

INTRUSION DETECTION SYSTEM USING DEEP LEARNING

*Project Review 3 for
CSE3501- Information Security Analysis and Audit (G2)*

By:
YUVRAJ KUMAR: 18BIT0276
(CNN Implementation)

Complete Team Members:

AMAN AGARWAL: 18BIT0256
SHIVANSH SRIVASTAVA: 18BIT0324
YUVRAJ KUMAR: 18BIT0276

Submitted to

PROF. SUMAIYA THASEEN I.
School of Information Technology and Engineering,
Vellore Institute of Technology, Vellore-632014



FALL 2020-21

1. Design and Description of the System

Network Intrusion Detection Systems (NIDSs) are essential tools for the network system administrators to detect various security breaches inside an organization's network. An NIDS monitors and analyzes the network traffic entering into or exiting from the network devices of an organization and raises alarms if an intrusion is observed. Based on the methods of intrusion detection, NIDSs are categorized into two classes:

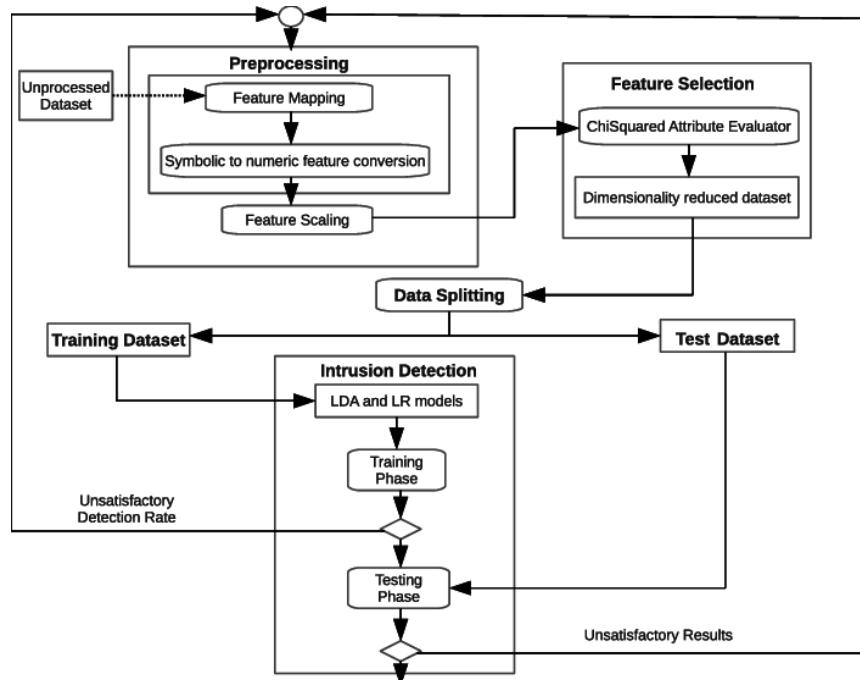
- i) signature (misuse) based NIDS (SNIDS),
- ii) anomaly detection-based NIDS (ADNIDS).

There are primarily two challenges that arise while developing an efficient and flexible NIDS for unknown future attacks. First, proper feature selections from the network traffic dataset for anomaly detection is difficult. The features selected for one class of attack may not work well for other categories of attacks due to continuously changing and evolving attack scenarios. Second, unavailability of labeled traffic dataset from real networks for developing an NIDS. Immense efforts are required to produce such a labeled dataset from the raw network traffic traces collected over a period or in real-time. Additionally, to preserve the confidentiality of the internal organizational network structure as well as the privacy of various users, network administrators are reluctant towards reporting any intrusion that might have occurred in their networks.

It is envisioned that the deep learning based approaches can help to overcome the challenges of developing an efficient NIDS . We can collect unlabeled network traffic data from different network sources and a good feature representation from these datasets using deep learning techniques can be obtained. These features can, then, be applied for supervised classification to a small, but labeled traffic dataset consisting of normal as well as anomalous traffic records. The traffic data for labeled dataset can be collected in a confined, isolated and private network environment. With this motivation, we use self-taught learning, a deep learning technique based on sparse autoencoder and soft-max regression, to develop an NIDS. We verify the usability of the self-taught learning based NIDS by applying on NSL-KDD intrusion dataset, an improved version of the benchmark dataset for various NIDS evaluations - KDD Cup 99. We provide a comparison of our current work with other techniques as well.

In the most widely used approach, the training data is used for both training and testing either using n-fold cross-validation or splitting the training data into training, cross-validation, and test sets. NIDSs based on this approach achieved very high accuracy and less false-alarm rates. The second approach uses the training and test data separately for the training and testing. Since the training and test data were collected in different environments, the accuracy obtained using the second approach is not as high as in the first approach. Therefore, we emphasize on the results of the second approach in our work for accurate evaluation of NIDS. However, we present the results of the first approach as well for completeness. We describe our NIDS implementation before discussing the results.

Architecture

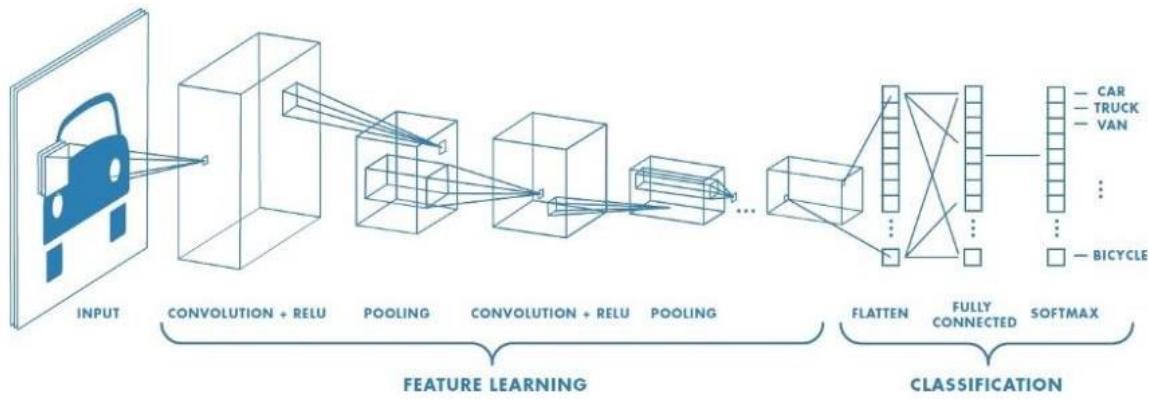


A Network Intrusion Detection System (NIDS) helps system administrators to detect network security breaches in their organizations. However, many challenges arise while developing a flexible and efficient NIDS for unforeseen and unpredictable attacks. We propose a deep learning based approach for developing such an efficient and flexible NIDS. We present the performance of our approach and compare it with a few previous work. Compared metrics include accuracy, precision, recall, and f-measure values.

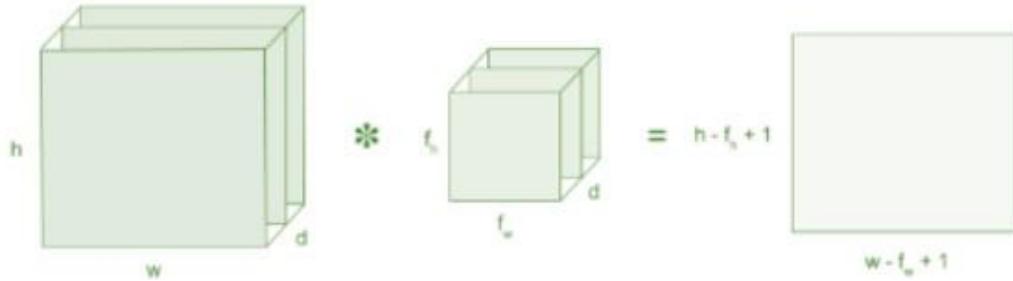
Convolutional Neural Network (CNN) — Deep Learning

In neural networks, Convolutional neural network (ConvNets or CNNs) is one of the main categories to do images recognition, images classifications. Objects detections, recognition faces etc., are some of the areas where CNNs are widely used.

Technically, deep learning CNN models to train and test, each input image will pass it through a series of convolution layers with filters (Kernels), Pooling, fully connected layers (FC) and apply Softmax function to classify an object with probabilistic values between 0 and 1. The below figure is a complete flow of CNN to process an input image and classifies the objects based on values.



- An image matrix (volume) of dimension $(h \times w \times d)$
- A filter $(f_h \times f_w \times d)$
- Outputs a volume dimension $(h - f_h + 1) \times (w - f_w + 1) \times 1$



Convolution of an image with different filters can perform operations such as edge detection, blur and sharpen by applying filters. The below example shows various convolution image after applying different types of filters (Kernels).

Operation	Filter	Convolved Image
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	
Sharpen	$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$	
Box blur (normalized)	$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$	
Gaussian blur (approximation)	$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$	

2. PRE-PROCESSING (Normalization, Sampling – Oversampling, Undersampling (technique used for sampling - atleast 70%balanced))

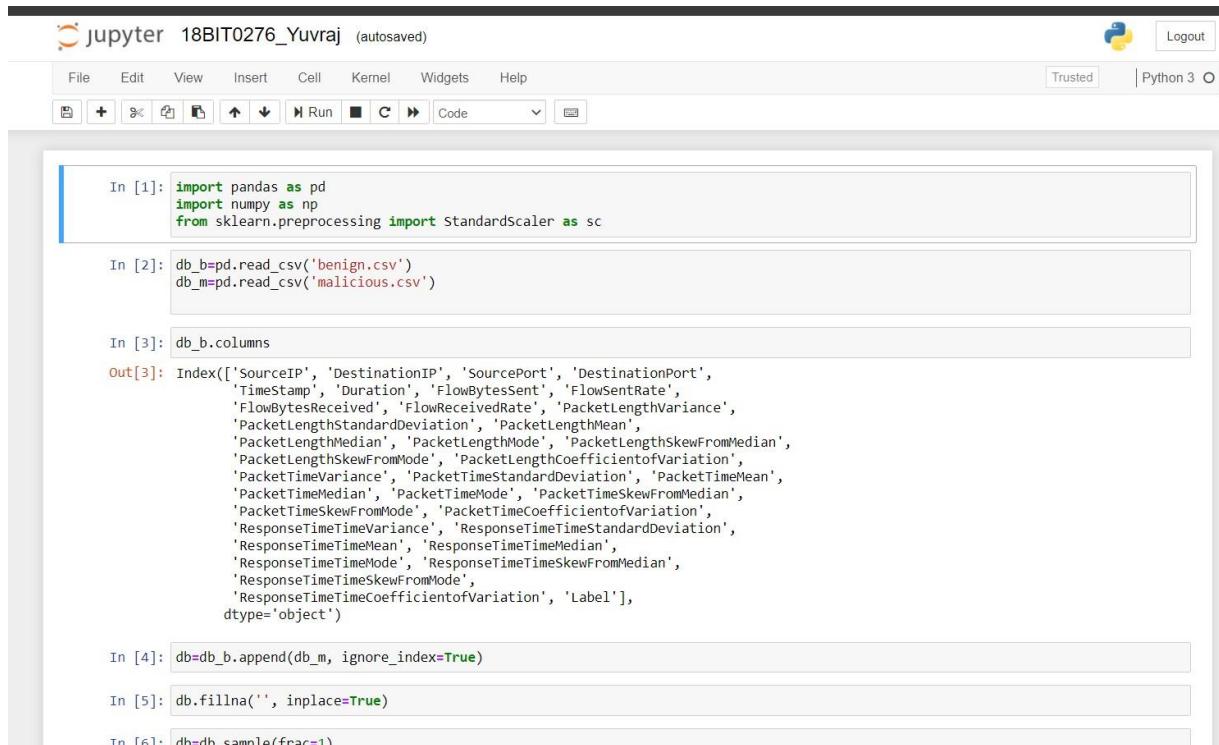
***Sampling is done before model creation. ->Pg.26*

There were 4 datasets in a Total CSVs ZIP file provided to us by the Faculty.

Link to the Dataset :

<https://www.unb.ca/cic/datasets/dohbrw-2020.html>

One having **12-benign** label and other containing **12-malicious** were assigned to me. So, first step of processing was to combine two datasets into one. The combined dataset will have all Benign entries on top and all Malicious entries at the bottom or vice versa. This could create a problem with splitting into training and testing data. To avoid such problem I randomly shuffled the dataset so that both the labels get mixed.



The screenshot shows a Jupyter Notebook interface with the following details:

- Header:** Jupyter 18BIT0276_Yuvraj (autosaved), Logout, Trusted, Python 3.
- Toolbar:** File, Edit, View, Insert, Cell, Kernel, Widgets, Help.
- Code Cells:**
 - In [1]:

```
import pandas as pd
import numpy as np
from sklearn.preprocessing import StandardScaler as sc
```
 - In [2]:

```
db_b=pd.read_csv('benign.csv')
db_m=pd.read_csv('malicious.csv')
```
 - In [3]:

```
db_b.columns
```
 - Out[3]:

```
Index(['SourceIP', 'DestinationIP', 'SourcePort', 'DestinationPort',
       'TimeStamp', 'Duration', 'FlowBytesSent', 'FlowSentRate',
       'FlowBytesReceived', 'FlowReceivedRate', 'PacketLengthVariance',
       'PacketLengthStandardDeviation', 'PacketLengthMean',
       'PacketLengthMedian', 'PacketLengthMode', 'PacketLengthSkewFromMedian',
       'PacketLengthCoefficientOfVariation',
       'PacketTimeVariance', 'PacketTimeStandardDeviation', 'PacketTimeMean',
       'PacketTimeMedian', 'PacketTimeMode', 'PacketTimeSkewFromMedian',
       'PacketTimeCoefficientOfVariation',
       'ResponseTimeVariance', 'ResponseTimeStandardDeviation',
       'ResponseTimeMean', 'ResponseTimeMedian',
       'ResponseTimeMode', 'ResponseTimeSkewFromMedian',
       'ResponseTimeCoefficientOfVariation', 'Label'],
      dtype='object')
```
 - In [4]:

```
db=db_b.append(db_m, ignore_index=True)
```
 - In [5]:

```
db.fillna(' ', inplace=True)
```
 - In [6]:

```
db=db.sample(frac=1)
```

jupyter 18BIT0276_Yuvraj (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 Logout

In [6]: db=db.sample(frac=1)

In [8]: db.head()

Out[8]:

	SourceIP	DestinationIP	SourcePort	DestinationPort	TimeStamp	Duration	FlowBytesSent	FlowSentRate	FlowBytesReceived	FlowReceivedR:
200825	9.9.9.11	192.168.20.212	443	60364	2020-03-25 01:16:29	120.058467	394025	3281.942622	201716	1680.1480
242379	1.1.1.1	192.168.20.207	443	51394	2020-03-20 20:31:33	120.047660	219795	1830.897828	41139	342.6888
229463	192.168.20.212	9.9.9.11	42458	443	2020-03-31 21:06:36	34.946188	1883	53.882844	4828	138.1552
245194	1.1.1.1	192.168.20.210	443	41550	2020-03-29 05:36:24	37.520967	485	12.926106	304	8.1021
93845	192.168.20.212	8.8.4.4	38332	443	2020-03-18 13:05:09	91.691830	493340	5380.413937	899119	9805.8791

5 rows × 35 columns

In [7]: db.to_csv(r'C:\Users\91742\Desktop\Nasscom_Final\Total-CSVs\dataset_db.csv')

In [10]: db_b.shape

Out[10]: (19807, 35)

In [11]: db_m.shape

Out[11]: (249836, 35)

In [8]: db.shape

Out[8]: (269613, 35)

