

INTRUSION DETECTION SYSTEM USING DEEP LEARNING

*Project Review 2 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:

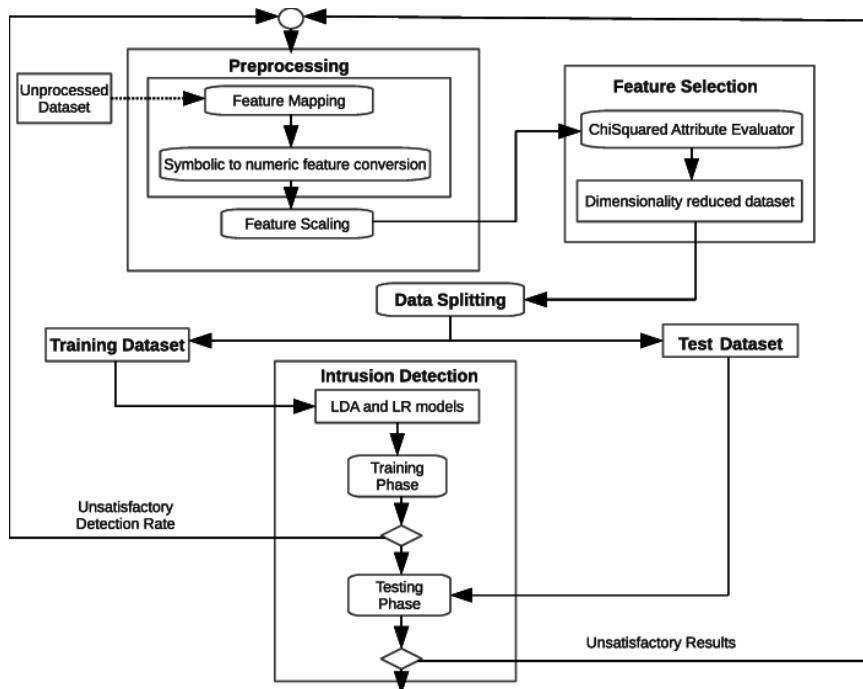
- i) signature (misuse) based NIDS (SNIDS),
- ii) anomaly detection-based NIDS (ADNIDS).

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

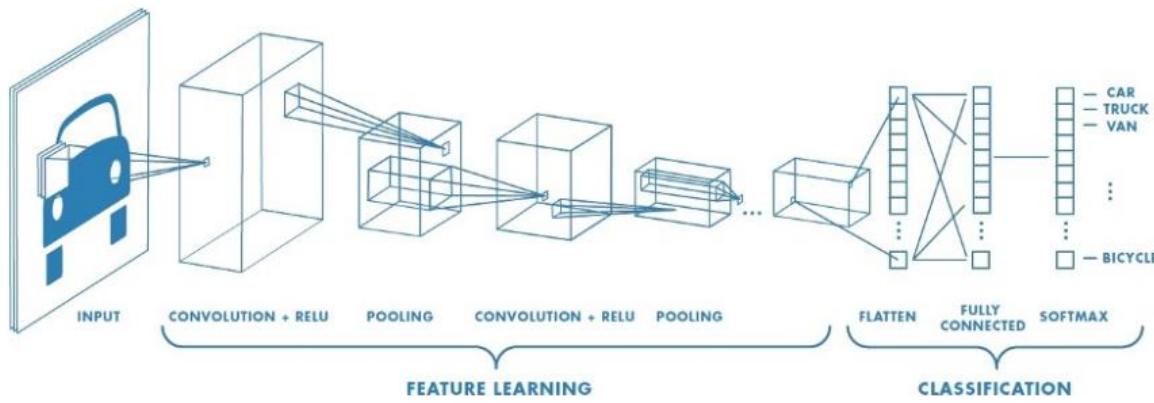


A Network Intrusion Detection System (NIDS) helps system administrators to detect network security breaches in their organizations. However, many challenges arise while developing a flexible and efficient NIDS for unforeseen and unpredictable attacks. We propose a deep learning based approach for developing 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.

