m a_i a_j \sigma_{ij} = \mathbf{A}^T \mathbf{Q} \mathbf{A}
\end{aligned} \tag{5}$$

where  $\mathbf{Q}$  is the covariance matrix.

The aim is to find out an attribution vector,  $\mathbf{A}$ , that the variance of this association  $\mathbf{A}^T \mathbf{r}$  is minimized, such that the non-trivial changes of the prediction scores can be obtained with the minimal fluctuations of the data features. The function of the optimization is given by:

$$\begin{aligned}
&\arg \min_{\mathbf{A}} \mathbf{A}^T \mathbf{Q} \mathbf{A} \\
&w.r.t. \sum_{i=1}^m a_i = 1 \text{ and } a_i \in [-1, 1]
\end{aligned} \tag{6}$$

The equation (6) considers both the individual features' contribution to the final predictions and the contribution made by the dependency among these features. Note that we argue that the features with negative attribution values do not mean that they have negative contribution to the prediction. The negative value means, instead, the associated feature has negligible affect on the final prediction.

## 6. Experiments

### 6.1. Experiment setup

The dataset utilised in this paper consists of 11,110 samples of Android apps drawn from DREBIN dataset, in which there are 5,555 benign apps and

5,555 malware apps, respectively. Further, we randomly separate this dataset into two dis-joint sub-sets: a training set with 7,442 apps and a test set with 3,668 apps. There are 179 malware families and one benign family in the whole dataset that the top four malware families (i.e., FakeInstaller (15.7%), Plankton (11.8%), DroidKungFu (11.7%), Opfake (11.0%)) with the largest number of samples account for more than 50% malware samples. The apps in the training set are from 160 malware families and the test set covers 124 malware families, respectively.

The two classifiers, i.e., SVM and BERT, are trained over the training set with the following settings: (1) **SVM** is with RBF kernel that  $\gamma = 1.0$ . (2) **BERT**. We use the uncased base model implementation of BERT from HuggingFace Transformers library [54] with the default configurations of parameters, e.g., the maximum sequence length of the text is 128, learning rate is  $4e-5$ . The batch size for training is set up as 8 for a trade-off between the performance and the memory limitation, and the model was trained for 5 epochs. Since BERT uses the word embedding technique to generate numeric representation of the text tokens that is not easy to perform the data perturbation directly by our method, we approximate its prediction behaviour by training a SVM, called *aprox\_SVM*. Specifically, BERT predicts a set of training data to have the predicted labels. The *aprox\_SVM* was trained on the same dataset and the “ground truth” are the predicted labels by BERT. By this way, the *aprox\_SVM* has a  $TPR=0.9984$  and  $FPR=0.0029$  that means it has a highly similar performance as the BERT. The *aprox\_SVM* will only be used to generate the predictions of perturbed data in the model explanation. Note that in the experiments, the features attribution for BERT is identified by *aprox\_SVM* with MPT explainer.

The classifiers’ performance were assessed on the test set and summarized in Table 1 with true positive rate (TPR), false positive rate (FPR) and Matthews correlation coefficient (MCC). The MCC combines true/false positives/negatives into a single metric range from  $[-1, 1]$  that 1 means a perfect prediction, 0 denotes to the prediction performance that is no better than the random prediction and -1 indicates a total disagreement between the prediction and ground truth.

It shows that the predictions by both classifiers match the ground truth quite well in terms of MCC, and SVM has a slightly better performance than BERT on Android malware detection, although SVM has approximate 0.0005 higher of FPR.

Table 1: **Classification performance for SVM and BERT.** Both classifiers obtained high true positive rate (TPR) and MCC, and relatively low false positive rate (FPR). It shows a high detection rate for both classifiers in the Android malware detection task.

	<b>SVM</b>	<b>BERT</b>
<b>True Positive Rate (TPR)</b>	<b>0.9684</b>	0.9591
<b>False Positive Rate (FPR)</b>	0.0425	<b>0.0420</b>
<b>Matthews correlation coefficient (MCC)</b>	<b>0.9259</b>	0.9171

The parameter of the perturbed data amount around the input data instance  $\mathbf{x}$  is set as 50 for the proposed MPT Explainer. For example, for an data instance  $\mathbf{x}$  with 10 features to describe the suspicious app, there are 500 perturbed data instances generated where each feature  $x_i$  was varied into 50 values.

