

First Name: Harish

Last Name: Marepalli

SJSU ID: 016707314

Professor: Dr. Younghee Park

Improving Accuracy Using Random Forest Algorithm

Improving accuracy using Random Forest Algorithm:

Here, we have been given two datasets, Train_data.csv and Test_data.csv. Train_data.csv is used here to train the model on how to detect whether an incoming traffic is normal or anomaly. In these datasets we there are so many columns. The Train_data.csv contains 42 columns and the Test_data.csv contains 41 columns. The difference of one column is the 'Class' column, which is used to detect normal and anomaly traffic.

Here, we were given the code to find the Random Forest algorithms output. By using the normal approach, the output achieved was 94% and we have to improve it to 98-99% by making modifications.

The procedure that I followed in this activity is I have installed jupyter labs and written the code in python notebook file. By using the normal code, I got the output as 94.7%.

First, to improve the accuracy, I have included 3 more columns of train dataset in the code and these columns are 'count', 'srv_count', and 'logged_in'. I have chosen only these columns is because I felt that these may contribute to the future prediction efficiently because remaining other column values are just 0's and 1's, which does not much contribute to the future prediction.

Some of the other columns might also contribute to the accurate prediction but we do not include more columns as it will be an overload/overfitting case. So, this is the reason I have included these columns.

Second, I have splitted the dataset into 70% and 30% to get the accurate result.

Third, in the modeling section, I have increased the max_depth from 2 to 7. The default value of number of trees is 100. The value of max_depth indicates the depth that we traverse for a tree. If the value of it is less, it means we are not deeply going through the random forest code. If the value is more, then we are diving deep into the process and can know how the process works and also the accurate output.

I have tried all the values of max depth from 2 to 7 and stopped at 7 because it is considered to be the limit. If the value is more than 7, then it will be considered as overfitting, which we don't want to do.

After making all these changes in the code, the output that I achieved is 0.9886213283937549 and that comes to be 98.8% which can be approximated to **99%**.

I have also written a few lines to display the confusion matrix and the matrix is as follows:

```
[[4017,    18],  
 [   68,   3455]]
```

The scheme of confusion matrix is:

		Predicted condition	
		Positive (PP)	Negative (PN)
Actual condition	Positive (P)	True positive (TP)	False negative (FN)
	Negative (N)	False positive (FP)	True negative (TN)

Here, the number of True positive (TP) values are 4017.

The number of False negative (FN) values are 18.

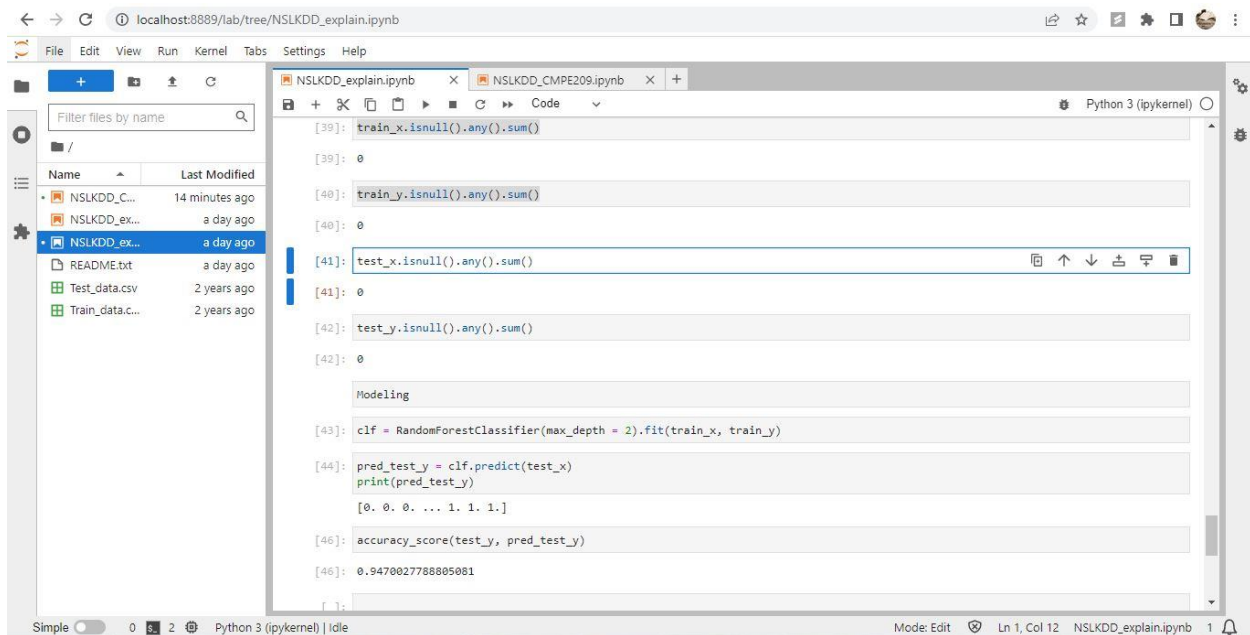
The number of False positive (FP) values are 68.

