# Configuration Manual

MSc Internship
Cyber Security

# Yash Shukla
Student ID: x18175104

School of Computing
National College of Ireland

Supervisor:     Mr.Vikas Sahni

**National College of Ireland**
**Project Submission Sheet**
**School of Computing**

| Student Name: | Yash Shukla |
|---|---|
| Student ID: | x18175104 |
| Programme: | Cyber Security |
| Year: | 2019 |
| Module: | MSc Internship |
| Supervisor: | Mr.Vikas Sahni |
| Submission Due Date: | 17/08/2020 |
| Project Title: | Configuration Manual |
| Word Count: | 610 |
| Page Count: | 8 |

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

**ALL** internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

I agree to an electronic copy of my thesis being made publicly available on NORMA the National College of Ireland's Institutional Repository for consultation.

| Signature: | Yash Shukla |
|---|---|
| Date: | 14th August 2020 |

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:**

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies). | ☐ |
| **Attach a Moodle submission receipt of the online project submission**, to each project (including multiple copies). | ☐ |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. | ☐ |

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

| Office Use Only | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# Configuration Manual

Yash Shukla

x18175104

Threat Hunting Using a Machine Learning Approach

# 1 Intended Audience

This document would be appropriate for researchers programmers or system administrators with a moderate degree of technical understanding of python. It explains how to configure a Linux system for executing python code and write adequate hypothesis required for threat hunting. [1].

Abbreviations:

1. IP= Internet Protocol

2. Pd= Pandas variable

3. Np= Numpy variable

4. plt= matplotlib.pyplot variable

5. cm= matplotlib.cm variable

6. NB= Naive Bayes

7. SVC= Support Vector Classification

# 2 Installing Python

Python is installed by default on some Linux distributions. (Note:This may not be the latest version or support all the functionality needed by the application)

Figure 1: Python Version Check

## 2.1 Steps for Ubuntu

- At first check the current installed version of python

  | Command :python —version |
  | --- |

- Get an update on your system

  | command : sudo apt update |
  | --- |

- Install python

  | command : sudo apt install python(req version) |
  | --- |

## 2.2 Data conversion steps

- The pcap file is processed by a program called Tshark. It is a network protocol analyzer, that has the ability to read saved network files. The program is responsible for the conversion of pcap files to text.

Figure 2: Tshark Installation

- Install Tshark using       sudo apt-get install -y tshark

- This step is essential for taking the pcap data into a python program using dataframe. A dataframe is a type of data structure in python which helps access data in the form of rows and columns.

- At this phase data cleaning is necessary, all the columns with missing values are checked and the missing values are either replaced or their mean is inserted.

- Dataset to be downloaded from https://www.netresec.com/?page=PcapFiles
  The pcap file has the following features available for selection, which are inserted into the dataframe.

The following are the columns that are derived from the pcap conversion:

1. 'Sr.No'

2. 'Date'

3. 'Time'

4. 'Source_IP'

5. 'Arrrow'

6. 'Dest_IP'

7. 'ProtocolUsed'

8. 'Length','Method'

9. 'link'

10. 'Version'

11. 'Extra'



Figure 3: Orignal Pcap file



Figure 4: Labelled data Frame

- Use the following command to generate text file from pcap file

  tshark -r /PATH/Filename. Pcap -t ad > /PATH/ Filename.txt

- The final step is to label the text file and use it for machine learning.

## 2.3 Libraries that have been imported for machine learning

1. from pandas import DataFrame [2]

2. import pandas as pd

3. import pandas

4. import numpy [3]

5. import numpy as np
   SKlearn [4]

6. from sklearn.externals import joblib

7. from sklearn.naive_bayes import MultinomialNB

8. from sklearn.feature_extraction.text import CountVectorizer

9. from sklearn.svm import LinearSVC

10. from sklearn.ensemble import VotingClassifier

11. from sklearn.linear_model import LogisticRegression

12. from optparse import OptionParser

13. from sklearn.cluster import KMeans

14. from sklearn import metrics

15. from sklearn.model_selection import train_test_split

16. from sklearn.model_selection import *

17. from sklearn import model_selection

18. from sklearn.linear_model import LogisticRegression

19. from sklearn import EnsembleVoteClassifier

20. import matplotlib [5]

