

# Configuration Manual

MSc Research Project Artificial Intelligence

# Shahna Shahul Hameed Student ID: x22235094

School of Computing National College of Ireland

Supervisor: Sheresh Zahoor

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



Student Name:	Shahna Shahul Hameed
Student ID:	x22235094
Programme:	Artificial Intelligence
Year:	2024
Module:	MSc Research Project
Supervisor:	Sheresh Zahoor
Submission Due Date:	12/12/2024
Project Title:	DeepDefend: Optimized Multi-Model Approach for Network
	Intrusion Detection Using Deep Learning and IoT Security
	Enhancement
Word Count:	835
Page Count:	54

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:	Shahna Shahul Hameed
Date:	11th December 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	
Signature:	
Date:	
Penalty Applied (if applicable):	

## Configuration Manual

Shahna Shahul Hameed x22235094

## 1 Introduction

This configuration manual gives a detailed instructions for replicating the setup, execution, and shows all the experimentation that is involved in the proposed project: "Deep-Defend: Optimized Multi-Model Approach for Network Intrusion Detection Using Deep Learning and IoT security Enhancement". The documentation describes the software tools, system requirements, and configuration steps that were required to demonstrate the experiments, train models and deploy the web application.

## 2 System Configuration

This section gives the details of the system's hardware requirements that were used to support this proposed study.

## 2.1 Hardware Requirements

- Operating System: Windows 11
- Processor: 13th Gen Intel(R) Core(TM) i5-1340P 1.90 GHz
- RAM: 16 GB
- Storage: 512 GB SSD
- System Type: 64-bit operating system, x64-based processor
- GPU: Intel(R) Iris(R) Xe Graphics

#### 2.2 Software Requirements

- Python Version: 3.11.1
- IDE/Platform Used: Jupyter Notebook (via Anaconda)

Python being the primary programming language that was used throughout the work, version 3.11.1 version is utilized. IDE platform that is used in this study is Jupyter Notebook for whole implementation and execution of this study, through a comprehensive package called Anaconda Navigator.



## 2.3 Required Python Libraries

- 1. **Numpy:** This provides support for numerical computations and matrix operations. Numpy is useful to handle array-based structures in the preprocessing steps. In this study, nu,py is sed for efficient data manipulation and computation in the models requiring mathematical operations on arrays.
- 2. **Pandas:** Pandas are required for versatile data manipulation and analysis library. It is used extensively for reading datasets like UNSW-NB15, NSL-KDD and preprocessing steps such as handling missing values, feature extraction and data splitting.
- 3. **Matplotlib:** This generates visualization like heatmaps, feature importance charts, and accuracy plots in the study to give more detailed insights into data and model performance.
- 4. **Seaborn:** Built on Matplotlib, this library gives high-level visualization capabilities for statistical graphics, and in the study it is used to create correlation heatmaps and distribution plots which helps in feature selection and explanatory data analysis.
- 5. **Scikit-learn:** This implements various ML algorithms like Random Forest, Gradient Boosting and evaluation metrics such as confusion matrices and classification reports.
- 6. **TensorFlow/Keras:** It is used for building, training, and evaluation CNN, LSTM and Hybrid CNN-LSTM models for intrusion detection.
- 7. Flask: This hosts the web interface that allows users to upload datasets and interact with the pre-trained intrusion detection models.
- 8. **Joblib:** This saves and loads trained machine learning models for use in the Flask web application.



## 3 Installation

## 3.1 Software Installation

- 1. **Install Anaconda:** Download and install Anaconda Navigator which is recommended to manage environments and to run Jupyter Notebooks.
- 2. Install Required Libraries:

## **Importing Modules and Libraries**

```
[1]: # importing required libraries
     import numpy as np
     import pandas as pd
     import seaborn as sns
     import matplotlib.pyplot as plt
     import pickle
     from os import path
     from sklearn.preprocessing import MinMaxScaler
     from sklearn.preprocessing import StandardScaler
     from sklearn.preprocessing import LabelEncoder
     from sklearn import metrics
     from sklearn import preprocessing
     from sklearn.metrics import accuracy_score
     from sklearn.model_selection import train_test_split
     from sklearn.metrics import classification_report
     from sklearn.svm import SVC
     from sklearn.linear model import LinearRegression
     from sklearn.linear_model import LogisticRegression
     from sklearn.tree import DecisionTreeClassifier
     from sklearn.ensemble import RandomForestClassifier
     from sklearn.neighbors import KNeighborsClassifier
     from sklearn.neural_network import MLPClassifier
```

Figure 1: Importing necessary libraries and modules in UNSW\_Intrusion notebook

## Imports and Data Setup

⊡ ↑ ↓ ≛ Ţ Ĩ

]:	import pandas as pd
	import numpy as np
	import seaborn as sns
	from scipy.stats import shapiro
	from sklearn.preprocessing import OneHotEncoder
	from sklearn.preprocessing import MinMaxScaler
	from keras.models import Sequential
	from keras.layers import Dense
	<pre>from sklearn.model_selection import train_test_split</pre>
	from sklearn.model_selection import KFold
	from sklearn.metrics import accuracy_score
	from sklearn.metrics import precision_score
	from sklearn.metrics import recall_score
	from sklearn.metrics import mean_squared_error
	from sklearn.metrics import classification_report
	<pre>import matplotlib.pyplot as plt</pre>
	import warnings
	warnings.filterwarnings('ignore')

Figure 2: Importing necessary libraries and modules in Cyber\_Intrusion notebook

[1]:	import pandas as pd
	import numpy as np
	<pre>import matplotlib.pyplot as plt</pre>
	import seaborn as sns
	<pre>import plotly.graph_objects as go</pre>
	%matplotlib inline
	import warnings
	warnings.filterwarnings("ignore")
	<pre>from sklearn.model_selection import train_test_split</pre>
	from sklearn.preprocessing import LabelEncoder
	from sklearn.tree import DecisionTreeClassifier
	from sklearn.metrics import classification report
	from sklearn.ensemble import RandomForestClassifier
	from sklearn.svm import SVC
	from sklearn.ensemble import GradientBoostingClassifier
	<pre>from sklearn.metrics import accuracy_score</pre>
	<pre>from sklearn.model_selection import train_test_split</pre>
	from sklearn.preprocessing import StandardScaler
	from sklearn.linear_model import LogisticRegression
	<pre>from sklearn.metrics import accuracy_score, classification_report</pre>
	from sklearn.neighbors import KNeighborsClassifier
	from sklearn.model selection import train test split
	from sklearn.preprocessing import StandardScaler
	from sklearn.naive_bayes import GaussianNB
	from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
	from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis
	from sklearn.metrics import confusion_matrix
	1

Figure 3: Importing necessary libraries and modules in ToN\_IoT notebook



Figure 4: Importing necessary libraries in Train\_Model notebook

## 4 Dataset Preparation

## 4.1 Description of Datasets

1. UNSW-NB15: Network traffic dataset for multi-class intrusion detection.

#### **Importing Datasets**

[2]:	da	<b>t</b> a .		coul!	/Usons	(chaha)	OneDed	vo (Doch	top (COP	E Totou	cion Dotoci	t i or	2024/datasats/	INCU NR15 cov	<b>`</b>		⊡ 个 ↓ ≛	<b>₽</b>
	ua	La	- putreau		/users/	Shanny	OneDri	ve/Desk	ccop/cob	c_incru	sion_bececi	101	1_2024/uacasets/	JNJW_NDIJ.CSV	)			
[3]:	da	ta.	head(n=5)															
[3]:		id	dur	proto	service	state	spkts	dpkts	sbytes	dbytes	rate		ct_dst_sport_ltm	ct_dst_src_ltm	is_ftp_login	ct_ftp_cmd	ct_flw_http_mthd	ct_src_ltn
	0	1	0.121478	tcp	-	FIN	6	4	258	172	74.087490		1	1	0	0	0	
	1	2	0.649902	tcp	-	FIN	14	38	734	42014	78.473372		1	2	0	0	0	
	2	3	1.623129	tcp	-	FIN	8	16	364	13186	14.170161		1	3	0	0	0	
	3	4	1.681642	tcp	ftp	FIN	12	12	628	770	13.677108		1	3	1	1	0	
	4	5	0.449454	tcp	-	FIN	10	6	534	268	33.373826		1	40	0	0	0	
	5 ro	ows	× 45 colun	nns														
	4		_	_	_			_	_	_								

[4]:	data	.info()		
	<cla< th=""><th>ss 'pandas.core.fra</th><th>me.DataFrame'&gt;</th><th></th></cla<>	ss 'pandas.core.fra	me.DataFrame'>	
	Rang	eIndex: 175341 entr	ies, 0 to 175340	
	Data	columns (total 45	columns):	- 1
	#	Column	Non-Null Count	Dtype
			475244 per pull	inter
	1	dun	175341 NON-NULL	float64
	2	proto	175341 non-null	object
	â	service	175341 non-null	object
	4	state	175341 non-null	object
	5	spkts	175341 non-null	int64
	6	dpkts	175341 non-null	int64
	7	sbytes	175341 non-null	int64
	8	dbytes	175341 non-null	int64
	9	rate	175341 non-null	float64
	10	sttl	175341 non-null	int64
	11	dttl	175341 non-null	int64
	12	sload	175341 non-null	float64
	13	dload	175341 non-null	float64
	14	sloss	175341 non-null	int64
	15	dloss	175341 non-null	int64
	16	sinpkt	175341 non-null	float64
	17	dinpkt	175341 non-null	float64
	18	sjit	175341 non-null	float64
	19	djit	175341 non-null	float64
	20	swin	175341 non-null	int64
	21	stcpb	175341 non-null	1nt64
	22	dtcpb	175341 non-null	1nt64
	23	dwin	1/5341 non-null	1Nt64
	24	sypack	175341 non-null	floatc4
	25	syndex	175341 Non-Null	flootc4
	20	dCKudL smean	175341 NON-NULL	int64
	28	dmean	175341 non-null	int64
	29	trans denth	175341 non-null	int64
	30	response body len	175341 non-null	int64
	31	ct srv src	175341 non-null	int64
	32	ct state ttl	175341 non-null	int64
	33	ct dst ltm	175341 non-null	int64
	34	ct src dport 1tm	175341 non-null	int64
	35	ct_dst_sport_ltm	175341 non-null	int64
	36	ct_dst_src_ltm	175341 non-null	int64
	37	is_ftp_login	175341 non-null	int64
	38	ct_ftp_cmd	175341 non-null	int64
	39	ct_flw_http_mthd	175341 non-null	int64
	40	ct_src_ltm	175341 non-null	int64
	41	ct_srv_dst	175341 non-null	int64
	42	is_sm_ips_ports	175341 non-null	int64
	43	attack_cat	175341 non-null	object
	44	label	175341 non-null	int64
	dtyp	es: float64(11), in	t64(30), object(	4)
	memo	ry usage: 60.2+ MB		

2. **NSL-KDD:** This is the improved version of existing KD'99 dataset which is widely used in various intrusion detection works.