The number of True negative (TN) values are 3455.

Code is provided in the Appendix.

Snippets:

1. Figure 1 shows the output before the improvement.

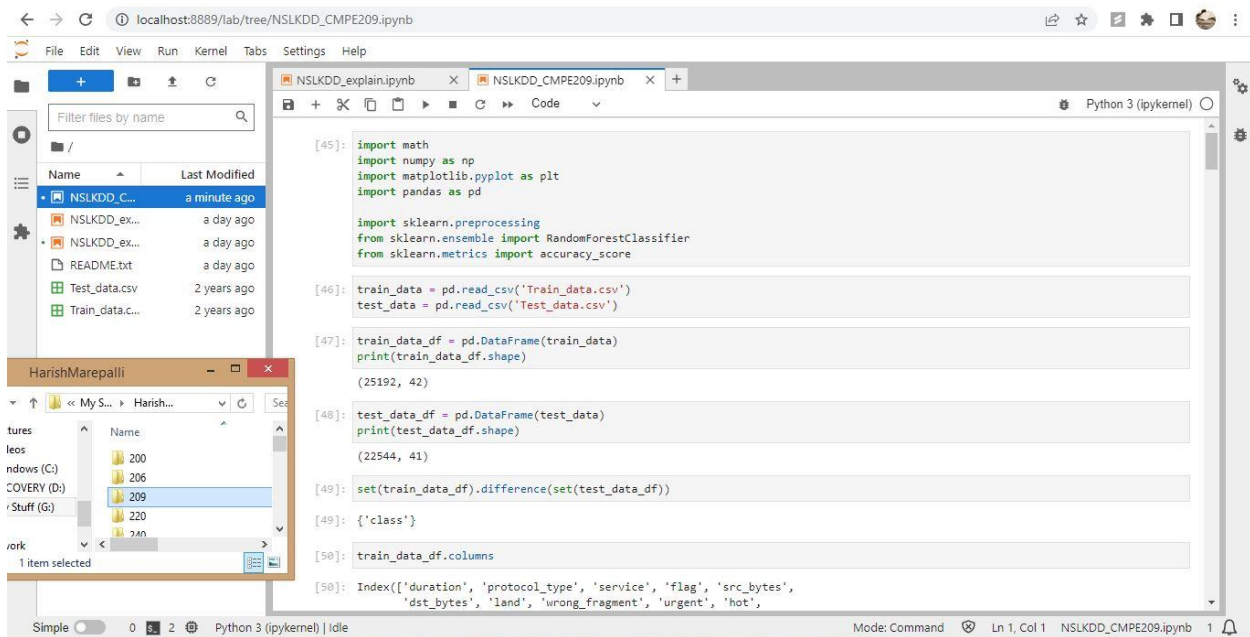


```
[39]: train_x.isnull().any().sum()
[39]: 0
[40]: train_y.isnull().any().sum()
[40]: 0
[41]: test_x.isnull().any().sum()
[41]: 0
[42]: test_y.isnull().any().sum()
[42]: 0

Modeling
[43]: clf = RandomForestClassifier(max_depth = 2).fit(train_x, train_y)
[44]: pred_test_y = clf.predict(test_x)
      print(pred_test_y)
      [0. 0. 0. ... 1. 1. 1.]
[46]: accuracy_score(test_y, pred_test_y)
[46]: 0.9470027788805081
```