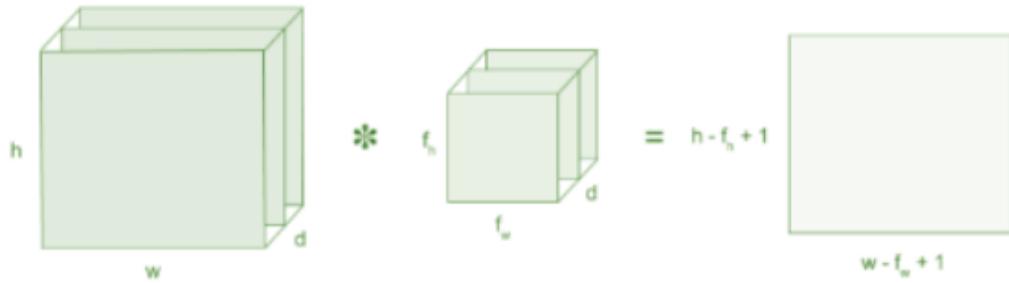
Convolutional Neural Network (CNN) — Deep Learning

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.

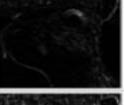
Technically, deep learning CNN models to train and test, each input image will pass it through a series of convolution layers with filters (Kernels), Pooling, fully connected layers (FC) and apply Softmax function to classify an object with probabilistic values between 0 and 1. The below figure is a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

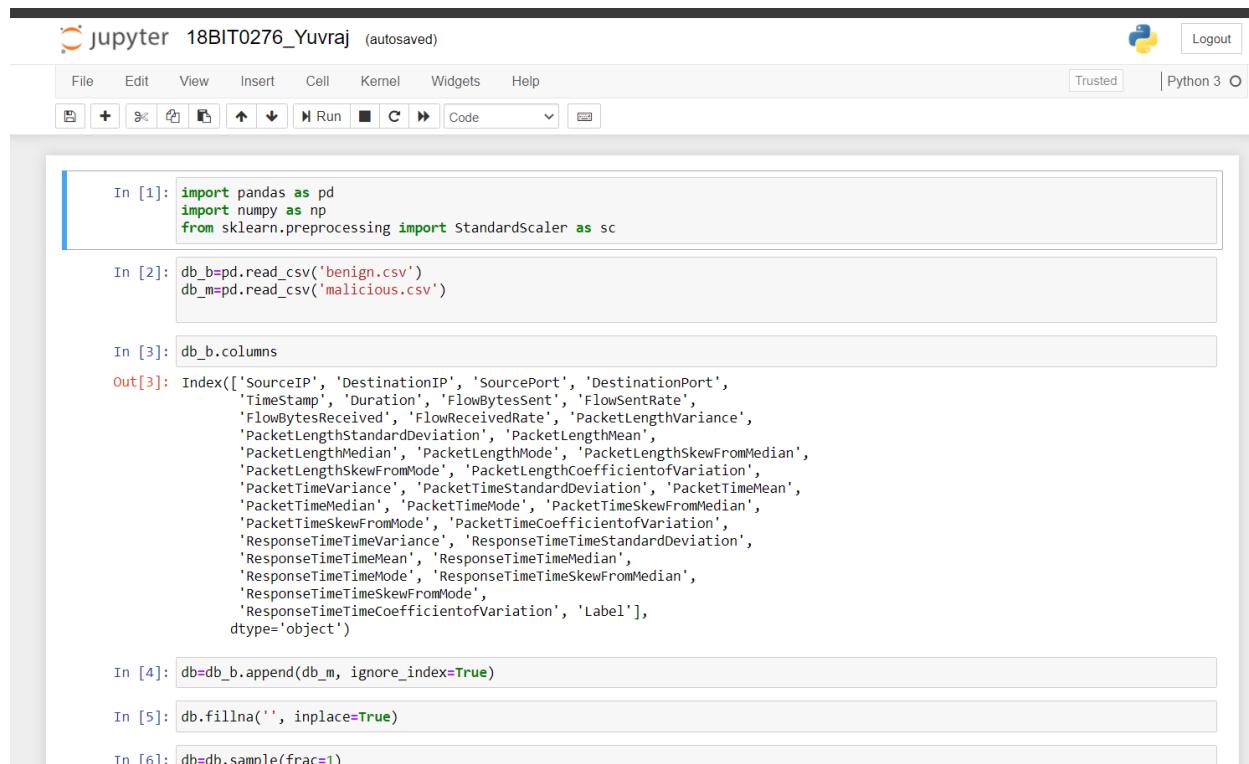
**Sampling is done before model creation. ->Pg.26