[3]:	<pre># open the df = pd.rea</pre>	.txt file here ad_csv('C:/User	e in the o rs/shahn/0	c <i>sv fe</i> OneDr:	<i>ormat</i> ive/Deskto	p/CODE_Int	rusio	n_Detection_2024	/Jupyter	_Cod	≘/NSL_KDD/KDDTrair	1+.csv',hea	der=None)		
	<pre># show all pd.set_opti</pre>	the columns ion('display.ma	ax_column:	5', N	one)										
[4]:	<pre># got from df.columns ,'num_fail@ ,'num_shell ,'srv_serrc ,'dst_host_ ,'dst_host_ print(df.sh df.head()</pre>	KaggLe NB = ['duration', d_logins','log is','num_access or_rate','rerroc same_srv_rate' srv_serror_rate hape)	,'protoco gged_in', s_files', pr_rate', ','dst_ho te','dst_l	l_type 'num_e 'srv_e st_di- host_e	e','servic compromise outbound_c rerror_rat ff_srv_rat rerror_rat	e','flag', d','root_s mds','is_f e','same_s e','dst_ho e','dst_ho	'src_ shell' nost_l srv_ra ost_sa ost_sr	bytes','dst_byte ,'su_attempted', ogin','is_guest_ te','diff_srv_ra me_src_port_rate v_rerror_rate','	s','lanc 'num_roc login',' te','srv ','dst_f outcome'	l','w ot',' coun (_dif nost_: ,'le	rong_fragment','ur num_file_creations t','srv_count','sst f_host_rate','dst srv_diff_host_rate vel']	rgent','hot error_rate' host_count ','dst_hos	','dst_host_srv_cc t_serror_rate'	unt'	
	(125973, 43	3)													
[4]:	duration	protocol_type	service	flag	src_bytes	dst_bytes	land	wrong_fragment	urgent	hot	num_failed_logins	logged_in	num_compromised	root_shell	su_att
	<b>0</b> 0	tcp	ftp_data	SF	491	0	0	0	0	0	0	0	0	0	
	1 0	udp	other	SF	146	0	0	0	0	0	0	0	0	0	
	<b>2</b> 0	tcp	private	S0	0	0	0	0	0	0	0	0	0	0	
	<b>3</b> 0	tcp	http	SF	232	8153	0	0	0	0	0	1	0	0	
	<b>4</b> 0	tcp	http	SF	199	420	0	0	0	0	0	1	0	0	
	4														۱.

#### 3. ToN\_IoT: This is an IoT\_specific dtatset with multiple attack types.

[5]: fridge = pd.read\_csv(r"C:\Users\shahn\OneDrive\Desktop\CODE\_Intrusion\_Detection\_2024\TON\_IOT\IoT\_Fridge.csv")
garage\_door = pd.read\_csv(r"C:\Users\shahn\OneDrive\Desktop\CODE\_Intrusion\_Detection\_2024\TON\_IOT\IoT\_Garage\_Door.csv")
gps = pd.read\_csv(r"C:\Users\shahn\OneDrive\Desktop\CODE\_Intrusion\_Detection\_2024\TON\_IOT\IoT\_GARAGE\_Door.csv")
modbus = pd.read\_csv(r"C:\Users\shahn\OneDrive\Desktop\CODE\_Intrusion\_Detection\_2024\TON\_IOT\IoT\_Modbus.csv")
motion\_light = pd.read\_csv(r"C:\Users\shahn\OneDrive\Desktop\CODE\_Intrusion\_Detection\_2024\TON\_IOT\IoT\_Modbus.csv")
thermostat = pd.read\_csv(r"C:\Users\shahn\OneDrive\Desktop\CODE\_Intrusion\_Detection\_2024\TON\_IOT\IoT\_Modbus.csv")
weather = pd.read\_csv(r"C:\Users\shahn\OneDrive\Desktop\CODE\_Intrusion\_Detection\_2024\TON\_IOT\IoT\_Weather.csv")

#### [4]: fridge.head()

4]:		date	time	fridge_temperature	temp_condition	label	type
	0	31-Mar-19	12:36:52	13.10	high	0	normal
	1	31-Mar-19	12:36:53	8.65	high	0	normal
	2	31-Mar-19	12:36:54	2.00	low	0	normal
	3	31-Mar-19	12:36:55	4.80	low	0	normal
	4	31-Mar-19	12:36:56	10.70	high	0	normal

[7]: fridge.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 587076 entries, 0 to 587075
Data columns (total 6 columns):
 #
     Column
                         Non-Null Count
                                           Dtype
     _____
                          -----
                                            _ _ _ _ _
 0
     date
                          587076 non-null
                                           object
     time
                         587076 non-null
                                           object
 1
                                           float64
 2
     fridge temperature
                         587076 non-null
 3
     temp_condition
                         587076 non-null
                                           object
                          587076 non-null
                                           int64
 4
     label
 5
                          587076 non-null
                                           object
     type
dtypes: float64(1), int64(1), object(4)
memory usage: 26.9+ MB
```

All the datasets to be placed in a "/datasets" directory.

## 5 Implementation

## 5.1 Implementation of Cyber\_Intrusion

- 1. The preprocessing of UNSW-NB15 dataset and generation of heatmaps for feature correlation is performed.
- 2. Load "Cyber\_Intrusion.ipynb" in Jupyter Notebook.
- 3. Execute preprocessing cells sequentially.
- 4. Heatmaps for correlation analysis.
- 5. Preprocessed dataset saved as processed\_unsw.csv

#### **Data Preprocessing**

[10]:	# save attribute and label strings
	all_attribute_names = intru_table.columns.drop('label')
	<pre>intru_label_names = ['Backdoor_Malware', 'BenignTraffic', 'BrowserHijacking', 'CommandInjection', 'DDoS-ACK_Fragmentation', 'DDoS-HTTP_Flood', 'DDoS-ICM# 'DDoS-PSHACK_Flood', 'DDoS-RSTFINFlood', 'DDoS-SYN_Flood', 'DDoS-SlowLoris', 'DDoS-SynonymousIP_Flood', 'DDoS-TCP_Flood', 'DDoS-UDP_ 'DNS_Spoofing', 'DictionaryBruteForce', 'DOS-HTTP_Flood', 'DOS-SYN_Flood', 'DOS-TCP_Flood', 'DOS-UDP_Flood', 'MITM-ArpSpoofing', 'Mi 'Mirai-udpplain', 'Recon-HostDiscovery', 'Recon-OSScan', 'Recon-PingSweep', 'Recon-PortScan', 'SqlInjection', 'Uploading_Attack', '\\</pre>
	4
[11]:	<pre># complete one-hot encoding on intrusion labels ohe = OneHotEncoder() ohe_intru_labels = pd.DataFrame(ohe.fit_transform(intru_table['label'].values.reshape(-1,1)).toarray(), columns=intru_label_names)</pre>
[12]:	# we decided to remove the following attributes since they provide no meaningful contributions to the model
	exclude = ['ece_flag_number', 'cwr_flag_number', 'SMTP', 'Telnet', 'IRC', 'Tot sum', 'Min', 'Max', 'AVG', 'Std', 'Tot size', 'Covariance', 'Variance']
	relevant = ['flow_duration', 'Header_Length', 'Protocol Type', 'Duration', 'Rate', 'Srate', 'Drate', 'syn_count', 'fin_count', 'urg_count', 'rst_count', 'IPv', 'LLC', 'DNS', 'SSH', 'TCP', 'UDP', 'DHCP', 'ARP', 'ICMP', 'IAT', 'Number', 'Mangitue', 'Radius', 'Weight']
	<pre># remove unwanted data attribute columns attribute_data = intru_labels_removed.drop(columns=exclude)</pre>
	# split data for all model runs X_train, X_test, y_train, y_test = train_test_split(attribute_data, ohe_intru_labels, train_size=0.9, random_state=5)

#### Figure 5: Data Preprocessing of Cyber\_Intrusion

for col_name in attribute_data:	
print(coi_name, snapiro(attribute_data[coi_name]), =\n\tmin:", np.min(attribute_data[coi_name]), "; max:", np.ma	x(attribute_data[coi_name]))
<pre>flow_duration ShapiroResult(statistic=0.003912689874606579, pvalue=3.287366402848467e-239)     min: 0.0 : max: 99435.76178</pre>	
Header_Length ShapiroResult(statistic=0.15595012723912527, pvalue=1.657367141034325e-232) min: 0.0 : max: 9815555.0	
Protocol Type ShapiroResult(statistic=0.6352051479259249, pvalue=4.925866299250941e-200) min: 0.0 : max: 47.0	
Duration ShapiroResult(statistic=0.18466471163268572, pvalue=4.0550473173632e-231) min: 0.0 : max: 255.0	
Rate ShapiroResult(statistic=0.05924111192856483, pvalue=6.924870220670144e-237) min: 0.0 : max: 7340032.0	
<pre>Srate ShapiroResult(statistic=0.05924111192858483, pvalue=6.924870220670144e-237)</pre>	
Drate ShapiroResult(statistic=9.673064369997153e-05, pvalue=2.2960230502648243e-239) min: 0.0 : max: 0.848465429	
<pre>fin_flag_number ShapiroResult(statistic=0.3146954186658635, pvalue=3.2703875468487175e-224)     min: 0 : max: 1</pre>	
<pre>syn_flag_number ShapiroResult(statistic=0.49781337233979617, pvalue=4.072253047902139e-212)     min: 0 : max: 1</pre>	
<pre>rst_flag_number ShapiroResult(statistic=0.3231803133049236, pvalue=1.0134578413708053e-223)     min: 0 : max: 1</pre>	
<pre>psh_flag_number ShapiroResult(statistic=0.31798828376526755, pvalue=5.064864595949719e-224)     min: 0 : max: 1</pre>	
<pre>ack_flag_number ShapiroResult(statistic=0.3855051669308578, pvalue=6.25496855367773e-220)     min: 0 : max: 1</pre>	
ack_count ShapiroResult(statistic=0.3320400435897358, pvalue=3.3476254818623606e-223) min: 0.0 : max: 4.0	
<pre>syn_count ShapiroResult(statistic=0.5663950769094286, pvalue=1.602892806221714e-206) min: 0.0 : max: 9.69</pre>	
<pre>fin_count ShapiroResult(statistic=0.32056961653993765, pvalue=7.146160098971559e-224)     min: 0.0 : max: 77.6</pre>	
urg_count ShapiroResult(statistic=0.05756365903671634, pvalue=5.863420209241335e-237)	
rst_count ShapiroResult(statistic=0.09728359740844761, pvalue=3.248066022733782e-235)	
HTTP ShapiroResult(statistic=0.2195085625971428, pvalue=2.260465979889006e-229)	
HTTPS ShapiroResult(statistic=0.23880164156491357, pvalue=2.2485706436126917e-228)	
min. V ; mun. 1 DNS (hereinenens)	

[14]: # attributes are NOT normally distributed (all p-values < 0.05) - normalize traffic attribute data in train and test sets scaler = MinMaxScaler()

X\_train = pd.DataFrame(scaler.fit\_transform(X\_train), columns=attribute\_data.columns) X\_test = pd.DataFrame(scaler.fit\_transform(X\_test), columns=attribute\_data.columns)

#### **Data Visualization**





## Fitting the Initial Model and Measuring Its Performance