Next, we present three quantitative experiments to assess the MPT Explainer’s capability of model explanation. In the first experiment (see Section 6.2), we test the capacity of features attribution by MPT explainer for adversarial samples that can evade the detection by the classifiers, i.e., SVM and BERT. We also compare the MPT explainer with another two state-of-the-art explainers (see Section 6.3), which are LIME [4] and SHAP [15], in the second experiment of “good” explanation. Finally, a fidelity test (Section 6.4) on the explanation is conducted for MPT explainer and SHAP.

### 6.2. Analysis of Activated Features in Adversarial Samples

In this experiment, we evaluate if the MPT explainer is able to identify the activated features in an adversarial example (see Section 4). It is expected that the MPT explainer can help to find out the reasons that the adversarial example evades the model’s detection in the features level. Specifically, we evaluate the

explanation performance in two ways: (1) the general performance assessment, and (2) the functional analysis in terms of the activated features number in an adversarial sample and the prediction scores, respectively.

The **first assessment** is to observe the general performance of MPT explainer on features attribution based on the percentage of the samples with “good” explanation. A sample has **“good” explanations**, only if a certain percentage of its activated features (i.e., camouflaged features) are attributed with positive values. The threshold for a sample to have “good” explanation varies from 0% to 90% with 10% intervals in our experiment. The increased thresholds of “good” explanation means less tolerance for MPT explainer to incorrectly attribute the activated features with negative values. In Table 2, it is clear that the amount of adversarial samples with “good” explanation by MPT Explainer is decreasing as the “good” explanation threshold is increasing. Specifically, in the **Configuration 1** that candidate activated features are from benign apps samples in the training set, the MPT explainer identifies a large number of the activated adversarial features as the positive contribution to the false negative prediction of benign class, when the “good” explanation threshold is less than 50%. In the usage of an analysis tool, such explanation of the classifier’s behavior on adversarial samples can reveal how the detector was evaded. The explanation performance, however, is deteriorated quickly once the “good” explanation threshold is larger than 50%. We argue that such performance deterioration is reasonable, because the model’s prediction for a sample is made by the cohesive contribution of all the features. That refers to a comprehensive consideration of all features by the classifier (i.e., over the distribution of the features vector), including the contributions made by individual features and that made by the complicated relationships among features.

Further, we verified the above argument by another experiment with the **Configuration 2** that the candidate activated features are from *malware samples* in the training set. In Table 2, it shows quite high percentages of adversarial samples with “good” explanation for both SVM and BERT, when the thresholds of “good” explanation are larger than 50%. Note that the activated features

Table 2: **Results on adversarial samples explanation in terms of the percentage of “good” explanation.** The varying thresholds of “good” explanation (top row) are from > 0% to >90% with 10% intervals. **C** for configuration. **S** for SVM and **B** for BERT.

		0%	10%	20%	30%	40%	50%	60%	70%	80%	90%
<b>C1</b>	<b>S</b>	1	1	0.992	0.908	0.581	0.309	0.216	0.188	0.188	0.182
<b>C1</b>	<b>B</b>	1	1	0.994	0.882	0.592	0.336	0.23	0.196	0.188	0.18
<b>C2</b>	<b>S</b>	1	1	1	0.99	0.936	0.679	0.293	0.164	0.142	0.128
<b>C2</b>	<b>B</b>	0.996	0.996	0.996	0.994	0.948	0.640	0.348	0.232	0.204	0.194

were used by malware samples and intuitively, all the features in malware samples should have been no contribution to the false negative predictions of benign class. Therefore, the features are attributed with positive values by MPT explainer imply that the classifiers consider all the features comprehensively. That is, a malware sample is camouflaged as a benign app sample to evade the detection successfully by: (1) the contribution from individual activated features; and (2) the contribution made by the activated feature and the existing features together, which is why some of the features in the original base sample are also attributed with positive values.

The **second assessment** is a *functional analysis* that measures the adversarial samples with “good” explanation, as a function of two variables: (1) the prediction score of benign class by the classifier; and (2) the number of activated features in each sample. A higher prediction score indicates the stronger confidence from the classifier on its prediction, when the classifier considers all the features, including the contributions made by individual features. In addition, more features are activated in an adversarial sample denotes to the significant changes of the features vector, which implies more complicated relationships among features, and thus the features attribution are more difficult than that on the adversarial samples with less amount of activated features.