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 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',
       'PacketLengthSkewFromMode', 'PacketLengthCoefficientOfVariation',
       'PacketTimeVariance', 'PacketTimeStandardDeviation', 'PacketTimeMean',
       'PacketTimeMedian', 'PacketTimeMode', 'PacketTimeSkewFromMedian',
       'PacketTimeSkewFromMode', 'PacketTimeCoefficientOfVariation',
       'ResponseTimeVariance', 'ResponseTimeStandardDeviation',
       'ResponseTimeMean', 'ResponseTimeMedian',
       'ResponseTimeMode', 'ResponseTimeSkewFromMedian',
       'ResponseTimeSkewFromMode',
       '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)
```

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

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\database_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]: (26963, 35)

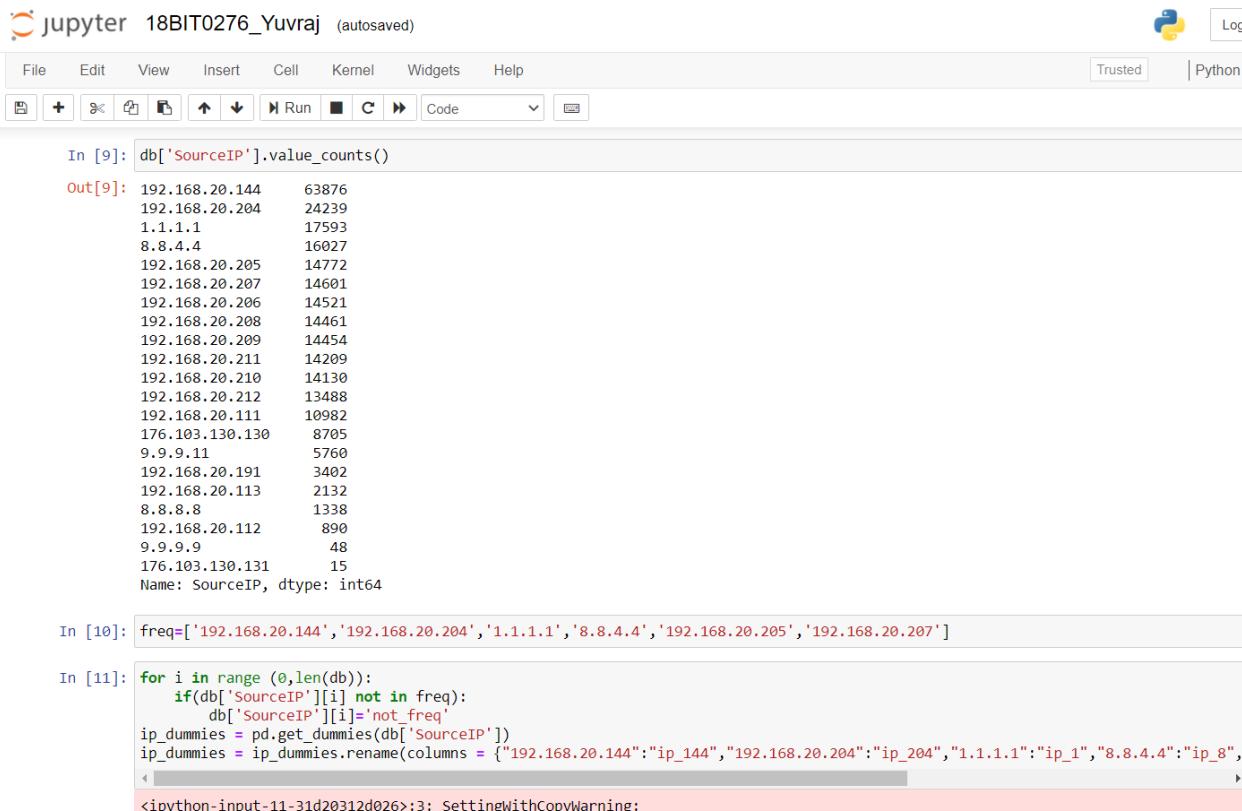
Now, the dataset basically had 2 kind of features, namely categorical and numeric. Categorical features comprised SourceIP, DestinationIP, SourcePort, DestinationPort etc. and numerical features included Duration, FlowSentRate, FlowReceived, FlowBytesRate etc. I have used different techniques to handle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

- Header:** jupyter 18BIT0276_Yuvraj (autosaved)
- Toolbar:** File, Edit, View, Insert, Cell, Kernel, Widgets, Help
- Cell 9:** `db['SourceIP'].value_counts()`
Output:
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
- Cell 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']`
- Cell 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",
 })`
In the status bar, it says: <ipython-input-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):
 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\q17a\anaconda3\lib\site-packages\pandas\core\indexing.py:671: SettingWithCopyWarning:
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

