

Project 4: Classify Android Goodware VS. Malware

This project is due on **November 08, 2024 at 11:59pm CST**. We strongly recommend that you get started early. You have a total of 3 grace days for the 5 MPs, which allow for a 24-hour extension on your MP deadlines. You can use up to 2 grace days per MP, but no more than 3 in total across all MPs. If you submit your assignment after the due date and beyond the grace day period, you will lose 20% of the total points for each late day.

This project consists of three parts. Checkpoint 1 is designed to help you familiarize yourself with the tools you'll be using throughout the project, as well as the process of extracting features from an Android APK. The recommended deadline for checkpoint 1 is October 23rd, 2024. We strongly recommend that you finish the first checkpoint before the recommended deadline. Please ensure that your code is submitted to your personal **GitHub** repository by **November 08, 2024, at 11:59 PM CST**.

In PrairieLearn, you are required to submit a JSON file that includes:

- "netid": Your NetID as a lowercase string.
- "github_token": The fine-grained GitHub access token for your CS 463 repository.

This JSON file will be part of your submission for each checkpoint. Please refer the instructions (<https://github.com/illinois-cs-coursework/generating-access-tokens>) on how to generate the token before making a submission.

This is an individual project; you **SHOULD** work **individually**.

The code and other answers you submit must be entirely your own work, and you are bound by the Student Code. You **MAY** consult with other students about the conceptualization of the project and the meaning of the questions, but you **MUST NOT** look at any part of someone else's solution or collaborate with anyone else. You may consult published references, provided that you appropriately cite them (e.g., with program comments), as you would in an academic paper.

This semester, our autograder includes anti-cheating checks. Plagiarism will be routinely monitored, and any form of cheating, whether giving or receiving help, is not tolerated. You will receive one warning if cheating is detected, and repeated offenses will result in a zero for the assignment and a report to the appropriate authorities.

Solutions **MUST** be submitted electronically on Github repository, following the submission checklist given at the end of each checkpoint. Details on the filename and submission guideline is listed at the end of the document.

Release date: Oct 18th, 2024

Checkpoint 1 Recommended Due date: Oct 23rd, 2024

Checkpoint 2 Due date: Nov 08th, 2024 at 11:59pm CST

Introduction

Android devices have been widely adopted and are capable of storing and accessing significant private and confidential information. Many malware developers started to target android devices and steal private information. In 2021, AV-Test (<https://www.av-test.org/en/statistics/malware/>) recorded almost 3.4 million new pieces of malware for Android alone and has already seen around 1 million malware in 2022. This has created an immediate need for security professionals to build detection frameworks and classify malware from all the applications in the market.

Please read the assignment carefully before you start implementing.

Objectives

- Know how to extract features from APKs with static analysis.
- Be able to build an ML classifier to distinguish malware from goodware.
- Create adversarial samples with the consideration of realizability of Android apps.

Checkpoints

This machine problem is split into 3 checkpoints. Checkpoint 1 will help you get familiar with what features are commonly extracted from Android apps and how to extract them. You will extract features for 2 apks with the apktool. In Checkpoint 2, you will be given an Android dataset (with feature vectors and ground-truth labels). You need to implement a LinearSVM classifier to predict whether a sample is benign or malicious. In Checkpoint 3, you will manipulate a few features (with and without realizability constraints) to create adversarial samples.

Implementations

To help you with the implementation, you are provided with a code skeleton which prototypes the main classes and includes some useful methods. You are free to modify the provided code as long as they generate correct results.

1 Checkpoint 1 (20 points)

This machine problem is split into 3 checkpoints. Checkpoint 1 will help you get familiar with the Android feature extraction. You don't need to write any code for this part, all you need to do is to set up the appropriate environment and run the tool.

1.1 Environment Setup

The tool we are using to extract features from Android apks is called apktool. This is a tool for reverse engineering 3rd party, closed, binary Android apps. You can install apktool using Docker by following section 4.1.1.1 (recommended), or alternatively install it on your local system by following instructions in section 4.1.1.2.

1.1.1 Set up Dockerfile Environment

1. We have provided a dockerfile, which you are not required to use, but we recommend doing so. If you are unfamiliar with docker, please read this overview: <https://docs.docker.com/get-started/>.
2. Docker can be installed using the instructions listed here: <https://docs.docker.com/get-docker/>.
3. To build the provided dockerfile, navigate to the location of the dockerfile in your command prompt and run the command `docker build -t apkdocker`. Note: you can replace apkdocker with any name you wish to give this image.
4. To run the built docker image, use `docker run -it apkdocker`. This will open a terminal in that image with the required setup.

1.1.2 APK Tool

The apktool will allow us to statically analyze the apk files that have been provided.