```
[17]: # setup neural net for predicting intrusion class
classifier = Sequential()
classifier.add(Dense(units = 36, activation = 'relu', input_dim = 33))
classifier.add(Dense(units = 35, activation = 'relu'))
classifier.add(Dense(units = 34, activation = 'softmax'))
classifier.compile(optimizer = "SGD", loss = 'CategoricalCrossentropy')
classifier.fit(X_train.astype(float), y_train, epochs = 50)
```



29492/29492 [====================================			101s 3ms/	/step
Model Training Accuracy (av	g.): 0.78	59835098869	682	
Model Training Precision (a	vg.): 0.6	10637616825	7833	
Model Training Recall (avg.	): 0.4603	96515253722	66	
Model Training MSE: 0.011884	404624541	9163		
pre	ecision	recall f	1-score	support
Backdoor_Malware	0.00	0.00	0.00	68
BenignTraffic	0.83	0.73	0.78	22033
BrowserHijacking	0.00	0.00	0.00	126
CommandInjection	0.00	0.00	0.00	88
DDoS-ACK_Fragmentation	0.98	0.98	0.98	5796
DDoS-HTTP_Flood	0.70	0.51	0.59	579
DDoS-ICMP_Flood	1.00	1.00	1.00	145241
DDoS-ICMP_Fragmentation	1.00	0.97	0.99	9235
DDoS-PSHACK_Flood	1.00	1.00	1.00	83181
DDoS-RSTFINFlood	1.00	1.00	1.00	81774
DDoS-SYN_Flood	0.67	0.94	0.78	82578
DDoS-SlowLoris	0.93	0.03	0.06	448
DDoS-SynonymousIP_Flood	0.93	0.53	0.67	72679
DDoS-TCP_Flood	0.64	0.97	0.77	90906
DDoS-UDP_Flood	0.71	0.95	0.82	109148
DDoS-UDP_Fragmentation	0.99	0.98	0.98	5782
DNS_Spoofing	0.48	0.21	0.30	3611
DictionaryBruteForce	0.00	0.00	0.00	294
DoS-HTTP_Flood	0.85	0.61	0.71	1482
DoS-SYN_Flood	0.55	0.44	0.49	40693
DoS-TCP_Flood	0.59	0.08	0.14	53793
DoS-UDP_Flood	0.83	0.38	0.52	67239
MITM-ArpSpoofing	0.82	0.51	0.63	6292
Mirai-greeth_flood	0.76	0.04	0.08	19923
Mirai-greip_flood	0.44	0.97	0.60	15262
Mirai-udpplain	1.00	0.99	0.99	18125
Recon-HostDiscovery	0.71	0.42	0.53	2682
Recon-OSScan	1.00	0.00	0.00	2010
Recon-PingSweep	0.00	0.00	0.00	36
Recon-PortScan	0.53	0.04	0.07	1673
SqlInjection	0.00	0.00	0.00	114
Uploading_Attack	0.00	0.00	0.00	20
VulnerabilityScan	0.81	0.39	0.52	746
XSS	0.00	0.00	0.00	60
micro avg	0.81	0.79	0.80	943717
macro avg	0.61	0.46	0.47	943717
weighted avg	0.82	0.79	0.76	943717
samples avg	0.79	0.79	0.79	943717

3277/3277 [===============			10s 3ms/ste	ep
Model Testing Accuracy (av	g.): 0.717	789772835	6444	
Model Testing Precision (a	vg.): 0.52	1458995034	43226	
Model Testing Recall (avg.	): 0.40873	783322410	05	
Model Testing MSE: 0.01591	0031830161	35		
p	recision	recall	f1-score	support
Backdoor_Malware	0.00	0.00	0.00	8
BenignTraffic	0.84	0.60	0.70	2443
BrowserHijacking	0.00	0.00	0.00	14
CommandInjection	0.00	0.00	0.00	17
DDoS-ACK_Fragmentation	0.94	0.99	0.96	635
DDoS-HTTP_Flood	0.68	0.40	0.51	47
DDoS-ICMP_Flood	1.00	1.00	1.00	16040
DDoS-ICMP_Fragmentation	1.00	0.97	0.98	988
DDoS-PSHACK_Flood	1.00	1.00	1.00	9214
DDoS-RSTFINFlood	1.00	1.00	1.00	9049
DDoS-SYN_Flood	0.89	0.01	0.01	9066
DDoS-SlowLoris	0.00	0.00	0.00	45
DDoS-SynonymousIP_Flood	0.38	0.94	0.54	8001
DDoS-TCP_Flood	0.64	0.99	0.78	10387
DDoS-UDP_Flood	0.71	0.96	0.82	12057
DDoS-UDP_Fragmentation	0.55	0.98	0.71	649
DNS_Spoofing	0.45	0.18	0.26	423
DictionaryBruteForce	0.00	0.00	0.00	30
DoS-HTTP_Flood	0.75	0.74	0.74	198
DoS-SYN_Flood	0.36	0.11	0.17	4514
DoS-TCP_Flood	0.71	0.02	0.04	6014
DoS-UDP_Flood	0.86	0.38	0.52	7548
MITM-ArpSpoofing	0.78	0.47	0.59	727
Mirai-greeth_flood	0.57	0.99	0.72	2192
Mirai-greip_flood	0.50	0.00	0.00	1690
Mirai-udpplain	0.99	0.74	0.85	2041
Recon-HostDiscovery	0.76	0.25	0.38	325
Recon-OSScan	0.29	0.01	0.02	215
Recon-PingSweep	0.00	0.00	0.00	5
Recon-PortScan	0.48	0.07	0.12	190
SqlInjection	0.00	0.00	0.00	8
Uploading_Attack	0.00	0.00	0.00	3
VulnerabilityScan	0.60	0.10	0.16	63
XSS	0.00	0.00	0.00	12
micro avg	0.74	0.72	0.73	104858
macro avg	0.52	0.41	0.40	104858
weighted avg	0.79	0.72	0.66	104858
samples avg	0.72	0.72	0.72	104858

```
[19]: # K-fold cross validation to generalize initial model performance (10 folds)
      classifier = Sequential()
      classifier.add(Dense(units = 36, activation = 'relu', input_dim = 33))
      classifier.add(Dense(units = 35, activation = 'relu'))
      classifier.add(Dense(units = 34, activation = 'softmax'))
      classifier.compile(optimizer = "SGD", loss = 'CategoricalCrossentropy')
      kfold_model = KFold(n_splits=10)
      # calculate accuracy and MSE scores for each fold iteration
      scaler = MinMaxScaler()
      X = scaler.fit_transform(attribute_data)
      y = ohe_intru_labels.values
      all_accuracy = []
      all_mse = []
      count = 0
      for train_index, test_index in kfold_model.split(X):
          count += 1
          classifier.fit(X[train_index].astype(float), y[train_index])
          yhat_pred_test = pd.DataFrame(((classifier.predict(X[test_index].astype(float))) > 0.5).astype(int))
          iter_accuracy = accuracy_score(y[test_index], yhat_pred_test)
          print("\tIteration", count, "Accuracy:", iter_accuracy)
          all_accuracy.append(iter_accuracy)
          iter_mse = mean_squared_error(y[test_index], yhat_pred_test)
          print("\tIteration", count, "MSE:", iter_mse)
          all_mse.append(iter_mse)
      # print results
      print("Overall Avg. K-Fold Accuracy:", np.mean(all_accuracy))
      print("Overall Avg. K-Fold MSE:", np.mean(all_mse))
```

```
29492/29492 [=======] - 107s 4ms/step - loss: 0.7231
3277/3277 [========] - 11s 3ms/step
        Iteration 1 Accuracy: 0.6926033302180091
        Iteration 1 MSE: 0.016283085360257514
```

```
29492/29492 [======] - 103s 3ms/step - loss: 0.5284
3277/3277 [=======] - 11s 3ms/step
       Iteration 2 Accuracy: 0.767609529077419
3277/3277 [=-----] - 11s 3ms/step
Iteration 3 Accuracy: 0.7688874477865304
3277/3277 [======] - 10s 3ms/step
        Iteration 4 Accuracy: 0.7648343474031547
       Iteration 4 MSE: 0.012716076531510963
29492/29492 [======] - 104s 4ms/step - loss: 0.4776
3277/3277 [===========] - 11s 3ms/step
       Iteration 5 Accuracy: 0.6997463235995346
       Iteration 5 MSE: 0.013544928547626877
29492/29492 [======] - 1025 3ms/step - loss: 0.4732
3277/3277 [========] - 115 3ms/step
       Iteration 6 Accuracy: 0.7747694479147792
Iteration 6 MSE: 0.012402324959089944
29492/29492 [------] - 103s 3ms/step - loss: 0.4700
3277/3277 [------] - 11s 3ms/step
3277/3277 [-----]
       Iteration 7 Accuracy: 0.7107584615237895
Iteration 7 MSE: 0.013231184879799887
29492/29492 [=------] - 1025 3ms/step - loss: 0.4673
3277/3277 [=-----] - 108 3ms/step
Iteration 8 Accuracy: 0.7731481922999895
        Iteration 8 MSE: 0.012627842176095288
Iteration & MSE: 0.02002/0042700052000
29492/29492 [=====================] - 1025 3ms/step - loss: 0.4655
3277/3277 [==============] - 115 3ms/step
       Iteration 9 Accuracy: 0.7824084228997588
       Iteration 9 MSE: 0.012148758336984432
Iteration 10 Accuracy: 0.7897708307504506
Iteration 10 MSE: 0.011722407379461888
Overall Avg. K-Fold Accuracy: 0.7524536333473415
```

```
Overall Avg. K-Fold MSE: 0.013002152841590894
```





14

Model Complexity

:

0.02 -

0.00

#### 5.2 Implementation of PPGO\_IoT

- 1. The evaluation of IoT-specific datasets using various ML classifiers are done here.
- 2. Load datasets such as GPS Tracker, Fridge, and Thermostat.
- 3. Train models such as Random Forest and Decision Trees.
- 4. Classification reports for each dataset.
- 5. Accuracy comparison for classifiers.



#### Figure 7: Optimization using PPGO algorithm

```
# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
# Encode target labels using one-hot encoding
label_encoder = LabelEncoder()
y_train_encoded = to_categorical(label_encoder.fit_transform(y_train))
y_test_encoded = to_categorical(label_encoder.transform(y_test))
```

Figure 8: Splitting the dataset into training and testing

