

Configuration Manual

MSc Research Project MSc Cyber Security

Jatinder Singh Saini Student ID: x21173656

School of Computing National College of Ireland

Supervisor: Diego Lugones

National College of Ireland



MSc Project Submission Sheet

School of Computing

Student Name:	Jatinder Singh Saini		
Student ID:	x21173656		
Programme:	MSc Cyber Security	Year:	2023/24
Module:	Research Project		
Lecturer:	Diego Lugones		
Date:	31/01/2024		
Project Title:	Enhancing Malware Detection in PE Files U Learning Techniques	sing Hyb	rid Ensemble

Word Count: 957 **Page Count:** 16

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.

<u>ALL</u> 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.

Signature: Jatinder Singh Saini

Date: 30/01/2024

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

Office Use Uniy	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Jatinder Singh Saini X21173656

1 Introduction

Using hybrid ensemble learning algorithms i.e. combining Random Forest, Gradient Boosting, and AdaBoost, into a stacking classifier, this configuration manual describes the suggested work and model used to detect malware in PE files. When it came to improving malware detection systems, the suggested technique worked well. This document contains the specifics of the setup that was used for this study. Hardware and software specs, as well as the environments and languages used, are detailed in this document.

2 System Configurations

2.1 Hardware Setup



Figure 1: System Config

Device	
Name:	Intel(R) UHD Graphics 620
Manufacturer:	Intel Corporation
Chip Type:	Intel(R) UHD Graphics Family
DAC Type:	Internal
Device Type:	Full Display Device
Approx. Total Memory:	8227 MB
Display Memory	128 MB
Shared Memory:	8099 MB

Figure 2: GPU Memory

2.2 Software Setup

The following Software Setup is needed to install and successfully execute the research project:

1. Python - 3.10.9: Python is an interpreted language that is both flexible and fast. It is great for scripting and for integrating components in a way that makes code clear and logical across projects ¹.



Figure 3: Python version used

2. Visual Studio Code: Microsoft VS code is quite a popular text editor and can be used on Windows, Linux, and macOS. It is an open-source and free editor, which was used for developing the application in this project

¹ <u>https://www.python.org/downloads/release/python-3109/</u>



Figure 4: Visual Studio code editor

3. Google Colab: Python coders may easily use GPUs through Google Colab, a free cloud-based tool for interactive and collaborative coding.

← -	C 🔒 colab.research.google.com/#scroll	To=Nma_JWh-W-IF E
cc	Welcome To Colaboratory File Edit View Insert Runtime Tools I	łelp
: = -	Table of contents \Box $ imes$	+ Code + Text 🔥 Copy to Drive
- ८ २ ८	Getting started Data science Machine learning More Resources Featured examples Section	Welcome to Colab! If you're already familiar with Colab, check out this video to learn about interactive tables, the executed code palette.
		What is Colab?

Figure 5: Interface of Google Colab

4. Libraries Configuration:

- Pandas
- Numpy
- Matplotlib
- Seaborn
- Plotly
- Scikit-learn
- Pefile
- Flask
- imblearn

3 Implementation

The section contains detailed instructions to run the project on any Windows system.

1. Browse the (Google Colab²) URL: <u>https://colab.research.google.com</u>, the following interface will be displayed, as shown in "Figure 6".



Figure 6: Google Colaboratory Interface

2. To run the code, the necessary libraries are imported into the project, as shown in "Figure 7".



Figure 7: Libraries

3. The dataset was downloaded, imported into the notebook, and read using the panda's library before the code could be run, as shown in "Figure 8".

² <u>https://colab.research.google.com/</u>



Figure 8: Dataset loading

4. Selecting only necessary columns based on the "Pefile" library helps to enhance the model efficiency and interpretability in malware classification, as shown in "Figure 9".



Figure 9: Taking on necessary columns

5. "Figure 10" shows, the dataset shape and information for data cleaning.

[]	#shape of dataset df.shape		
	(19611, 53)		
0	#information of datatype df.info()		
Ð	Data columns (total 53 columns): # Column 	Non-Null Count 	Dtype int64 int64 int64 int64 int64 int64 int64 int64 int64 int64

Figure 10: Info of dataset

6. Verifying that the dataset does not include any null values. The dataset used in this research does not contain any null values. This is shown in the below "Figure 11".

0	<pre>#checking null values i df.isna().sum()</pre>	n the dataset	
(†)	e_magic e_cblp e_cp e_crlc e_cparhdr e_minalloc e_maxalloc e_ss e_sp e_cs e_ip e_cs e_lfarlc e_oemid e_oemid e_oeminfo e_lfanew Machine	0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	

Figure 11: Checking null values in the dataset

7. Displaying if any unique values are present in each column, as shown in "Figure 12".



Figure 12: Displaying unique values in column

8. "Figure 13" the data set is imbalanced which might not result in better performance of the models. The below figure shows the distribution of samples i.e. (malware and no malware).



Figure 13: Distribution of samples in the dataset

9. "Figure 14" shows the data set is imbalanced, so to balance the dataset SMOTE () function is used. This results in removing unwanted noise in the data set and making it easier for a model to perform well once the training is completed.