```

0      1807

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
<ipython-input-47-67f2d0160206> in <module>
      3     db['FlowBytesSent'][i] = 'fs_1807'
      4     elif(db['FlowBytesSent'][i] >= 1810):
----> 5         df['FlowBytesSent'][i] = 'fsless'
      6     else:
      7         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 3

In [39]: `rsus-a-copy`
db['FlowBytesSent'][i] = 'fs_1807'

In [40]: `fs_dummy = pd.get_dummies(db['FlowBytesSent'])`

In [40]: `destination_port_dummy.head()`

Out[40]:

Wrong_port	destination_port_443
251170	0
250984	0
90316	0
6366	0
135318	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]: `frec = db[['FlowBytesReceived']]`

In [49]: `frec.mean()`

Out[49]: `FlowBytesReceived 67362.488946`
dtype: float64

In [50]: `frec.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

Trusted Python 3

```
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

Trusted Python 3

```
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)`

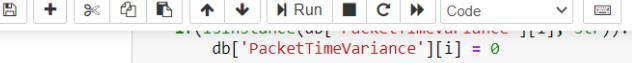
In [85]: `for i in range (0,len(db)):`
`if(isinstance(db['PacketLengthCoefficientofVariation'][i], str)):`
`db['PacketLengthCoefficientofVariation'][i] = 0`

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

In [91]: ptv = ptv.rename(columns = {0:"PTV"})

In [92]: db_create = pd.concat([db_create,ptv],axis=1)

In [93]: for i in range (0,len(db)):
    if(isinstance(db['PacketTimeStandardDeviation'][i], str)):
        db['PacketTimeStandardDeviation'][i] = 0

In [94]: pt_std = db[['PacketTimeStandardDeviation']]
pt_std = sc.fit_transform(pt_std)
pt_std = pd.DataFrame(pt_std)

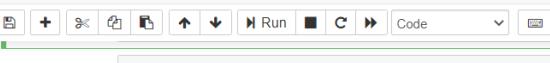
In [95]: pt_std = pt_std.rename(columns = {0:"PT_std"})
db_create = pd.concat([db_create,pt_std],axis=1)

In [96]: 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 [97]: 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 Logout

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  
Logou  
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  
  
In [115]: r_time_mean = db[['ResponseTimeTimeMean']]  
r_time_mean = sc.fit_transform(r_time_mean)  
r_time_mean = pd.DataFrame(r_time_mean)  
  
In [116]: r_time_mean = r_time_mean.rename(columns = {0:"r_time_mean"})  
  
In [117]: db_create = pd.concat([db_create,r_time_mean],axis=1)  
  
In [118]: for i in range (0,len(db)):  
    if(isinstance(db['ResponseTimeTimeMedian'][i], str)):  
        db['ResponseTimeTimeMedian'][i] = 0  
  
    r_time_med = db[['ResponseTimeTimeMedian']]
```

jupyter 10110270_YUVRDj (autosaved)

Trusted

File Edit View Insert Cell Kernel Widgets Help

Code

```
r_time_med = r_time[["ResponseTimeTimeMedian"]]
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-or-slice
    db['ResponseTimeTimeMedian'][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['ResponseTimeTimeMode'][i], str)):
        db['ResponseTimeTimeMode'][i] = 0

    r_time_mode = db[['ResponseTimeTimeMode']]
    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['ResponseTimeTimeSkewFromMedian'][i], str)):
        db['ResponseTimeTimeSkewFromMedian'][i] = 0

    r_s_med_time = db[['ResponseTimeTimeSkewFromMedian']]
    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



```
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)
```

File Edit View Insert Cell Kernel Widgets Help

Trusted | Python 3 O

Run Cell | Code

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_iqr
0	0	0	0	0	0	0	1	0	0	0	...	0.754077	-0.178941	-0.244670	-0.239715	-0.221598	-0.159623	
1	0	0	0	0	0	0	1	0	0	0	...	-0.814219	-0.132047	0.277059	-0.183655	-0.225091	-0.159623	
2	0	0	0	0	0	0	1	0	0	0	...	-0.817715	-0.178940	-0.244316	-0.238336	-0.221617	-0.159623	
3	0	0	0	0	0	0	1	0	0	0	...	-0.316040	-0.178946	-0.250146	-0.238243	-0.221523	-0.159623	
4	0	0	0	0	0	0	1	0	0	0	...	-0.770578	-0.178937	-0.242955	-0.240378	-0.221754	-0.159623	
...	
269638	1	0	0	0	0	0	0	0	0	0	...	-0.997728	-0.178940	-0.244337	-0.240169	-0.221636	-0.159623	
269639	1	0	0	0	0	0	0	0	0	0	...	0.947832	-0.178940	-0.244243	-0.240186	-0.221621	-0.159623	
269640	1	0	0	0	0	0	0	0	0	0	...	-0.858274	-0.178946	-0.250175	-0.243940	-0.226451	-0.159623	
269641	1	0	0	0	0	0	0	0	0	0	...	0.898332	-0.178940	-0.244267	-0.240184	-0.221634	-0.159623	
269642	1	0	0	0	0	0	0	0	0	0	...	-0.588891	-0.178785	-0.219341	-0.231216	-0.218044	-0.159623	

