

Selina Narain, Neelam Boywah, Zoya Haq
DTSC 870 - Masters Project - Fall 2023
Advisor: Professor Dr. Wenjia Li

Progress Report 5

Timeline: November 15th, 2023 - November 29th, 2023

Accomplishments: What did you accomplish?

Research Topic Idea:

- Comparing machine learning and deep learning algorithms for accuracy and efficiency in detecting malware in android applications.
- Applying adversarial attack on Random Forest machine learning model.

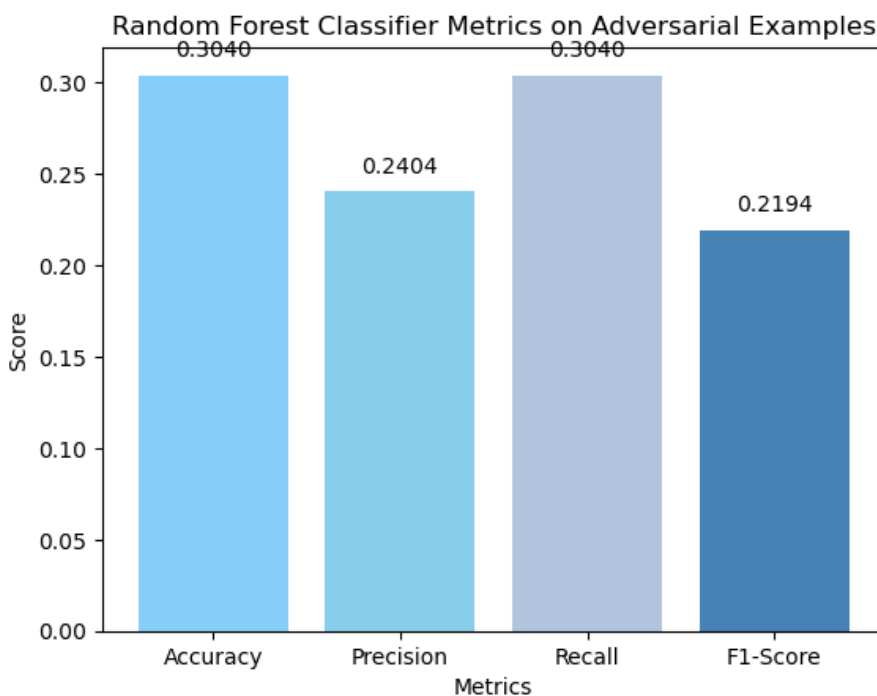
Google Site

- Contains links to Research Paper, Final Report and Presentation Slides.
- Team Members and education are displayed.

Implementation for CICMaldroid 2020 Dataset:

- Trained and tested SVM Models with 2 kernels
- SVM parameters: kernel = 'rbf', C=1.0 and random_state = 42
- SVM Statistics (Kernel - 'rbf'):
 - Accuracy: 0.7979
 - Precision: 0.8000
 - Recall: 0.7979
 - F1-Score: 0.7942
- SVM parameters: kernel = 'linear', C=1.0 and random_state = 42
- SVM Statistics (Kernel - 'linear'):
 - Accuracy: 0.8317
 - Precision: 0.8323
 - Recall: 0.8317
 - F1-Score: 0.8285
- Deep Learning Model: Connected Neural Network (CNN)
- Built a Deep Learning CNN model and trained it. The model architecture uses the Sequential function and is built using 2 Convolutional 2D layers, 2 Max Pooling 2D layers, 1 Flatten layer and a Dense layer.