We show the analysis results for SVM and BERT in Table 3. The number in each cell is calculated by the number of adversarial samples with “good” ex-

Table 3: **Function analysis for model explanation (SVM left, BERT right), as a function of the model’s prediction scores and the number of activated features in each sample.** The threshold of an adversarial sample with the “good” explanation  $> 0.4$ . The candidate activated features are by Configuration 1. The intervals of the classifier’s prediction scores are shown in the first column, and the intervals of the number of activated features in a sample is shown in the first row.

SVM	(0, 30]	(30,50]	50+	BERT	(0, 30]	(30,50]	50+
(0.5,0.6]	0.448	0.278	0.286	(0.5,0.6]	0.579	0.333	0.417
(0.6,0.7]	0.392	0.500	0.111	(0.6,0.7]	0.444	0.545	0.706
(0.7,0.8]	0.661	0.636	0.185	(0.7,0.8]	0.600	0.167	0.786
(0.8,0.9]	0.737	0.636	0.490	(0.8,0.9]	0.553	0.222	0.667
(0.9,1.0]	1.000	1.000	1.000	(0.9,1.0]	0.667	0.667	0.909

planation divides the number of the adversarial samples that has the prediction scores in a certain interval (i.e., a row in the table) and the amount of activated features (i.e., a column in the table). The functional analysis reveals the features attributions for both SVM and BERT by MPT explainer have better performance when the classifier has a higher confidence on its false negative prediction of benign class (i.e., prediction score close to 1) than those sample with less confidence. In addition, the performance of features attribution by MPT explainer for SVM is deteriorated if the number of activated features in an adversarial sample is increased. And more activated features in the adversarial samples are identified by the MPT explainer for BERT as the number of activated features in a sample increased. The difference reveals the following findings: (1) SVM and BERT made the false negative prediction by considering different features; and (2) it requires to activate less amount of features to evade the detection by SVM than that by BERT, which maybe helpful for the analysts to find out the solution to a more secure malware detector.

### 6.3. Explanation Capability Comparison

In this experiment, we compared the explanation capabilities of MPT explainer with another two state-of-the-art explainable AI methods: LIME [4] and SHAP [15], when they are used to explain SVM and BERT’s false negative predictions on the adversarial samples (see Fig. 2). Totally, the 500 adversarial samples generated with the Configuration 1 are used. The performance of each explainer is evaluated in terms of the “good” explanation metric (see Section 6.2).

Fig. 2a shows the percent of samples with “good” explanations for SVM’s prediction behaviors by MPT explainer, SHAP and LIME. It seems these three explainers explain the model’s behavior accurately when the threshold of “good” explanation is small. LIME shows the most stable percentage of “good” explanation, as it keeps relatively high percentage of “good” explanation in most time even if the threshold is increasing. MPT explainer shows competitive performance when the threshold is small, but its percentage of “good” explanation is decreased when the threshold is increasing. The percentage of samples with “good” explanations by SHAP is deteriorated quickly to zero, as the the threshold of a “good” explanation is becoming larger. This is because SHAP only explains a small number of features that have the top attribution values. In summary, we can have a few findings: (1) Most of the features with high attributed values generated by the three explainers are from the activated features of the adversarial samples. (2) LIME shows the relatively stable performance, regardless of the thresholds; (3) MPT explainer has competitive explanation ability with LIME, when the threshold is small; and (4) SHAP’s performance is not stable and deteriorated quickly as the threshold is increasing.

### 6.4. Fidelity Test

We also compare the explanation fidelity between MPT explainer and kernel SHAP, when they are used to explain the PDF malware detectors (e.g., SVM and Random Forest classifiers) for the samples from the PDF malware dataset [57]. The PDF malware dataset has 4,999 malicious samples and 5,000 normal PDF

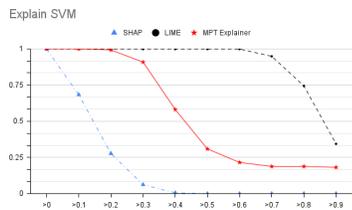
files. We follow [40] to use the 135 features, which values have been encoded into binary values. The SVM and Random Forest are trained on the training set, and we observe the explanation performance when MPT explainer and SHAP are used to explain the classification of the testing samples.

