National
College of
Ireland

# Configuration Manual for A Hybrid Ensemble Model using XGBoost and AdaBoost to detect and distinguish zero-day attacks

MSc Research Project
MSc Cyber Security

Ajay Krishna Edakkat Parambil
Student ID: X22110674

School of Computing
National College of Ireland

Supervisor: Vikas Sahni

| | |
|---|---|
| **Student Name:** | Ajay krishna Edakkat Parambil |
| **Student ID:** | X22110674 |
| **Programme:** | MSc Cyber Security |
| **Module:** | MSc Academic Internship |
| **Lecturer:** | Vikas Sahni |
| **Submission Due Date:** | January 31 |
| **Project Title:** | A Hybrid Ensemble Model using XGBoost and AdaBoost to detect and distinguish zero-day attacks |

**Year:** 2023-2024

**Word Count:** …………445…………… **Page Count:** ………………4…………………..……

**Signature:** AJAY KRISHNA EDAKKAT PARAMBIL……………………………

**Date:** …31-01-2024……………………………………………………………………

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST**

# Configuration Manual

Ajay Krishna Edakkat Parambil
x22110674

# 1  Introduction

This document discusses about the Hybrid Ensemble Model using base classifiers Random Forest, Decision Tree and AdaBoost as well as a meta-classifier called XGBoost was implemented and executed. The project is run in python programming language in google collab.

# 2  System Specifications

The Application is run in the following specifications
Code Editor : Google Collab
Python Version : 3.6.3
Operating System : MacOs Sonoma Version 14.1.2

## 2.2 Software Requirements

Google Colaboratory: cloud-based jupyter notebook, python version 3.6.3

# 3  Package Details

NumPy and Pandas for data manipulation, Seaborn for data visualisation, imbalanced-learn (imblearn) for addressing class imbalance using SMOTE, scikit-learn for various machine learning functionalities, matplotlib for plotting and charting, XGBoost for gradient boosting, mlxtend for stacking classifier implementation, and other specific modules for tasks such as classification reports, train-test splitting, decision tree and random forest classifiers label.

| Package | Version |
|---------|---------|
| numpy | 1.23.5 |
| pandas | 1.5.3 |
| seaborn | 0.12.2 |
| matplotlib | 3.7.1 |
| xgboost | 2.0.2 |
| mlxtend | 0.22.0 |

numpy : It is used for working with arrays. It also has functions for working in domain of linear algebra, fourier transform, and matrices

pandas : It offers data structures and operations for manipulating numerical tables and time series

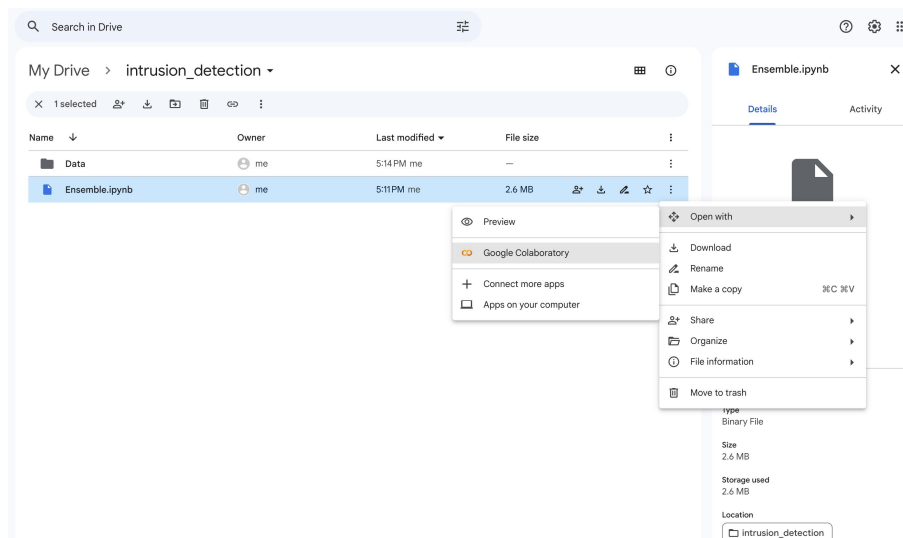matplotlib: It is used for plotting various graphical visualisations

# 4 DataSet

The dataset is from Kaggle[1], UNSW-NB15 a network intrusion dataset that contains raw network packets. It was created by the IXIA PerfectStorm tool in the Cyber Range Lab of the Australian Centre for Cyber Security at the University of New South Wales (UNSW) Canberra for generating a hybrid of real modern normal activities and synthetic contemporary attack behaviors. A data training set named "UNSW_NB15_training-set.csv" was taken from this dataset.

# 5 Execution

## Environment Setup

Google Collaboratory is used for the process of developing models. It is necessary to have a working Gmail account in order to access the Google Collaboratory. Python version 3.6.3 is used during the whole Model creation process.

1) After uploading the project file to your google repository, Now select the project file "open with -> Google Colaboratory



---

1 https://www.kaggle.com/datasets/mrwellsdavid/unsw-nb15/

2) After opening the project file with Colaboratory



3) Now change the location of the dataset to the correct location in your drive and click run



4) Now click on Run time and select run all

5) Now the dataset will be balanced with Smote(Juanjuan et al., 2007) all the models will start running and you can find the precision, accuracy, f1 score and recall of each model.

6) Additionaly the explanation for the code is provided in the link given below:

[x22110674_academic_thesis_presentation&demo.mp4](x22110674_academic_thesis_presentation&demo.mp4)

# References

Juanjuan, W., Mantao, X., Hui, W. and Jiwu, Z. (2007) 'Classification of imbalanced data by using the SMOTE algorithm and locally linear embedding', *In International Conference on Signal Processing Proceedings, ICSP*.