

Security Vulnerability Detection with Enhanced Privacy Preservation for Edge Computing Using Hybrid Machine Learning Approach - Configuration Manual

MSc Research Project
Cloud Computing

Shubham Prashant Patil
Student ID: 21139261

School of Computing
National College of Ireland

Supervisor: Sean Heeney

National College of Ireland
Project Submission Sheet
School of Computing



Student Name:	Shubham Prashant Patil
Student ID:	21139261
Programme:	Cloud Computing
Year:	2022
Module:	MSc Research Project
Supervisor:	Sean Heeney
Submission Due Date:	15/12/2022
Project Title:	Security Vulnerability Detection with Enhanced Privacy Preservation for Edge Computing Using Hybrid Machine Learning Approach - Configuration Manual
Word Count:	573
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.

Signature:	
Date:	14th December 2022

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).	<input type="checkbox"/>
Attach a Moodle submission receipt of the online project submission , to each project (including multiple copies).	<input type="checkbox"/>
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.	<input type="checkbox"/>

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):	

Security Vulnerability Detection with Enhanced Privacy Preservation for Edge Computing Using Hybrid Machine Learning Approach - Configuration Manual

Shubham Prashant Patil
21139261

1 Introduction

This configuration manual/handbook outlines the process that can be used to recreate the development of a hybrid machine-learning algorithm with infinite feature selection and PCA. In addition to this, it includes the basic configuration that is required to set up the tools that are essential for the investigation. The overarching goal of this paper is to provide the instructions and background information necessary to successfully run the code that was included as part of the research project submission. Jupyter notebook is used throughout the whole of the project in order to carry out the coding.

2 System configuration requirements

This section refers to a technical document that serves as a guideline for creating software. It describes the program or application in general, such as its purpose, its parameters, how it will communicate with its environment as well as its users, and its necessary hardware and software.

2.1 Hardware Requirements

2.1.1 The experiment was carried out on the following hardware:

- Processor: AMD Ryzen 5 5600U with Radeon Graphics 2.30 GHz
- RAM: 16.0 GB
- System type: 64-bit operating system, x64-based processor

2.1.2 Minimum hardware requirements are:

- Modern Operating System:
 1. Windows 7 or 10
 2. Mac OS X 10.11 or higher 64-bit

- 3. Linux: RHEL 6/7, 64-bit (almost all libraries also work in Ubuntu)
- x86 64-bit CPU (Intel / AMD architecture)
- 4 GB RAM
- 5 GB free disk space

2.2 Software Requirements

- Google Colaboratory: cloud-based jupyter notebook, python version 3.8.
- Email: Gmail account needed for accessing the drive.
- Browser: Any web browser.
- Other Software: Word.

3 Project Recreation

Within this area, instructions are given on how to construct the infrastructure and where to get information.

3.1 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 3.10.5 was used during the whole Model creation process.

- Step 1: Launch the Google Collaboratory by clicking on the link. Link: <https://research.google.com/colaboratory/>
- Step 2: Navigate to the file area, and then launch a new notebook.