```
# Define LSTM model
model = Sequential()
model.add(LSTM(units=50, return_sequences=True, input_shape=(X_train_lstm.shape[1], X_train_lstm.shape[2])))
model.add(LSTM(units=50))
model.add(Dense(units=num_classes, activation='softmax'))
model.compile(optimizer='adam', loss='categorical crossentropy', metrics=['accuracy'])
# Get the shapes of model weights
shapes = [w.shape for w in model.get_weights()]
# Run PPGO optimization
num_parameters = sum(np.prod(shape) for shape in shapes) # Number of parameters in the model
initial_parameters = ppgo_initialize(num_parameters)
def objective_function(parameters):
    # Reshape parameters to match the shapes of model weights
    weights = []
    start = 0
    for shape in shapes:
       weight size = np.prod(shape)
       weight = parameters[start:start+weight_size].reshape(shape)
       weights.append(weight)
       start += weight_size
    model.set_weights(weights)
    loss, _ = model.evaluate(X_test_lstm, y_test_encoded, verbose=0)
    return loss
perpetual_pigeon_galvanized_optimization(objective_function, initial_parameters, num_iterations=100)
# Evaluate the model with the optimized weights
loss, accuracy = model.evaluate(X_test_lstm, y_test_encoded)
print("Test Loss:", loss)
print("Test Accuracy:", accuracy)
```

```
Figure 9: Deep Learning with LSTM
```

## 5.3 Implementation of ToN\_IoT

- 1. Feature Selection and ML model evaluation for ToN\_IoT dataset is done.
- 2. Used correlation matrices for selection of key features.
- 3. Train models like Logistic Regression and SVM.
- 4. Feature importance visualizations.
- 5. Confusion matrices for model predictions.

[5]: fridge = pd.read\_csv(r"C:\Users\shahn\OneDrive\Desktop\CODE\_Intrusion\_Detection\_2024\TON\_IOT\IoT\_Fridge.csv")
garage\_door = pd.read\_csv(r"C:\Users\shahn\OneDrive\Desktop\CODE\_Intrusion\_Detection\_2024\TON\_IOT\IoT\_Garage\_Door.csv")
gps = pd.read\_csv(r"C:\Users\shahn\OneDrive\Desktop\CODE\_Intrusion\_Detection\_2024\TON\_IOT\IoT\_Motbus.csv")
modbus = pd.read\_csv(r"C:\Users\shahn\OneDrive\Desktop\CODE\_Intrusion\_Detection\_2024\TON\_IOT\IoT\_Motbus.csv")
motion\_light = pd.read\_csv(r"C:\Users\shahn\OneDrive\Desktop\CODE\_Intrusion\_Detection\_2024\TON\_IOT\IoT\_Motbus.csv")
thermostat = pd.read\_csv(r"C:\Users\shahn\OneDrive\Desktop\CODE\_Intrusion\_Detection\_2024\TON\_IOT\IoT\_Motion\_Light.csv")
weather = pd.read\_csv(r"C:\Users\shahn\OneDrive\Desktop\CODE\_Intrusion\_Detection\_2024\TON\_IOT\IoT\_Weather.csv")

[4]:	fr	idge.head(	)				
[4]:		date	time	fridge_temperature	temp_condition	label	type
	0	31-Mar-19	12:36:52	13.10	high	0	normal
	1	31-Mar-19	12:36:53	8.65	high	0	normal
	2	31-Mar-19	12:36:54	2.00	low	0	normal
	3	31-Mar-19	12:36:55	4.80	low	0	normal
	4	31-Mar-19	12:36:56	10.70	high	0	normal

Figure 10: Preprocessing of ToN\_IoT notebook

```
fridge.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 587076 entries, 0 to 587075
Data columns (total 6 columns):
 #
     Column
                           Non-Null Count
                                            Dtype
     _ _ _ _ _ _
                                             _ _ _ _ _
                           _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _
_ _ _
                           587076 non-null object
 0
     date
 1
     time
                           587076 non-null object
    fridge_temperature 587076 non-null float64
 2
     temp condition
 3
                           587076 non-null object
     label
                           587076 non-null int64
 4
 5
                           587076 non-null object
     type
dtypes: float64(1), int64(1), object(4)
memory usage: 26.9+ MB
```

<pre>type_counts = fridge['type'].value_counts() fig = go.Figure(data=[go.Pie(labels=type_counts.in fig.update_traces(textinfo='percent+label', pull=[ fig.update_layout(title='Fridge Type Distribution' fig.show()</pre>	dex, values=type_counts)]) 0.1, 0.1, 0.1, 0.1]) )	[□ ↑ ↓ 告 ⊋ ■
Fridge Type Distribution	backdoor 6.06% 4.84% ddos 1.74% injection 1.2.1% ransomware 0.494% xss 0.348%	normal backdoor password ddos injection ransomware

[11]: fridge['temp\_condition'] = fridge['temp\_condition'].apply(lambda x:x.strip())
fridge['time'] = fridge['time'].apply(lambda x: x.strip())

```
[12]: # number of seconds that have passed since midnight
def time_to_seconds(time_str):
    h, m, s = map(int, time_str.split(':'))
    return h * 3600 + m * 60 + s
fridge['time_seconds'] = fridge['time'].apply(time_to_seconds)
# target variable + features
X = fridge.drop(['date', 'time', 'label', 'type'], axis=1)
y = fridge['type']
```

[13]: # encoding categorical variables label\_encoder = LabelEncoder() X['temp\_condition'] = label\_encoder.fit\_transform(X['temp\_condition']) # training + testing X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.2, random\_state=42) X train

[13]:		fridge_temperature	temp_condition	time_seconds
	189639	7.65	0	15495
	304660	1.00	1	35883
	408687	4.95	1	74706
	460624	11.25	0	45642
	197265	5.05	1	15525
	110268	6.95	0	15108
	259178	5.85	1	76677
365838 131932	365838	8.55	0	10669
	131932	7.10	0	15229
	121958	10.55	0	15175

469660 rows × 3 columns

[20]: scaler = StandardScaler() X\_train\_scaled = scaler.fit\_transform(X\_train) X\_test\_scaled = scaler.transform(X\_test) k = 5 # Choose an appropriate value for K knn\_classifier = KNeighborsClassifier(n\_neighbors=k, metric='euclidean') knn\_classifier.fit(X\_train\_scaled, y\_train) y\_pred = knn\_classifier.predict(X\_test\_scaled)

```
[14]: dt_classifier = DecisionTreeClassifier(random_state=42)
      dt_classifier.fit(X_train, y_train)
      dt_predictions = dt_classifier.predict(X_test)
      dt_accuracy = accuracy_score(y_test, dt_predictions)
      print("Decision Tree Accuracy:", dt_accuracy)
      Decision Tree Accuracy: 0.7756694147305307
      # Get classification report
[15]:
      classification_rep = classification_report(y_test, dt_predictions)
      # Print classification report
      print("Decision Tree Classification Report:")
      print(classification_rep)
      Decision Tree Classification Report:
                    precision
                               recall f1-score
                                                    support
          backdoor
                         0.12
                                   0.12
                                             0.12
                                                       7192
              ddos
                         0.08
                                   0.07
                                             0.08
                                                       2090
         injection
                         0.22
                                   0.24
                                             0.23
                                                       1415
                         0.88
                                   0.89
                                             0.88
            normal
                                                     100027
                         0.19
                                   0.17
                                             0.18
          password
                                                       5707
         ransomware
                         0.12
                                   0.13
                                             0.13
                                                        577
                         0.14
                                   0.14
                                             0.14
                                                        408
               XSS
          accuracy
                                             0.78
                                                     117416
         macro avg
                         0.25
                                   0.25
                                             0.25
                                                     117416
      weighted avg
                         0.77
                                   0.78
                                             0.77
                                                     117416
```

Figure 11: Decision Tree of fridge dataset

```
[16]: rf_classifier = RandomForestClassifier(random_state=42)
rf_classifier.fit(X_train, y_train)
rf_predictions = rf_classifier.predict(X_test)
rf_accuracy = accuracy_score(y_test, rf_predictions)
print("Random Forest Accuracy:", rf_accuracy)
```

```
Random Forest Accuracy: 0.7855317844246099
```

```
[17]: classification_rep = classification_report(y_test, rf_predictions)
```

```
# Print classification report
print("Random Forest Classifier Report:")
print(classification_rep)
```

Random	Forest	Classifier	Report:		
		precision	recall	f1-score	support
bac	kdoor	0.13	0.11	0.12	7192
	ddos	0.07	0.06	0.06	2090
inje	ction	0.23	0.25	0.24	1415
n	ormal	0.88	0.90	0.89	100027
pas	sword	0.18	0.15	0.16	5707
ranso	mware	0.13	0.11	0.12	577
	XSS	0.15	0.13	0.14	408
acc	uracy			0.79	117416
macr	o avg	0.25	0.24	0.25	117416
weighte	d avg	0.77	0.79	0.78	117416

Figure 12: Random Forest of fridge dataset

```
[18]: gb_classifier = GradientBoostingClassifier(random_state=42)
      gb_classifier.fit(X_train, y_train)
      gb_predictions = gb_classifier.predict(X_test)
      gb_accuracy = accuracy_score(y_test, gb_predictions)
      print("Gradient Boosting Accuracy:", gb_accuracy)
      Gradient Boosting Accuracy: 0.8518941200517817
      classification_rep = classification_report(y_test, gb_predictions)
[19]:
      # Print classification report
      print("Gradient Boosting Classifier Report:")
      print(classification rep)
      Gradient Boosting Classifier Report:
                     precision
                                  recall f1-score
                                                     support
          backdoor
                          0.00
                                    0.00
                                              0.00
                                                        7192
               ddos
                          1.00
                                    0.00
                                              0.00
                                                        2090
         injection
                          0.00
                                    0.00
                                              0.00
                                                        1415
            normal
                         0.85
                                    1.00
                                              0.92
                                                      100027
          password
                         1.00
                                    0.00
                                              0.00
                                                        5707
                          0.00
        ransomware
                                    0.00
                                              0.00
                                                         577
                          0.00
                                    0.00
                                              0.00
                                                         408
               xss
                                              0.85
                                                      117416
          accuracy
                          0.41
                                    0.14
                                              0.13
                                                      117416
         macro avg
```

Figure 13: Gradient Boosting of fridge dataset

0.85

0.78

117416

0.79

weighted avg