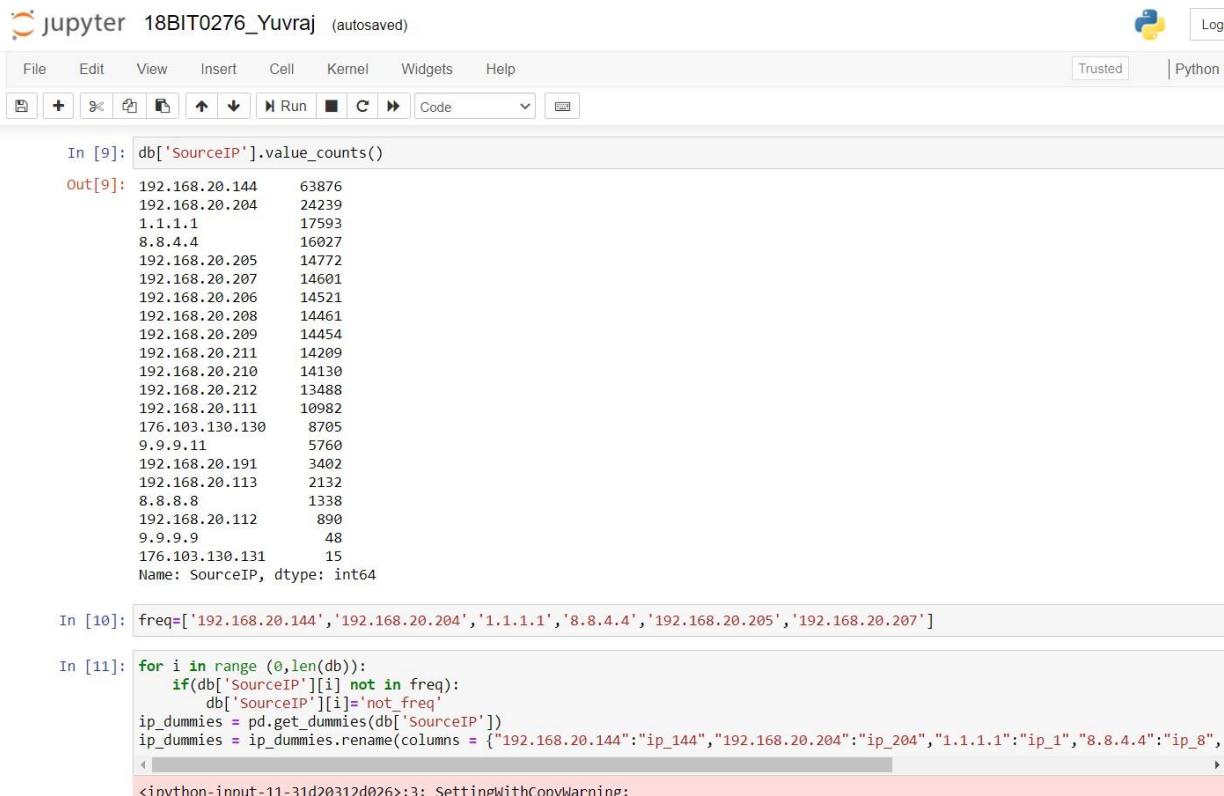
Now, the dataset basically had 2 kind of features, namely categorical and numeric. Categorical features comprised SourceIP, DestinationIP, SourcePort, DestinationPort etc. and numerical features included Duration, FlowSentRate, FlowReceived, FlowBytesRate etc. I have used different techniques to handle both kind of data.

HANDLING CATEGORICAL DATA

Since total number of distinct values in each of the column are large, therefore encoding all them would have created a great sparsity in the dataset. So, in order to overcome the problem of sparsity, I analyzed each feature(column) of the dataset and choose a suitable threshold for the same. If the value is not in the threshold or given range, then I categorized it as other's. Different threshold criteria was chosen for each of the feature based on it's analysis. Then I one-hot encoded all of them so that I can feed them into my neural network. Below is the demonstration for the same.

This was done on the following column:

- **SOURCE IP COLUMN**
- **DESTINATION IP COLUMN**



The screenshot shows a Jupyter Notebook interface with the following content:

```
In [9]: db['SourceIP'].value_counts()
```

```
Out[9]:
```

SourceIP	Count
192.168.20.144	63876
192.168.20.204	24239
1.1.1.1	17593
8.8.4.4	16027
192.168.20.205	14772
192.168.20.207	14601
192.168.20.206	14521
192.168.20.208	14461
192.168.20.209	14454
192.168.20.211	14209
192.168.20.210	14130
192.168.20.212	13488
192.168.20.111	10982
176.103.130.130	8705
9.9.9.11	5760
192.168.20.191	3402
192.168.20.113	2132
8.8.8.8	1338
192.168.20.112	890
9.9.9.9	48
176.103.130.131	15

```
Name: SourceIP, dtype: int64
```

```
In [10]: freq=['192.168.20.144','192.168.20.204','1.1.1.1','8.8.4.4','192.168.20.205','192.168.20.207']
```

```
In [11]: for i in range (0,len(db)):
```

```
    if(db['SourceIP'][i] not in freq):
```

```
        db['SourceIP'][i]='not_freq'
```

```
ip_dummies = pd.get_dummies(db['SourceIP'])
```

```
ip_dummies = ip_dummies.rename(columns = {"192.168.20.144":"ip_144","192.168.20.204":"ip_204","1.1.1.1":"ip_1","8.8.4.4":"ip_8",
```

```
<ipython-inout-11-31d20312d026>:3: SettingWithCopyWarning:
```

jupyter 18BIT0276_Yuvraj (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 Logout

In [14]: `db['DestinationIP'].value_counts()`

Out[14]:

9.9.9.11	150932
8.8.4.4	26044
176.103.130.130	16290
1.1.1.1	13633
8.8.8.8	10177
192.168.20.144	9855
192.168.20.206	4536
192.168.20.205	4461
192.168.20.204	4457
192.168.20.207	4322
192.168.20.208	4011
192.168.20.211	3944
192.168.20.209	3942
192.168.20.210	3870
192.168.20.212	3687
9.9.9.9	2947
192.168.20.111	1353
192.168.20.113	459
192.168.20.112	457
176.103.130.131	134
192.168.20.191	132

Name: DestinationIP, dtype: int64

In [15]: `freq_dest=['9.9.9.11', '8.8.4.4']`

In [16]: `freq_dest_2=['176.103.130.130', '1.1.1.1','8.8.8.8']`

In [17]: `for i in range (0,len(db)):
 if(db['DestinationIP'][i] in freq_dest_2):
 db['DestinationIP'][i] = 'dest_ip_130_1_8'
 elif(db['DestinationIP'][i] not in freq_dest):
 db['DestinationIP'][i] = 'not_dest_freq'`

SIMILARLY, AS ABOVE, THE SOURCE PORT AND DESTINATION PORT FEATURES HAVE BEEN PRE-PROCESSED AND CONCATENATING PRE-PROCESSED DATA INTO A NEW DATAFRAME

```
In [18]: destination_dummy=pd.get_dummies(db['DestinationIP'])

In [19]: destination_dummy = destination_dummy.rename(columns = {"9.9.9.11":"ip_11","8.8.4.4":"ip_4"})

In [20]: db_create=pd.concat([db_create,destination_dummy], axis = 1)

In [21]: db['SourcePort'].value_counts()

Out[21]: 443      49486
58759      73
41618      60
44886      60
40378      53
...
63314       1
64686       1
61443       1
60121       1
59999       1
Name: SourcePort, Length: 15446, dtype: int64

In [22]: for i in range (0,len(db)):
    if(db['SourcePort'][i] != 443):
        db['SourcePort'][i] = 'Wrong_port'
    else:
        db['SourcePort'][i] = 'port_is_443'

<ipython-input-22-522c55c71581>:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
db['SourcePort'][i] = 'Wrong_port'
c:\users\9170\anaconda\lib\site-packages\pandas\core\indexing.py:671: SettingWithCopyWarning:
db['SourcePort'][i] = 'port_is_443'

In [23]: db['SourcePort'][i] = 'port_is_443'

In [23]: for i in range (0,len(db)):
    if(db['DestinationPort'][i] != 443):
        db['DestinationPort'][i] = 'Wrong_port'
    else:
        db['DestinationPort'][i] = 'port_443'

<ipython-input-23-ec0a6933b03d>:5: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
db['DestinationPort'][i] = 'port_443'
<ipython-input-23-ec0a6933b03d>:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
db['DestinationPort'][i] = 'Wrong_port'

In [24]: source_port_dummy = pd.get_dummies(db['SourcePort'])

In [25]: destination_port_dummy = pd.get_dummies(db['DestinationPort'])

In [26]: destination_port_dummy = destination_port_dummy.rename(columns = {"port_443":"destination_port_443"})

In [27]: db_create = pd.concat([db_create,source_port_dummy],axis=1)
db_create = pd.concat([db_create,destination_port_dummy],axis=1)
```

HANDLING NUMERIC FEATURES / NORMALIZATION

Some values were really large and some values were small. To avoid biasness and create a balanced dataset, I used Standard Scaler to normalize the numeric data.

I WILL BE TAKING CARE OF ALL THE NUMERIC FEATURE BELOW LIKE, FlowSentRate, ByteRate, Duration etc. using different Methodologies like Median, Standard Scaling, Sc_Fit_Transform etc.

- **DURATION COLUMN**
- **FLOW BYTES SENT**
- **FLOW SENT RATE**
- **FLOW BYTES RECEIVED**
- **PACKET LENGTH (USING MEDIAN METHOD)**
- **PACKET LENGTH CoV**
- **PACKET TIME VARIANCE**
- **PACKET LENGTH MEAN**
- **OTHER SPECIAL NUMERIC FEATURES**

```
In [28]: from sklearn.preprocessing import StandardScaler
In [29]: sc=StandardScaler()
In [30]: drn_db = db[["Duration"]]
drn_db = sc.fit_transform(drn_db)
drn_db = pd.DataFrame(drn_db)
In [31]: drn_db = drn_db.rename(columns = {0:"Duration"})
In [32]: db_create = pd.concat([db_create,drn_db],axis=1)
In [33]: db['FlowBytesSent'].value_counts()
Out[33]: 1807      44384
1875      15189
1806      7627
1085      7474
1739      7471
...
211692      1
199402      1
201449      1
217825      1
421554      1
Name: FlowBytesSent, Length: 72805, dtype: int64
In [34]: fs = db[['FlowBytesSent']]
In [35]: fs.mean()
Out[35]: FlowBytesSent    65757.889977
```

jupyter 18BIT0276_Yuvraj (autosaved)  Logout

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

In [47]:

```
for i in range (0,len(db)):
    if(db['FlowBytesSent'][i] == 1807):
        db['FlowBytesSent'][i] = 'fs_1807'
    elif(db['FlowBytesSent'][i] >= 1810):
        db['FlowBytesSent'][i] = 'fsless'
    else:
        db['FlowBytesSent'][i] = 'fs_more'
```

NameError: name 'df' is not defined

In [38]:

```
for i in range (0,len(db)):
    if(db['FlowBytesSent'][i] == 1807):
        db['FlowBytesSent'][i] = 'fs_1807'
    elif(db['FlowBytesSent'][i] >= 1810):
        db['FlowBytesSent'][i] = 'fsless'
    else:
        db['FlowBytesSent'][i] = 'fs_more'
```

<ipython-input-38-4a8b68076bc7>:5: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

jupyter 18BIT0276_Yuvraj (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Trusted | Python

```
rsus-a-copy
db['FlowBytesSent'][i] = 'fs_1807'

In [39]: fs_dummy = pd.get_dummies(db['FlowBytesSent'])

In [40]: destination_port_dummy.head()

Out[40]:
Wrong_port  destination_port_443
281170          0          1
250984          0          1
90316           0          1
6366            0          1
135318          0          1

In [41]: db_create = pd.concat([db_create,fs_dummy],axis=1)

In [42]: fs_r = db[["FlowSentRate"]]

In [43]: fs_r = sc.fit_transform(fs_r)

In [44]: fs_r = pd.DataFrame(fs_r)

In [45]: fs_r = fs_r.rename(columns = {0:"FlowSentRate"})

In [46]: db_create = pd.concat([db_create,fs_r],axis=1)
```

jupyter 18BIT0276_Yuvraj (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Trusted | Python 3

```
see the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
db["FlowBytesReceived"][115] = 0

In [48]: freq = db[['FlowBytesReceived']]

In [49]: freq.mean()

Out[49]: FlowBytesReceived    67362.488946
dtype: float64

In [50]: freq.median()

Out[50]: FlowBytesReceived    4896.0
dtype: float64

In [51]: for i in range (0,len(db)):
    if(isinstance(db['FlowBytesReceived'][i], str)):
        db['FlowBytesReceived'][i] = 0

    if(db['FlowBytesReceived'][i] == 98):
        db['FlowBytesReceived'][i] = 'fr_98'
    elif(db['FlowBytesReceived'][i] > 4900):
        db['FlowBytesReceived'][i] = 'fr_more'
    else:
        db['FlowBytesReceived'][i] = 'fr_less'

ipython-input-51-b979bb7d5b2d>:8: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

see the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
```

jupyter 18BIT0276_Yuvraj (autosaved)

File Edit View Insert Cell Kernel Widgets Help

In [52]: `freq_dummies = pd.get_dummies(db['FlowBytesReceived'])`

In [64]: `db_new = pd.concat([db_new,freq_dummies],axis=1)`

NameError Traceback (most recent call last)
<ipython-input-64-34d07a978212> in <module>
----> 1 db_new = pd.concat([db_new,freq_dummies],axis=1)

NameError: name 'db_new' is not defined

In [53]: `db_create = pd.concat([db_create,freq_dummies],axis=1)`

In [54]: `frc = db[["FlowReceivedRate"]]
frc = sc.fit_transform(frc)
frc = pd.DataFrame(frc)`

In [55]: `frc = frc.rename(columns = {0:"FlowRecvRate"})`

In [56]: `db_create = pd.concat([db_create,frc],axis=1)`

In [57]: `plv = db[["PacketLengthVariance"]]
plv = sc.fit_transform(plv)
plv = pd.DataFrame(plv)`

In [58]: `plv = plv.rename(columns = {0:"Pcket_lenVar"})`

In [59]: `db_create = pd.concat([db_create,plv],axis=1)`

In [60]: `plsd = db[["PacketLengthStandardDeviation"]]`

jupyter 18BIT0276_Yuvraj (autosaved)

File Edit View Insert Cell Kernel Widgets Help

In [62]: `db_create = pd.concat([db_create,plsd],axis=1)`

In [63]: `for i in range (0,len(db)):
 if(isinstance(db['PacketLengthMean'][i], str)):
 df['PacketLengthMean'][i] = 0`

In [64]: `plm = db[['PacketLengthMean']]
plm = sc.fit_transform(plm)
plm = pd.DataFrame(plm)`

In [65]: `plm = plm.rename(columns = {0:"Length_mean"})`

In [66]: `db_create = pd.concat([db_create,plm],axis=1)`

In [67]: `db_temp_half = db_create`

In [68]: `for i in range (0,len(db)):
 if(isinstance(db['PacketLengthMedian'][i], str)):
 db['PacketLengthMedian'][i] = 0`

In [81]: `pl_med = df[['PacketLengthMedian']]`

NameError Traceback (most recent call last)
<ipython-input-81-c6e68f44c154> in <module>
----> 1 pl_med = df[['PacketLengthMedian']]

jupyter 18BIT0276_Yuvraj (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Py