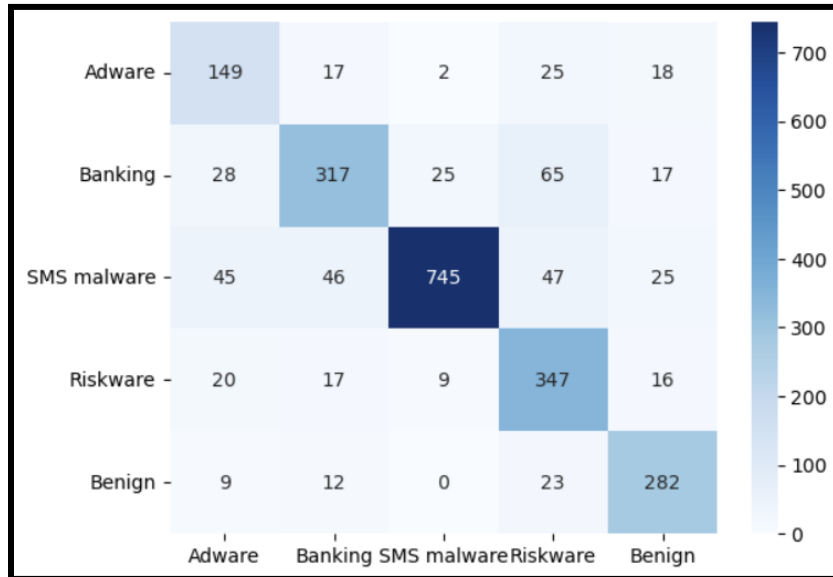
- Parameters for the Convolutional 2D layers: number of filters and filter size, and activation function - 'relu'
- Parameter for the Max Pooling 2D layers: strides = 2, moving 2 filters at a time.
- Parameter for the Dense layer: softmax activation function.
- Comparison of models: Display 6 models using a bar graph to show their metrics.
- Adversarial Attack on Random Forest Model
- The adversarial attack has been defined and integrated with the Random Forest model.
- Bar Graph Visualization was created to show the Random Forest Classifier metrics reacted to the adversarial attack.
- Random Forest Classifier Metrics on Adversarial Attack:
 - Adversarial Random Forest Accuracy: 0.3040
 - Adversarial Random Forest Precision: 0.2404
 - Adversarial Random Forest Recall: 0.3040
 - Adversarial Random Forest F1-Score: 0.2194



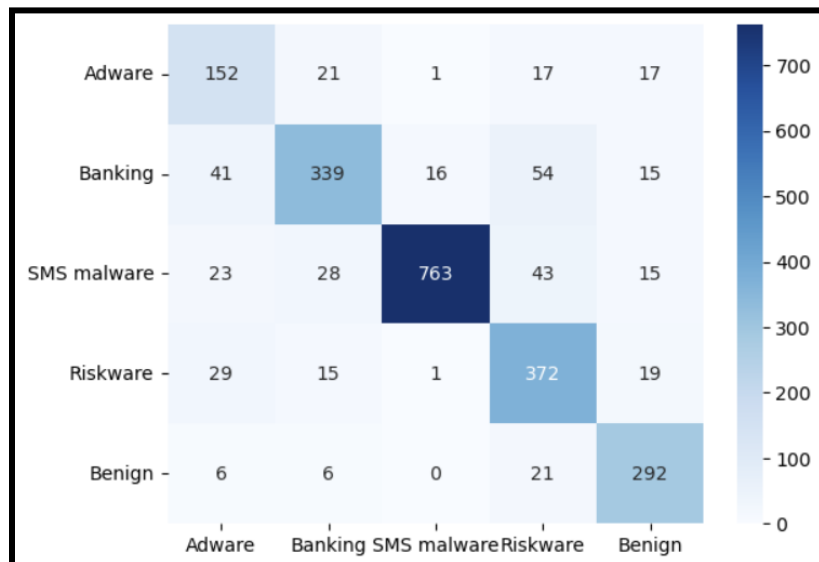
Analysis:

- The SVM models did not outperform the Random Forest Model and K-Nearest Neighbor model. However, they did do better than the Naive Bayes model and Logistic Regression.
- Heatmaps were created for the SVM models to display the amount of features used for Adware, Banking, SMS, Malware, Riskware and Benign.