```
[24]: from sklearn.naive_bayes import GaussianNB
      nb_classifier = GaussianNB()
      nb_classifier.fit(X_train, y_train)
      nb_predictions = nb_classifier.predict(X_test)
      nb_accuracy = accuracy_score(y_test, nb_predictions)
      print("Naive Bayes Accuracy:", nb_accuracy)
      Naive Bayes Accuracy: 0.8519026367786332
      classification_rep = classification_report(y_test, nb_predictions)
[25]:
      # Print classification report
      print("Naive Bayes Classifier Report:")
      print(classification_rep)
      Naive Bayes Classifier Report:
                     precision
                                  recall f1-score
                                                     support
          backdoor
                          0.00
                                    0.00
                                              0.00
                                                        7192
               ddos
                          0.00
                                    0.00
                                              0.00
                                                        2090
         injection
                          0.00
                                    0.00
                                              0.00
                                                        1415
             normal
                          0.85
                                    1.00
                                              0.92
                                                      100027
                          0.00
                                    0.00
                                              0.00
                                                        5707
          password
        ransomware
                          0.00
                                    0.00
                                              0.00
                                                         577
                          0.00
                                    0.00
                                              0.00
                                                         408
               XSS
          accuracy
                                              0.85
                                                      117416
         macro avg
                          0.12
                                    0.14
                                              0.13
                                                      117416
```

Figure 14: Naive Bayes of fridge dataset

0.85

0.73

0.78

117416

weighted avg

[26]: from sklearn.discriminant\_analysis import LinearDiscriminantAnalysis lda\_classifier = LinearDiscriminantAnalysis() lda\_classifier.fit(X\_train, y\_train) lda\_predictions = lda\_classifier.predict(X\_test) lda\_accuracy = accuracy\_score(y\_test, lda\_predictions) print("Linear Discriminant Analysis Accuracy:", lda\_accuracy)

```
Linear Discriminant Analysis Accuracy: 0.8519026367786332
```

[27]: classification\_rep = classification\_report(y\_test, lda\_predictions)

```
# Print classification report
print("Linear Discriminant Analysis Report:")
print(classification_rep)
```

Linear Discri	minant Analy	ysis Repor	٠t:	
	precision	recall	f1-score	support
backdoon	0 00	0 00	0 00	7100
Dackdoor	0.00	0.00	0.00	7192
ddos	0.00	0.00	0.00	2090
injection	0.00	0.00	0.00	1415
normal	0.85	1.00	0.92	100027
password	0.00	0.00	0.00	5707
ransomware	0.00	0.00	0.00	577
XSS	0.00	0.00	0.00	408
accuracy			0.85	117416
macro avg	0.12	0.14	0.13	117416
weighted avg	0.73	0.85	0.78	117416

Figure 15: LDA of fridge dataset

- [28]: **from** sklearn.discriminant\_analysis **import** QuadraticDiscriminantAnalysis qda\_classifier = QuadraticDiscriminantAnalysis() qda\_classifier.fit(X\_train, y\_train) qda\_predictions = qda\_classifier.predict(X\_test) qda\_accuracy = accuracy\_score(y\_test, qda\_predictions) print("Quadratic Discriminant Analysis Accuracy:", qda\_accuracy) Quadratic Discriminant Analysis Accuracy: 0.8519026367786332

[29]: classification\_rep = classification\_report(y\_test, qda\_predictions)

```
# Print classification report
print("Quadratic Discriminant Analysis Report:")
print(classification_rep)
```

Quadratic Discriminant Analysis Report:					
	precision	recall	f1-score	support	
backdoor	0.00	0.00	0.00	7192	
ddos	0.00	0.00	0.00	2090	
injection	0.00	0.00	0.00	1415	
normal	0.85	1.00	0.92	100027	
password	0.00	0.00	0.00	5707	
ransomware	0.00	0.00	0.00	577	
XSS	0.00	0.00	0.00	408	
			0.95	117410	
accuracy			0.85	11/416	
macro avg	0.12	0.14	0.13	117416	
weighted avg	0.73	0.85	0.78	117416	

	Figure	16:	QDA	of fridge	dataset
--	--------	-----	-----	-----------	---------

[30]: garage\_door.head()

[30]:		date	time	door_state	sphone_signal	label	type
	0	1-Apr-19	20:53:44	open	true	0	normal
	1	1-Apr-19	20:53:49	closed	false	0	normal
	2	1-Apr-19	20:53:49	open	true	0	normal
	3	1-Apr-19	20:53:54	closed	false	0	normal
	4	1-Apr-19	20:53:54	open	true	0	normal

```
[31]:
```

garage\_door.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 591446 entries, 0 to 591445
Data columns (total 6 columns):
 #
    Column
                   Non-Null Count
                                    Dtype
_ _ _
    _ _ _ _ _ _
                    -----
                                     _ _ _ _ _
    date
                   570111 non-null object
0
    time
                   570111 non-null object
1
 2
    door_state
                 541177 non-null object
 3
    sphone_signal 541177 non-null object
4
    label
                   591446 non-null int64
5
                   591446 non-null object
    type
dtypes: int64(1), object(5)
memory usage: 27.1+ MB
```

Figure 17: Preprocessing of garage\_door dataset





Figure 18: Distribution of garage\_door dataset

```
[37]: dt_classifier = DecisionTreeClassifier(random_state=42)
      dt_classifier.fit(X_train, y_train)
      dt_predictions = dt_classifier.predict(X_test)
      dt_accuracy = accuracy_score(y_test, dt_predictions)
      print("Decision Tree Accuracy:", dt_accuracy)
      Decision Tree Accuracy: 0.8651365288697269
      # Get classification report
[38]:
      classification_rep = classification_report(y_test, dt_predictions)
      # Print classification report
      print("Decision Tree Classification Report:")
      print(classification rep)
      Decision Tree Classification Report:
                    precision recall f1-score
                                                     support
          backdoor
                         0.36
                                    0.37
                                              0.36
                                                        7228
              ddos
                         0.31
                                    0.36
                                              0.34
                                                        2032
         injection
                         0.19
                                   0.21
                                              0.20
                                                        1199
                         0.94
                                   0.94
                                              0.94
            normal
                                                      103123
                         0.42
                                   0.38
                                              0.40
                                                        3792
          password
                         0.15
                                   0.12
                                              0.14
                                                         577
        ransomware
                         0.00
                                   0.00
                                              0.00
                                                         114
          scanning
               xss
                         0.44
                                   0.48
                                              0.46
                                                         225
                                              0.87
                                                     118290
          accuracy
         macro avg
                         0.35
                                   0.36
                                              0.35
                                                      118290
      weighted avg
                         0.87
                                    0.87
                                              0.87
                                                      118290
```

Figure 19: Decision Tree of garage\_door dataset

```
rf_classifier = RandomForestClassifier(random_state=42)
[39]:
      rf_classifier.fit(X_train, y_train)
      rf_predictions = rf_classifier.predict(X_test)
      rf_accuracy = accuracy_score(y_test, rf_predictions)
      print("Random Forest Accuracy:", rf_accuracy)
      Random Forest Accuracy: 0.8735818750528362
      # Get classification report
[40]:
      classification_rep = classification_report(y_test, rf_predictions)
      # Print classification report
      print("Random Forest Classification Report:")
      print(classification_rep)
      Random Forest Classification Report:
                               recall f1-score
                    precision
                                                   support
          backdoor
                        0.44
                                  0.34
                                            0.38
                                                      7228
              ddos
                       0.31
                                  0.31
                                            0.31
                                                     2032
         injection
                        0.19
                                  0.19
                                            0.19
                                                     1199
                        0.94
            normal
                                 0.95
                                            0.94
                                                   103123
                        0.44
                                0.51
                                           0.47
                                                     3792
          password
                       0.19
                                 0.16
                                           0.18
                                                      577
        ransomware
          scanning
                        0.00
                                  0.00
                                            0.00
                                                      114
                        Q 11
                                  0 56
                                            0 10
                                                       225
```

XSS	0.44	0.56	0.49	225
accuracy			0.87	118290
macro avg	0.37	0.38	0.37	118290
weighted avg	0.87	0.87	0.87	118290

Figure 20: Random Forest of garage\_door dataset

[41]: gb\_classifier = GradientBoostingClassifier(random\_state=42) gb\_classifier.fit(X\_train, y\_train) gb\_predictions = gb\_classifier.predict(X\_test) gb\_accuracy = accuracy\_score(y\_test, gb\_predictions) print("Gradient Boosting Accuracy:", gb\_accuracy)

```
Gradient Boosting Accuracy: 0.8926198326147603
```

```
[42]: # Get classification report
classification_rep = classification_report(y_test, gb_predictions)
```

```
# Print classification report
print("Random Forest Classification Report:")
print(classification_rep)
```

```
Random Forest Classification Report:
```

	precision	recall	f1-score	support
backdoor	0.90	0.28	0.43	7228
ddos	0.33	0.00	0.00	2032
injection	0.41	0.07	0.11	1199
normal	0.93	0.98	0.95	103123
password	0.37	0.74	0.49	3792
ransomware	0.00	0.00	0.00	577
scanning	0.00	0.00	0.00	114
XSS	0.50	0.01	0.03	225
accuracy			0.89	118290
macro avg	0.43	0.26	0.25	118290
weighted avg	0.89	0.89	0.87	118290

Figure 21: Gradient Boosting of garage\_door dataset

```
[44]: scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
k = 5 # Choose an appropriate value for K
knn_classifier = KNeighborsClassifier(n_neighbors=k, metric='euclidean')
knn_classifier.fit(X_train_scaled, y_train)
y_pred = knn_classifier.predict(X_test_scaled)
knn_accuracy = accuracy_score(y_test, y_pred)
print("KNNs Accuracy:", knn_accuracy)
```

KNNs Accuracy: 0.8611294276777411

```
[45]: classification_rep = classification_report(y_test, y_pred)
print("KNNs Classifier Report:")
print(classification_rep)
```

KNNs Classifi	er Report:			
	precision	recall	f1-score	support
backdoor	0.42	0.38	0.40	7228
ddos	0.02	0.01	0.02	2032
injection	0.21	0.24	0.22	1199
normal	0.93	0.95	0.94	103123
password	0.29	0.25	0.27	3792
ransomware	0.20	0.12	0.15	577
scanning	0.00	0.00	0.00	114
XSS	0.26	0.21	0.23	225
accuracy			0.86	118290
macro avg	0.29	0.27	0.28	118290
weighted avg	0.85	0.86	0.86	118290

[47]:	<pre>y_pred = logistic_model.predict(X_test_scaled) ll_accuracy = accuracy_score(y_test, y_pred) print("Logistic Regression Accuracy:", ll_accuracy)</pre>											
	Logistic Regre	ession Accura	cy: 0.872	73649505452	27							
[48]:	<pre>report = class print("Classic</pre>	sification_re fication Repo	port(y_te rt:\n", r	st, y_pred) eport)	)							
	Classification	n Report:										
		precision	recall	f1-score	support							
	backdoor	0.65	0.23	0.34	7228							
	ddos	0.00	0.00	0.00	2032							
	injection	0.00	0.00	0.00	1199							
	normal	0.89	0.98	0.93	103123							
	password	0.17	0.05	0.08	3792							
	ransomware	0.00	0.00	0.00	577							
	scanning	0.00	0.00	0.00	114							
	XSS	0.00	0.00	0.00	225							
	accuracy			0.87	118290							
	macro avg	0.21	0.16	0.17	118290							
	weighted avg	0.82	0.87	0.84	118290							