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_iqr
0	0	0	0	0	0	0	1	0	0	0	...	-0.178941	-0.244670	-0.239715	-0.221598	-0.159623	-0.159623	
1	0	0	0	0	0	0	1	0	0	0	...	-0.132047	0.277059	-0.183655	-0.225091	-0.159623	-0.159623	

File Edit View Insert Cell Kernel Widgets Help

Trusted | Python 3 O

Run Cell | Code

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_iqr
0	0	0	0	0	0	0	1	0	0	0	...	0.754077	-0.178941	-0.244670	-0.239715	-0.221598	-0.159623	
1	0	0	0	0	0	0	1	0	0	0	...	-0.814219	-0.132047	0.277059	-0.183655	-0.225091	-0.159623	
2	0	0	0	0	0	0	1	0	0	0	...	-0.817715	-0.178940	-0.244316	-0.238336	-0.221617	-0.159623	
3	0	0	0	0	0	0	1	0	0	0	...	-0.316040	-0.178946	-0.250146	-0.238243	-0.221523	-0.159623	
4	0	0	0	0	0	0	1	0	0	0	...	-0.770578	-0.178937	-0.242955	-0.240378	-0.221754	-0.159623	

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_iqr
0	0	0	0	0	0	0	1	0	0	0	...	0.754077	-0.178941	-0.244670	-0.239715	-0.221598	-0.159623	
1	0	0	0	0	0	0	1	0	0	0	...	-0.814219	-0.132047	0.277059	-0.183655	-0.225091	-0.159623	
2	0	0	0	0	0	0	1	0	0	0	...	-0.817715	-0.178940	-0.244316	-0.238336	-0.221617	-0.159623	
3	0	0	0	0	0	0	1	0	0	0	...	-0.316040	-0.178946	-0.250146	-0.238243	-0.221523	-0.159623	
4	0	0	0	0	0	0	1	0	0	0	...	-0.770578	-0.178937	-0.242955	-0.240378	-0.221754	-0.159623	

jupyter 18BIT0276_Yuvraj (autosaved)

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

5 rows x 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_u
0	0	0	0	0	0	0	1	0	0	0	...	0.754077	-0.178941	-0.244670	-0.239715	-0.221598	-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
2	0	0	0	0	0	0	1	0	0	0	...	-0.817715	-0.178940	-0.244316	-0.238336	-0.221617	-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
4	0	0	0	0	0	0	1	0	0	0	...	-0.770578	-0.178937	-0.242955	-0.240378	-0.221754	-0.11

5 rows x 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:

12-benign

SourceIP	Destinatio	SourcePor	Destinatio	TimeStamp	Duration	FlowBytes	FlowSentr	FlowBytes	FlowRecei	PacketLen	PacketTim												
192.168.2.176.103.1	50749	443	#####	95.08155	62311	655.3427	65358	687.3889	7447.674	68.45621	135.6738	102	54	1168647	944483	637326	670.5858	25.8567	54.0565	48.81129	1.49506	4	
192.168.2.176.103.1	50749	443	#####	12.3093	5928	767.137	101232	827.672	10458.12	102.2649	141.2455	114	54	0.799261	0.85132	0.724023	708.4659	26.61702	52.2879	48.83031	31.7196	0	
192.168.2.176.103.1	50749	443	#####	120.9584	38784	320.6391	38236	316.108	7300.294	85.44176	133.7153	89	54	15.70027	9.932978	0.638893	158.911	36.86341	0.51361	39.77705	0.417528	0	
192.168.2.176.103.1	50749	443	#####	110.5011	5611.061	1571	69757	631.2789	8498.293	92.19155	139.1235	114	54	0.817544	0.92333	0.62666	111.138	33.43835	51.69373	34.8825	13.28093	0	
16.103.1.2	192.168.2.1	443	#####	54.22989	83641	154.3241	76804	146.2167	8025.746	89.73709	138.9134	114	54	0.83288	0.277627	0.645993	341.9666	18.48504	36.43562	49.82256	7.342519	0	
7.168.2.176.103.1	52491	443	#####	145.4607	54084	371.8117	63843	438.902	16074.98	126.7871	141.061	89	54	1.231852	0.686671	0.898811	151.785	38.9842	58.2639	47.49244	33.80841	0	
192.168.2.176.103.1	52742	443	#####	0.15541	1718	11054.63	7411	76876.6	234216	483.8657	551.1154	129.5	54	1.37403	0.614045	1.378082	0.01068	0.032679	0.062809	0.0609	0.087941	0	
9.168.2.176.103.1	52742	443	#####	0.02700	55	2036.92	66	244.354	30.25	5.5	60.5	60.5	55	0	1	0.090994	0.00182	0.13501	0.13501	0	0	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	0.937181	0.654724	19.93652	13.92708	23.25506	23.03529	0.642233	0	
11.168.2.176.103.1	52742	443	#####	0.046455	55	118.9341	66	1420.73	30.25	5.5	60.5	60.5	55	0	1	0.09099	0.00054	0.023226	0.023228	0	0	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	0.195566	377.2307	19.42423	33.6514	44.86306	44.86306	0	
13.168.2.176.103.1	52491	443	#####	91.33806	8090	88.57206	7805	85.45178	6158.02	78.47307	125.1575	89	60	1.382289	0.830216	0.626995	140.5086	11.85363	78.94149	81.96176	0	0	
14.168.2.176.103.1	52491	443	#####	121.3329	53598	411.7434	59634	491.4906	8309.003	91.15374	136.2599	89	54	1.555392	0.902431	0.66869	107.683	31.90114	57.2504	55.49805	78.45495	0	
192.168.2.176.103.1	52491	443	#####	12.3599	56840	505.8745	60236	536.0988	8218.83	90.65759	135.819	89	54	1.549311	0.902504	0.66749	117.2673	33.58009	34.36958	25.3958	72.24556	0	
16.103.1.2	192.168.2.1	443	#####	135.2442	69610	514.6985	65101	481.358	8353.12	91.39542	136.6237	89	54	1.56322	0.904025	0.668957	90.743	30.0623	55.51041	44.10633	20.5055	0	
192.168.2.176.103.1	52491	443	#####	121.1859	35194	290.4132	37434	308.8972	7686.456	87.67244	136.5188	114	54	0.770554	0.941217	0.6422	116.5663	34.15502	56.07565	62.02651	76.86132	0	
19.168.2.176.103.1	52491	443	#####	98.55257	78613	285.1738	7709.703	80.78401	135.4259	89	54	1.586218	0.92735	0.648361	707.6148	26.60103	54.10346	54.5677	48.8673	29.66674	0		
16.103.1.2	192.168.2.1	443	#####	33.05635	13565	40.3599	11354	343.4741	9409.972	97.00501	139.9944	89	54	1.57064	0.886494	0.692291	102.984	10.14811	20.12085	7.854821	33.02868	0	
19.168.2.176.103.1	54640	443	#####	114.9067	57076	496.7161	68354	594.8653	17032.2	130.5114	146.5304	114	54	0.74776	0.070983	0.890678	117.686	34.332	46.67889	49.60508	0.054005	0	
21.168.2.176.103.1	52491	443	#####	96.75011	42044	434.5628	44920	464.2889	7899.639	88.87992	133.7908	89	54	1.511841	0.87737	0.66432	488.3844	22.09942	39.3314	33.98518	0.145493	0	
192.168.2.176.103.1	54640	443	#####	105.5193	61248	580.4434	63497	601.757	8029.011	89.60476	135.1517	89	54	1.545175	0.905663	0.662994	670.7841	25.8995	53.89672	56.66998	26.27162	0	
192.168.2.176.103.1	54640	443	#####	110.6149	2797	247.6675	29436	266.1124	7085.691	88.34982	135.0786	89	54	1.564641	0.917699	0.654062	147.843	38.3769	55.2808	41.25153	22.90654	1	
192.168.2.176.103.1	54640	443	#####	121.773	67491	554.2362	73742	605.569	9569.62	97.82444	141.3744	114	54	0.839945	0.893175	0.619593	107.6993	32.85869	53.07471	52.74703	20.96734	0	
25.162.1.2	192.168.2.1	443	#####	121.6762	38891	320.4489	41627	342.1113	8195.517	90.5291	135.4924	89	54	1.540669	0.90018	0.668149	185.212	34.4269	47.1759	46.7256	8.89319	0	
192.168.2.176.103.1	54640	443	#####	121.2499	7213	955.5916	80570	664.4953	10426.65	102.2871	139.9112	89	54	1.493185	0.839902	0.713086	107.711	32.7371	50.99465	45.52109	15.50502	0	
27.162.1.2	192.168.2.1	443	#####	74.03461	30733	415.1167	32409	437.7547	7758.135	88.08027	132.6513	89	54	1.486755	0.89295	0.663999	650.7667	25.51013	33.52663	32.25093	1.671356	0	
28.168.2.176.103.1	56611	443	#####	122.0805	35590	295.5291	45875	365.2264	21960.19	148.1897	146.3084	89	54	1.160169	0.622907	1.012859	552.312	23.01032	45.36756	41.00397	38.21969	0	
29.168.2.176.103.1	56611	443	#####	96.36561	41539	431.0563	45477	462.5821	8937.76	94.59381	137.5655	89	54	1.541112	0.513704	0.687235	875.4395	29.58783	53.88762	40.85792	23.00198	0	

