

# COMP3217 Assignment 2: ML For Security

By Gherman Nicolisin (gn2g21)

## Contents

1. Part A .....	2
1.1. Problem .....	2
1.2. Solution .....	2
1.3. Results .....	2
2. Part B .....	3
2.1. Problem .....	3
2.2. Solution .....	3
2.3. Results .....	3
3. Code .....	4
3.1. Usage .....	4
3.2. Local settings .....	4

## 1. Part A

## 1.1. Problem

## Binary classification of 6000 system traces traces

## 1.2. Solution

The training data was split into 80/20 ratio. With 20% of data used for testing. After the model was selected, it is trained with 100% of the training data to further improve the accuracy.

**Logistic Regression Classifier** was initially used and produced results of decent accuracy (~90%). Further action was taken to optimise a hyperparameter  $C$ , which resulted in the boost of ~3%. **Decision Tree Classifier** was then used as it does not require scaling and normalisation of data and produced a model with ~96% accuracy. **Support Vector Machine** was also tried but did not yield accurate model even with tuning. **Random Forest Classifier** was also tried since it operates well on continuous data which is present in the dataset. The model yielded ~98% accuracy.

Finally, during the part B classification, the **XGBoost** showed remarkable performance over the Random Forest Classifier when performing multiclass classification. Therefore, it was also tested in the part A with binary classification and showed a significant improvement in accuracy up to **~99%**.

### 1.3. Results

- Model: **XGBoost**
- Accuracy: **0.99 (99.00%)**
- F1 score: **0.9899959983993598 (99.00%)**
- Labels

```
[1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0  
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0]
```

## 2. Part B

### 2.1. Problem

Multiclass classification of 6000 system traces

### 2.2. Solution

Similarly, as in part A, the data was split into 80/20 ration for training and testing purposes, and 100% of data for final training.

**Random Forest Classifier** was again tested and produced a consistent good accuracy of **~95%**. **Logic Regression Classifier** performed even worse than in the part A even with hyperparameter tuning with the accuracy of **~67%**, and significant downgrade in runtime speed. The **XGBoost** was discovered due to its incorporated regularization while working with this dataset and better runtime performance. The raw model perfomed faster and slightly more accurately with the score of **~96%** than the Random Tree Forest. Further tuning took place using RandomizedSearchCV to find optimum values for `colsample_bytree`, `gamma`, `learning_rate`, `max_depth`, `n_estimators`, `subsample` as well as running on scaled features. The calculated hyperparameters was then used in the refined model and gave further improvement in accuracy up to **~97%**.

### 2.3. Results

- Model: **XGBoost (tuned)**
- Accuracy: **0.97 (97.00%)**
- F1 score: **0.9699129151738807 (96.99%)**
- Labels:

```
[0 0 0 0 0 0 2 2 2 2 2 2 0 0 0 0 1 2 2 2 0 0 0 0 0 0 0 0 1 1 1 1 1 1 2 2
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 0 0 0 0 0 2
0 0 0 0 0 0 0 0 2 2 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0]
```

## 3. Code

The code can be access in Google Collab and is published on GitHub: [https://github.com/SkymanOne/COMP3217\\_CW2](https://github.com/SkymanOne/COMP3217_CW2).

### 3.1. Usage

- Navigate to `COMP3217_CW_2.ipynb` to view the model submissions.
- Execute code blocks as instructed at the beginning of the section A and the Section B
- Navigate to `COMP3217_CW_2_Experiments.ipynb` to view other models.

### 3.2. Local settings

- Update paths for `train_data` and `test_data` CSV files in parts A and B.
- The notebook has been configured to use the local running image for better performance, therefore, it is important to change remote instance if you do not have the local execution environment running.
- If you run hyperparameter fitting for verifying the results, adjust `n_jobs` according to your environment computational capabilities.
- Finally, adjust the paths to the datasets' CSVs in the prelude blocks depending on your configuration of the environment.