- SVM Heatmap (Kernel - 'RBF')
 - Adware Features: 149
 - Banking Features: 317
 - SMS Malware Features: 745
 - Riskware Features: 347
 - Benign Features: 282



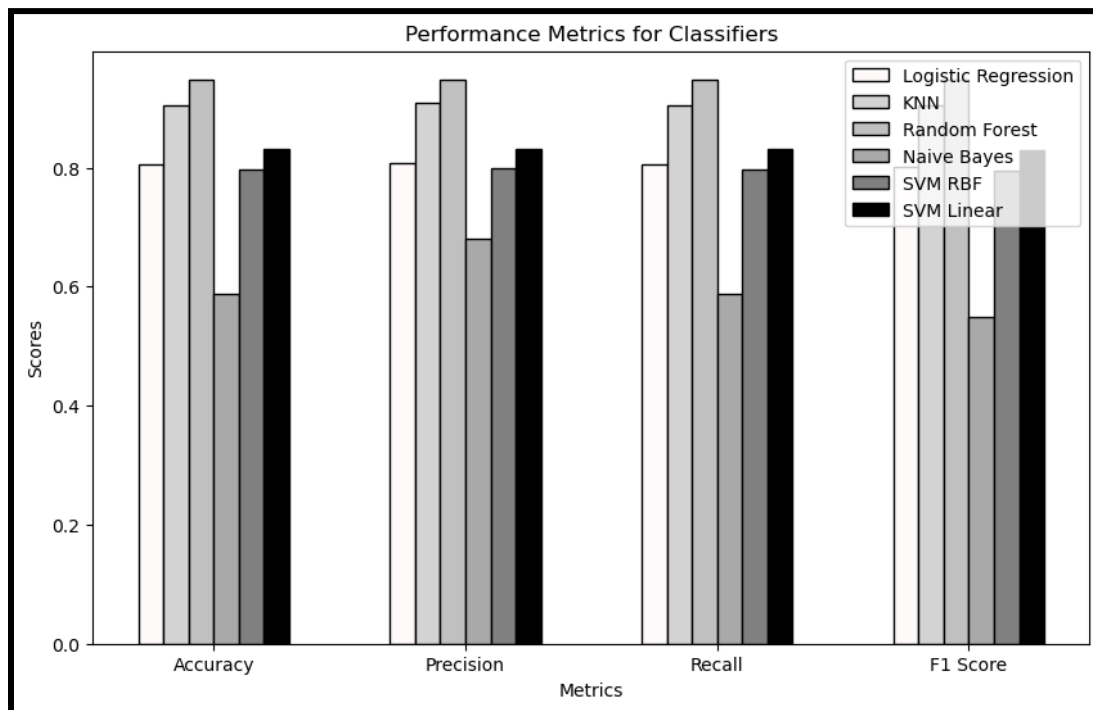
- SVM Heatmap (Kernel - 'Linear')
 - Adware Features: 152
 - Banking Features: 339
 - SMS Malware Features: 763
 - Riskware Features: 372
 - Benign Features: 292



- Adversarial Attack on Random Forest Model
- We can see that after applying the adversarial attack on the Random Forest Model, we see that the attack caused a significant change in Random Forest models performance.
- From the Random Forest model being able to perform and have an accuracy score of 0.9422 to being attacked and currently has an accuracy of only 0.3040.

Performance Metrics for Classifiers Bar Graph Visualization

- The graph displays 6 machine learning models: Logistic Regression, KNN, Random Forest, Naive Bayes, SVM RBF and SVM Linear.
- For each model, the values for accuracy, precision, recall and f1-score are shown.
- The visualization shows that the Random Forest Model has the highest values for all metrics and the Naive Bayes Model has the lowest values for all metrics.