21. import matplotlib.pyplot as plt

22. import matplotlib.cm as cm

# 3 Execution steps

1. python file.py

   Execute the file using [ python -w file.py ] to convert the pcap text file to a dataframe.
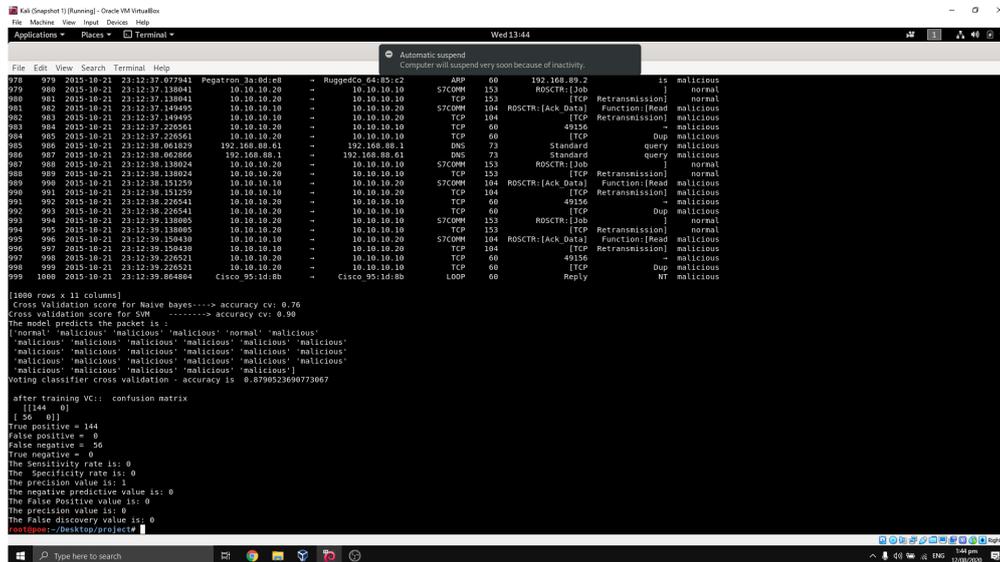


Figure 5: Python Code for machine learning models

2. python check.py

   Execute the file using [ python -w check.py ] to convert the pcap text file to a dataframe.
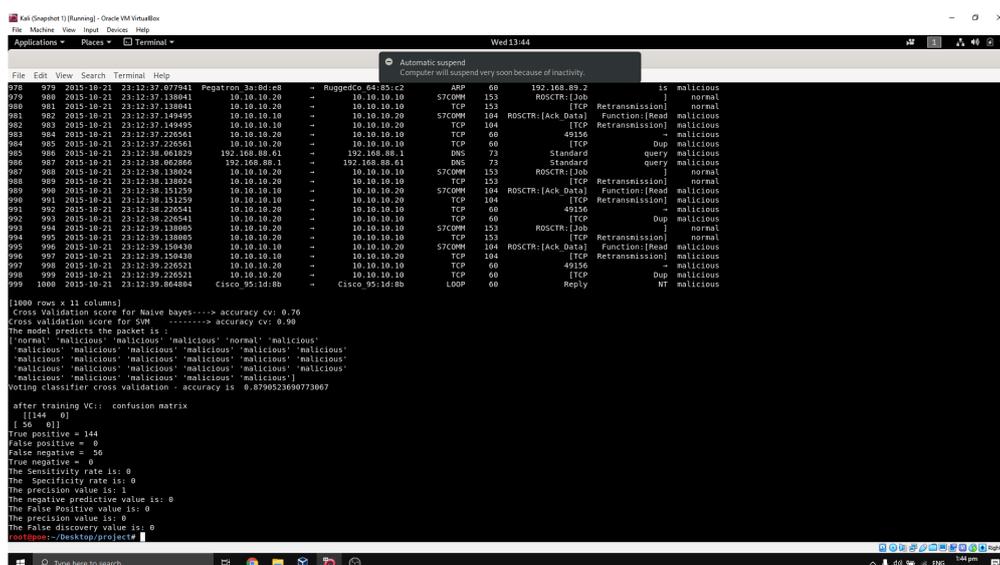