The fidelity tests include *deduction test* and *augmentation test*. The **deduction test** assumes that the AI model’s prediction of a manipulated sample, in which a certain number of features with high attribution values (so called “important” features) are removed, will towards the opposite class of the original sample, if the explanation (i.e., the attribution values) by the explainer is correct. The **augmentation test** refers to adding a certain number of features with high attribution values from a malware sample to a benign sample (i.e., the sample in the opposite class) may make the model’s prediction towards the class of malware, if the explanation is correct. We follow [40] to evaluate the explainer’s performance in the fidelity tests by **positive classification rate (PCR)** that refers to the percentage of samples remains its original class label after the manipulation of deduction or augmentation. Therefore, the PCR for a “good” explainer is as low as possible in a deduction test. A “good” explainer has a high PCR in augmentation.

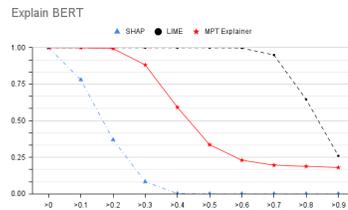
Fig. 3 shows the fidelity of the MPT explainer and SHAP. In deduction test (see Fig. 3a), removing a small number of “important” features (less than 3 features) identified by both SHAP and MPT explainer can reduce the PCR to around 98.92% for SVM classifier. The PCR curve of SHAP then fluctuates sharply when more and more “important” features are used, and it will finally increase to nearly 100% PCR that is higher than MPT explainer’s PCR, when more than 48 features are used (135 features totally). The similar case happens when SHAP explains Random Forest. In comparison, the PCR of the MPT explainer keeps relatively lower than 100% in a stable way (around 98.92% PCR). The unstable PCR curve of SHAP increases the uncertainty of accurately identified “important” features towards the model’s prediction. We analyzed the explanations generated by SHAP and MPT explainer, and found that SHAP usually only focus on a small amount of features (6% features for SVM, and

10% features for Random Forest), which means most of the remaining features are attributed with zeros. This may cause the problem that if one “important” feature is not identified or an unimportant feature is incorrectly identified in the explanation, the accuracy of the whole explanation will be decreased significantly. In comparison, MPT explainer attributes the features by not only considering their individual contributions, but also the contributions made by the dependencies among them. Therefore, in both SVM and RF classifiers, MPT explainer attributes more than 60% features with non-zero values, which forms a large pool of candidate features that have contribution (positive or negative) to the model’s prediction, and these features keeps the stable explanation regardless the number of features used in the deduction tests.

In Fig. 3b, it shows that in the augmentation test, adding a small number of “important” features identified by SHAP from the malware samples to a non-malware sample can make the model’s prediction towards malware (i.e., the opposite class of non-malware samples) with a high PCR value (>80%). However, as the instability issue for SHAP in the deduction test, as long as more “important” features are used, the PCR of SHAP becomes not stable and intends to decrease significantly. The MPT explainer shows a low PCR when a small number of features are used, and a rapid growth of the PCR as the increasing of the “important” features. And finally the PCR of MPT explainer is getting close to that of SHAP. We argue that (1) SHAP is only able to accurately identify the contributions of a small number of features, and therefore, the explanation is not complete and the fluctuating PCR curves along the increasing number of features implies the remaining features that are attributed with low or zero attribution values by SHAP can also affect the model’s behavior significantly and likely to mean inaccurate attributions. (2) MPT explainer, however, explains the model’s prediction on a sample by attributing more features’ contribution/importance, and the rapid growth of PCR curves in the augmentation test implies that the attribution values for the remaining features by MPT explainer are still accurate. This is because MPT explainer takes care of both the individual features’ contribution to the model’s

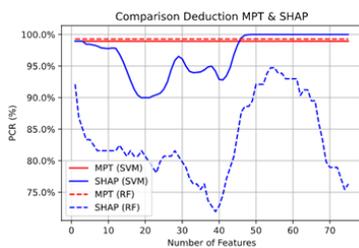


(a) Percentage of “good” explanation identified by LIME, SHAP & MPT Explainer for SVM

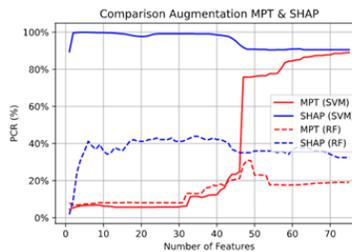


(b) Percentage of “good” explanation identified by LIME, SHAP & MPT Explainer for BERT

Figure 2: “Good” segmentation comparison



(a) Deduction Test



(b) Augmentation Test

Figure 3: **Fidelity test.** (a) Deduction test: SHAP shows unstable PCRs that fluctuate at the beginning for SVM and increased to nearly 100% when more than 50 features used that is higher than the PCR of MPT explainer. (b) Augmentation test: MPT explainer has a low PCR when a small number of features used, and increased quickly to close to the PCR of SHAP, which is decreasing, as the number of features is increasing.

prediction and the dependencies among features.