1. To install apktool follow the instructions provided here <https://ibotpeaches.github.io/Apktool/install/>. If you are using the dockerfile environment, use the instructions provided for the Linux environment.
2. Alternatively, you can download Homebrew using this command: `/bin/bash -c "$(curl -fsSL https://raw.githubusercontent.com/Homebrew/install/HEAD/install.sh)"`
Note: There will be two required updates at the end of the comments, `(echo; echo 'eval "$(/home/linuxbrew/.linuxbrew/bin/brew shellenv)"') » /root/.profile` and `eval "$(/home/linuxbrew/.linuxbrew/bin/brew shellenv)"`. Once you have Homebrew installed, use the command `brew install apktool`
3. Once the apktool has been installed, run the analysis on each of the apk files using the command `apktool d <apkfile>` where <apkfile> should be replaced with the name of the apk file (with the .apk extension) to disassemble.

What to submit

- AndroidManifest_49875A.xml
- AndroidManifest_5E06B7.xml

2 Checkpoint 2 (40 points (35 + 5))

For this checkpoint, you'll need to implement a LinearSVM classification model to differentiate malware and goodware. We have provided a .json file for both apks with the features that were extracted during Checkpoint 1 analysis that you will use as baseline samples.

2.1 Environment Setup

Due to GitHub storage limit, you need to unzip “data/apg-X.json.zip” first, the unzipped file should be “data/apg-X.json”.

You can follow the steps to install these libraries (If you already set it up during MP1, you only need to do the last two steps:

1. Install Anaconda here: <https://www.anaconda.com/products/distribution>.
2. Create a virtual environment named as cs463 with Python 3.8.8: `conda create -n cs463 python==3.8.8`
3. Activate the virtual environment: `conda activate cs463`
4. `pip install -r requirements.txt` (the one in the root folder of MP4 with numpy and scikit-learn only.)

2.2 Train an SVM model (20 points)

SVM is very popular in Android malware detection. You need to implement the “fit()” function of class SVM in “mp4_model.py” and return the classifier. Since our “perform_feature_selection()” is pretty generic, you should also implement half of it to deal with the case that we actually don't do any feature selection. You should carefully read how the performance is calculated and F_1 , precision, and recall on the testing set would be returned. Another function you need to implement is “stat()” in “main.py”. You should calculate how many samples were used in training; how many samples were from testing; how many malware and goodware in the training set, and how many malware and goodware in the testing set. Think about if the goodware VS malware ratio matches with your common sense.

Tips

- Please refer to Scikit-learn documentation on LinearSVC here: <https://scikit-learn.org/stable/modules/generated/sklearn.svm.LinearSVC.html>.

What to submit

- “main.py”.
- You can submit “mp4_model.py” together with the next subsection.

2.3 Feature Selection (15 points)

If you take a look at the shape of `svm_model.X_train` in `main.py`, you will notice that the original dataset is very sparse and could easily go over to more than 200 thousand features. This is neither efficient nor effective. So you need to implement the `perform_feature_selection()` function in `mp4_model.py`. To achieve slightly better performance, you can use Scikit-learn LinearSVC module to perform a feature selection to reduce the feature dimension to 5,000. Please feel free to try other dimensions. The final code would be graded based on 5,000.

What to submit

- `mp4_model.py`
- Implementations for this section.
- `data/all_feature_names_5000.json`, generated by your model.
- Please put your thought as code comments to compare the performance of with and without feature selection and why do you think one could achieve a little bit better results.

Comments are worth 5 points.

3 Checkpoint 3 (40 points)

As part of this checkpoint you'll manually create some adversarial samples to make your trained classifier misclassify the sample after your perturbation. We have provided a `.json` file for both apks with the features that were extracted during Checkpoint 1 analysis that you will use as baseline samples.

3.1 Create Adversarial Samples

You should have learned about adversarial samples from the lecture. But how could we generate them? There are many optimization-based methods but here we want to focus on the simplicity. Given the two apks provided in Checkpoint 1 in the `cp1_feature_extract/mp4_apks` folder, 49875AXXX is malware while 5E06B7XXX is goodware. Both of them are correctly predicted as their corresponding label given the classifier trained using 5000 features (already verified). Your goal is to change some feature values to artificially create a new sample and it would be misclassified as "benign/malicious". For example, 49875AXXX should be predicted as benign by the original classifier while 5E06B7XXX was predicted as malicious. To make it easier, I have provided the name of the 5000 features in `data/all_feature_names_5000_example.json`. This file corresponds to the feature vector. You can use the first few lines of `train_model()` to generate your feature names. Theoretically, they should be the same as the one I provided.

To create an adversarial example, usually you should find some features and change their values to make them misclassified. But in the context of security, there are more constraints compared to common Machine Learning domains such as images.

1. To protect the integrity of these apks, usually we can only add features instead of removing features because the latter can easily break the dependencies.
2. Imagine that we need the adversarial sample to be effective in the real world, we need to make sure that it still keep the original malicious / benign functionality, this is usually called problem space attack compared to commonly used feature space attack (more details can be found here: <https://arxiv.org/abs/1911.02142>). In reality, it's usually difficult to add some features like the activity or API call or other features that need to add code. So similar to <https://www.usenix.org/conference/usenixsecurity21/presentation/severi>, we only allow certain types of features can be added, i.e., the features from the Android Manifest file which does not affect the code logic at all.
3. You should not add too many features which would make the sample suspicious.