Figure 22: Logistic Regression of garage\_door dataset

```
[49]: nb_classifier = GaussianNB()
      nb_classifier.fit(X_train, y_train)
      nb predictions = nb classifier.predict(X test)
      nb_accuracy = accuracy_score(y_test, nb_predictions)
      print("Naive Bayes Accuracy:", nb_accuracy)
      Naive Bayes Accuracy: 0.8526164510947671
      classification_rep = classification_report(y_test, nb_predictions)
[50]:
      # Print classification report
      print("Naive Bayes Classifier Report:")
      print(classification_rep)
      Naive Bayes Classifier Report:
                    precision recall f1-score
                                                    support
          backdoor
                         0.22
                                   0.10
                                             0.14
                                                       7228
              ddos
                                             0.00
                         0.00
                                   0.00
                                                       2032
         injection
                         0.27
                                   0.25
                                             0.26
                                                       1199
                                             0.94
            normal
                         0.92
                                   0.95
                                                     103123
                                             0.35
          password
                         0.26
                                   0.57
                                                       3792
                                   0.00
                                             0.00
                                                        577
                         0.00
        ransomware
                         0.00
                                   0.00
                                             0.00
          scanning
                                                        114
                         0.00
                                   0.00
                                             0.00
                                                        225
               xss
                                             0.85
                                                     118290
          accuracy
         macro avg
                                   0.23
                                             0.21
                                                     118290
                         0.21
      weighted avg
                         0.83
                                   0.85
                                             0.84
                                                     118290
```

Figure 23: Naive Bayes of garage\_door dataset

```
[37]: dt_classifier = DecisionTreeClassifier(random_state=42)
      dt_classifier.fit(X_train, y_train)
      dt_predictions = dt_classifier.predict(X_test)
      dt_accuracy = accuracy_score(y_test, dt_predictions)
      print("Decision Tree Accuracy:", dt_accuracy)
      Decision Tree Accuracy: 0.8651365288697269
      # Get classification report
[38]:
      classification_rep = classification_report(y_test, dt_predictions)
      # Print classification report
      print("Decision Tree Classification Report:")
      print(classification_rep)
      Decision Tree Classification Report:
                    precision recall f1-score
                                                     support
          backdoor
                         0.36
                                   0.37
                                             0.36
                                                        7228
              ddos
                         0.31
                                   0.36
                                             0.34
                                                        2032
         injection
                         0.19
                                   0.21
                                             0.20
                                                       1199
                         0.94
                                   0.94
                                             0.94
            normal
                                                     103123
                         0.42
                                   0.38
                                             0.40
                                                       3792
          password
                         0.15
                                   0.12
                                             0.14
                                                        577
        ransomware
                         0.00
                                   0.00
                                             0.00
                                                         114
          scanning
               xss
                         0.44
                                   0.48
                                             0.46
                                                         225
                                             0.87
                                                     118290
          accuracy
         macro avg
                         0.35
                                   0.36
                                             0.35
                                                     118290
      weighted avg
                         0.87
                                   0.87
                                             0.87
                                                     118290
```

Figure 24: LDA of garage\_door dataset

```
dt_classifier = DecisionTreeClassifier(random_state=42)
[37]:
      dt_classifier.fit(X_train, y_train)
      dt_predictions = dt_classifier.predict(X_test)
      dt_accuracy = accuracy_score(y_test, dt_predictions)
      print("Decision Tree Accuracy:", dt_accuracy)
      Decision Tree Accuracy: 0.8651365288697269
[38]:
      # Get classification report
      classification_rep = classification_report(y_test, dt_predictions)
      # Print classification report
      print("Decision Tree Classification Report:")
      print(classification rep)
      Decision Tree Classification Report:
                                  recall f1-score
                     precision
                                                      support
           backdoor
                          0.36
                                    0.37
                                               0.36
                                                         7228
               ddos
                          0.31
                                    0.36
                                               0.34
                                                         2032
          injection
                          0.19
                                    0.21
                                               0.20
                                                         1199
             normal
                          0.94
                                    0.94
                                               0.94
                                                       103123
                          0.42
                                    0.38
                                                         3792
           password
                                               0.40
         ransomware
                          0.15
                                    0.12
                                               0.14
                                                          577
                          0.00
                                    0.00
                                               0.00
           scanning
                                                          114
                xss
                          0.44
                                    0.48
                                               0.46
                                                          225
                                               0.87
                                                       118290
           accuracy
         macro avg
                          0.35
                                    0.36
                                               0.35
                                                       118290
      weighted avg
                          0.87
                                    0.87
                                               0.87
                                                       118290
```

Figure 25: QDA of garage\_door dataset

#### 5.4 Implementation of UNSW-NB15

- 1. This shows the comprehensive evaluation using UNSW-NB15 dataset.
- 2. Preprocess the dataset by handling missing values.
- 3. Train ML models and evaluate their performance.
- 4. Accuracy, precision, recall, and F1-score for each model.
- 5. Comparison charts for UNSW-specific models.

## **Binary Classification**

```
[26]: plt.figure(figsize=(8,8))
    plt.pie(data.label.value_counts(),labels=['normal','abnormal'],autopct='%0.2f%%')
    plt.title("Pie chart distribution of normal and abnormal labels",fontsize=16)
    plt.legend()
    plt.savefig('C:/Users/shahn/OneDrive/Desktop/CODE_Intrusion_Detection_2024/plots_UNSW/Pie_chart_binary.png')
    plt.show()
```





Figure 26: Visualization of Binary classification

#### **Multi-class Classification**

```
[27]: plt.figure(figsize=(8,8))
plt.pie(data.attack_cat.value_counts(),labels=data.attack_cat.unique(),autopct='%0.2f%%')
plt.title('Pie chart distribution of multi-class labels')
plt.legend(loc='best')
plt.savefig('C:/Users/shahn/OneDrive/Desktop/CODE_Intrusion_Detection_2024/plots_UNSW/Pie_chart_multi.png')
plt.show()
```

Pie chart distribution of multi-class labels



Figure 27: Visualization of Multi-class classification









Figure 29: Correlation Matrix for Multi-class labels

	А	В	С	D	E	F	G	Н	
1		dttl	swin	dwin	tcprtt	synack	ackdat	label	
2	3	0.99213	1	1	0	0	0	6	
3	11	0.99213	1	1	0.08797	0.103	0.05498	6	
4	15	0	0	0	0	0	0	6	
5	17	0.99213	1	1	0.08071	0.11537	0.04291	6	
5	21	0.99213	1	1	0.1036	0.1074	0.06975	6	
7	22	0.99213	1	1	0.07542	0.10249	0.04202	6	
3	28	0.99213	1	1	0.05695	0.08548	0.02881	6	
9	30	0.99213	1	1	0.07084	0.07229	0.04811	6	
0	31	0.99213	1	1	0.06751	0.09819	0.03529	6	
1	32	0.99213	1	1	0.06119	0.05654	0.04368	6	
2	34	0.99213	1	1	0.07875	0.05629	0.06215	6	
3	35	0.99213	1	1	0.0586	0.09039	0.02877	6	
4	39	0.99213	1	1	0.06008	0.06972	0.03777	6	
5	41	0.99213	1	1	0.07036	0.06592	0.0499	6	
6	43	0.99213	1	1	0.05152	0.04037	0.03938	6	
7	44	0.99213	1	1	0.07816	0.10302	0.0447	6	
8	45	0.99213	1	1	0.05965	0.05856	0.04134	6	
9	47	0.99213	1	1	0.0491	0.03991	0.03702	6	
0	48	0.99213	1	1	0.05433	0.04198	0.04174	6	
1	56	0.99213	1	1	0.04607	0.03297	0.03635	6	
2	60	0.99213	1	1	0.04081	0.02268	0.03455	6	
3	63	0.99213	1	1	0.05121	0.02529	0.04449	6	
4	64	0.99213	1	1	0.04022	0.02041	0.03475	6	
5	83	0 00012	1	1	0 06388	0 05190	0.0/817	6	

Figure 30: Features saved as "multi\_data.csv"

Α	В	С	D	E	F	G	Н	l I	J	K	L	Μ	N	0
	rate	sttl	sload	dload	ct_srv_src	ct_state_t	ct_dst_ltm	ct_src_dpd	ct_dst_spo	ct_dst_src	ct_src_ltm	ct_srv_dst	label	
3	1.3677107	0.1383928	1.1893137	0.00015	0	0.1666666	0.02	0	0	0.039216	0.016949	0	1	
11	4.2520966	0.1383928	9.1937988	0.000364	0.019608	0.1666666	0	0	0	0.019608	0	0	1	
15	0.4999999	0.9955357	0.1197916	0	0	0.3333333	0	0	0	0.0588235	0.016949	0	1	
17	4.3195885	0.1383928	6.8287612	0.000869	0.039216	0.1666666	0	0	0	0.019608	0	0.039216	1	
21	4.4376467	0.9955357	9.2348251	0.000236	0.098039	0.1666666	0	0	0	0	0.016949	0.039216	1	
22	2.5915973	0.1383928	2.3072416	0.000322	0	0.5	0	0	0	0.019608	0.016949	0	1	
28	3.1111314	0.9955357	6.6848661	0.000594	0.098039	0.1666666	0	0	0	0.019608	0	0.1372549	1	
30	3.4435450	0.1383928	3.5391991	0.00563	0.098039	0.1666666	0	0	0	0.019608	0.016949	0	1	
31	5.4623912	0.1383928	1.3984733	0.000291	0.078431	0.1666666	0	0	0	0.0588235	0	0.1372549	1	
32	5.0376939	0.1383928	7.4083735	0.002091	0.039216	0.1666666	0	0	0	0	0	0.1372549	1	
34	5.0563613	0.9955357	1.0417040	0.000269	0.039216	0.1666666	0	0	0	0	0	0.039216	1	
35	3.8610096	0.1383928	8.8535096	0.001499	0.098039	0.1666666	0	0	0	0.039216	0	0.019608	1	
39	5.1932942	0.9955357	8.1626086	0.000277	0.098039	0.1666666	0	0	0	0.039216	0	0	1	
41	3.1912768	0.1383928	4.4206648	0.000784	0.098039	0.1666666	0	0	0	0.019608	0	0.1372549	1	
43	4.9882333	0.1383928	7.3763910	0.002486	0.078431	0.1666666	0	0	0	0.0588235	0	0.1372549	1	
44	3.1976033	0.1383928	4.1664678	0.001668	0.078431	0.1666666	0	0	0	0.0588235	0	0.1372549	1	
45	8.4473813	0.1383928	9.3352839	0.0260275	0.098039	0.1666666	0	0	0	0.7647058	0	0.1372549	1	
47	5.7355892	0.1383928	1.0182749	0.004155	0.098039	0.1666666	0	0	0	0.019608	0	0	1	
48	4.4047123	0.1383928	4.7018093	0.000587	0.019608	0.1666666	0.02	0	0	0.019608	0	0.019608	1	
56	5.5814011	0.1383928	8.6411395	0.002538	0.078431	0.1666666	0	0	0	0.0588235	0	0.039216	1	
60	6.5728021	0.1383928	4.1512434	0.001434	0.098039	0.1666666	0	0	0	0	0	0.019608	1	
63	3.7938214	0.9955357	6.5709978	0.000247	0.098039	0.1666666	0	0	0	0.039216	0	0	1	
64	6.5099927	0.1383928	9.2082138	0.005248	0.078431	0.1666666	0	0	0	0.0588235	0	0.1372549	1	
68	6.7309893	0.1383928	5.3019896	0.0150552	0.098039	0.1666666	0	0	0	0	0	0.039216	1	
91	5.7614107	0.1383928	8.4743557	0.00061	0	0.1666666	0.02	0	0	0.039216	0	0.019608	1	
100	9.1749499	0	1.6072597	0.000101	0.039216	0	0.16	0.02	0	0.019608	0.050847	0.078431	1	
	F	-	· · · · · · · · · ·			-			-	-				

Figure 31: Features saved as "bin\_data.csv"

#### 5.5 Train\_Model

- 1. Training CNN and hybrid CNN-LSTM models for intrusion detection.
- 2. Import libraries and load preprocessed datasets.
- 3. Train CNN and CNN-LSTM models.
- 4. Model Evaluation metrics such as accuracy and F1-score
- 5. Plots of training vs. validation accuracy.