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_port_is_44	Wrong_po_destination	Duration	fs_1807	fs_more	fsless	FlowSent	FlowRecv	Pcket	Jenk	Le		
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	1.
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
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	-0.78239	-4
5	3	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	2.703618	0	0	1	-0.0473	-0.06989	-0.91759	-4
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	0.
7	5	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	1.18219	0	0	1	0.016064	-0.00313	-0.91769	-4
8	6	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	1.453431	0	1	0	-0.04301	-0.05703	-0.89671	-4
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	0.
10	8	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	-0.4348	0	0	1	-0.04642	-0.06694	0.378454	0.
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	0.
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	0.
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	0.
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	0.
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	0.
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	0.
17	15	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	-1.14745	0	0	1	-0.04391	-0.04546	0.675613	0.
18	16	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	1.465147	0	0	1	-0.04405	-0.06645	-0.9041	-4
19	17	0	0	0	0	0	0	1	0	0	0	1	0	1	1	0	1.467265	0	0	1	-0.04678	-0.06784	-0.87492	-4
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	0.192244	0.
21	19	0	0	0	0	0	0	1	0	0	0	1	1	0	0	1	1.446425	0	0	1	-0.0471	-0.06984	-0.91934	-4
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	0.
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	0.
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	2.
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	-4
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	-4
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	0.
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	1.
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	-4

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

Comment Share

+ Code + Text

Connect Editing

```
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	1	0	0	0	1	1	0	0	0	1	-0.440208
1	1	0	0	0	0	0	1	0	0	0	1	1	0	0	0	1	1.171644
2	2	0	0	0	0	0	1	0	0	0	1	1	0	0	0	1	1.145800
3	3	0	0	0	0	0	1	0	0	0	1	1	0	0	0	1	2.703618
4	4	0	0	0	0	0	1	0	0	0	1	0	1	1	1	0	1.152860

FINAL_CNN_MODEL.ipynb

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

Comment Share

+ Code + Text

Connect Editing

```
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 np_utils
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