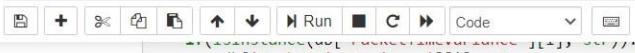
```
In [69]: pl_med = db[['PacketLengthMedian']]  
  
In [70]: pl_med = sc.fit_transform(pl_med)  
pl_med = pd.DataFrame(pl_med)  
  
In [71]: pl_med = pl_med.rename(columns = {0:"L_median"})  
  
In [72]: db_create = pd.concat([db_create,pl_med],axis=1)  
  
In [73]: for i in range (0,len(db)):  
    if(isinstance(db['PacketLengthMode'][i], str)):  
        db['PacketLengthMode'][i] = 0  
  
In [74]: pl_mode = db[['PacketLengthMode']]  
pl_mode = sc.fit_transform(pl_mode)  
pl_mode = pd.DataFrame(pl_mode)  
  
In [75]: pl_mode = pl_mode.rename(columns = {0:"L_mode"})  
  
In [76]: db_create = pd.concat([db_create,pl_mode],axis=1)  
  
In [77]: for i in range (0,len(db)):  
    if(isinstance(db['PacketLengthSkewFromMedian'][i], str)):  
        db['PacketLengthSkewFromMedian'][i] = 0  
  
In [78]: s_med = db[['PacketLengthSkewFromMedian']]  
s_med = sc.fit_transform(s_med)  
s_med = pd.DataFrame(s_med)
```

jupyter 18BIT0276_Yuvraj (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Py

```
In [79]: s_med = s_med.rename(columns = {0:"s_median"})  
  
In [80]: db_create = pd.concat([db_create,s_med],axis=1)  
  
In [81]: for i in range (0,len(db)):  
    if(isinstance(db['PacketLengthSkewFromMode'][i], str)):  
        db['PacketLengthSkewFromMode'][i] = 0  
  
In [82]: s_mode = db[['PacketLengthSkewFromMode']]  
s_mode = sc.fit_transform(s_mode)  
s_mode = pd.DataFrame(s_mode)  
  
In [83]: s_mode = s_mode.rename(columns = {0:"s_mode"})  
  
In [84]: db_create = pd.concat([db_create,s_mode],axis=1)  
  
In [85]: for i in range (0,len(db)):  
    if(isinstance(db['PacketLengthCoefficientofVariation'][i], str)):  
        db['PacketLengthCoefficientofVariation'][i] = 0  
  
In [86]: plcv = db[['PacketLengthCoefficientofVariation']]  
plcv = sc.fit_transform(plcv)  
plcv = pd.DataFrame(plcv)  
  
In [87]: plcv = plcv.rename(columns = {0:"PLCV"})  
  
In [88]: db_create = pd.concat([db_create,plcv],axis=1)
```

jupyter 18BIT0276_Yuvraj (autosaved)

```
File Edit View Insert Cell Kernel Widgets Help Trusted

In [90]: db['PacketTimeVariance'][i] = 0

In [91]: ptv = db[['PacketTimeVariance']]
ptv = sc.fit_transform(ptv)
ptv = pd.DataFrame(ptv)

In [92]: ptv = ptv.rename(columns = {0:"PTV"})

In [93]: db_create = pd.concat([db_create,ptv],axis=1)

In [94]: for i in range (0,len(db)):
    if(isinstance(db['PacketTimeStandardDeviation'][i], str)):
        db['PacketTimeStandardDeviation'][i] = 0

In [95]: pt_std = db[['PacketTimeStandardDeviation']]
pt_std = sc.fit_transform(pt_std)
pt_std = pd.DataFrame(pt_std)

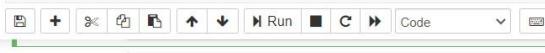
In [96]: pt_std = pt_std.rename(columns = {0:"PT_std"})
db_create = pd.concat([db_create,pt_std],axis=1)

In [97]: for i in range (0,len(db)):
    if(isinstance(db['PacketTimeMean'][i], str)):
        db['PacketTimeMean'][i] = 0

    pt_mean = db[['PacketTimeMean']]
    pt_mean = sc.fit_transform(pt_mean)
    pt_mean = pd.DataFrame(pt_mean)

In [98]: pt_mean = pt_mean.rename(columns = {0:"PT_mean"})
db_create = pd.concat([db_create,pt_mean],axis=1)
```

jupyter 18BIT0276_Yuvraj (autosaved)

```
File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

In [98]: for i in range (0,len(db)):
    if(isinstance(db['PacketTimeMedian'][i], str)):
        db['PacketTimeMedian'][i] = 0

    pt_med = db[['PacketTimeMedian']]
    pt_med = sc.fit_transform(pt_med)
    pt_med = pd.DataFrame(pt_med)

In [99]: pt_med = pt_med.rename(columns = {0:"PT_median"})
db_create = pd.concat([db_create,pt_med],axis=1)

In [100]: for i in range (0,len(db)):
    if(isinstance(db['PacketTimeMode'][i], str)):
        db['PacketTimeMode'][i] = 0

    packt_time_mode = db[['PacketTimeMode']]
    packt_time_mode = sc.fit_transform(packt_time_mode)
    packt_time_mode = pd.DataFrame(packt_time_mode)

In [101]: packt_time_mode = packt_time_mode.rename(columns = {0:"Packt_Time_mode"})
db_create = pd.concat([db_create,packt_time_mode],axis=1)

In [102]: for i in range (0,len(db)):
    if(isinstance(db['PacketTimeSkewFromMedian'][i], str)):
        db['PacketTimeSkewFromMedian'][i] = 0

    s_med_time = db[['PacketTimeSkewFromMedian']]
    s_med_time = sc.fit_transform(s_med_time)
    s_med_time = pd.DataFrame(s_med_time)

In [103]: s_med_time = s_med_time.rename(columns = {0:"s_med_time"})
dfb_create = pd.concat([db_create,s_med_time],axis=1)
```

jupyter 18BIT0276_Yuvraj (autosaved)

```
File Edit View Insert Cell Kernel Widgets Help
Trusted | Python 3
In [104]: for i in range (0,len(db)):
    if(isinstance(db['PacketTimeSkewFromMode'][i], str)):
        db['PacketTimeSkewFromMode'][i] = 0

    s_mode_time = db[['PacketTimeSkewFromMode']]
    s_mode_time = sc.fit_transform(s_mode_time)
    s_mode_time = pd.DataFrame(s_mode_time)

In [105]: s_mode_time = s_mode_time.rename(columns = {0:"s_mode_time"})

In [106]: db_create = pd.concat([db_create,s_mode_time],axis=1)

In [107]: for i in range (0,len(db)):
    if(isinstance(db['PacketTimeCoefficientofVariation'][i], str)):
        db['PacketTimeCoefficientofVariation'][i] = 0

    pkt_cov_time = db[['PacketTimeCoefficientofVariation']]
    pkt_cov_time = sc.fit_transform(pkt_cov_time)
    pkt_cov_time = pd.DataFrame(pkt_cov_time)

In [108]: pkt_cov_time = pkt_cov_time.rename(columns = {0:"pkt_cov_time"})
db_create = pd.concat([db_create,pkt_cov_time],axis=1)

In [109]: for i in range (0,len(db)):
    if(isinstance(db['ResponseTimeTimeVariance'][i], str)):
        db['ResponseTimeTimeVariance'][i] = 0

    r_time_var = db[['ResponseTimeTimeVariance']]
    r_time_var = sc.fit_transform(r_time_var)
    r_time_var = pd.DataFrame(r_time_var)
```

jupyter 18BIT0276_Yuvraj (autosaved)

```
File Edit View Insert Cell Kernel Widgets Help
Trusted | Python 3
In [110]: r_time_var = r_time_var.rename(columns = {0:"r_time_var"})
db_create = pd.concat([db_create,r_time_var],axis=1)

In [111]: for i in range (0,len(db)):
    if(isinstance(db['ResponseTimeTimeStandardDeviation'][i], str)):
        db['ResponseTimeTimeStandardDeviation'][i] = 0

    r_time_std = db[['ResponseTimeTimeStandardDeviation']]
    r_time_std = sc.fit_transform(r_time_std)
    r_time_std = pd.DataFrame(r_time_std)

In [112]: r_time_std = r_time_std.rename(columns = {0:"r_time_std"})

In [113]: db_create = pd.concat([db_create,r_time_std],axis=1)

In [114]: for i in range (0,len(db)):
    if(isinstance(db['ResponseTimeTimeMean'][i], str)):
        db['ResponseTimeTimeMean'][i] = 0

    r_time_mean = db[['ResponseTimeTimeMean']]
    r_time_mean = sc.fit_transform(r_time_mean)
    r_time_mean = pd.DataFrame(r_time_mean)

In [115]: r_time_mean = r_time_mean.rename(columns = {0:"r_time_mean"})

In [116]: db_create = pd.concat([db_create,r_time_mean],axis=1)

In [117]: for i in range (0,len(db)):
    if(isinstance(db['ResponseTimeTimeMedian'][i], str)):
        db['ResponseTimeTimeMedian'][i] = 0

    r_time_med = db[['ResponseTimeTimeMedian']]
```

File Edit View Insert Cell Kernel Widgets Help

```
r_time_med = db[['ResponseTimeMedian']]
r_time_med = sc.fit_transform(r_time_med)
r_time_med = pd.DataFrame(r_time_med)

<ipython-input-118-4b739f030769>:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
    db['ResponseTimeMedian'][i] = 0

In [119]: r_time_med = r_time_med.rename(columns = {0:"r_time_med"})
db_create = pd.concat([db_create,r_time_med],axis=1)

In [120]: for i in range (0,len(db)):
    if(isinstance(db['ResponseTimeMode'][i], str)):
        db['ResponseTimeMode'][i] = 0

    r_time_mode = db[['ResponseTimeMode']]
    r_time_mode = sc.fit_transform(r_time_mode)
    r_time_mode = pd.DataFrame(r_time_mode)

In [121]: r_time_mode = r_time_mode.rename(columns = {0:"r_time_mode"})
db_create = pd.concat([db_create,r_time_mode],axis=1)

In [122]: for i in range (0,len(db)):
    if(isinstance(db['ResponseTimeSkewFromMedian'][i], str)):
        db['ResponseTimeSkewFromMedian'][i] = 0

    r_s_med_time = db[['ResponseTimeSkewFromMedian']]
    r_s_med_time = sc.fit_transform(r_s_med_time)
    r_s_med_time = pd.DataFrame(r_s_med_time)

<ipython-input-122-1abd97a78dfc>:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame
```

jupyter 18BIT0276_Yuvraj (autosaved)

File Edit View Insert Cell Kernel Widgets Help

A horizontal toolbar with various icons: a floppy disk, a plus sign, a scissor, a copy, a paste, an up arrow, a down arrow, a play button, a run button, a stop button, a clear button, a right arrow, a code editor, and a dropdown arrow.

```
In [123]: r_s_med_time = r_s_med_time.rename(columns = {0:"r_s_med_time"})
db_create = pd.concat([db_create,r_s_med_time],axis=1)
```

```
In [124]: for i in range (0,len(db)):
    if(isinstance(db['ResponseTimeTimeSkewFromMode'][i], str)):
        db['ResponseTimeTimeSkewFromMode'][i] = 0
```

```
In [125]: r_s_mode_time = db[['ResponseTimeTimeSkewFromMode']]
r_s_mode_time = sc.fit_transform(r_s_mode_time)
r_s_mode_time = pd.DataFrame(r_s_mode_time)
```

```
In [126]: r_s_mode_time = r_s_mode_time.rename(columns = {0:"r_s_mode_time"})
db_create = pd.concat([db_create,r_s_mode_time],axis=1)
```

```
In [127]: for i in range (0,len(db)):
    if(isinstance(db['ResponseTimeTimeCoefficientofVariation'][i], str)):
        db['ResponseTimeTimeCoefficientofVariation'][i] = 0
```

```
In [128]: r_c_time = db[['ResponseTimeTimeCoefficientofVariation']]
r_c_time = sc.fit_transform(r_c_time)
r_c_time = pd.DataFrame(r_c_time)
```

```
In [129]: r_c_time = r_c_time.rename(columns = {0:"r_c_time"})
db_create = pd.concat([db_create,r_c_time],axis=1)
```

```
In [142]: labels = pd.get_dummies(db['Label'])
```

```
In [143]: db_create = pd.concat([db_create,labels],axis=1)
```

jupyter 18BIT0276_Yuvraj (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

```
In [146]: db_create.drop('Malicious', axis=1)
```

```
Out[146]:
```

	ip_1	ip_144	ip_204	ip_205	ip_207	ip_8	not_freq	ip_4	ip_11	dest_ip_130_1_8	...	0	r_time_var	r_time_std	r_time_mean	r_time_med	r_time_mode	r_time_min	r_time_max
0	0	0	0	0	0	0	1	0	0	0	...	0.754077	-0.178941	-0.244670	-0.239715	-0.221598	-0.11	-0.11	
1	0	0	0	0	0	0	1	0	0	0	...	-0.814219	-0.132047	0.277059	-0.183655	-0.225091	-0.11	-0.11	
2	0	0	0	0	0	0	1	0	0	0	...	-0.817715	-0.178940	-0.244316	-0.238336	-0.221617	-0.11	-0.11	
3	0	0	0	0	0	0	1	0	0	0	...	-0.316040	-0.178946	-0.250146	-0.238243	-0.221523	-0.11	-0.11	
4	0	0	0	0	0	0	1	0	0	0	...	-0.770578	-0.178937	-0.242955	-0.240378	-0.221754	-0.11	-0.11	
...	
269638	1	0	0	0	0	0	0	0	0	0	...	-0.997728	-0.178940	-0.244337	-0.240169	-0.221636	-0.11	-0.11	
269639	1	0	0	0	0	0	0	0	0	0	...	0.947832	-0.178940	-0.244243	-0.240186	-0.221621	-0.11	-0.11	
269640	1	0	0	0	0	0	0	0	0	0	...	-0.858274	-0.178946	-0.250175	-0.243940	-0.226451	-0.11	-0.11	
269641	1	0	0	0	0	0	0	0	0	0	...	0.898332	-0.178940	-0.244267	-0.240184	-0.221634	-0.11	-0.11	
269642	1	0	0	0	0	0	0	0	0	0	...	-0.588891	-0.178785	-0.219341	-0.231216	-0.218044	-0.11	-0.11	

269643 rows × 45 columns

```
In [147]: db_create['Benign'] = db_create['Benign'].replace(['1'], 'Benign')
db_create['Benign'] = db_create['Benign'].replace(['0'], 'Malicious')
```

```
In [148]: db_create.head()
```

```
Out[148]:
```

	ip_1	ip_144	ip_204	ip_205	ip_207	ip_8	not_freq	ip_4	ip_11	dest_ip_130_1_8	...	0	r_time_var	r_time_std	r_time_mean	r_time_med	r_time_mode	r_time_min	r_time_max
0	0	0	0	0	0	0	1	0	0	0	...	-0.178941	-0.244670	-0.239715	-0.221598	-0.159623	-0.11	-0.11	
1	0	0	0	0	0	0	1	0	0	0	...	-0.132047	0.277059	-0.183655	-0.225091	-0.11	-0.11	-0.11	

jupyter 18BIT0276_Yuvraj (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

```
In [150]: db_create.drop('Malicious', inplace=True, axis=1)
```

```
In [151]: db_create.head()
```

```
Out[151]:
```

	ip_1	ip_144	ip_204	ip_205	ip_207	ip_8	not_freq	ip_4	ip_11	dest_ip_130_1_8	...	0	r_time_var	r_time_std	r_time_mean	r_time_med	r_time_mode	r_time_min	r_time_max
0	0	0	0	0	0	0	1	0	0	0	...	0.754077	-0.178941	-0.244670	-0.239715	-0.221598	-0.11	-0.11	
1	0	0	0	0	0	0	1	0	0	0	...	-0.814219	-0.132047	0.277059	-0.183655	-0.225091	-0.11	-0.11	
2	0	0	0	0	0	0	1	0	0	0	...	-0.817715	-0.178940	-0.244316	-0.238336	-0.221617	-0.11	-0.11	
3	0	0	0	0	0	0	1	0	0	0	...	-0.316040	-0.178946	-0.250146	-0.238243	-0.221523	-0.11	-0.11	
4	0	0	0	0	0	0	1	0	0	0	...	-0.770578	-0.178937	-0.242955	-0.240378	-0.221754	-0.11	-0.11	

5 rows × 45 columns

```
In [152]: db_create.to_csv(r'C:\Users\91742\Desktop\Nasscom_Final\Total-CSVs\final_data_input_FINAL_Semifinal.csv')
```

```
In [153]: db_create['Benign'] = db_create['Benign'].replace(['1'], 'Benign')
db_create['Benign'] = db_create['Benign'].replace(['0'], 'Malicious')
```

```
In [154]: db_create.head()
```

```
Out[154]:
```

	ip_1	ip_144	ip_204	ip_205	ip_207	ip_8	not_freq	ip_4	ip_11	dest_ip_130_1_8	...	0	r_time_var	r_time_std	r_time_mean	r_time_med	r_time_mode	r_time_min	r_time_max
0	0	0	0	0	0	0	1	0	0	0	...	0.754077	-0.178941	-0.244670	-0.239715	-0.221598	-0.11	-0.11	
1	0	0	0	0	0	0	1	0	0	0	...	-0.814219	-0.132047	0.277059	-0.183655	-0.225091	-0.11	-0.11	
2	0	0	0	0	0	0	1	0	0	0	...	-0.817715	-0.178940	-0.244316	-0.238336	-0.221617	-0.11	-0.11	
3	0	0	0	0	0	0	1	0	0	0	...	-0.316040	-0.178946	-0.250146	-0.238243	-0.221523	-0.11	-0.11	

jupyter 18BIT0276_Yuvraj (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

```
5 rows × 45 columns
In [157]: db_create['Benign'] = db_create['Benign'].replace([1], 'Benign')
db_create['Benign'] = db_create['Benign'].replace([0], 'Malicious')

In [158]: db_create.head()

Out[158]:
   ip_1  ip_144  ip_204  ip_205  ip_207  ip_8  not_freq  ip_4  ip_11  dest_ip_130_1_8 ...  0  r_time_var  r_time_std  r_time_mean  r_time_med  r_time_
0  0  0  0  0  0  0  1  0  0  0  ...  0.754077 -0.178941 -0.244670 -0.239715 -0.221598 -0.1
1  0  0  0  0  0  0  1  0  0  0  ... -0.814219 -0.132047 0.277059 -0.183655 -0.225091 -0.1
2  0  0  0  0  0  0  1  0  0  0  ... -0.817715 -0.178940 -0.244316 -0.238336 -0.221617 -0.1
3  0  0  0  0  0  0  1  0  0  0  ... -0.316040 -0.178946 -0.250146 -0.238243 -0.221523 -0.1
4  0  0  0  0  0  0  1  0  0  0  ... -0.770578 -0.178937 -0.242955 -0.240378 -0.221754 -0.1

5 rows × 45 columns
In [159]: db_create.to_csv(r'C:\Users\91742\Desktop\Nasscom_Final\Total-CSVs\final_data_input_LAST_FINAL.csv')

In [ ]:
```

As I had 269643 records in DoH and only 19807 records in benign, I tried Under Sampling and went Ahead with 19807 Records and upon Random Sampling, Each Dataset Picked 19807 records to form 39614 Record Data Together.