Fig. 1: Output before improvement

2. Figure 2 shows the improved code.



```
[45]: import math
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

import sklearn.preprocessing
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score

[46]: train_data = pd.read_csv('Train_data.csv')
test_data = pd.read_csv('Test_data.csv')

[47]: train_data_df = pd.DataFrame(train_data)
print(train_data_df.shape)

(25192, 42)

[48]: test_data_df = pd.DataFrame(test_data)
print(test_data_df.shape)

(22544, 41)

[49]: set(train_data_df).difference(set(test_data_df))

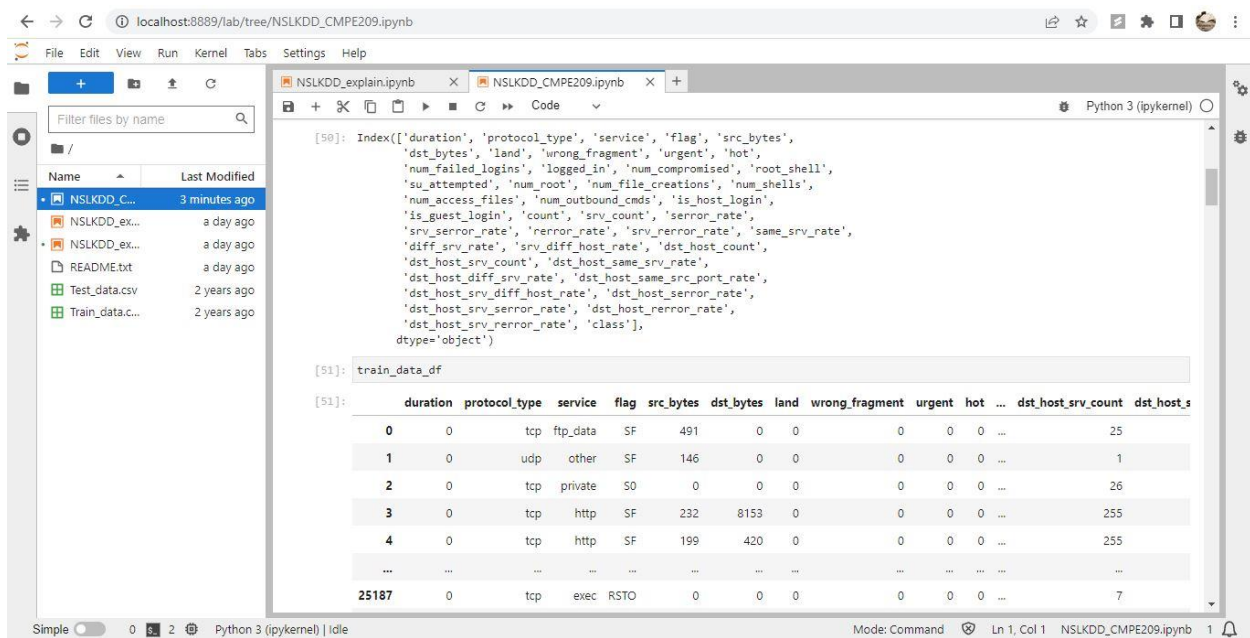
[49]: {'class'}

[50]: train_data_df.columns

[50]: Index(['duration', 'protocol_type', 'service', 'flag', 'src_bytes',
'dst_bytes', 'land', 'wrong_fragment', 'urgent', 'hot',
'dst_host_srv_count', 'dst_host_same_src_port_rate',
'dst_host_srv_diff_host_rate', 'dst_host_error_rate',
'dst_host_srv_error_rate', 'class'],
dtype='object')
```