Therefore, we suggest (1) MPT explainer can be used to generate a more complete explanation of a model’s behavior on a given sample. Users can understand the AI model’s behavior through the contributions made by most of the features in a sample. (2) SHAP’s explanation on a model’s behavior usually focuses on the contribution made by a small number of features. Therefore, it may be a good practice to use SHAP to have a compact explanation (abstract) of a model’s prediction on a given sample.

### 6.5. Running Time Performance

The average running time that MPT explainer explains the Android malware detection by SVM for a single data sample is around 15.44 seconds, and 15.20 seconds for BERT, respectively. The explanation was running through an un-optimized Python code on a PC with GeForce RTX 2070 GPU, 3.60 GHz  $\times$  8 Intel i7-9700K CPU and 62.8 GB memory. Most of the time were spent on the prediction process by the classifiers over the hundreds of perturbed data.

## 7. Conclusion

In this article, we presented a novel explainable AI method that addresses the problem of features attribution for machine learning classifiers used in Android malware detection. This method is inspired by the modern portfolio theory (MPT) that minimizes the variance of the association of the prediction scores changes and the attribution values. By this way, a higher value of a feature’s attribution implies that the small change of this feature will cause a non-trivial change of the model’s prediction score. The effectiveness of the proposed method is assessed by three experiments. The first experiment presents a comprehensive analysis on the capacity of features attribution by the MPT explainer for the adversarial samples, where the malware samples are camouflaged as the benign apps samples and can evade the detection by the SVM and BERT. In the second experiment, we compare the explanation capacity between the MPT explainer and the state-of-the-arts methods, e.g., SHAP and LIME. The results for these two experiments prove the MPT explainer is helpful for the security analysts to find out the reasons that the classifiers are fooled by the adversarial samples, such that a machine learning model for malware detection that is more resistant against such attacks can be developed. The third experiment is to test the explanation fidelity by the MPT explainer, where the comparison with SHAP shows the MPT explainer can explain the model’s behavior in a high fidelity. These results hold a promise that the proposed model explanation method is helpful for both AI and cyber security practitioners.

## References

- [1] A. L. Buczak, E. Guven, A survey of data mining and machine learning methods for cyber security intrusion detection, *IEEE Communications surveys & tutorials* 18 (2) (2015) 1153–1176.
- [2] Y. Xin, L. Kong, Z. Liu, Y. Chen, Y. Li, H. Zhu, M. Gao, H. Hou, C. Wang, Machine learning and deep learning methods for cybersecurity, *IEEE Access* 6 (2018) 35365–35381.
- [3] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno, D. Song, Robust physical-world attacks on deep learning visual classification, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1625–1634.
- [4] M. T. Ribeiro, S. Singh, C. Guestrin, “why should i trust you?” explaining the predictions of any classifier, in: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 2016, pp. 1135–1144.
- [5] M. Sundararajan, A. Taly, Q. Yan, Axiomatic attribution for deep networks, in: *Proceedings of the 34th International Conference on Machine Learning-Volume 70, JMLR. org*, 2017, pp. 3319–3328.
- [6] A. Demontis, M. Melis, B. Biggio, D. Maiorca, D. Arp, K. Rieck, I. Corona, G. Giacinto, F. Roli, Yes, machine learning can be more secure! a case study on android malware detection, *IEEE Transactions on Dependable and Secure Computing*.
- [7] N. McLaughlin, J. Martinez del Rincon, B. Kang, S. Yerima, P. Miller, S. Sezer, Y. Safaei, E. Trickel, Z. Zhao, A. Doupé, et al., Deep android malware detection, in: *Proceedings of the Seventh ACM on Conference on Data and Application Security and Privacy*, 2017, pp. 301–308.
- [8] J. Yan, Y. Qi, Q. Rao, Lstm-based hierarchical denoising network for android malware detection, *Security and Communication Networks* 2018.