DATASETS:

l2-benign

SourceIP	DestIP	SourcePort	DestPort	TimeStamp	Duration	FlowBytes	FlowSentR	FlowBytes	FlowRecei	PacketLen	PacketTim	PacketTim	PacketTim	PacketTim	PacketTim	PacketTim							
2.192.16.2.176.103.1.	50749	443	#####	95.08155	62311	655.3427	65358	687.3889	7447.677	86.4562	135.6738	102	54	1.168467	0.944683	0.673236	70.5858	25.89567	45.05262	48.81129	1.49506		
192.168.2.176.103.1.	50749	443	#####	122.3093	5987	76.137	101232	827.672	10458.12	102.264	141.2455	114	54	1.29968	0.85312	0.742042	70.48569	26.6170	52.2879	48.83031	31.71996		
4.192.16.2.176.103.1.	50749	443	#####	120.9584	38784	320.6391	38236	316.1086	7300.299	85.44176	133.7153	89	54	1.570027	0.932978	0.638893	135.8911	36.86341	50.31611	37.7075	0.4715728		
192.168.2.176.103.1.	50749	443	#####	110.5011	5611.0617	69757	631.2789	8498.283	92.19155	139.1235	114	54	0.817544	0.92333	0.66266	11.1833	33.43835	51.69373	34.88825	13.28093			
176.103.1.292.176.103.1.	443	#####	54.22989	83641	154.34241	76804	1416.267	8052.746	89.73709	138.9134	114	54	0.83288	0.277627	0.645993	341.6966	18.48504	36.43562	49.82526	7.342519			
7.192.168.2.176.103.1.	52491	443	#####	145.4607	50408	371.8117	63843	438.902	16074.98	126.7871	141.061	89	54	1.231852	0.686671	0.988811	151.785	38.98442	76.2529	47.49244	50.88041		
192.168.2.176.103.1.	52742	443	#####	1.15541	1713	11054.63	7411	4768.76	234126	483.8657	351.1154	129.5	54	1.37403	0.61405	1.37802	0.01068	0.032679	0.06289	0.0609	0.087941		
9.192.16.2.176.103.1.	52742	443	#####	0.027001	55	2036.962	66	244.354	30.25	5.5	60.5	60.5	55	0	1	0.09009	0.000182	0.01350	0.01350	0.01350	0		
192.168.2.176.103.1.	52491	443	#####	64.64955	32285	723.0755	36193	810.6016	8371.95	91.49836	139.751	114	54	0.844311	0.937186	0.654724	19.93524	13.97209	23.2550	23.03529	0.644233		
11.192.16.2.176.103.1.	52742	443	#####	0.046455	55	1183.941	66	1420.73	30.25	5.5	60.5	60.5	55	0	1	0.09009	0.00054	0.023228	0.023228	0.023228	0		
192.168.2.176.103.1.	52742	443	#####	44.89041	163	3.631065	357	7.9527	168.25	12.97112	65	60	54	1.156415	0.848038	0.199556	37.73707	19.42243	33.65414	44.86306	44.86306		
13.192.16.2.176.103.1.	52491	443	#####	91.33806	8090	88.57206	7805	85.45178	6158.022	78.47307	135.1575	89	60	1.382289	0.830316	0.626995	140.5086	11.85363	78.94149	81.96176	0		
14.192.16.2.176.103.1.	52491	443	#####	121.3332	53598	441.7434	59634	491.4908	8309.039	91.15374	136.2599	89	54	1.555393	0.92431	0.686969	101.7683	31.90114	57.2504	55.49805	78.54958		
192.168.2.176.103.1.	52491	443	#####	12.3599	56840	505.8745	60236	536.0988	8218.835	90.65779	135.819	89	54	1.545931	0.90250	0.66749	117.1273	33.58008	34.36958	25.3958	72.4556		
176.103.1.192.16.2.176.103.1.	443	52491	#####	135.2442	69616	514.9088	65101	481.3588	835.123	91.39542	136.6237	89	54	1.561232	0.904025	0.686957	90.743	30.06253	50.1431	44.10633	20.5056		
192.168.2.176.103.1.	52491	443	#####	121.1859	290.201	41324	37434	308.8972	766.456	87.67244	136.5188	114	54	0.705505	0.94122	0.64172	116.1656	34.15105	56.05765	62.02651	76.86132		
192.168.2.176.103.1.	52491	443	#####	98.55257	7815	767.9016	81233	825.1738	7709.703	80.84091	135.4259	89	54	1.586218	0.92735	0.648361	70.6148	26.60103	54.34677	56.88736	29.66267		
176.103.1.192.16.2.176.103.1.	443	52491	#####	33.053653	13565	410.3599	11354	343.4741	940.9792	97.00501	139.9944	89	54	1.570664	0.886494	0.69291	102.984	10.14811	12.20385	7.858421	33.02868		
20.192.16.2.176.103.1.	56460	443	#####	114.9067	57076	495.7611	68354	594.8653	1703.22	130.5114	145.5304	114	54	0.74776	0.709883	0.890678	117.8.686	34.332	46.67889	49.60508	0.054805		
21.192.16.2.176.103.1.	56460	443	#####	96.75071	4024	434.5628	44920	464.28849	7899.639	88.87992	133.7908	89	54	1.511841	0.897737	0.66432	48.3844	22.09942	39.3142	55.38985	1.445943		
192.168.2.176.103.1.	56460	443	#####	105.5193	61248	580.4434	63497	601.757	802.091	89.60647	135.1517	89	54	1.54575	0.90563	0.662994	7.670.844	25.3597	45.8976	56.66998	26.27152		
23.192.16.2.176.103.1.	56460	443	#####	110.6149	2747	246.775	29436	266.1124	7805.691	88.34982	135.0786	89	54	1.56464	0.917699	0.654062	147.2847	38.37769	55.58280	55.12513	20.50564		
192.168.2.176.103.1.	56460	443	#####	121.773	67491	554.2362	73742	605.5694	959.662	97.82444	141.3744	114	54	0.83495	0.893175	0.691953	107.93	32.85886	53.07471	52.74703	20.96734		
25.192.16.2.176.103.1.	56460	443	#####	121.6762	38920	320.4489	41627	42.312	819.5157	90.5291	135.4924	89	54	1.54069	0.90018	0.668149	118.212	34.4269	47.17519	70.62596	89.83199		
192.168.2.176.103.1.	56460	443	#####	121.2499	72213	505.5516	80570	664.4953	1046.625	102.2871	139.9112	89	54	1.493185	0.88902	0.730126	107.178	23.7371	32.7371	50.99465	45.52109		
27.192.16.2.176.103.1.	56460	443	#####	74.03461	30733	415.1167	32409	437.7547	7758.135	88.08027	132.6513	89	54	1.486755	0.8925	0.663999	650.7667	25.51013	33.52663	32.25093	1.671356		
29.192.16.2.176.103.1.	56611	443	#####	122.0805	35590	291.2591	4587	365.2264	21960.19	148.1897	146.3084	89	54	1.160169	0.622907	1.012859	552.312	23.50125	45.36756	41.00397	38.21969		
192.168.2.176.103.1.	56611	443	#####	96.36561	41539	431.0563	45477	462.5821	8937.76	94.5981	137.5655	89	54	1.541112	0.513704	0.687235	78.74395	29.58783	53.87872	48.80572	22.30518		

12-malicious

Combined Final dataset

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W
1		ip_1	ip_144	ip_204	ip_205	ip_207	ip_8	not_freq	ip_4	ip_11	dest_ip_1	not_dest_1	Wrong_po	po_port_is_44	Wrong_po	destination	Duration	fs_1807	fs_more	fsless	FlowSentR	FlowRecvF	Pcket_JenLe
2	0	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	-0.44021	0	0	1	-0.04745	-0.06851	1.373156
3	1	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	1.171644	0	0	1	-0.04652	-0.06607	-0.80391
4	2	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	1.1458	0	0	1	-0.04638	-0.06589	-0.78339
5	3	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	2.703618	0	0	1	-0.04773	-0.06989	-0.91759
6	4	0	0	0	0	0	0	1	0	0	0	1	0	1	1	0	1.15286	0	0	1	-0.04657	-0.06179	0.003366
7	5	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	-1.8219	0	0	1	0.016064	-0.00319	-0.91769
8	6	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	1.453431	0	1	0	-0.04301	-0.05705	-0.89671
9	7	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	-0.4519	0	1	0	-0.04742	-0.06843	-0.04242
10	8	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	-0.4348	0	0	1	-0.04642	-0.06692	0.378454
11	9	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	0.44384	0	1	0	-0.04744	-0.06849	0.114244
12	10	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	-0.4345	0	1	0	-0.04744	-0.06848	0.087027
13	11	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	0.45441	0	0	1	-0.04743	-0.06844	0.086357
14	12	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	1.14569	0	0	1	-0.03766	-0.03764	-0.19232
15	13	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	0.43936	0	0	1	-0.04744	-0.06847	0.087027
16	14	0	0	0	0	0	0	1	0	0	0	1	0	1	1	0	-1.1667	0	0	1	-0.03899	-0.0145	0.678708
17	15	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	-1.14745	0	0	1	-0.04391	-0.04548	0.675613
18	16	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	0.465147	0	0	1	-0.04405	-0.06645	-0.9041
19	17	0	0	0	0	0	0	1	0	0	0	1	0	1	1	0	1.467265	0	0	1	-0.04678	-0.06789	-0.87492
20	18	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	1.479826	0	0	1	-0.04624	-0.06172	0.192224
21	19	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	1.446425	0	0	1	-0.04771	-0.06984	-0.191934
22	20	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	-0.45159	0	0	1	-0.04743	-0.06847	0.114933
23	21	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	-0.42489	0	0	1	-0.04743	-0.06847	0.054202
24	22	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	-1.08575	0	0	1	-0.04641	-0.06085	2.646019
25	23	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	1.468775	0	0	1	-0.0442	-0.06597	-0.78779
26	24	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	1.439916	0	0	1	-0.04771	-0.06984	-0.91931
27	25	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	-0.45422	0	0	1	-0.04743	-0.06847	0.114244
28	26	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	-0.43477	0	0	1	-0.04745	-0.06852	1.371659
29	27	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	-0.22303	0	0	1	-0.0477	-0.0698	-0.91799

3. TRAINING AND CROSS VALIDATION

IMPORTING NECESSARY LIBRARIES AND PACKAGES ON GOOGLE COLAB

FINAL_CNN_MODEL.ipynb

File Edit View Insert Runtime Tools Help Last edited on October 21

+ Code + Text

```
[ ] import pandas as pd
import numpy as np
from google.colab import drive
#from imblearn.under_sampling import TomekLinks

[ ] drive.mount('/content/drive')
Mounted at /content/drive

[ ] db = pd.read_csv('/content/drive/My Drive/final_data_input_LAST_FINAL.csv')

[ ] db.head()
```

Unnamed: 0	ip_1	ip_144	ip_204	ip_205	ip_207	ip_8	not_freq	ip_4	ip_11	dest_ip_130_1_8	not_dest_freq	Wrong_port	port_is_443	Wrong_port.1	destination_port_443	Duration	fs
0	0	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	-0.440208
1	1	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	1.171644
2	2	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	1.145800
3	3	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	2.703618
4	4	0	0	0	0	0	0	1	0	0	0	1	0	1	1	0	1.152860

FINAL_CNN_MODEL.ipynb

File Edit View Insert Runtime Tools Help Last edited on October 21

+ Code + Text

```
[ ] from keras.preprocessing import sequence
from keras.models import Sequential
from keras.layers import Dense, Dropout, Activation, Lambda
from keras.layers import Embedding
from keras.layers import Convolution1D, MaxPooling1D, Flatten
from keras.datasets import imdb
from keras import backend as K
from sklearn.model_selection import train_test_split
import pandas as pd
from keras.utils import to_categorical

from sklearn.preprocessing import Normalizer
from keras.models import Sequential
from keras.layers import Convolution1D, Dense, Dropout, Flatten, MaxPooling1D
from keras.utils import np_utils
import numpy as np
import h5py
from keras import callbacks
from keras.layers import LSTM, GRU, SimpleRNN
from keras.callbacks import CSVLogger
from keras.callbacks import ModelCheckpoint, EarlyStopping, ReduceLROnPlateau, CSVLogger
```

FINAL_CNN_MODEL.ipynb

File Edit View Insert Runtime Tools Help Last edited on October 21

+ Code + Text

```
[ ] from keras.layers import LSTM, GRU, SimpleRNN
from keras.callbacks import CSVLogger
from keras.callbacks import ModelCheckpoint, EarlyStopping, ReduceLROnPlateau, CSVLogger

[ ] db = db.iloc[:,1:46]

[ ] db = db.iloc[:,0:45]

[ ] db
```

	ip_1	ip_144	ip_204	ip_205	ip_207	ip_8	not_freq	ip_4	ip_11	dest_ip_130_1_8	not_dest_freq	Wrong_port	port_is_443	Wrong_port.1	destination_port_443	Duration	fs_1807	
0	0	0	0	0	0	0	1	0	0	0	0	1	1	0	0	1	-0.440208	0
1	0	0	0	0	0	0	1	0	0	0	0	1	1	0	0	1	1.171644	0
2	0	0	0	0	0	0	1	0	0	0	0	1	1	0	0	1	1.145800	0
3	0	0	0	0	0	0	1	0	0	0	0	1	1	0	0	1	2.703618	0
4	0	0	0	0	0	0	1	0	0	0	0	1	0	1	1	0	1.152860	0
...
269638	1	0	0	0	0	0	0	0	0	0	0	1	0	1	1	0	-1.180485	0
269639	1	0	0	0	0	0	0	0	0	0	0	1	0	1	1	0	-0.434267	0
269640	1	0	0	0	0	0	0	0	0	0	0	1	0	1	1	0	1.466744	0
269641	1	0	0	0	0	0	0	0	0	0	0	1	0	1	1	0	-0.446174	0
269642	1	0	0	0	0	0	0	0	0	0	0	1	0	1	1	0	-0.853504	0