Figure 6: Python Code for machine learning models
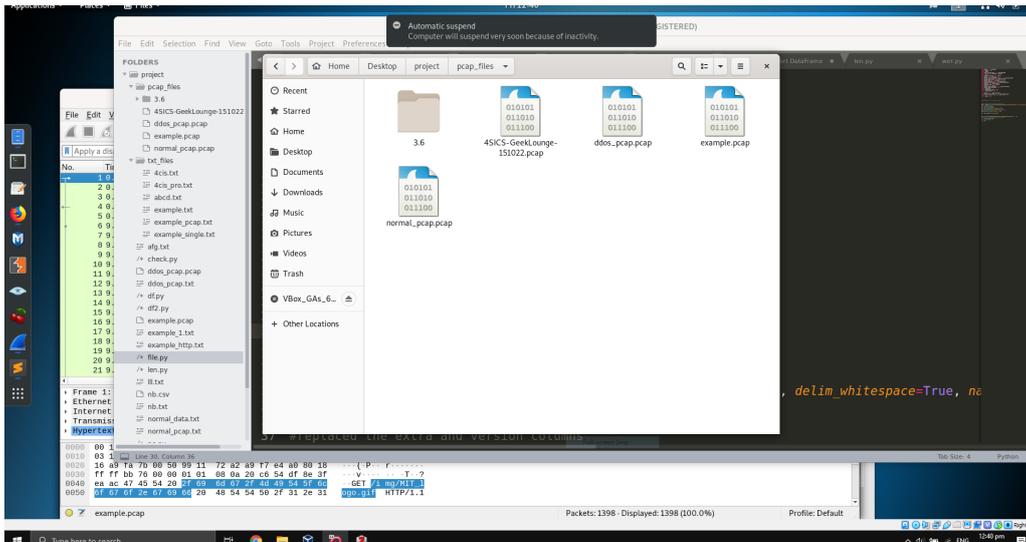
3. file structure for pcap files



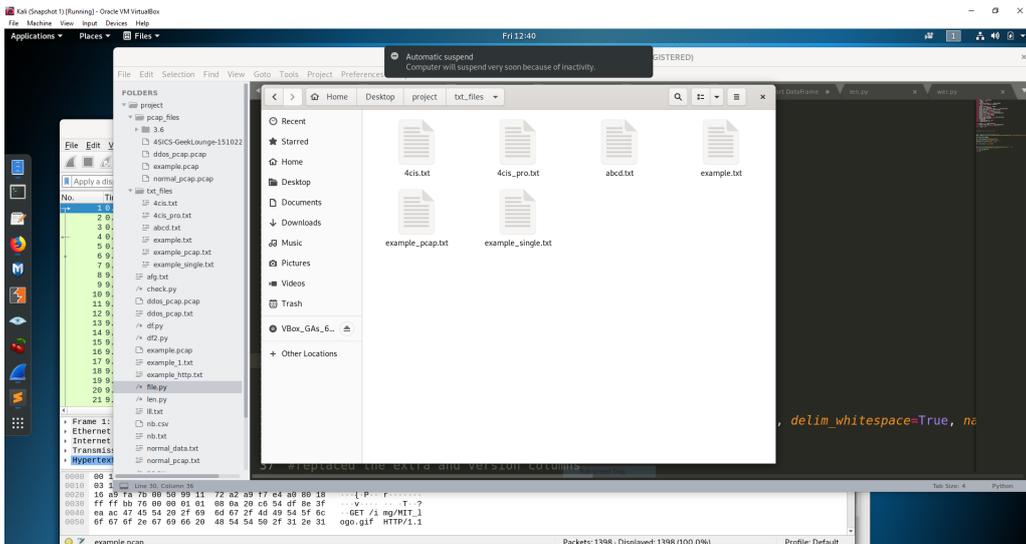Figure 7: Pcap file Repository

4. file structure for text files



Figure 8: Text file Repository

# References

[1] "Welcome to Python.org," library Catalog: www.python.org. [Online]. Available: https://www.python.org/

[2] "pandas - Python Data Analysis Library." [Online]. Available: https://pandas. pydata.org/

[3] "NumPy." [Online]. Available: https://numpy.org/

[4] "scikit-learn: machine learning in Python — scikit-learn 0.23.2 documentation." [Online]. Available: https://scikit-learn.org/stable/

[5] "Matplotlib: Python plotting — Matplotlib 3.3.0 documentation." [Online]. Available: https://matplotlib.org/