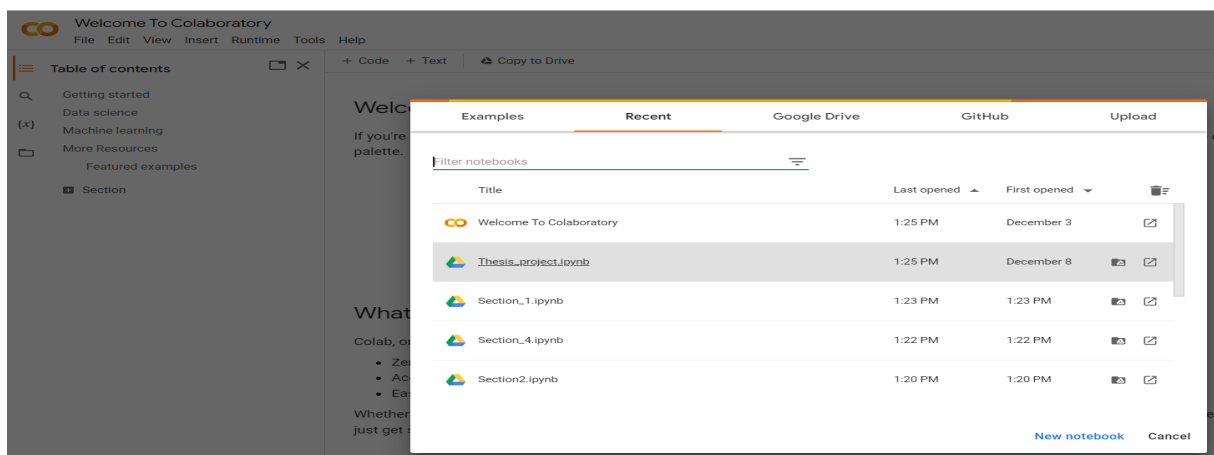


Figure 1: Creation of new notebook

- Step 3: Navigate to the area labeled "upload," then choose the file to upload.

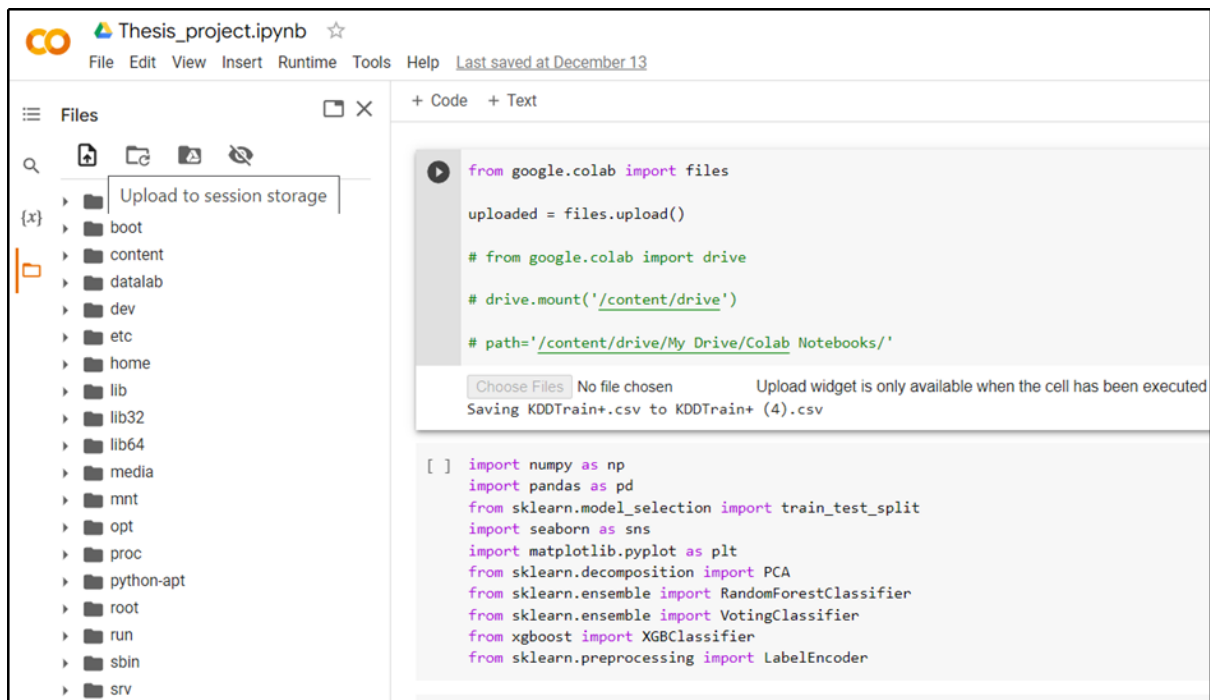


Figure 2: Uploading the files

- Step 4: There are two options to import the dataset.
 1. Importing dataset file via Drive.