269643 rows x 45 columns

FINAL_CNN_MODEL.ipynb

File Edit View Insert Runtime Tools Help Last edited on October 21

+ Code + Text

```
[ ] X1 = db.iloc[:,0:44]

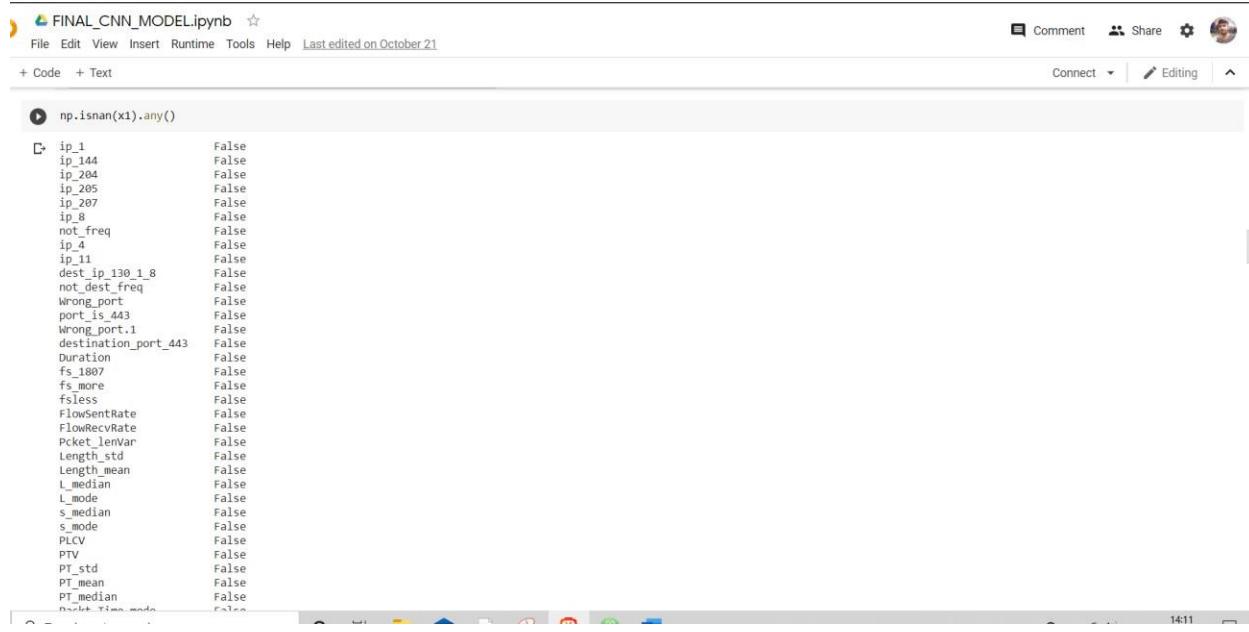
[ ] y1 = db.iloc[:,44]

[ ] x1
```

	ip_1	ip_144	ip_204	ip_205	ip_207	ip_8	not_freq	ip_4	ip_11	dest_ip_130_1_8	not_dest_freq	Wrong_port	port_is_443	Wrong_port.1	destination_port_443	Duration	fs_1807	
0	0	0	0	0	0	0	1	0	0	0	0	1	1	0	0	1	-0.440208	0
1	0	0	0	0	0	0	1	0	0	0	0	1	1	0	0	1	1.171644	0
2	0	0	0	0	0	0	1	0	0	0	0	1	1	0	0	1	1.145800	0
3	0	0	0	0	0	0	1	0	0	0	0	1	1	0	0	1	2.703618	0
4	0	0	0	0	0	0	1	0	0	0	0	1	0	1	1	0	1.152860	0
...
269638	1	0	0	0	0	0	0	0	0	0	0	1	0	1	1	0	-1.180485	0
269639	1	0	0	0	0	0	0	0	0	0	0	1	0	1	1	0	-0.434267	0
269640	1	0	0	0	0	0	0	0	0	0	0	1	0	1	1	0	1.466744	0
269641	1	0	0	0	0	0	0	0	0	0	0	1	0	1	1	0	-0.446174	0
269642	1	0	0	0	0	0	0	0	0	0	0	1	0	1	1	0	-0.853504	0

269643 rows x 44 columns

Now I checked any invalid python data type or any number is present in my data which I will be using or if array contains at least one non-numeric value.



```
FINAL_CNN_MODEL.ipynb
File Edit View Insert Runtime Tools Help Last edited on October 21
+ Code + Text
np.isnan(x1).any()
ip_1        False
ip_144      False
ip_204      False
ip_205      False
ip_207      False
ip_8        False
not_freq    False
ip_4        False
ip_11       False
dest_ip_130_1_8  False
not_dest_freq  False
Wrong_port  False
port_is_443  False
Wrong_port_1 False
destination_port_443  False
Duration    False
fs_1807     False
fs_more     False
fsless      False
FlowSentRate False
FlowRecvRate False
Pcket_lenVar False
Length_std   False
Length_mean  False
l_median    False
l_mode      False
s_median    False
s_mode      False
PLCV        False
PTV         False
PT_std      False
PT_mean    False
PT_median   False
Duration_min  False
Duration_max False
14:11
```

Sampling Of the Dataset

Dataset given is highly unbalanced. So, to prevent overfitting, sampling was done. Also since this was network data, therefore creating dummy data might not be feasible in real time. Therefore I preferred undersampling instead of over-sampling. To overcome this, I used technique of random sampling



```
[ ] r_s_med_time      False
r_s_mode_time      False
r_c_time          False
dtype: bool

[ ] y1.head()
0    Benign
1    Benign
2    Benign
3    Benign
4    Benign
Name: Benign, dtype: object

[ ] from collections import Counter
print(Counter(y1))
Counter({'Malicious': 249836, 'Benign': 19807})

[ ] c_class_0, c_class_1 = db.Benign.value_counts()

[ ] c_class_0, c_class_1 = db.Benign.value_counts()

[ ] print(c_class_0)
print(c_class_1)
249836
19807
```

As I had 249836 records in malicious and only 19807 records in benign, I tried Under Sampling and went Ahead with 19807 Records and upon Random

Sampling, Each Dataset Picked 19807 records to form 39614 Record Data Together.

FINAL_CNN_MODEL.ipynb

File Edit View Insert Runtime Tools Help Last edited on October 21

Comment Share

+ Code + Text

Connect Editing

```
[ ] df_c_0 = db[db['Benign'] == 'Malicious']
df_c_1 = db[db['Benign'] == 'Benign']

[ ] df_c_0_under = df_c_0.sample(c_class_1)

[ ] df_under = pd.concat([df_c_0_under, df_c_1], axis=0)

[ ] print(df_under.Benign.value_counts())

Benign    19807
Malicious 19807
Name: Benign, dtype: int64

[ ] df_Totol = df_under

[ ] df_Totol.head()

ip_1 ip_144 ip_204 ip_205 ip_207 ip_8 not_freq ip_4 ip_11 dest_ip_130_1_8 not_dest_freq Wrong_port port_is_443 Wrong_port.1 destination_port_443 Duration fs_1807
253816 0 1 0 0 0 0 0 0 1 0 0 1 0 0 0 1 1 0 0 1 1.443449 0
131414 0 1 0 0 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 -1.135324 0
144159 0 0 0 0 0 0 1 0 0 1 0 1 0 0 0 1 0 0 0 0 1 -0.440212 0
167495 0 0 0 1 0 0 0 0 1 0 0 0 1 0 0 1 0 0 0 0 1 -1.150709 0
253272 0 0 0 0 0 0 1 0 0 1 0 0 1 0 0 1 0 0 0 0 1 1.140970 0
```

TRAINING AND TESTING ACCURACY

```
[ ] x_train,x_test,y_train,y_test = train_test_split(x,y,train_size=0.8,random_state=42)

[ ] from sklearn.model_selection import StratifiedKFold
kfold = StratifiedKFold(n_splits=10, shuffle=True, random_state=42)
```

MODEL CREATION

I also added two layers in my CNN model which are:

Convolution1D and MaxPooling1D

And the optimizer used here is “adam”



```
File Edit View Insert Runtime Tools Help Last edited on October 21
Comment Share Connect Editing ↑

[ ] cvscores = []
i=0
for train, test in kfold.split(x_train, y_train):
    cnn = Sequential()
    cnn.add(Convolution1D(64, 3, padding="same", activation="relu", input_shape=(44,1)))
    cnn.add(MaxPooling1D(pool_size=(2)))
    cnn.add(Flatten())
    cnn.add(Dense(128, activation="relu"))
    cnn.add(Dropout(0.5))
    cnn.add(Dense(2, activation="sigmoid"))
    # define optimizer and objective, compile cnn

    cnn.compile(loss="binary_crossentropy", optimizer="adam", metrics=['accuracy'])

    x_tn = x_train.iloc[train]
    y_tn = y_train.iloc[train]
    y_tn = pd.get_dummies(y_tn)
    x_ts = x_train.iloc[test]
    y_ts = y_train.iloc[test]
    y_ts = pd.get_dummies(y_ts)

    x_tn1 = x_tn.to_numpy()
    x_tn1 = np.reshape(x_tn1, (x_tn1.shape[0],x_tn1.shape[1],1))

    x_ts1 = x_ts.to_numpy()
    x_ts1 = np.reshape(x_ts1, (x_ts1.shape[0],x_ts1.shape[1],1))

    cnn.fit(x_tn1, y_tn, epochs=50, batch_size=64, verbose=1)
```

+ Code

+ Text

```

[ ] x_ts1 = x_ts.to_numpy()
x_ts1 = np.reshape(x_ts1, (x_ts1.shape[0],x_ts1.shape[1],1))

cnn.fit(x_ts1, y_ts, epochs=50,batch_size=64,verbose=1)
scores = cnn.evaluate(x_ts1, y_ts, verbose=1)
print(str(i)+"th Fold :")
print("%s: %.2f%%" % (cnn.metrics_names[1], scores[1]*100))
cvscores.append(scores[1] * 100)
i = i+1
print("-----")

print("Average validation accuracy : ")
print("%.2f%% (+/- %.2f%%)" % (np.mean(cvscores), np.std(cvscores)))

Epoch 1/50
446/446 [=====] - 3s 7ms/step - loss: 0.2065 - accuracy: 0.9240
Epoch 2/50
446/446 [=====] - 3s 7ms/step - loss: 0.1566 - accuracy: 0.9440
Epoch 3/50
446/446 [=====] - 3s 7ms/step - loss: 0.1512 - accuracy: 0.9442
Epoch 4/50
446/446 [=====] - 3s 7ms/step - loss: 0.1473 - accuracy: 0.9446
Epoch 5/50
446/446 [=====] - 3s 6ms/step - loss: 0.1482 - accuracy: 0.9444
Epoch 6/50
446/446 [=====] - 3s 7ms/step - loss: 0.1464 - accuracy: 0.9451
Epoch 7/50
446/446 [=====] - 3s 7ms/step - loss: 0.1441 - accuracy: 0.9473
Epoch 8/50

```

[] Epoch 36/50

```

446/446 [=====] - 3s 7ms/step - loss: 0.1385 - accuracy: 0.9480
Epoch 37/50
446/446 [=====] - 3s 7ms/step - loss: 0.1382 - accuracy: 0.9483
Epoch 38/50
446/446 [=====] - 3s 7ms/step - loss: 0.1378 - accuracy: 0.9480
Epoch 39/50
446/446 [=====] - 3s 7ms/step - loss: 0.1375 - accuracy: 0.9484
Epoch 40/50
446/446 [=====] - 3s 7ms/step - loss: 0.1376 - accuracy: 0.9487
Epoch 41/50
446/446 [=====] - 3s 7ms/step - loss: 0.1371 - accuracy: 0.9475
Epoch 42/50
446/446 [=====] - 3s 7ms/step - loss: 0.1371 - accuracy: 0.9486
Epoch 43/50
446/446 [=====] - 3s 7ms/step - loss: 0.1384 - accuracy: 0.9483
Epoch 44/50
446/446 [=====] - 3s 7ms/step - loss: 0.1390 - accuracy: 0.9484
Epoch 45/50
446/446 [=====] - 3s 7ms/step - loss: 0.1373 - accuracy: 0.9485
Epoch 46/50
446/446 [=====] - 3s 7ms/step - loss: 0.1363 - accuracy: 0.9486
Epoch 47/50
446/446 [=====] - 3s 7ms/step - loss: 0.1378 - accuracy: 0.9481
Epoch 48/50
446/446 [=====] - 3s 7ms/step - loss: 0.1376 - accuracy: 0.9488
Epoch 49/50
446/446 [=====] - 3s 7ms/step - loss: 0.1372 - accuracy: 0.9486
Epoch 50/50
446/446 [=====] - 3s 7ms/step - loss: 0.1377 - accuracy: 0.9487
100/100 [=====] - 0s 2ms/step - loss: 0.1527 - accuracy: 0.9382
9th Fold :
accuracy: 93.82%
-----
```

Average validation accuracy :
94.63% (+/- 0.47%)

Validation Accuracy is : **94.63% (+/-0.47%)**:-

Result of 10-fold cross validation with 60 epochs each:

94.63% (+/-0.47%)

Training accuracy: 94.83449459075928

Test accuracy: 94.44654583930969

HYPER PARAMETER TUNING

Parameters for the model like number of hidden units, number of stacked units, choice of optimizer, loss function chosen, value of Dropouts, number of epochs, folds variations and choice of metrics has been tuned to get the final model with better accuracy.

```

[ ] x_tr = x_train.to_numpy()
x_tr = np.reshape(x_tr, (x_tr.shape[0], x_tr.shape[1],1))

[ ] x_ts = x_test.to_numpy()
x_ts = np.reshape(x_ts, (x_ts.shape[0], x_ts.shape[1],1))

[ ] y_tr = pd.get_dummies(y_train)
y_ts = pd.get_dummies(y_test)

[ ] _, train_acc = cnn.evaluate(x_tr, y_tr, verbose=1)
_, test_acc = cnn.evaluate(x_ts, y_ts, verbose=1)

991/991 [=====] - 1s 2ms/step - loss: 0.1343 - accuracy: 0.9483
248/248 [=====] - 0s 2ms/step - loss: 0.1407 - accuracy: 0.9445

[ ] print("Training accuracy: "+str(train_acc*100))
print("Test accuracy: "+str(test_acc*100))

Training accuracy: 94.83449459075928
Test accuracy: 94.44654583930969

[ ] y_probs = cnn.predict(x_ts, verbose=1).ravel()
y_classes = cnn.predict_classes(x_ts, verbose=1)

248/248 [=====] - 0s 2ms/step
WARNING:tensorflow:From <ipython-input-41-ed2a250cf2a8>:2: Sequential.predict_classes (from tensorflow.python.keras.engine.sequential) is deprecated and will be removed after 2021-01-01.
Instructions for updating:
Please use instead: `np.argmax(model.predict(x), axis=-1)`, if your model does multi-class classification (e.g. if it uses a `softmax` last-layer activation).` (model.predict(x)
248/248 [=====] - 0s 2ms/step

```

PERFORMANCE METRICS (Precision, Recall, F1-Score, Support)

```

[ ] from sklearn.metrics import classification_report
print(classification_report(y_ts,y_pred))

precision    recall  f1-score   support

          0       0.92      0.97      0.95      3926
          1       0.97      0.92      0.94      3997

   micro avg       0.94      0.94      0.94      7923
   macro avg       0.95      0.94      0.94      7923
weighted avg       0.95      0.94      0.94      7923
samples avg       0.94      0.94      0.94      7923

[ ] from sklearn.preprocessing import LabelEncoder

[ ] le = LabelEncoder()

[ ] y_test_le = le.fit_transform(y_test)

[ ] from sklearn.metrics import roc_curve
fpr_keras, tpr_keras, thresholds_keras = roc_curve(y_test_le, y_classes)

[ ] from sklearn.metrics import auc
auc_keras = auc(fpr_keras, tpr_keras)

[ ] from matplotlib import pyplot

```

NOTE : THE F-SCORE VALUE NEAR 1 (~0.95) PROVES THAT THE MODEL IS BEST TRAINED, NEITHER UNDER-FITTED OR OVER-FITTED.