Figure 14: Balanced Dataset

10. Splitting dataset for model's training and testing phase, so that class distribution in the target variable is maintained in both sets. In this project, the dataset is split into an 80:20 ratio as shown in the below "Figure 15".



Figure 15: Splitting the dataset

11. The code for model building and evaluation metrics for the classification of samples (malware or benign) are shown below in "Figure 16,17,18,19":

AdaBoost Classifier







Figure 17: Model building and evaluation metrices for Gradient Boosting



Figure 18: Model building and evaluation metrices for Random Forest



Figure 19: Model building and evaluation metrices for Stacking Classifier

Correlation Matrix: A correlation is a measure of the degree to which two variables are directly related to one another. We can predict one variable from another through relationships. Below "Figure 20" shows the correlation matrix.



Figure 20: Correlation Matrix

Feature Importance Bar Chart for Data Visualization Only: This is used in the project to get an idea of which feature, or column is contributing more to enhancement of the model performance. The feature importance is shown in below "Figure 21".



Figure 21: Feature Importance

4 Confusion Matrix, Classification Report and Web Application (Running on Local Machine)

The below results show how the model has performed:





	precision	recall	f1-score	support
0	0.97 0.94	0.93 0.97	0.95 0.95	2920 2920
accuracy	0104	0.57	0.95	5840
macro avg	0.95	0.95	0.95	5840
weighted avg	0.95	0.95	0.95	5840

Figure 23: Classification Report (CLF) for AdaBoost



Figure 24: CM for Gradient Boosting

	nrecision	recall	f1_score	support
	precision	Tecarr	11-30016	Suppor C
0	0.99	0.94	0.96	2920
1	0.94	0.99	0.96	2920
accuracy			0.96	5840
macro avg	0.96	0.96	0.96	5840
weighted avg	0.96	0.96	0.96	5840

Figure 25: CLF for Gradient Boosting



Figure 26: CM for Random Forest (RF)

	precision	recall	f1-score	support	
0	0.95	0.95	0.95	2920	
1	0.95	0.95	0.95	2920	
accuracy			0.95	5840	
macro avg	0.95	0.95	0.95	5840	
weighted avg	0.95	0.95	0.95	5840	

Figure 27: CLF for RF



Figure 28: CM for Stacking Classifier

	precision	recall	f1-score	support	
0	1.00 0.99	0.99 1.00	0.99 0.99	2920 2920	
accuracy		1.00	0.99	5840	
macro avg weighted avg	0.99 0.99	0.99 0.99	0.99 0.99	5840 5840	
		0.00			

Figure 29: CLF for Stacking Classifier

The stacking classifier achieves a 99% accuracy rate by combining Gradient Boosting, Random Forest, and AdaBoost classifiers.

Web Application:

The "Figure 30" shows the code for the application. The models after the implementation were integrated into the web application and run using the VS code editor.

C	EXPLORER ***	🔹 app.py X
0	X 🍨 app.py	85 det miprediction():
~	V WEB APPLICATION PE MALWARE OR BENI	86 file = None
0.0	> pvcache	87 file = request.files['file']
5	> vscode	88 print(Tile)
	> Modele	
) statis	90 if file and allowed file/file.filename):
	> static	92 #prediction
H0	> templates	93 file.save(save path)
ш.	> uploadednies	94 content = pefile.PE(save_path)
π	app.py	95 img = None
A	irequirements.bdt	96 dataframe = createDataframeFromPEdump(content)
		97 prediction = model.predict(dataframe)
		98 output = prediction[0]
		99 print(output)
		100 content.c.tose() 101 if output = 1:
		102 return render template('Classifier.html'.
		103 results 'UPLOADED PE FILE DETECTED AS : ', positive='MALMARE', res2='Risk is HIGH', image = malwareime pat
		105 return render_template('Classifier.html',
		<pre>106 result='UPLOADED PE FILE DETECTED AS : ', positive='BENIGN', res2='Risk is LOW', image = safeimg_path)</pre>
		109 return render_template('Classifier.ntml', result='UPLOADED FILE EXTENSION IS NOT CORRECT YOU NEED TO UPLOAD EXE FILE', image =
		112 if name == ' main ':
		113 # Run the application
		114 app.run(debug=False)

Figure 30: Web Application Code

"Figure 31" shows the interface of the application. For the classification of the model whether it is malware or safe file, click on the "Pe File Classifier" button.



Figure 31: Web application interface

Now, the classification can be done using the PE file sample by providing the sample file to the application to check the risk of the file. If the file that is provided to the application is malicious, then the application provides the following output, as shown in "Figure 32" and if the file is safe, the application displays the output, as shown in "Figure 33".



Figure 32: The file is detected as malware

PE MALWARE CLASSIFICATION SYSTEM	
	PE MALWARE CLASSIFICATION
	Select EXE File & Click on Predict
	Select EXE File
	Predict Uploaded Exe File
	UPLOADED PE FILE DETECTED AS : BENIGN Risk is LOW

Figure 33: The file is detected as benign