Comment Share

+ 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 × 45 columns

FINAL_CNN_MODEL.ipynb

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

Comment Share

+ Code + Text

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

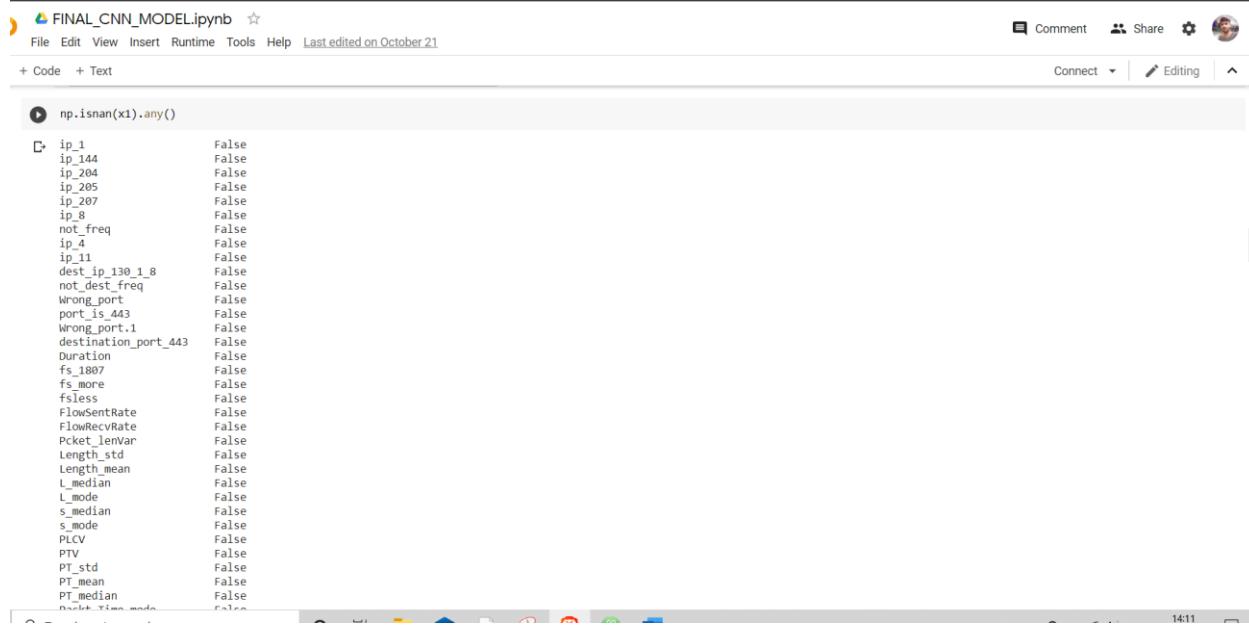
```
[ ] y1 = db.iloc[:,44]
```

```
[ ] 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 × 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_time_mean  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_mean_time      False
r_s_med_time       False
r_s_mode_time      False
r_s_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      0      1      0      0      1      1.443449      0
131414  0      1      0      0      0      0      0      1      0      0      0      0      1      0      0      0      -1.135324      0
144159  0      0      0      0      0      0      1      0      0      1      0      1      0      1      0      0      -0.440212      0
167495  0      0      0      1      0      0      0      0      1      0      0      0      1      0      0      0      -1.150709      0
253272  0      0      0      0      0      0      1      0      0      1      0      0      1      0      0      0      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”



```
[ ] 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)
```

FINAL_CNN_MODEL.ipynb

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

Comment Share

Connect Editing

```
+ 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("Xs: %.2F%% (%.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
446/446 [=====] - 3s 7ms/step - loss: 0.1441 - accuracy: 0.9473
accuracy: 93.82%
-----
```

```
[ ] 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%)

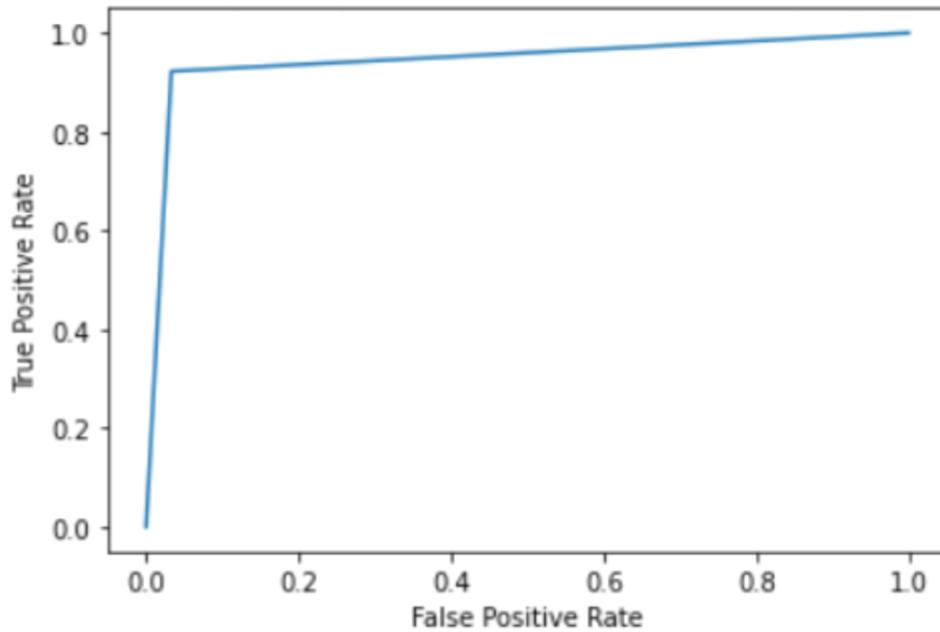
Training accuracy: 94.83449459075928
Test accuracy: 94.44654583930969

	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()
```



Yuvraj Kumar
(18BIT0276)

Student's Signature