[14]:

```
# find the correlation coefficient between the numerical variables with class label and make it absolute
corr = abs(df_numerical.corr()['class_label']).sort_values(ascending=False)
# print(corr)
# pick only correlation values greater than 0.5
corr = corr[corr > 0.5]
print(corr)
class_label 1.000000
same_srv_rate 0.751913
dst_host_srv_count 0.722535
dst_host_same_srv_rate 0.693803
logged in 0.690171
```

doc_nooc_balle_briv_race	0.000000
logged_in	0.690171
dst_host_srv_serror_rate	0.654985
dst_host_serror_rate	0.651842
serror_rate	0.650652
<pre>srv_serror_rate</pre>	0.648289
count	0.576444
Name: class_label, dtype:	float64

[15]:

# now save the above correlation values in a list and check for direct and inverse relation between the variables wrt class label corr\_list = corr\_list)

# # now plot the correlation matrix # plt.figure(figsize(15,15)) sns.heatmap(df\_numerical[corr\_list].corr(), annot=True, cmap='viridis') plt.show()

['class\_label', 'same\_srv\_rate', 'dst\_host\_srv\_count', 'dst\_host\_same\_srv\_rate', 'logged\_in', 'dst\_host\_srv\_serror\_rate', 'dst\_host\_serror\_rate', 'serror\_rate', 'srv\_serror\_rate', 'count']

class_label -	- 1	-0.75	-0.72	-0.69	-0.69	0.65	0.65	0.65	0.65	0.58			- 1.0
same_srv_rate -	-0.75	1	0.71	0.79	0.6	-0.77	-0.76	-0.76	-0.76	-0.63			- 0.8
dst_host_srv_count -	-0.72	0.71	1	0.9	0.62	-0.57	-0.57	-0.56	-0.56	-0.4		-	- 0.6
dst_host_same_srv_rate -	-0.69	0.79	0.9	1	0.6	-0.63	-0.64	-0.62	-0.62	-0.47		-	- 0.4
logged_in -	-0.69	0.6	0.62	0.6	1	-0.49	-0.49	-0.49	-0.49	-0.54		-	- 0.2
dst_host_srv_serror_rate -	0.65	-0.77	-0.57	-0.63	-0.49	1	0.99	0.98	0.99	0.46			- 0.0
dst_host_serror_rate -	0.65	-0.76	-0.57	-0.64	-0.49	0.99	1	0.98	0.98	0.46			0.2
serror_rate -	0.65	-0.76	-0.56	-0.62	-0.49	0.98	0.98	1	0.99	0.46			0.4
srv_serror_rate -	0.65	-0.76	-0.56	-0.62	-0.49	0.99	0.98	0.99	1	0.45			-0.4
count -	0.58	-0.63	-0.4	-0.47	-0.54	0.46	0.46	0.46	0.45	1			0.6
	class_label -	same_srv_rate -	dst_host_srv_count -	dst_host_same_srv_rate -	logged_in -	dst_host_srv_serror_rate -	dst_host_serror_rate -	serror_rate -	srv_serror_rate -	count -	-	-	

[10].	: # do one hot encoding for the categorical variables and check its correlation with class label df_categorical = pd.get_dummies(df_categorical, drop_first=True) df_categorical.head()															
[16]:	protocol	_type_tcp	protocol_typ	pe_udp	service_X11	service_Z39_	50	service_aol serv	ce_auth	service_bgp	service_couri	er service	_csnet_ns s	ervice_ctf	service	_daytim
	0	True		False	False	Fa	lse	False	False	False	Fal	se	False	False		Fals
	1	False		True	False	Fa	lse	False	False	False	Fal	se	False	False		Fals
	2	True		False	False	Fa	lse	False	False	False	Fal	se	False	False		Fals
	3	True		False	False	Fa	lse	False	False	False	Fal	se	False	False		Fals
	4	True		False	False	Fa	lse	False	False	False	Fal	se	False	False		Fals
	•															۱.
[17]:	<pre># now conce df_new = pe df_new.head</pre>	atenate ti d.concat( d()	he numerica [df_numeric	L and al, df	categorical v _categorical	variables  , axis=1)	h-4	num feiled leei								<i>6</i> 1.
[1/].	duration	src_byte:	ast_bytes	land	wrong_fragm	ent urgent	not	num_tailed_logi	ns logge	a_in_num_c	ompromised	root_snell	su_attempt	ea num_r		Im_nie_(
	• •	49	0	0		0 0	0		0	0	0	0		0	0	
	1 0	140		0		0 0	0		0	0	0	0		0	0	
	2 0	(	0 0	0		0 0	0		0	0	0	0		0	0	
	-		0152	<u>_</u>			- 0			1	0	0			0	
	<b>3</b> 0	232	8153	0		0 0	•		0					0	0	
	<b>3</b> 0 <b>4</b> 0	232	8153 420	0		0 0	0		0	1	0	0		0	0	

[18]: # now find the correlation between the new dataframe and class label corr = abs(df\_new.corr()['class\_label']).sort\_values(ascending=False) corr = corr[corr > 0.5] # corr\_list = corr.index.tolist() # print(corr\_list) corr

[18]:	attack_class_normal	1.000000
	attack_type_normal	1.000000
	class_label	1.000000
	flag_SF	0.756286
	same_srv_rate	0.751913
	attack_type_neptune	0.747336
	dst_host_srv_count	0.722535
	dst_host_same_srv_rate	0.693803
	logged_in	0.690171
	dst_host_srv_serror_rate	0.654985
	dst_host_serror_rate	0.651842
	serror_rate	0.650652
	flag_S0	0.650206
	srv_serror_rate	0.648289
	count	0.576444
	service_http	0.562312
	Name: class_label, dtype:	float64

[19]:

# now plot the correlation matrix
# plt.figure(figsize=(15,15))
sns.heatmap(df\_new[corr\_list].corr(), annot=True, cmap='viridis') plt.show()

class_label -	1	-0.75	-0.72	-0.69	-0.69	0.65	0.65	0.65	0.65	0.58	- 1.0
same_srv_rate -	-0.75	1	0.71	0.79	0.6	-0.77	-0.76	-0.76	-0.76	-0.63	- 0.8
dst_host_srv_count -	-0.72	0.71	1	0.9	0.62	-0.57	-0.57	-0.56	-0.56	-0.4	- 0.6
dst_host_same_srv_rate -	-0.69	0.79	0.9	1	0.6	-0.63	-0.64	-0.62	-0.62	-0.47	- 0.4
logged_in -	-0.69	0.6	0.62	0.6	1	-0.49	-0.49	-0.49	-0.49	-0.54	- 0.2
dst_host_srv_serror_rate -	0.65	-0.77	-0.57	-0.63	-0.49	1	0.99	0.98	0.99	0.46	- 0.0
dst_host_serror_rate -	0.65	-0.76	-0.57	-0.64	-0.49	0.99	1	0.98	0.98	0.46	0.2
serror_rate -	0.65	-0.76	-0.56	-0.62	-0.49	0.98	0.98	1	0.99	0.46	0.2
srv_serror_rate -	0.65	-0.76	-0.56	-0.62	-0.49	0.99	0.98	0.99	1	0.45	0.4
count -	0.58	-0.63	-0.4	-0.47	-0.54	0.46	0.46	0.46	0.45	1	0.6
	class_label -	same_srv_rate -	dst_host_srv_count -	dst_host_same_srv_rate -	logged_in -	dst_host_srv_serror_rate -	dst_host_serror_rate -	serror_rate -	srv_serror_rate -	count -	-

[20]:	df_nume df_cate	rica gori	l = df.select cal = df.selec	_dtypes(includ ct_dtypes(excl	de=np.u Lude=n	number) p.number)									
	<pre># label from sk le = La df_e_ca # df_e_</pre>	lear belE tego	oding for the n.preprocessin ncoder() rical = df_cat gorical.head(,	categorical w ng <b>import</b> Labe tegorical.app] )	variab elEnco Ly(le.	les in df data der fit_transform)									
	<pre># norma from sk scaler df_n_nu</pre>	lize lear = Min merio	the numerican n.preprocessin nMaxScaler() cal = pd.Data	l variables ng <b>import</b> MinM Frame(scaler.1	MaxScal Fit_tr	ler ansform(df_numer	ical), d	colum	ns=df_numerical.co	olumns)					
	<pre># conca df_ne = # df_ne # drop df_ne.d df_ne.h</pre>	atta atta nop(	e encoded and concat([df_n_i d() ck_class and d ['attack_class )	normalized da numerical, df_ attack_type co s', 'attack_ty	ptaset: _e_cato 	s egorical], axis= axis=1, inplace	1) =True)								
[20]:	dura	tion	src_bytes	dst_bytes	land	wrong_fragment	urgent	hot	num_failed_logins	logged_in	num_compromised	root_shell	su_attempted	num_root	n
	0	0.0	3.558064e-07	0.000000e+00	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
	1	0.0	1.057999e-07	0.000000e+00	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
	2	0.0	0.000000e+00	0.000000e+00	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
	3	0.0	1.681203e-07	6.223962e-06	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	
	4	0.0	1.442067e-07	3.206260e-07	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	
	4	_													

```
[22]: # perfrom random forest feature selection
      from sklearn.ensemble import RandomForestClassifier
      from sklearn.feature_selection import SelectFromModel
      from sklearn.model_selection import train_test_split
      X = df_ne.drop(['class_label'], axis=1)
      y = df_ne['class_label']
      X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
      print(X_train.shape, X_test.shape, y_train.shape, y_test.shape)
      # fit the modeL
      rf = RandomForestClassifier(n_estimators=100, random_state=42)
      rf.fit(X_train, y_train)
      # select the features
      sel = SelectFromModel(rf)
      sel.fit(X_train, y_train)
      # make a list of the selected features
      selected_feat= X_train.columns[(sel.get_support())]
      print(len(selected_feat))
      print(selected_feat)
      # plot all the selected features
      plt.figure(figsize=(10,10))
      sns.barplot(x=sel.estimator_.feature_importances_, y=X_train.columns)
      plt.show()
```



```
[23]: # best features by gradient boosting
      from sklearn.ensemble import GradientBoostingClassifier # or GradientBoostingRegressor
      from sklearn.feature_selection import SelectFromModel
      from sklearn.model_selection import train_test_split
      X = df_ne.drop(['class_label'], axis=1)
      y = df_ne['class_label']
      X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
      print(X_train.shape, X_test.shape, y_train.shape, y_test.shape)
      # fit the modeL
      gb = GradientBoostingClassifier(n_estimators=100, random_state=42)
      gb.fit(X_train, y_train)
      # select the features
      sel = SelectFromModel(gb)
      sel.fit(X_train, y_train)
      # make a list of the selected features
      selected_feat= X_train.columns[(sel.get_support())]
      print(len(selected_feat))
      print(selected_feat)
      # plot the feature importance
      plt.figure(figsize=(12,8))
      feat_importances = pd.Series(gb.feature_importances_, index=X_train.columns)
      feat_importances.nlargest().plot(kind='barh')
      plt.show()
```