Implementing AndroZoo Dataset

We used the API documentation provided by AndroZoo to test how we can obtain the APK files and once obtained, how we can reverse engineer the files. AndroZoo has a CSV file with all the hash values that we can use with the API key they provided to us in order to obtain the APK file. We were able to successfully obtain the APK and then use apktool to obtain the manifest.xml file. From the image below, you can see the android permissions from this specific application.

```
APKtool — -zsh — 89x22
Last login: Mon Nov 27 18:05:54 on ttys001
(base) selinanarain@Selinas-MacBook-Pro APKtool % apktool --version
2.9.0
(base) selinanarain@Selinas-MacBook-Pro APKtool % apktool d 038847E28B4E70D8DA5BF49B36C16]
616A517964CF9F81058F28DD237D1186D46.apk
I: Using Apktool 2.9.0 on 038847E28B4E70D8DA5BF49B36C16616A517964CF9F81058F28DD237D1186D4
6.apk
I: Loading resource table...
I: Decoding file-resources...
I: Loading resource table from file: /Users/selinanarain/Library/apktool/framework/1.apk
I: Decoding values */* XMLs...
I: Decoding AndroidManifest.xml with resources...
I: Regular manifest package...
I: Baksmaling classes.dex...
I: Copying assets and libs...
I: Copying unknown files...
I: Copying original files...
(base) selinanarain@Selinas-MacBook-Pro APKtool % █

▼<manifest xmlns:android="http://schemas.android.com/apk/res/android" package="com.asish.megapiayer.mp3downloader.views">
  <uses-permission android:name="android.permission.INTERNET"/>
  <uses-permission android:name="android.permission.ACCESS_NETWORK_STATE"/>
  <uses-permission android:name="android.permission.READ_PHONE_STATE"/>
  <uses-permission android:name="android.permission.SEND_SMS"/>
  <uses-permission android:name="android.permission.WRITE_EXTERNAL_STORAGE"/>
  <uses-permission android:name="android.permission.WRITE_SETTINGS"/>
  <uses-permission android:name="android.permission.RECEIVE_SMS"/>
  <uses-permission android:name="android.permission.READ_SMS"/>
  <uses-permission android:name="android.permission.WRITE_SMS"/>
  <uses-permission android:name="android.permission.READ_PHONE_STATE"/>
  <uses-permission android:name="android.permission.ACCESS_WIFI_STATE"/>
  <uses-permission android:name="android.permission.CHANGE_WIFI_STATE"/>
```

PowerBI Visualizations Report

(Based on the performance of ML models for 2 Datasets)

- An excel spreadsheet was created with the following categories and statistical metrics: Datasets, Machine Learning Models, Accuracy Score, Precision, Recall and F1-Score.
- The Android Malware Detection Visualizations Report displays six filters at the top that can be manipulated based on the user's preference of what they would like to see.
- The two datasets we are trying to compare metrics for are the New Brunswick: CICMaldroid 2020 Dataset and AndroZoo Dataset.
- The table displays the outputs of the filter inputs. The default is set to all datasets, models and metrics showing at once.
- The four visualizations created are Comparisons between Datasets and machine learning models for each of the metrics (Accuracy Score, Precision, F1-Score, and Recall).

Upcoming Plan: What do you plan to do in upcoming weeks?

- Continue to edit and fine tune the adversarial attack applied on the Random Forest Model.
- Work on achieving metrics for the Deep Learning - Convolutional Neural Network (CNN) model.
- Gather apks from AndroZoo API and extract permissions.
- Edit the SVM models parameters to try and improve the metrics. Utilize verbose and max_iter parameters.

Obstacles & Concerns: Were there any obstacles or barriers that prevented you from getting things done?

- We are getting an error when compiling the Deep Learning CNN model, so we are currently trying to fix the error.
- We are still tweaking the Adversarial Attack Implementation integrated with the Random Forest model.
- Clarification: Any Specific Requirements for Google Site?
- Finalize date for presentation ?
- Review/Edit Final Report Outline
- Review/Edit Research Paper Outline