To perform the problem space attack, you should follow these guidelines:

1. You can only add features. That means, you can choose some features and mark their values as 1 and add to the original feature vector of the sample, see if they can be misclassified by your SVM model (trained using 5000 features).
2. You can start with 49875AXXX.json and 5E06B7XXX.json as the base features for the two samples and use added-features-xxx.txt (please see the boilerplate code) to add features to make it adversarial.
3. You can at most add 30 features. We need to stay stealthy.
4. After adding features, you will need to merge and create one feature vector to run the classification (more details in the verify_adv.py)
5. You should only choose feature names start with "app_permission" or "api_permission".
6. You should complete script "verify_adv.py" on how to merge your added features with the original feature vector and use your trained SVM model to make a prediction (output the predicted label in the way mentioned in the comments of verify_adv.py). You only need to find one solution for each sample.

Tips

- Think about which features are benign-prone or malicious-prone. You may use SVM to make an explanation on the feature weights.

What to submit

- mp4_model.py from previous checkpoint
- data/all_feature_names_5000.json

- added-features-5E06B7.txt
- added-features-49875A.txt
- svmf5000.p that will be saved from checkpoint 2 automatically
- verify_adv.py with code comments why these features are helpful.

You would not need to delete any files or submit selected files to the autograder, the above is a minimum checklist the autograde requires to grade the checkpoint. You can submit the repository as it is after making sure these files are present

Autograder

We have set up an autograder on PrairieLearn. To use it, you must first enroll in the course CS 463: Computer Security II by visiting <https://us.prairielearn.com/pl>.

After enrolling, navigate to the **Assessments** section within CS 463: Computer Security II, where you will find **MP4** here:

https://us.prairielearn.com/pl/course_instance/163006/assessment/2468409.

In PrairieLearn, you are required to submit a JSON file containing the following key-value pairs:

- "netid": Your NetID as a lowercase string.
- "github_token": The fine-grained GitHub access token for your CS 463 repository.

This JSON file is essential for each checkpoint submission, as it verifies your identity and access rights.

Before submitting, please refer to the instructions on how to generate the access token at the following link: (<https://github.com/illinois-cs-coursework/generating-access-tokens>). Ensure that you create the token correctly to avoid any submission issues.

Important Note

1. Your highest grade and grading report will be automatically recorded in our server
2. You must submit through **PrairieLearn** to our auto-grader at least once before the deadline to get a grade. We don't grade your code manually, the only way is to use our auto-grader.
3. Try auto-grader as early as possible to avoid the heavy traffic before the deadline. **DON'T DO EVERYTHING IN THE LAST MINUTE!**
4. If the auto-grader is experiencing downtime or behaving unexpectedly, don't worry. Reach out to us on Campuswire, and we will address the issue and provide an extension if needed.

5. We have anti-cheating mechanisms, NEVER CHEAT OR HELP OTHERS CHEAT!
6. Please use Python version 3.7 for checkpoint 2. This checkpoint execution may take some time as you are building a model. Please be patient and allow the process to complete.
7. Never abuse the auto-grader. You should only submit your own code. Please do not abuse the auto-grader in any form. Abusing the autograder is considered a form of violation of the student code of conduct, and the corresponding evidence will be gathered and reported

Submission Instructions

You need to complete two tasks to submit your work:

1. Upload to GitHub

Check your repository to confirm that all files have been uploaded successfully at the end of each checkpoint. We will also perform anti-cheating checks to detect any plagiarism. Please note that any form of code plagiarism is strictly prohibited, including copying parts of someone else's code or code from previous students.

Folder Structure

- Make sure there is a folder "mp4" inside the root folder of your Github repository in the main branch. Your score will be deducted if you fail to do so.

Below is the folder structure for MP4:

```
mp4
├── data
│   └── all_feature_names_5000.json
├── AndroidManifest_49875A.xml
├── AndroidManifest_5E06B7.xml
├── added-features-49875A.txt
├── added-features-5E06B7.txt
├── main.py
├── mp4_model.py
├── svm-f5000.p
└── verify_adv.py
```

2. Upload to PrairieLearn

To facilitate the grading of your submissions, it is essential that you submit a JSON file containing: your NetID and your GitHub token.

Make sure you follow the below steps to ensure that your submission process runs smoothly.

1. JSON File Requirements:

- Your submission must be in JSON format, which should include the following keys:
 - "netid": Your NetID as a lowercase string.
 - "github_token": The GitHub access token generated for your account.

2. Token Visibility:

- Please remember to save the generated token securely, as it will only be displayed temporarily after creation. Once you navigate away from the page, you will not be able to retrieve it again.

3. Access Permissions:

- Ensure that you have provided Read-Only access to the contents of your repository.

Submission Limit

Please note that you are limited to a maximum of 30 submissions for each checkpoint on PrairieLearn. Be mindful of each submission you make, as your submission count is limited.

Git Upload

Use the following commands to push your submission for grading (it is advisable to do this frequently for better version control). Grading will be based on the latest submission to the repository before the deadline. Please push to your main branch rather than creating a separate branch.

```
git add -A
git commit -m "REPLACE THIS WITH YOUR COMMIT MESSAGE"
git push origin main
```