AUC:

0.9446646242960366

```
print(auc_keras)

0.9446646242960366
```

ROC CURVE:



TOTAL COMPUTATION TIME :

ON LOCAL MACHINE (JUPYTER NOTEBOOK) : ~ 9-10 Hours (TOTAL TIME)

ON GOOGLE COLAB: <1 Hours

Overview of results:

(10 folds and 50 epochs)

Average validation accuracy :
94.63% (+/- 0.47%)

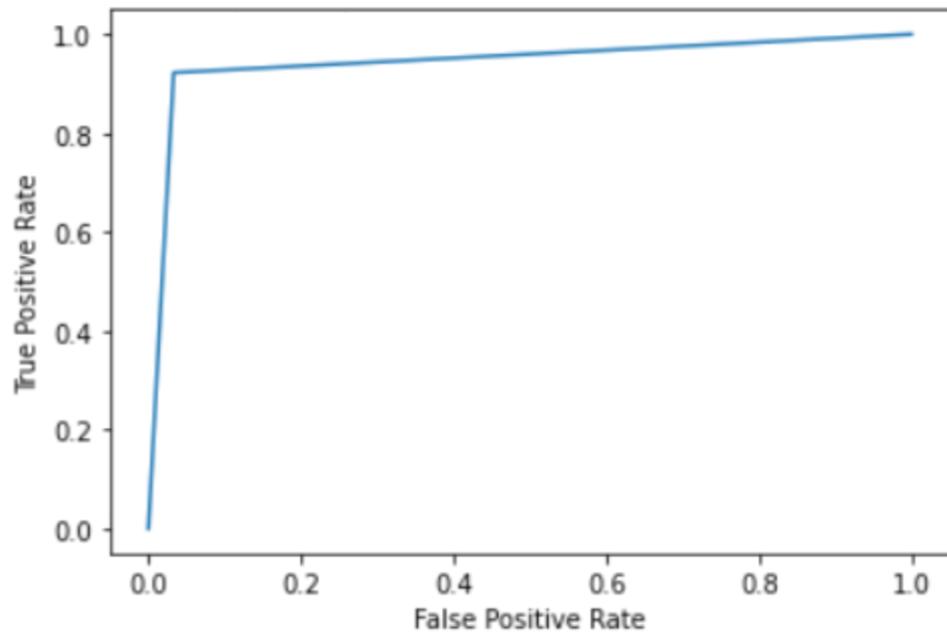
Training accuracy: 94.83449459075928
Test accuracy: 94.44654583930969

	precision	recall	f1-score	support
0	0.92	0.97	0.95	3926
1	0.97	0.92	0.94	3997
micro avg	0.94	0.94	0.94	7923
macro avg	0.95	0.94	0.94	7923
weighted avg	0.95	0.94	0.94	7923
samples avg	0.94	0.94	0.94	7923

AUC:

0.9446646242960366

```
pyplot.plot(fpr_keras, tpr_keras)
pyplot.xlabel('False Positive Rate')
pyplot.ylabel('True Positive Rate')
pyplot.show()
```



Changing the type of optimizer used along with number of epochs:

Opimizer="RMSprop"

Epochs=50

Model:-

```
cvscores = []
i=0
for train, test in kfold.split(x_train, y_train):
    cnn = Sequential()
    cnn.add(Convolution1D(64, 3, padding="same",activation="relu",input_shape=(44,1)))
    cnn.add(MaxPooling1D(pool_size=(2)))
    cnn.add(Flatten())
    cnn.add(Dense(128, activation="relu"))
    cnn.add(Dropout(0.5))
    cnn.add(Dense(2, activation="sigmoid"))
    # define optimizer and objective, compile cnn

    cnn.compile(loss="binary_crossentropy", optimizer="RMSprop",metrics=['accuracy'])

    x_tn = x_train.iloc[train]
    y_tn = y_train.iloc[train]
    y_tn = pd.get_dummies(y_tn)
    x_ts = x_train.iloc[test]
    y_ts = y_train.iloc[test]
    y_ts = pd.get_dummies(y_ts)

    x_tn1 = x_tn.to_numpy()
    x_tn1 = np.reshape(x_tn1, (x_tn1.shape[0],x_tn1.shape[1],1))

    x_ts1 = x_ts.to_numpy()
    x_ts1 = np.reshape(x_ts1, (x_ts1.shape[0],x_ts1.shape[1],1))

    cnn.fit(x_tn1, y_tn, epochs=50,batch_size=64,verbose=1)
    scores = cnn.evaluate(x_ts1, y_ts, verbose=1)
    print(str(i)+"th Fold :")
    print("%s: %.2f%%" % (cnn.metrics_names[1], scores[1]*100))
    cvscores.append(scores[1] * 100)
    i = i+1
    print("-----")

print("Average validation accuracy : ")
print("%.2f%% (+/- %.2f%%)" % (np.mean(cvscores), np.std(cvscores)))
```

Training and Testing accuracy:

```
446/446 [=====] - 3s 7ms/step - loss: 0.1443 - ac
Epoch 49/50
446/446 [=====] - 3s 7ms/step - loss: 0.1439 - ac
Epoch 50/50
446/446 [=====] - 3s 7ms/step - loss: 0.1438 - ac
100/100 [=====] - 0s 2ms/step - loss: 0.1422 - ac
9th Fold :
accuracy: 94.32%
-----
Average validation accuracy :
94.37% (+/- 0.36%)
```

```
37] print("Training accuracy: "+str(train_acc*100))
    print("Test accuracy: "+str(test_acc*100))

Training accuracy: 94.57574486732483
Test accuracy: 94.34557557106018
```

PERFORMANCE METRICS (Precision, Recall, F1-Score, Support)

```
] from sklearn.metrics import classification_report
print(classification_report(y_ts,y_pred))

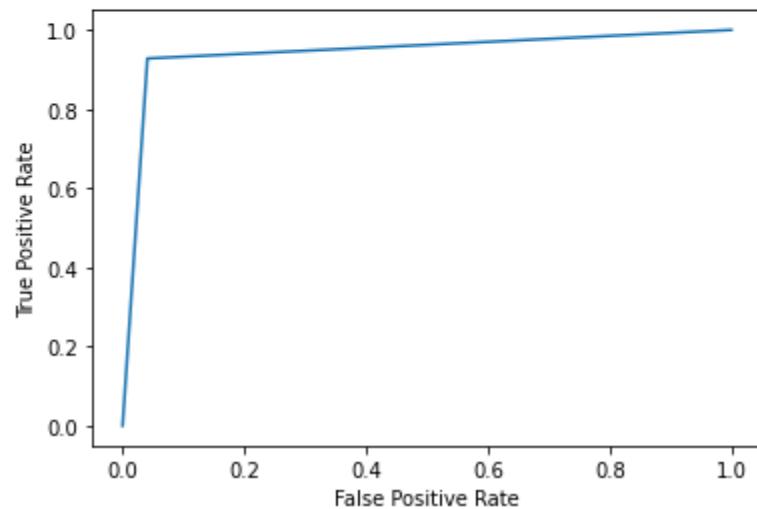
precision    recall  f1-score   support

          0       0.93      0.96      0.94     3926
          1       0.96      0.93      0.94     3997

      micro avg       0.94      0.94      0.94     7923
      macro avg       0.94      0.94      0.94     7923
  weighted avg       0.94      0.94      0.94     7923
   samples avg       0.94      0.94      0.94     7923
```

AUC and ROC CURVE:

```
| print(auc_keras)  
  
0.943596005715443  
  
| pyplot.plot(fpr_keras, tpr_keras)  
| pyplot.xlabel('False Positive Rate')  
| pyplot.ylabel('True Positive Rate')  
| pyplot.show()
```



Optimizer="RMSprop" Epochs=65

Model:-

```
cvscores = []
i=0
for train, test in kfold.split(x_train, y_train):
    cnn = Sequential()
    cnn.add(Convolution1D(64, 3, padding="same",activation="relu",input_shape=(44,1)))
    cnn.add(MaxPooling1D(pool_size=(2)))
    cnn.add(Flatten())
    cnn.add(Dense(128, activation="relu"))
    cnn.add(Dropout(0.5))
    cnn.add(Dense(2, activation="sigmoid"))
    # define optimizer and objective, compile cnn

    cnn.compile(loss="binary_crossentropy", optimizer="RMSprop",metrics=['accuracy'])

    x_tn = x_train.iloc[train]
    y_tn = y_train.iloc[train]
    y_tn = pd.get_dummies(y_tn)
    x_ts = x_train.iloc[test]
    y_ts = y_train.iloc[test]
    y_ts = pd.get_dummies(y_ts)

    x_tn1 = x_tn.to_numpy()
    x_tn1 = np.reshape(x_tn1, (x_tn1.shape[0],x_tn1.shape[1],1))

    x_ts1 = x_ts.to_numpy()
    x_ts1 = np.reshape(x_ts1, (x_ts1.shape[0],x_ts1.shape[1],1))

    cnn.fit(x_tn1, y_tn, epochs=65, batch_size=64,verbose=1)
    scores = cnn.evaluate(x_ts1, y_ts, verbose=1)
    print(str(i)+"th Fold :")
    print("%s: %.2f%%" % (cnn.metrics_names[1], scores[1]*100))
    cvscores.append(scores[1] * 100)
    i = i+1
    print("-----")

print("Average validation accuracy : ")
print("%.2f%% (+/- %.2f%%)" % (np.mean(cvscores), np.std(cvscores)))
```

Training and Testing accuracy:

```
9th Fold :  
accuracy: 93.88%  
-----  
Average validation accuracy :  
94.60% (+/- 0.41%)
```

```
| print("Training accuracy: "+str(train_acc*100))  
| print("Test accuracy: "+str(test_acc*100))
```

```
Training accuracy: 94.59152221679688  
Test accuracy: 94.59800720214844
```

PERFORMANCE METRICS (Precision, Recall, F1-Score, Support)

```
| from sklearn.metrics import classification_report  
| print(classification_report(y_ts,y_pred))
```

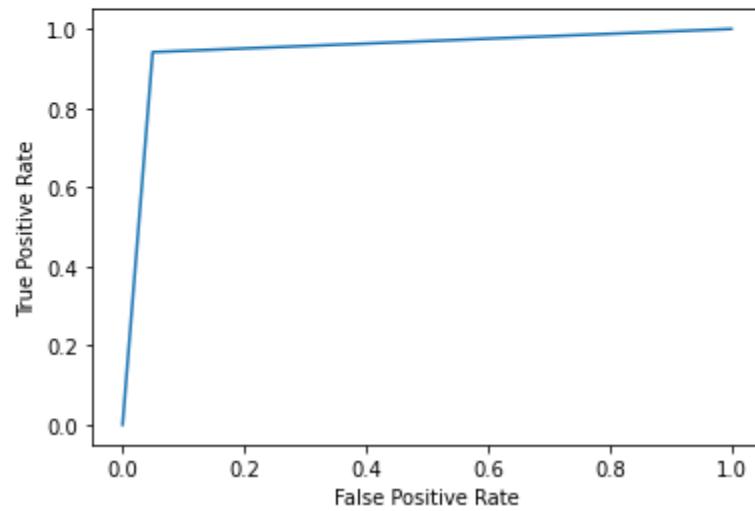
	precision	recall	f1-score	support
0	0.94	0.95	0.95	3926
1	0.95	0.94	0.95	3997
micro avg	0.95	0.95	0.95	7923
macro avg	0.95	0.95	0.95	7923
weighted avg	0.95	0.95	0.95	7923
samples avg	0.95	0.95	0.95	7923

AUC and ROC CURVE:

```
print(auc_keras)
```

```
0.946020965036054
```

```
pyplot.plot(fpr_keras, tpr_keras)
pyplot.xlabel('False Positive Rate')
pyplot.ylabel('True Positive Rate')
pyplot.show()
```



Optimizer="Adagrad"

Epochs=50

Model:-

```
cvscores = []
i=0
for train, test in kfold.split(x_train, y_train):
    cnn = Sequential()
    cnn.add(Convolution1D(64, 3, padding="same",activation="relu",input_shape=(44,1)))
    cnn.add(MaxPooling1D(pool_size=(2)))
    cnn.add(Flatten())
    cnn.add(Dense(128, activation="relu"))
    cnn.add(Dropout(0.5))
    cnn.add(Dense(2, activation="sigmoid"))
    # define optimizer and objective, compile cnn

    cnn.compile(loss="binary_crossentropy", optimizer="Adagrad",metrics=['accuracy'])

    x_tn = x_train.iloc[train]
    y_tn = y_train.iloc[train]
    y_tn = pd.get_dummies(y_tn)
    x_ts = x_train.iloc[test]
    y_ts = y_train.iloc[test]
    y_ts = pd.get_dummies(y_ts)

    x_tn1 = x_tn.to_numpy()
    x_tn1 = np.reshape(x_tn1, (x_tn1.shape[0],x_tn1.shape[1],1))

    x_ts1 = x_ts.to_numpy()
    x_ts1 = np.reshape(x_ts1, (x_ts1.shape[0],x_ts1.shape[1],1))

    cnn.fit(x_tn1, y_tn, epochs=50,batch_size=64,verbose=1)
    scores = cnn.evaluate(x_ts1, y_ts, verbose=1)
    print(str(i)+"th Fold :")
    print("%s: %.2f%%" % (cnn.metrics_names[1], scores[1]*100))
    cvscores.append(scores[1] * 100)
    i = i+1
    print("-----")

print("Average validation accuracy : ")
print("%.2f%% (+/- %.2f%%)" % (np.mean(cvscores), np.std(cvscores)))
```

Training and Testing accuracy:

```
9th Fold :  
accuracy: 94.32%  
-----  
Average validation accuracy :  
94.23% (+/- 0.58%)
```

```
| print("Training accuracy: "+str(train_acc*100))  
| print("Test accuracy: "+str(test_acc*100))
```

```
Training accuracy: 94.40219402313232  
Test accuracy: 94.53489780426025
```

PERFORMANCE METRICS (Precision, Recall, F1-Score, Support)

```
from sklearn.metrics import classification_report  
print(classification_report(y_ts,y_pred))
```

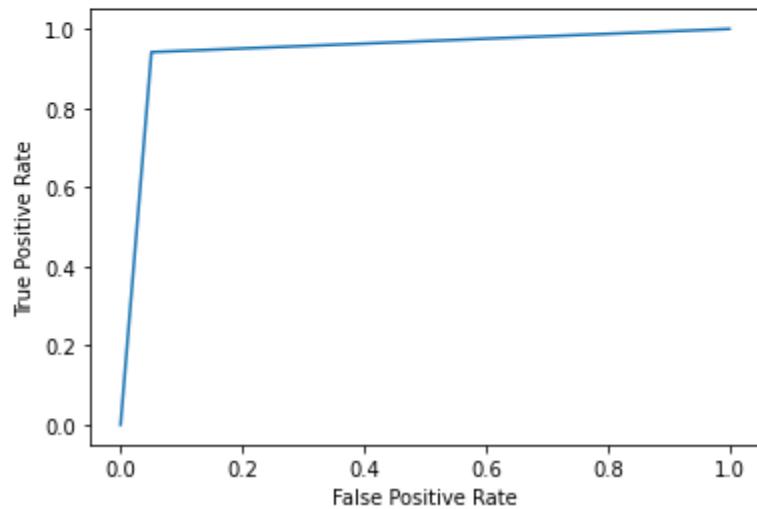
	precision	recall	f1-score	support
0	0.94	0.95	0.95	3926
1	0.95	0.94	0.95	3997
micro avg	0.95	0.95	0.95	7923
macro avg	0.95	0.95	0.95	7923
weighted avg	0.95	0.95	0.95	7923
samples avg	0.95	0.95	0.95	7923

AUC and ROC CURVE:

```
print(auc_keras)
```

```
0.9453841845979493
```

```
pyplot.plot(fpr_keras, tpr_keras)
pyplot.xlabel('False Positive Rate')
pyplot.ylabel('True Positive Rate')
pyplot.show()
```



Opimizer="Adagrad" Epochs=65

Model:-

```
cvscores = []
i=0
for train, test in kfold.split(x_train, y_train):
    cnn = Sequential()
    cnn.add(Convolution1D(64, 3, padding="same",activation="relu",input_shape=(44,1)))
    cnn.add(MaxPooling1D(pool_size=(2)))
    cnn.add(Flatten())
    cnn.add(Dense(128, activation="relu"))
    cnn.add(Dropout(0.5))
    cnn.add(Dense(2, activation="sigmoid"))
    # define optimizer and objective, compile cnn

    cnn.compile(loss="binary_crossentropy", optimizer="Adagrad",metrics=[ 'accuracy'])

    x_tn = x_train.iloc[train]
    y_tn = y_train.iloc[train]
    y_tn = pd.get_dummies(y_tn)
    x_ts = x_train.iloc[test]
    y_ts = y_train.iloc[test]
    y_ts = pd.get_dummies(y_ts)

    x_tn1 = x_tn.to_numpy()
    x_tn1 = np.reshape(x_tn1, (x_tn1.shape[0],x_tn1.shape[1],1))

    x_ts1 = x_ts.to_numpy()
    x_ts1 = np.reshape(x_ts1, (x_ts1.shape[0],x_ts1.shape[1],1))

    cnn.fit(x_tn1, y_tn, epochs=65, batch_size=64,verbose=1)
    scores = cnn.evaluate(x_ts1, y_ts, verbose=1)
    print(str(i)+"th Fold :")
    print("%s: %.2f%%" % (cnn.metrics_names[1], scores[1]*100))
    cvscores.append(scores[1] * 100)
    i = i+1
    print("-----")

print("Average validation accuracy : ")
print("%.2f%% (+/- %.2f%%)" % (np.mean(cvscores), np.std(cvscores)))
```

Training and Testing accuracy:

```
... ...
9th Fold :
accuracy: 94.76%
-----
Average validation accuracy :
94.57% (+/- 0.42%)

] print("Training accuracy: "+str(train_acc*100))
print("Test accuracy: "+str(test_acc*100))

Training accuracy: 94.55366134643555
Test accuracy: 95.01451253890991
```