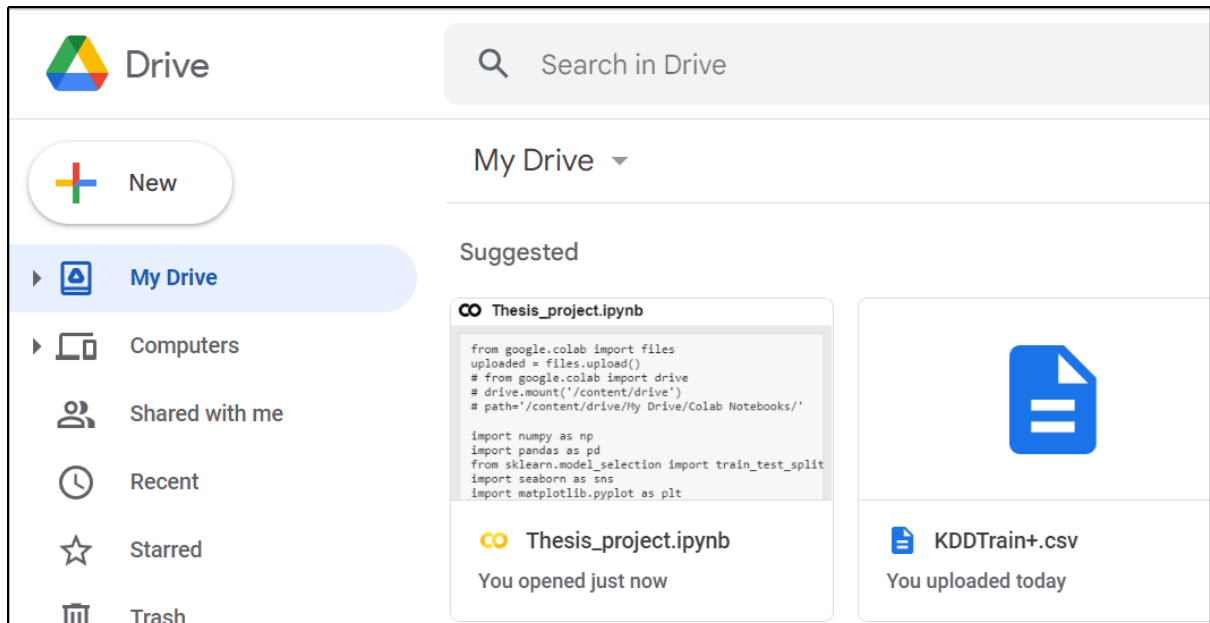


Figure 3: Importing dataset

2. Importing dataset file via Code.

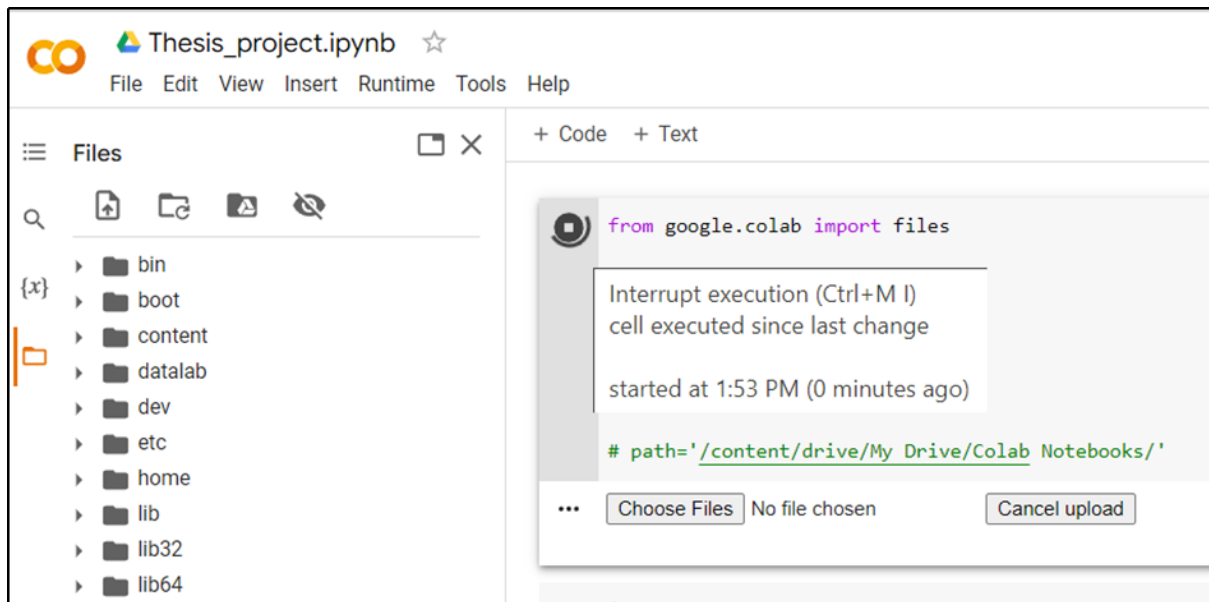


Figure 4: Imported via browse option

3.2 Packages and libraries

After the dataset has been successfully uploaded, check that the following libraries have been imported before continuing with the code implementation.