- [9] X. Xiao, S. Zhang, F. Mercaldo, G. Hu, A. K. Sangaiah, Android malware detection based on system call sequences and lstm, *Multimedia Tools and Applications* 78 (4) (2019) 3979–3999.
- [10] D. Arp, M. Spreitzenbarth, M. Hubner, H. Gascon, K. Rieck, C. Siemens, Drebin: Effective and explainable detection of android malware in your pocket., in: *Ndss*, Vol. 14, 2014, pp. 23–26.
- [11] I. D. C. (IDC), Smart phone market share, accessed: 2020-02-26.  
URL <https://www.idc.com/promo/smartphone-market-share/os>
- [12] Google play, accessed: 2020-02-25.  
URL <https://play.google.com/>
- [13] Av-test, accessed: 2020-02-25.  
URL <https://www.av-test.org/en/statistics/malware/>
- [14] L. Zhang, V. L. Thing, Y. Cheng, A scalable and extensible framework for android malware detection and family attribution, *Computers & Security* 80 (2019) 120–133.
- [15] S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, in: *Advances in neural information processing systems*, 2017, pp. 4765–4774.
- [16] B. Efron, T. Hastie, I. Johnstone, R. Tibshirani, et al., Least angle regression, *The Annals of statistics* 32 (2) (2004) 407–499.
- [17] H. Markowitz, Portfolio selection, *Journal of Finance* 7 (1) (1952) 77–91.
- [18] M. Zhao, F. Ge, T. Zhang, Z. Yuan, Antimaldroid: An efficient svm-based malware detection framework for android, in: *International Conference on Information Computing and Applications*, Springer, 2011, pp. 158–166.
- [19] W. Li, J. Ge, G. Dai, Detecting malware for android platform: An svm-based approach, in: *2015 IEEE 2nd International Conference on Cyber Security and Cloud Computing*, IEEE, 2015, pp. 464–469.

- [20] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, Bert: Pre-training of deep bidirectional transformers for language understanding, arXiv preprint arXiv:1810.04805.
- [21] H. Meng, V. L. Thing, Y. Cheng, Z. Dai, L. Zhang, A survey of android exploits in the wild, *Computers & Security* 76 (2018) 71–91.
- [22] K. Bakour, H. M. Ünver, R. Ghanem, The android malware detection systems between hope and reality, *SN Applied Sciences* 1 (9) (2019) 1120.
- [23] A. A. A. Samra, Kangbin Yim, O. A. Ghanem, Analysis of clustering technique in android malware detection, in: 2013 Seventh International Conference on Innovative Mobile and Internet Services in Ubiquitous Computing, 2013, pp. 729–733. doi:10.1109/IMIS.2013.111.
- [24] A. Feizollah, N. B. Anuar, R. Salleh, G. Suarez-Tangil, S. Furnell, Androdialysis: Analysis of android intent effectiveness in malware detection, *computers & security* 65 (2017) 121–134.
- [25] M. Zhang, Y. Duan, H. Yin, Z. Zhao, Semantics-aware android malware classification using weighted contextual api dependency graphs, in: Proceedings of the 2014 ACM SIGSAC conference on computer and communications security, 2014, pp. 1105–1116.
- [26] Y. Aafer, W. Du, 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.
- [27] H.-J. Zhu, Z.-H. You, Z.-X. Zhu, W.-L. Shi, X. Chen, L. Cheng, Droiddet: Effective and robust detection of android malware using static analysis along with rotation forest model, *Neurocomputing* 272 (2018) 638 – 646. doi:<https://doi.org/10.1016/j.neucom.2017.07.030>.  
URL <http://www.sciencedirect.com/science/article/pii/S0925231217312870>

- [28] H. Lu, L. Yang, K. Yan, Y. Xue, Z. Gao, A cost-sensitive rotation forest algorithm for gene expression data classification, *Neurocomputing* 228 (2017) 270–276.
- [29] R. Vinayakumar, K. Soman, P. Poornachandran, S. Sachin Kumar, Detecting android malware using long short-term memory (lstm), *Journal of Intelligent & Fuzzy Systems* 34 (3) (2018) 1277–1288.
- [30] F. K. Došilović, M. Brčić, N. Hlupić, Explainable artificial intelligence: A survey, in: 2018 41st International convention on information and communication technology, electronics and microelectronics (MIPRO), IEEE, 2018, pp. 0210–0215.
- [31] Z. C. Lipton, The mythos of model interpretability, *Queue* 16 (3) (2018) 31–57.
- [32] C. Molnar, *Interpretable Machine Learning*, 2019, <https://christophm.github.io/interpretable-ml-book/>.
- [33] M. D. Zeiler, R. Fergus, Visualizing and understanding convolutional networks, in: *European conference on computer vision*, Springer, 2014, pp. 818–833.
- [34] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, Y. Bengio, Show, attend and tell: Neural image caption generation with visual attention, in: *International conference on machine learning*, 2015, pp. 2048–2057.
- [35] W. Samek, T. Wiegand, K.-R. Müller, Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models, *arXiv preprint arXiv:1708.08296*.
- [36] R. Aumann, L. Shapley, Values of non-atomic games. 1974.
- [37] P. K. Mudrakarta, A. Taly, M. Sundararajan, K. Dhamdhere, Did the model understand the question?, *arXiv preprint arXiv:1805.05492*.