PERFORMANCE METRICS (Precision, Recall, F1-Score, Support)

```
from sklearn.metrics import classification_report
print(classification_report(y_ts,y_pred))

precision    recall    f1-score    support
0            0.95     0.95     0.95     3926
1            0.95     0.95     0.95     3997

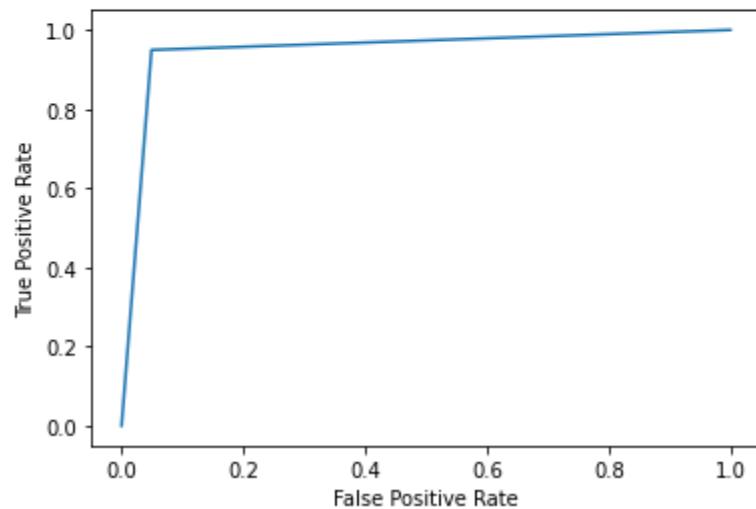
    micro avg     0.95     0.95     0.95     7923
    macro avg     0.95     0.95     0.95     7923
  weighted avg     0.95     0.95     0.95     7923
  samples avg     0.95     0.95     0.95     7923
```

AUC and ROC CURVE:

```
print(auc_keras)
```

```
0.9501513233753638
```

```
pyplot.plot(fpr_keras, tpr_keras)
pyplot.xlabel('False Positive Rate')
pyplot.ylabel('True Positive Rate')
pyplot.show()
```



Running CNN on Aman's and Shivansh's dataset with optimizer="Adam" and epoch=50:

Aman(NonDoH, Benign):

Dataset:

Z	AA	AB	AC	AD	AE	AF	AG	AH	AI	AJ	AK	AL	AM	AN	AO	AP	AQ	AR	AS	AT	AU	AV
3.450375	1.33275	1.387531	0.133807	-1.48866	-0.1352	-0.28203	-0.35867	0.23029	-0.10095	0.456427	-1.06931	-1.42116	-0.82871	-0.11313	-0.12664	-0.05018	0.140587	-0.04807	0.179427	0.691027	0.091464	
-0.1633	-0.65795	-0.1806	-1.95845	-1.13494	-0.85858	-0.28469	-0.45267	-0.40968	-0.31316	-0.12911	-1.50409	0.328504	-0.46039	-0.11348	-0.1828	0.166193	0.322024	0.506174	-1.44439	-1.4897	-0.5275	NonDoH
-0.2898	0.392213	-0.23887	0.916266	0.305363	0.355579	-0.28473	-0.4607	-0.41677	-0.32717	-0.12911	-0.54841	0.867792	-0.80344	-0.11348	-0.18273	-0.133	-0.0559	0.009551	0.838991	0.692224	-0.52518	NonDoH
-0.2898	-0.42536	-0.19392	0.360076	0.238742	0.107735	-0.27802	-0.29691	-0.34145	-0.32198	-0.0774	0.919694	-0.57546	1.265228	-0.11347	-0.17298	-0.14424	-0.05677	-0.04806	0.264739	0.727377	0.010519	NonDoH
-0.2898	0.018507	-0.25385	0.308179	-0.00277	1.328569	-0.28474	-0.46255	-0.42152	-0.33332	-0.12911	0.176719	0.465554	-0.57207	-0.11348	-0.17881	-0.15147	-0.0749	-0.03969	0.349409	0.778978	-0.21312	NonDoH
-0.2898	-0.02211	-0.16729	0.547834	0.201478	0.49324	-0.28473	-0.46205	-0.4194	-0.33061	-0.12911	-0.20837	0.81025	-0.77665	-0.11347	-0.17493	-0.14505	-0.05945	-0.04783	0.226929	0.782825	-0.08087	NonDoH
-0.2898	0.155428	-0.28714	0.775609	0.217716	0.524299	-0.26629	-0.18548	-0.28201	-0.32899	-0.12911	0.1099388	-0.43139	1.19843	-0.11343	-0.16253	-0.13129	-0.05567	-0.04814	0.706986	0.657305	0.190419	NonDoH
-0.2568	-0.68215	-0.2938	0.924183	0.255064	-0.92773	-0.28473	-0.46134	-0.42231	-0.33374	-0.12911	-0.0383	-0.04398	0.010963	-0.11348	-0.1828	-0.05285	0.045362	0.142751	-1.44439	-1.4897	-0.5275	NonDoH
-0.28705	-0.72956	-0.31211	-0.48773	0.577172	-1.07042	-0.28474	-0.46544	-0.42597	-0.33731	-0.12911	-0.03522	-0.05055	0.02271	-0.11348	-0.1828	-0.14708	-0.07366	-0.01365	-1.44439	-1.4897	-0.5275	NonDoH
0.14955	0.636262	-0.16396	-1.13494	-0.93075	-0.28474	-0.46496	-0.42503	-0.33555	-0.12911	-1.50372	0.328504	-0.46039	-0.11348	-0.1828	-0.13401	-0.05715	0.00809	-1.44439	-1.4897	-0.5275	NonDoH	
-0.2568	-0.2448	-0.2938	0.323738	-0.07835	1.424362	6.985	5.108939	4.491384	2.843157	-0.12911	0.660287	-0.05121	0.023903	-0.11346	-0.17235	-0.1368	-0.05363	-0.04814	0.472266	0.794226	-0.09182	NonDoH
-0.2568	1.844555	-0.2938	1.0992	0.339457	0.374174	2.012322	2.667616	0.688238	-0.32061	-0.12911	0.788583	-0.59685	2.465835	-0.11346	-0.17179	-0.15657	-0.09685	-0.04813	0.987668	0.566348	0.731486	NonDoH
-0.1633	-0.49101	-0.1804	0.134242	0.12585	0.059821	0.942791	1.824645	2.982883	1.772787	-0.12911	1.238119	0.577147	-0.64835	-0.11347	-0.17375	-0.13484	-0.0528	-0.04797	0.495669	0.87214	-0.17263	NonDoH
-0.28705	-0.72956	-0.31211	-0.48773	0.577172	-1.07042	-0.28474	-0.46481	-0.42541	-0.33676	-0.12911	-0.03522	-0.05055	0.02271	-0.11348	-0.1828	-0.13251	-0.05526	0.01058	-1.44439	-1.4897	-0.5275	NonDoH
-0.2733	-0.226608	-0.20724	0.303247	0.043836	0.838338	-0.28473	-0.46116	-0.41941	-0.33222	-0.10304	0.565121	-1.165555	-0.59865	-0.11345	-0.16693	-0.10976	-0.03766	0.009236	0.891124	0.619569	-0.17336	NonDoH
-0.2568	-0.01889	-0.26717	0.191587	-0.08661	1.771481	-0.28471	-0.4546	-0.40538	-0.3192	-0.12911	0.36024	0.921608	-0.82691	-0.1124	-0.15575	-0.07223	0.04492	-0.0376	0.406871	0.823194	-0.16072	NonDoH
-0.28705	-0.72956	-0.31211	-0.48773	0.577172	-1.07042	-0.28474	-0.46324	-0.42402	-0.33542	-0.12911	-0.03522	-0.05055	0.02271	-0.11348	-0.1828	-0.13401	-0.05715	0.00809	-1.44439	-1.4897	-0.5275	NonDoH
-0.28705	-0.72956	-0.31211	-0.48773	0.577172	-1.07042	-0.28474	-0.46544	-0.42597	-0.33731	-0.12911	-0.03522	-0.05055	0.02271	-0.11348	-0.1828	-0.1471	-0.07368	-0.01363	-1.44439	-1.4897	-0.5275	NonDoH
-0.2568	1.858754	-0.2938	0.413568	0.008846	1.562845	1.355003	2.181515	0.816208	-0.30867	-0.12911	1.042655	-0.47904	1.475064	-0.11345	-0.16849	-0.14263	-0.04525	-0.04812	0.192615	0.652523	0.206642	NonDoH
-0.2898	-0.72183	-0.31378	0.02051	0.291957	-0.92831	-0.20052	0.133753	0.565388	0.902131	-0.12911	-1.14447	0.748218	-0.746	-0.11347	-0.1748	-0.14916	-0.06447	-0.04809	0.227816	0.719673	0.025777	NonDoH
-0.2898	0.637818	-0.20724	0.558092	0.12311	0.91355	-0.25902	-0.13476	-0.22242	-0.27876	-0.12911	0.982247	-0.32683	0.736022	-0.11244	-0.08621	-0.00812	-0.06438	-0.04812	1.201676	0.639396	0.256862	NonDoH
-0.2898	-0.13268	-0.17394	0.333959	0.089386	0.794669	0.207114	0.983843	0.234223	-0.33452	-0.12911	1.031507	-0.49317	1.567826	8.210831	8.476549	5.704928	-0.07304	-0.0481	0.937475	0.525254	1.385799	NonDoH
-0.2568	-0.09346	-0.26717	0.41863	0.028668	0.964085	-0.28467	-0.44953	-0.40464	-0.306	-0.12911	-1.55512	0.39805	-0.51987	-0.11346	-0.17166	-0.14935	-0.07687	-0.04814	0.686859	0.640792	0.251057	NonDoH
-0.2733	0.058149	-0.1806	0.525207	0.149009	0.638995	-0.28464	-0.44564	-0.41142	-0.32979	-0.11615	0.739657	-0.46763	0.338204	-0.11347	-0.17511	-0.13707	-0.05697	-0.04772	0.517364	0.911086	-0.204	NonDoH
-0.2898	-0.21518	-0.2938	0.591899	0.132005	0.586714	-0.28463	-0.44484	-0.39287	-0.3091	-0.12911	0.412681	0.668155	-0.70284	-0.11348	-0.18276	0.014311	0.130211	0.253733	0.486308	0.756486	-0.52727	NonDoH
-0.2898	0.463518	-0.31378	0.21739	-0.08069	2.105162	2.916401	3.233298	2.039795	-0.27051	-0.12911	1.48961	-0.27453	0.557258	-0.11338	-0.15333	-0.06256	0.061539	-0.04814	0.383663	0.864095	-0.16458	NonDoH
-0.2568	-0.32143	-0.2938	0.275877	-0.05755	1.142643	4.770838	4.183012	3.507186	2.259715	-0.12911	0.604228	-0.08854	0.094142	-0.11347	-0.17839	-0.14087	-0.06439	-0.04808	0.579454	1.162102	-0.31594	NonDoH
-0.28705	-0.72956	-0.31211	-0.48773	0.577172	-1.07042	-0.28474	-0.46474	-0.42534	-0.3367	-0.12911	-0.03522	-0.05055	0.02271	-0.11348	-0.1828	-0.13085	-0.05316	0.013333	-1.44439	-1.4897	-0.5275	NonDoH

```

cvscores = []
i=0
for train, test in kfold.split(x_train, y_train):
    cnn = Sequential()
    cnn.add(Convolution1D(64, 3, padding="same",activation="relu",input_shape=(44,1)))
    cnn.add(MaxPooling1D(pool_size=(2)))
    cnn.add(Flatten())
    cnn.add(Dense(128, activation="relu"))
    cnn.add(Dropout(0.5))
    cnn.add(Dense(2, activation="sigmoid"))
    # define optimizer and objective, compile cnn

    cnn.compile(loss="binary_crossentropy", optimizer="adam",metrics=['accuracy'])

    x_tn = x_train.iloc[train]
    y_tn = y_train.iloc[train]
    y_tn = pd.get_dummies(y_tn)
    x_ts = x_train.iloc[test]
    y_ts = y_train.iloc[test]
    y_ts = pd.get_dummies(y_ts)

    x_tn1 = x_tn.to_numpy()
    x_tn1 = np.reshape(x_tn1, (x_tn1.shape[0],x_tn1.shape[1],1))

    x_ts1 = x_ts.to_numpy()
    x_ts1 = np.reshape(x_ts1, (x_ts1.shape[0],x_ts1.shape[1],1))

    cnn.fit(x_tn1, y_tn, epochs=50,batch_size=64,verbose=1)
    scores = cnn.evaluate(x_ts1, y_ts, verbose=1)
    print(str(i)+"th Fold :")
    print("%s: %.2f%%" % (cnn.metrics_names[1], scores[1]*100))
    cvscores.append(scores[1] * 100)
    i = i+1
    print("-----")

print("Average validation accuracy : ")

```

Training and Testing accuracy:

```
446/446 [=====] - 3s 6ms/step - loss: 0.3365 - accuracy: 0.83/3
Epoch 47/50
446/446 [=====] - 3s 6ms/step - loss: 0.3369 - accuracy: 0.8369
Epoch 48/50
446/446 [=====] - 3s 6ms/step - loss: 0.3373 - accuracy: 0.8374
Epoch 49/50
446/446 [=====] - 3s 6ms/step - loss: 0.3370 - accuracy: 0.8367
Epoch 50/50
446/446 [=====] - 3s 7ms/step - loss: 0.3348 - accuracy: 0.8369
100/100 [=====] - 0s 2ms/step - loss: 0.3113 - accuracy: 0.8571
9th Fold :
accuracy: 85.71%
-----
Average validation accuracy :
83.93% (+/- 1.00%)
```

```
[48] print("Training accuracy: "+str(train_acc*100))
     print("Test accuracy: "+str(test_acc*100))
```

```
↳ Training accuracy: 84.1185212135315
    Test accuracy: 83.99596214294434
```

PERFORMANCE METRICS (Precision, Recall, F1-Score, Support)

```
[1] from sklearn.metrics import classification_report
    print(classification_report(y_ts,y_pred))
```

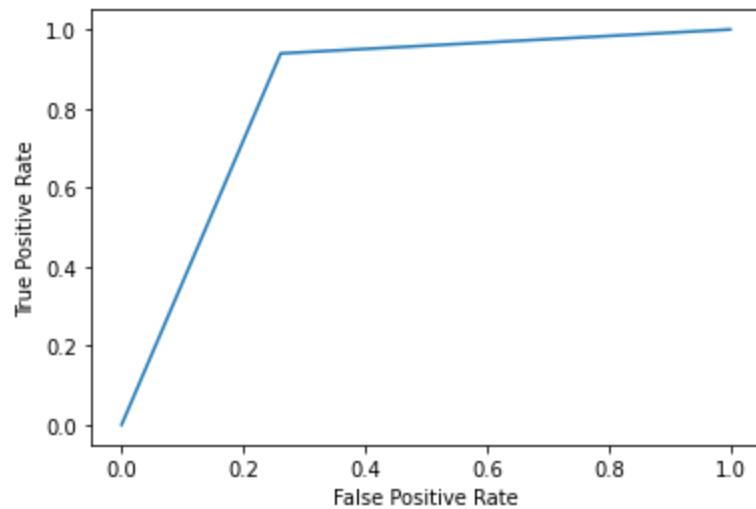
	precision	recall	f1-score	support
0	0.92	0.74	0.82	3926
1	0.79	0.94	0.86	3997
micro avg	0.84	0.84	0.84	7923
macro avg	0.85	0.84	0.84	7923
weighted avg	0.85	0.84	0.84	7923
samples avg	0.84	0.84	0.84	7923

AUC and ROC CURVE:

```
[58] print(auc_keras)
```

```
0.8390599495724697
```

```
[59] pyplot.plot(fpr_keras, tpr_keras)
pyplot.xlabel('False Positive Rate')
pyplot.ylabel('True Positive Rate')
pyplot.show()
```



Shivansh(DoH, Benign):

Dataset:

U	V	W	X	Y	Z	AA	AB	AC	AD	AE	AF	AG	AH	AI	AJ	AK	AL	AM	AN	AO	AP	AQ
<code>length_miLength_miLength_mi skew_mi skew_mi med skew_mocLength_co Time_var Time_std Time_med Time_mod skew_mi med skew_moc cov_time skew_mi med res_time res_time res_time res_time res_time res_skew res_cov ti Label</code>																						
.062383	-0.60254	-0.19077	0.717028	0.187329	0.794178	-0.55644	-0.85171	-0.64247	-0.181	1.254661	-0.40539	0.327724	1.254661	-0.14549	-0.24009	-0.20009	-0.1529	-0.11249	0.246676	0.488486	-0.25187 DoH	
.682527	-0.60254	-0.19077	0.494456	0.091134	1.044114	-0.31955	-0.00296	-0.64133	-0.181	1.17284	-0.70034	1.218502	1.17284	-0.1455	-0.242321	-0.20296	-0.15359	-0.1125	-0.32153	0.497271	-0.26058 DoH	
.417151	-0.72591	-0.19077	0.349595	0.012761	2.044237	-0.55651	-0.86352	-0.6416	-0.16704	0.226964	2.81537	-0.80369	0.226964	-0.14549	-0.24035	-0.19891	-0.14972	-0.11239	-0.16168	0.564216	-0.31047 DoH	
.682527	-0.60254	-0.19077	0.49471	0.091254	1.043262	-0.31391	0.007251	-0.64126	-0.181	1.178104	-0.69748	1.206346	1.178104	-0.1455	-0.24277	-0.20291	-0.15381	-0.11257	-0.20427	0.472583	-0.23483 DoH	
0.58539	0.353578	-0.19077	0.04301	0.258015	-0.4018	-0.2336	0.141517	0.252861	-0.181	-0.04316	0.255948	-0.52549	0.04316	-0.1455	-0.24277	-0.20304	-0.15383	-0.11258	-0.23061	0.46301	-0.22495 DoH	
.820781	-0.60254	-0.19077	0.289844	-0.00741	2.076458	-0.31163	0.011332	-0.64184	-0.181	1.240541	-0.66568	1.077417	1.240541	-0.14549	-0.24083	-0.20312	-0.15383	-0.11257	0.209601	0.461406	-0.22383 DoH	
.682124	-0.60254	-0.19077	0.494789	0.091296	1.042808	-0.32449	-0.01201	-0.64197	-0.181	1.167933	-0.70442	1.235156	1.167933	-0.1455	-0.24206	-0.20171	-0.15367	-0.11256	0.070198	0.51039	-0.2676 DoH	
.935822	5.01081	-0.19077	-0.31709	0.491806	-0.23653	-0.08648	0.349619	0.702179	-0.181	-0.12868	0.602667	-0.74915	-0.12868	-0.14258	-0.10296	-0.15625	-0.14879	-0.09646	0.689661	0.165507	0.400294 DoH	
0.75004	-0.60254	-0.19077	0.506804	0.096371	1.033396	-0.31916	-0.00226	-0.64197	-0.181	1.201185	-0.68665	1.161194	1.201185	-0.1455	-0.24317	-0.20327	-0.15384	-0.11257	-0.3252	0.472969	-0.23505 DoH	
.570402	-0.60254	-0.19077	0.468103	0.0797	1.079515	-0.30044	0.031174	-0.64076	-0.181	1.223954	-0.67235	1.10352	1.223954	-0.14549	-0.23792	-0.20147	-0.15388	-0.11257	0.22364	0.342012	-0.05943 DoH	
.918549	-0.32496	-1.07705	0.517806	0.149749	0.934837	-0.37334	-0.10727	-0.64272	-0.17758	0.91514	-0.84255	1.95586	0.91514	-0.1455	-0.24396	-0.20274	-0.15374	-0.11198	-0.30244	0.548701	-0.32485 DoH	
.735221	5.01081	-0.19077	-0.60053	0.805161	-0.61599	1.005191	1.350061	1.659632	-0.181	-0.06169	0.515481	-0.70047	-0.06169	-0.1455	-0.24425	-0.20728	-0.16045	-0.11235	0.608367	0.140659	0.817417 DoH	
.592327	5.01081	-0.19077	-0.83545	0.839773	-0.66055	-0.3516	-0.06348	0.237313	-0.181	-0.16668	0.566034	-0.72922	-0.16668	-0.14549	-0.24007	-0.19984	-0.15375	-0.10329	0.447304	0.090111	-0.26279 DoH	
0.27561	0.862481	-0.19077	0.26016	0.517399	-0.58544	1.008678	1.352534	1.738546	-0.181	-0.17916	0.515319	-0.70038	-0.17916	-0.1288	0.099871	-0.0601	-0.16062	-0.11257	0.917709	0.240561	0.235049 DoH	
0.51907	0.353578	-0.19077	0.168241	0.31502	-0.44966	1.001326	1.347317	1.710988	0.957554	-0.10698	0.173744	-0.71214	-0.10698	-0.14548	-0.23589	-0.20249	-0.16061	-0.11256	0.96594	0.25681	0.167614 DoH	
.619923	-0.60254	-0.19077	0.47875	0.084267	1.067583	-0.32422	-0.01153	-0.64194	-0.181	1.142958	-0.71742	1.293208	1.142958	-0.1455	-0.2432	-0.20327	-0.15383	-0.11252	-0.32965	0.472794	-0.23783 DoH	
0.92906	0.291893	2.15155	-1.65098	-0.89916	-1.14124	3.822928	3.845217	1.707772	-0.181	0.036843	-0.16507	-0.08688	0.036843	-0.1455	-0.24877	-0.20067	-0.10319	0.373743	-0.63931 DoH			
0.63261	0.353578	-0.19077	0.066902	0.485785	-0.7536	1.056085	1.385886	1.801238	2.613076	-0.42655	-0.46734	-0.63253	-0.42655	-0.14461	-0.16797	-0.20097	-0.16061	-0.11256	0.440022	0.08319	3.384907 DoH	
.682527	-0.60254	-0.19077	0.49471	0.091254	1.043262	-0.31455	0.006093	-0.64186	-0.181	1.177784	-0.69887	1.212274	1.177784	-0.1455	-0.24305	-0.20301	-0.15382	-0.11256	-0.25065	0.48556	-0.24672 DoH	
.506814	-0.60254	-0.19077	0.479138	0.085711	1.007474	-0.32753	-0.01763	-0.6394	-0.181	1.091652	-0.73917	1.393831	1.091652	-0.14548	-0.23492	-0.20058	-0.1538	-0.11257	0.324046	0.30711	0.017023 DoH	
.605377	-0.60254	-0.19077	0.466235	0.068871	1.682257	-0.54864	-0.709	-0.64073	-0.181	1.61088	-0.43082	0.383461	1.61088	-0.1455	-0.24349	-0.19696	-0.15213	-0.10219	0.82119	0.297117	-0.46812 DoH	
1.07465	-0.60254	-0.19077	0.598605	0.136955	0.886715	-0.55577	-0.81796	-0.59805	-0.181	-0.13805	0.334659	-0.58405	-0.13805	-0.14063	-0.06044	-0.12978	-0.14882	-0.11258	0.805138	0.236885	0.251437 DoH	
0.77985	-0.72591	-0.19077	0.564582	0.113024	1.027083	-0.55648	-0.8563	-0.64229	-0.181	1.555399	-0.09555	-0.17862	1.555399	-0.14549	-0.23939	-0.20354	-0.15868	-0.11257	0.720013	0.285742	0.073505 DoH	
.506814	-0.60254	-0.19077	0.479138	0.085711	1.007474	-0.33295	-0.02774	-0.642	-0.181	1.088073	-0.74653	1.429361	1.088073	-0.1455	-0.2429	-0.20375	-0.15381	-0.11258	-0.42187	0.415211	-0.17219 DoH	
.560982	-0.60254	-0.19077	0.493838	0.092126	0.987025	-0.3357	-0.0329	-0.64183	-0.181	1.106607	-0.73629	1.380174	1.106607	-0.1455	-0.24277	-0.20366	-0.15382	-0.11257	-0.38225	0.414074	-0.17155 DoH	
0.11987	2.358347	-0.19077	-0.56053	0.831372	-0.78912	0.996312	1.343752	1.700547	-0.181	-0.03272	0.569635	-0.73121	-0.03272	-0.1455	-0.2438	-0.19458	-0.14881	-0.09645	0.46043	0.090235	-0.50625 DoH	
.729418	0.106838	-0.19077	0.800895	0.260773	0.546152	0.346465	0.819007	1.342227	-0.181	-0.2749	0.653874	-0.77582	-0.2749	-0.14393	-0.14193	-0.17883	-0.14888	-0.11257	0.527366	0.170463	0.703526 DoH	
.062383	-0.60254	-0.19077	0.717028	0.187329	0.794178	-0.55617	-0.83379	-0.64154	-0.181	1.334147	-0.45392	0.436659	1.334147	-0.14548	-0.2371	-0.1964	-0.14847	-0.11251	0.10414	0.523771	-0.27859 DoH	

```

| cvscores = []
i=0
for train, test in kfold.split(x_train, y_train):
    cnn = Sequential()
    cnn.add(Convolution1D(64, 3, padding="same",activation="relu",input_shape=(41,1)))
    cnn.add(MaxPooling1D(pool_size=(2)))
    cnn.add(Flatten())
    cnn.add(Dense(128, activation="relu"))
    cnn.add(Dropout(0.5))
    cnn.add(Dense(2, activation="sigmoid"))
    # define optimizer and objective, compile cnn

    cnn.compile(loss="binary_crossentropy", optimizer="adam",metrics=['accuracy'])

    x_tn = x_train.iloc[train]
    y_tn = y_train.iloc[train]
    y_tn = pd.get_dummies(y_tn)
    x_ts = x_train.iloc[test]
    y_ts = y_train.iloc[test]
    y_ts = pd.get_dummies(y_ts)

    x_tn1 = x_tn.to_numpy()
    x_tn1 = np.reshape(x_tn1, (x_tn1.shape[0],x_tn1.shape[1],1))

    x_ts1 = x_ts.to_numpy()
    x_ts1 = np.reshape(x_ts1, (x_ts1.shape[0],x_ts1.shape[1],1))

    cnn.fit(x_tn1, y_tn, epochs=50,batch_size=64,verbose=1)
    scores = cnn.evaluate(x_ts1, y_ts, verbose=1)
    print(str(i)+"th Fold :")
    print("%s: %.2f%%" % (cnn.metrics_names[1], scores[1]*100))
    cvscores.append(scores[1] * 100)
    i = i+1
    print("-----")

print("Average validation accuracy : ")
print("%.2f%% (+/- %.2f%%)" % (np.mean(cvscores), np.std(cvscores)))

```

Training and Testing accuracy:

```
446/446 [=====] - 3s 7ms/step - loss: 0.0000 - acc: 1.0000
Epoch 50/50
446/446 [=====] - 3s 7ms/step - loss: 0.0000 - acc: 1.0000
100/100 [=====] - 0s 2ms/step - loss: 0.0000 - acc: 1.0000
9th Fold :
accuracy: 87.85%
-----
Average validation accuracy :
88.71% (+/- 0.55%)
```

```
print("Training accuracy: "+str(train_acc*100))
print("Test accuracy: "+str(test_acc*100))
```

```
Training accuracy: 88.70972990989685
Test accuracy: 88.24940323829651
```

PERFORMANCE METRICS (Precision, Recall, F1-Score, Support)

```
from sklearn.metrics import classification_report
print(classification_report(y_ts,y_pred))
```

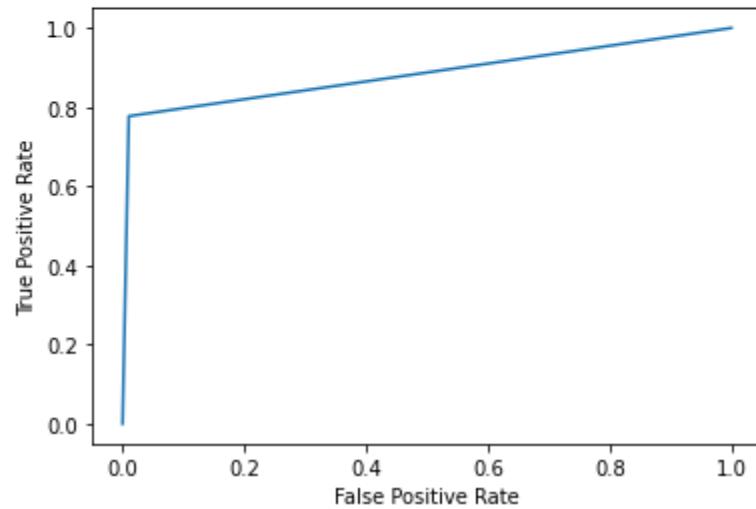
	precision	recall	f1-score	support
0	0.81	0.99	0.89	3926
1	0.99	0.78	0.87	3997
micro avg	0.88	0.88	0.88	7923
macro avg	0.90	0.88	0.88	7923
weighted avg	0.90	0.88	0.88	7923
samples avg	0.88	0.88	0.88	7923

AUC and ROC CURVE:

```
print(auc_keras)
```

```
0.8834471625497013
```

```
pyplot.plot(fpr_keras, tpr_keras)
pyplot.xlabel('False Positive Rate')
pyplot.ylabel('True Positive Rate')
pyplot.show()
```



Results and comparison :

1) Dataset used here for comparison consists of Malicious and Benign data:

<u>OPTIMIZER</u>	<u>EPOCH</u>	<u>TRAINING ACCURACY (%)</u>	<u>TEST ACCURACY (%)</u>
Adam	50	94.83449459075928	94.44654583930969
RMSprop	50	94.57574486732483	94.34557557106018
RMSprop	65	94.59152221679688	94.59800720214844
Adagrad	50	94.40219402313232	94.53489780426025
Adagrad	65	94.55366134643555	95.01451253890991

Feature/ Name	OPTIMIZER Adam	OPTIMIZER RMSprop	OPTIMIZER RMSprop	OPTIMIZER Adagrad	OPTIMIZER Adagrad
Precision	0.92	0.93	0.94	0.94	0.95
Recall	0.97	0.96	0.95	0.95	0.95
f-Score	0.95	0.94	0.95	0.95	0.95
AUC Score	0.9446646 242960366	0.9435960 05715443	0.9460209 65036054	0.9453841 845979493	0.9501513 233753638
Avg. Validation Accuracy	94.63% (+/- 0.47%)	94.37% (+/- 0.36%)	94.60% (+/- 0.41%)	94.23% (+/- 0.58%)	94.57% (+/- 0.42%)
Training Accuracy	94.834494 59075928	94.575744 86732483	94.591522 21679688	94.402194 02313232	94.55366 134643555
Testing Accuracy	94.446545 83930969	94.345575 57106018	94.59800 720214844	94.53489 780426025	95.014512 53890991
Computation Time	1533.566s	1513.191s	2026.284s	1627.167s	1996.546s

- 2) Dataset here used are all four datasets i.e. Malicious, Benign, NonDoH and NonDoH on CNN model.

Feature/Name	Shivansh (DoH + Benign)	Yuvraj (Malicious + Benign)	Aman (Non-DoH + Benign)
Precision	0.81	0.92	0.92
Recall	0.99	0.97	0.74
f-Score	0.89	0.95	0.82
AUC Score	0.8834471625497013	0.9446646242960366	0.8390599495724697
Avg. Validation Accuracy	88.71% (+/- 0.55%)	94.63% (+/- 0.47%)	83.93% (+/- 1.00%)
Training Accuracy	88.70972990989685	94.83449459075928	84.1185212135315
Testing Accuracy	88.24940323829651	94.44654583930969	83.99596214294434
Computation Time	1530.099s	1533.566s	1466.379s

CONCLUSION:

I have successfully made a deep neural network consisting of stacked Convolution1D and MaxPooling1D layers along with appropriate change in optimizers and epochs. Appropriate preprocessing & sampling was done to avoid biasness in the model & fed the data to algorithm for training. The model has been validated using 10-fold cross validation and hyper-parameter tunning was done to improve the accuracy.

I found that Adagrad optimizer gave the best results on the same dataset and same number of epoch which is 50. Also when the number of epoch increased to 65, the accuracy shown a great increase in accuracy which even crossed 95% accuracy.

Apart from trying combinations on my dataset (malicious and benign), I tried CNN model on dataset of Aman (Non-DoH and benign) and Shivansh (DoH and Benign). I found that CNN model performed exceptionally well on dataset for Malicious and Benign.

Link for video demonstration:

<https://drive.google.com/file/d/1oMzIOdPpv9-bDKQuQuxsSugWBMRM143Z/view?usp=sharing>

Link to my Github project repository
(Network-Intrusion-Detection-using-
Deep-Learning):

<https://github.com/yuvrajkumarViT/Network-Intrusion-Detection-using-Deep-Learning>



Yuvrajkumar
(18BIT0276)

Student's Signature