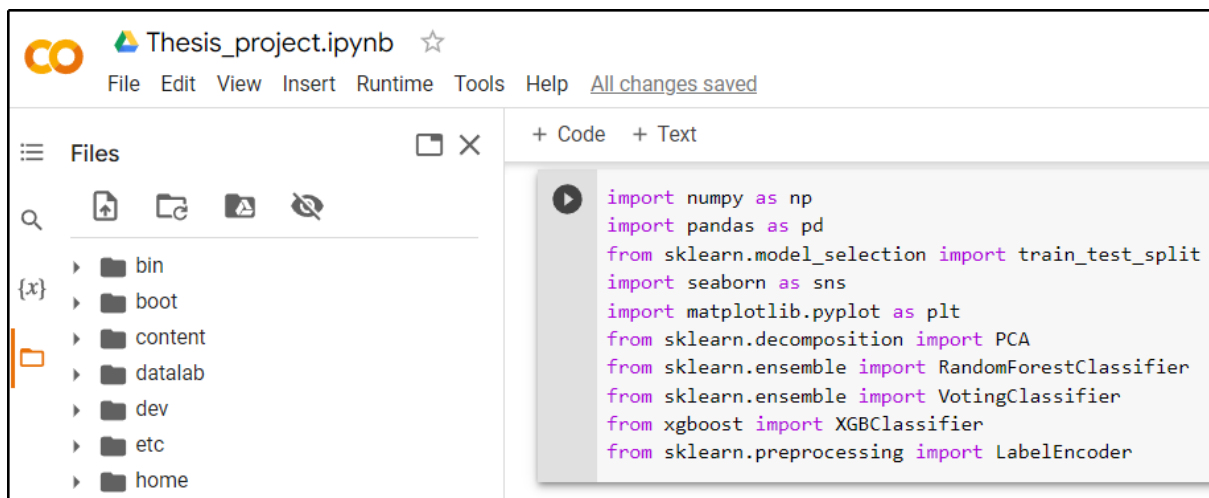


Figure 5: Required Libraries

The following is a list of the needed libraries for execution:

- Matplotlib
- Scipy
- Numpy
- Pandas

- Sci-kit learn
- Seaborn
- Plotly
- Xgboost
- Math
- Stats

4 Phases

The complete methodology of the research process, including data collecting, Exploratory Data Analysis is presented in the following paragraphs.

4.1 Data gathering

- Reading dataset using pandas.

The screenshot shows a Jupyter Notebook interface with the following code and output:

```
df_train=pd.read_csv('KDDTrain+.csv', header=0, na_values='')
df_train = df_train.fillna(df_train.mean())
df_train
```

The output is a DataFrame with 125973 rows and 41 columns. The columns are labeled f2 through f41, plus a 'Label' column. The data is as follows:

	f2	f3	f4	f5	f6	f7	f8	f9	f10	f11	...	f33	f34	f35	f36	f37	f38	f39	f40	f41	Label
0	tcp	ftp_data	SF	491	0	0	0	0	0	0	...	25	0.17	0.03	0.17	0.00	0.00	0.00	0.05	0.00	normal
1	udp	other	SF	146	0	0	0	0	0	0	...	1	0.00	0.60	0.88	0.00	0.00	0.00	0.00	0.00	normal
2	tcp	private	S0	0	0	0	0	0	0	0	...	26	0.10	0.05	0.00	0.00	1.00	1.00	0.00	0.00	dos
3	tcp	http	SF	232	8153	0	0	0	0	0	...	255	1.00	0.00	0.03	0.04	0.03	0.01	0.00	0.01	normal
4	tcp	http	SF	199	420	0	0	0	0	0	...	255	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	normal
...
125968	tcp	private	S0	0	0	0	0	0	0	0	...	25	0.10	0.06	0.00	0.00	1.00	1.00	0.00	0.00	dos
125969	udp	private	SF	105	145	0	0	0	0	0	...	244	0.96	0.01	0.01	0.00	0.00	0.00	0.00	0.00	normal
125970	tcp	smtp	SF	2231	384	0	0	0	0	0	...	30	0.12	0.06	0.00	0.00	0.72	0.00	0.01	0.00	normal
125971	tcp	klogin	S0	0	0	0	0	0	0	0	...	8	0.03	0.05	0.00	0.00	1.00	1.00	0.00	0.00	dos
125972	tcp	ftp_data	SF	151	0	0	0	0	0	0	...	77	0.30	0.03	0.30	0.00	0.00	0.00	0.00	0.00	normal