- [38] M. Wu, M. C. Hughes, S. Parbhoo, M. Zazzi, V. Roth, F. Doshi-Velez, Beyond sparsity: Tree regularization of deep models for interpretability, in: Thirty-Second AAAI Conference on Artificial Intelligence, 2018.
- [39] N. F. Rajani, B. McCann, C. Xiong, R. Socher, Explain yourself! leveraging language models for commonsense reasoning, in: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, Florence, Italy, 2019, pp. 4932–4942. doi:10.18653/v1/P19-1487.  
URL <https://www.aclweb.org/anthology/P19-1487>
- [40] W. Guo, D. Mu, J. Xu, P. Su, G. Wang, X. Xing, Lemna: Explaining deep learning based security applications, in: Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security, 2018, pp. 364–379.
- [41] E. C. R. Shin, D. Song, R. Moazzezi, Recognizing functions in binaries with neural networks, in: 24th {USENIX} Security Symposium ( {USENIX} Security 15), 2015, pp. 611–626.
- [42] R. Tibshirani, M. Saunders, S. Rosset, J. Zhu, K. Knight, Sparsity and smoothness via the fused lasso, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 67 (1) (2005) 91–108.
- [43] A. Khalili, J. Chen, Variable selection in finite mixture of regression models, *Journal of the American Statistical Association* 102 (479) (2007) 1025–1038.
- [44] What-if tool, accessed: 2020-03-30.  
URL <https://pair-code.github.io/what-if-tool/index.html>
- [45] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke,

- V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, X. Zheng, TensorFlow: Large-scale machine learning on heterogeneous systems, software available from tensorflow.org (2015).  
URL <https://www.tensorflow.org/>
- [46] Captum, accessed: 2020-03-20.  
URL <https://captum.ai/>
- [47] A. Shrikumar, P. Greenside, A. Kundaje, Learning important features through propagating activation differences, in: Proceedings of the 34th International Conference on Machine Learning-Volume 70, JMLR. org, 2017, pp. 3145–3153.
- [48] P.-J. Kindermans, K. Schütt, K.-R. Müller, S. Dähne, Investigating the influence of noise and distractors on the interpretation of neural networks, arXiv preprint arXiv:1611.07270.
- [49] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, A. Lerer, Automatic differentiation in pytorch.
- [50] M. Melis, A. Demontis, M. Pintor, A. Sotgiu, B. Biggio, secml: A python library for secure and explainable machine learning, arXiv preprint arXiv:1912.10013.
- [51] C. D. Manning, P. Raghavan, H. Schütze, Scoring, term weighting, and the vector space model, Cambridge University Press, 2008, p. 100–123. doi:10.1017/CB09780511809071.007.
- [52] Virustotal, accessed: 2020-02-26.  
URL <https://www.virustotal.com/>
- [53] J.-P. Vert, K. Tsuda, B. Schölkopf, A primer on kernel methods, Kernel methods in computational biology 47 (2004) 35–70.
- [54] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, J. Brew, Huggingface’s transformers: State-of-the-art natural language processing, ArXiv abs/1910.03771.

- [55] M. Barreno, B. Nelson, A. D. Joseph, J. D. Tygar, The security of machine learning, *Machine Learning* 81 (2) (2010) 121–148.
- [56] X. Liu, X. Du, X. Zhang, Q. Zhu, H. Wang, M. Guizani, Adversarial samples on android malware detection systems for iot systems, *Sensors* 19 (4) (2019) 974.
- [57] C. Smutz, A. Stavrou, Malicious pdf detection using metadata and structural features, in: *Proceedings of the 28th annual computer security applications conference*, 2012, pp. 239–248.