[24]:	<pre># drop att df.drop(['d</pre>	ack_class and attack_class',	attack_type co 'attack_type'	oLumns '], ax	is=1, inplace=Tr	ue)									
	<pre>df_numerical = df.select_dtypes(include=np.number) df_categorical = df.select_dtypes(exclude=np.number)</pre>														
	<pre># one hot encoding of the categorical variables df_he_categorical = pd.get_dummies(df_categorical, drop_first=True) df_he_categorical.head() # normalize the numerical variables from sklearn.preprocessing import MinMaxScaler scaler = MinMaxScaler() df_n_numerical = pd.DataFrame(scaler.fit_transform(df_numerical), columns=df_numerical.columns) # concat the encoded and normalized datasets df_he_n = pd.concat[[df_n_numerical, df_he_categorical], axis=1) df he no head()</pre>														
[24]:	duration	src_bytes	dst_bytes	land	wrong_fragment	urgent	hot	num_failed_logins	logged_in	num_compromised	root_shell	su_attempted	num_root		
	0 0.0	3.558064e-07	0.000000e+00	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0		
	1 0.0	1.057999e-07	0.000000e+00	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0		
	2 0.0	0.000000e+00	0.000000e+00	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0		
	3 0.0	1.681203e-07	6.223962e-06	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0		
	4 0.0	1.442067e-07	3.206260e-07	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0		
	•												•		

[25]: # fit the model gb = GradientBoostingClassifier(n\_estimators=100, random\_state=42) gb.fit(X\_train, y\_train) # select the features sel = SelectFromModel(gb) sel.fit(X\_train, y\_train) # make a list of the selected features selected\_feat= X\_train.columns[(sel.get\_support())] print(len(selected\_feat)) print(selected\_feat) # plot the feature importance plt.figure(figsize=(12,8)) feat\_importances = pd.Series(gb.feature\_importances\_, index=X\_train.columns) feat\_importances.nlargest().plot(kind='barh') plt.tight\_layout()

5



[26]:

```
from xgboost import XGBClassifier
# fit the model
xgb = XGBClassifier(n_estimators=100, random_state=42)
xgb.fit(X_train, y_train)
# select the features
sel = SelectFromModel(xgb)
sel.fit(X_train, y_train)
# make a list of the selected features
selected_feat= X_train.columns[(sel.get_support())]
print(len(selected_feat))
print(selected_feat)
# plot the feature importance
plt.figure(figsize=(12,8))
feat_importances = pd.Series(gb.feature_importances_, index=X_train.columns)
feat_importances.nlargest().plot(kind='barh')
plt.tight_layout()
plt.show()
```

8



[39]: # Train the GAW for epoch in range(epochs + 1): idx = np.random.randIn(e, df.shape[0], batch\_size) real\_data = df.iloc[idx] noise = np.random.normal(0, 1, (batch\_size, latent\_dim)) generated\_data = generator.predict(noise, verbose=0) real\_labels = np.cens((batch\_size, 1)) fake\_labels = np.zeros((batch\_size, 1)) d\_loss\_real = discriminator.train\_on\_batch(real\_data, real\_labels) d\_loss\_fake = discriminator.train\_on\_batch(real\_data, fake\_labels) d\_loss = 0.5 \* np.add(d\_loss\_real, d\_loss\_fake) noise = np.random.normal(0, 1, (batch\_size, latent\_dim)) valid\_labels = np.ones((batch\_size, 1)) g\_loss = gan.train\_on\_batch(noise, valid\_labels) # Print the progress if epoch % 100 == 0: print(f\*Epoch (epoch), D\_Loss: (d\_loss), G\_Loss: (g\_loss)") WANNING:tensorFlow:From C:\Program Files\Anaconda3\Lib\site-packages\keras\src\utils\tf\_utils.py:492: The name tf.ragged.RaggedTensorValue is deprecate d. Ploss: 0.7589017152786255, G\_Loss: 0.57380211353302

Epoch 0, D Loss: 0.7580017152786255, G Loss: 0.57380211353302 Epoch 100, D Loss: 0.7977718412876129, G Loss: 0.4618147015571594 Epoch 200, D Loss: 0.6338235139846802, G Loss: 0.60539380139237 Epoch 300, D Loss: 0.607880361366272, G Loss: 0.756239302859277 Epoch 400, D Loss: 0.6138206720352173, G Loss: 0.756289585494995 Epoch 500, D Loss: 0.658779263496398, G Loss: 0.7668085252334595 Epoch 600, D Loss: 0.6657924263362885, G Loss: 0.7583783864974976 Epoch 700, D Loss: 0.665239426332184, G Loss: 0.7568085252334595 Epoch 900, D Loss: 0.66523631627559662, G Loss: 0.7583783864974976 Epoch 900, D Loss: 0.66523934232184, G Loss: 0.758408512619019 Epoch 900, D Loss: 0.6993010342121124, G Loss: 0.6383166578292847 Epoch 1000, D Loss: 0.7098850011825562, G Loss: 0.674918293952919

```
[41]: import matplotlib.pyplot as plt
      # Defining subplot format
      features_per_row = 4
      num_rows = len(df.columns) // features_per_row
      fig, axs = plt.subplots(num_rows, features_per_row, figsize=(15, 5 * num_rows))
      # Visualize the generated samples
      for row in range(num_rows):
          for col in range(features_per_row):
              feature_idx = row * features_per_row + col
              if feature_idx < len(df.columns):</pre>
                  axs[row, col].hist(generated_data[:, feature_idx], bins=50, color='blue', alpha=0.7)
                  axs[row, col].set_title(df.columns[feature_idx])
                  axs[row, col].set_xlabel('Value')
                  axs[row, col].set_ylabel('Frequency')
      plt.suptitle('Generated Data Distribution', y=1.02)
      plt.tight_layout()
      plt.show()
```







0.0 Value

0.5

0



5

service\_aol





0.75

0

8

6

Frequency

2

-0.75



-0.25 0.00 Value

0.50

0.75

0.25



service\_bgp

Frequency

0



## 6 Web Application

## 6.1 Purpose

This web application is just a demonstration for uploading datasets and predicting intrusion attack types' categories using pre-trained models.

## 6.2 Running the Application

- 1. Navigate to the folder that contains "app.py".
- 2. Copy the path of the folder.
- 3. Open Anaconda Prompt and change directory to this folder, by pasting the path.
- 4. Run the web application by entering the following command:



Figure 32: Running the web application in Anaconda Prompt

5. Open the web interface by using a browser.

S Intrusion Detection System - Ha					
< > e	127.0.0.1:5000/predict			ବ ୯୮୭ 🖉	🔲 🖬 🎄 🔍 VPN 🚍
Intru	usion Detection System View Data				*
		Check Intrusio	on Attack Class		
		Choos	e Model		
		Select Model	~		
		Attack Class	hould be DOS		
	Attack Type		Status of the connection – Normal or Erro	or	
	Other	~	Other		×
	Number of connections to the same dest	tination host in the past two seconds	Last Flag		
	Count		Last Flag		
	Percentage of connections to different se	ervices	1 if successfully logged in; 0 otherwise		
	Diff. Service Rate		Logged In		
	Percentage of connections to the same s	ource port	Percentage of connections to the same se	ervice (count-based)	
	Same Source Port Rate		Same Service Rate (Count-based)		
	Percentage of connections to the same s	ervice	Percentage of connections with flag (4) s0	), s1, s2, or s3 (count-based)	
	Same Service Rate		Serror Rate (Count-based)		
	Number of connections having the same	port number	Destination network service used http or	not	
	Same Port Number Count		No		×
		Pro	edict		

Figure 33: Website homepage

## 6.3 Directory structure

#### 6.4 Inputs and Outputs

- Input: CSV file with network traffic data
- Output: Predicted attack classes displayed on the web interface

Intrusion Detection System - H x +			• - 0 ×
< > C ↓ C ↓ D ↓ 127.0.0		ର ୯୮୦ 🦉	□ □ � � • VPN =
Intrusion Detection System View Data Report Downloads *			•
	Attack Class should be DOS		
	Predicted Input		
	Selected Model: Logistic Regression		
	Attack Type: Neptune		
	Diff. Service Rate: 0.2		
	Same Source Port Rate: 0.6		
	Same Service Rate: 0.7		
	Same Port Number Count: 25		
	Status of the Connection: S1		
	Last Flag: 0		
	Logged In: 1		
	Same Service Rate (Count-based): 0.9		
	Serror Rate (Count-based): 0.1		
	Destination Network Service Used HTTP: Yes		
	Attack Class should be DOS		

Figure 34: Attack Classes predictions - Result 1

S Intrusion Detection System - H × +			~ - 0 ×
< > C □ 127.0		ର ଅ।ତୃ 🥂	🔲 🖬 🎄 🔹 VPN 🚍
Intrusion Detection S	ystem View Data Report Downloads 👻		-
	Attack Class should be PROBE		
	Predicted Input		
	Selected Model: Bernoulli NB		
	Attack Type: Satan		
	Diff. Service Rate: 0.8		
	Same Source Port Rate: 0.9		
	Same Service Rate: 0.6		
	Same Port Number Count: 75		
	Status of the Connection: S1		
	Last Flag: 0		
	Logged In: 1		
	Same Service Rate (Count-based): 0.6		
	Serror Rate (Count-based): 0.8		
	Destination Network Service Used HTTP: Yes		
	Attack Class should be PROBE		

Figure 35: Attack Classes predictions - Result 2