Figure 6: Reading

- Renaming the column names.



Figure 7: Renaming

- Display information of the dataset.

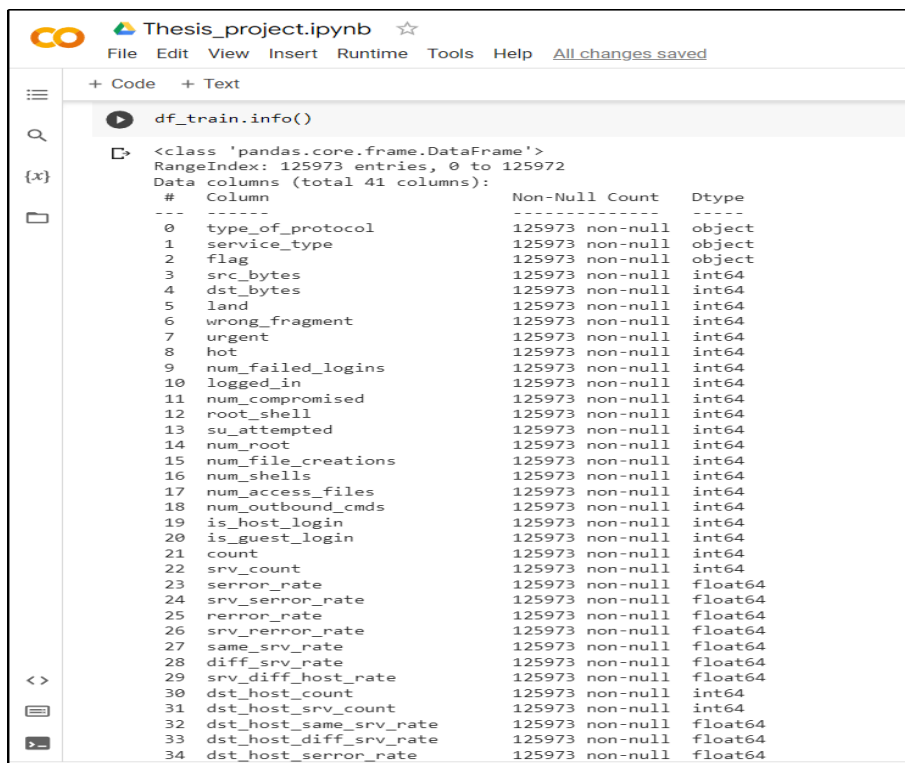


Figure 8: Content of dataset

4.2 Exploratory Data Analysis

- Displaying attack categorization bar graph:

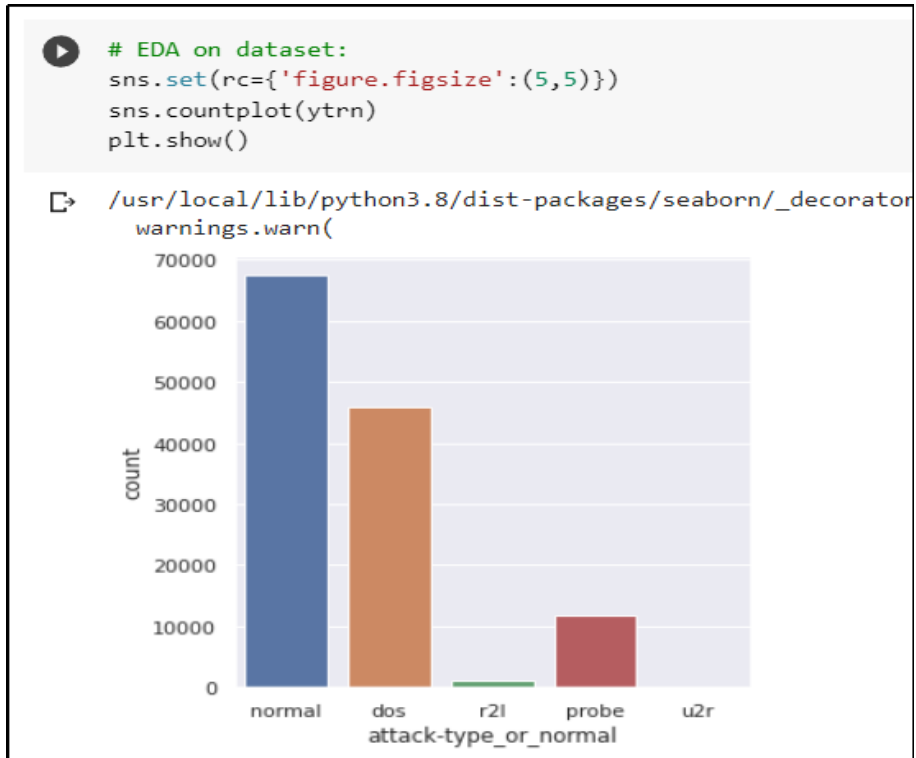


Figure 9: Attack Count

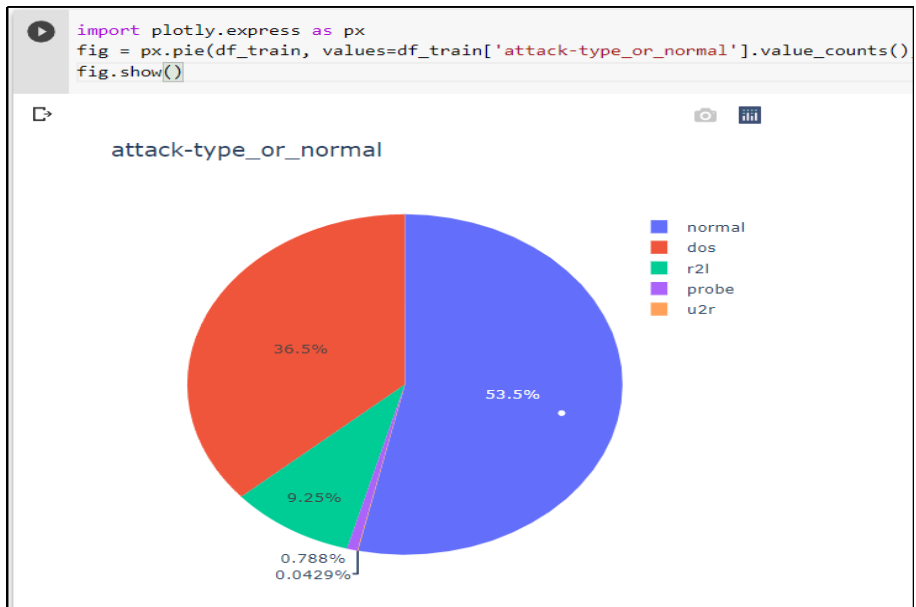


Figure 10: Attack percentage

- Displaying types of a protocol bar graph:



Figure 11: Protocol Type

5 Data Pre-processing

We have used label encoding for the pre-processing stage.

```

# Pre-processing Dataset
lbEn = LabelEncoder()
x['type_of_protocol']=lbEn.fit_transform(x['type_of_protocol'])
x['service_type']=lbEn.fit_transform(x['service_type'])
x['flag']=lbEn.fit_transform(x['flag'])
x['attack-type_or_normal']=lbEn.fit_transform(x['attack-type_or_normal'])

ytrn=lbEn.fit_transform(ytrn)

```

Figure 12: Label Encoding

6 Results

The accuracy score for the test data is 99.33 percent.

References

Radovanovic, I. (2022), 'Google colab - a step-by-step guide - algo trading101 blog'.

URL: <https://algotrading101.com/learn/google-colab-guide/>

The python tutorial (n.d.).

URL: <https://docs.python.org/3/tutorial/>