Fig. 2: Improved code 1

3. Figure 3 shows the continued code.



```
[50]: Index(['duration', 'protocol_type', 'service', 'flag', 'src_bytes',
'dst_bytes', 'land', 'wrong_fragment', 'urgent', 'hot',
'dst_host_srv_count', 'dst_host_same_src_port_rate',
'dst_host_srv_diff_host_rate', 'dst_host_error_rate',
'dst_host_srv_error_rate', 'class'],
dtype='object')

[51]: train_data_df

[51]:
```

	duration	protocol_type	service	flag	src_bytes	dst_bytes	land	wrong_fragment	urgent	hot	...	dst_host_srv_count	dst_host_s
0	0	tcp	ftp_data	SF	491	0	0	0	0	0	...	25	
1	0	udp	other	SF	146	0	0	0	0	0	...	1	
2	0	tcp	private	SO	0	0	0	0	0	0	...	26	
3	0	tcp	http	SF	232	8153	0	0	0	0	...	255	
4	0	tcp	http	SF	199	420	0	0	0	0	...	255	
...	
25187	0	tcp	exec	RSTO	0	0	0	0	0	0	...	7	

Fig. 3: Improved code 2

4. Figure 4 shows the continued code.

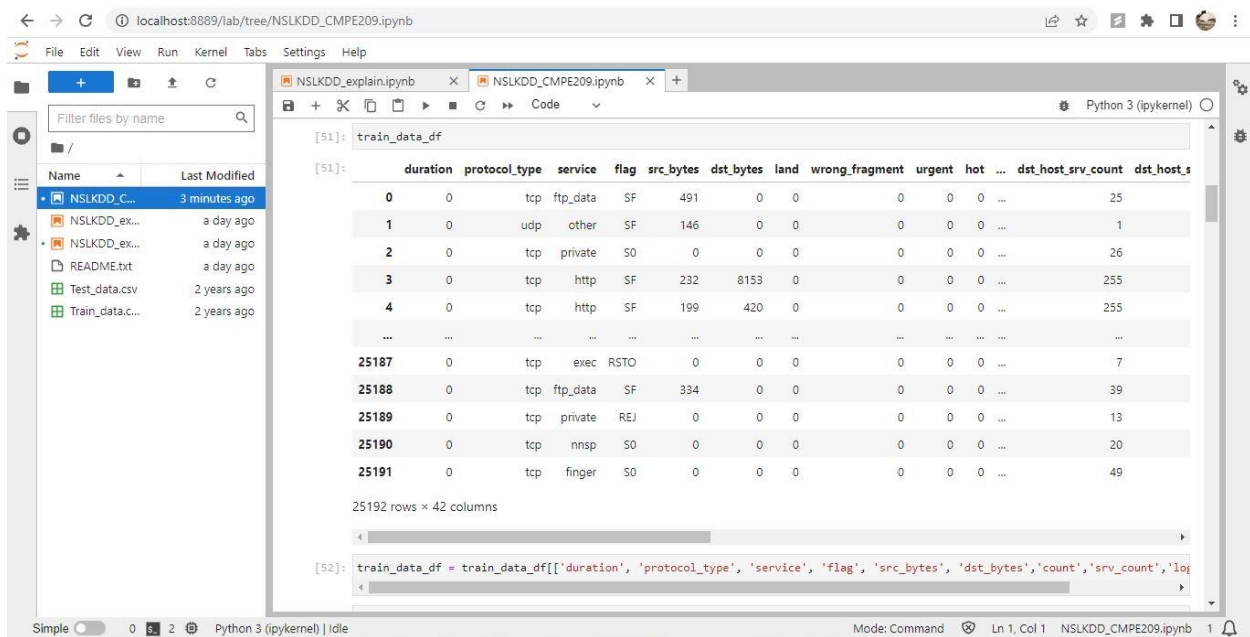


Fig. 4: Improved code 3

5. Figure 5 shows the added columns.



Fig. 5: Added columns

6. Figure 6 shows the continued code.

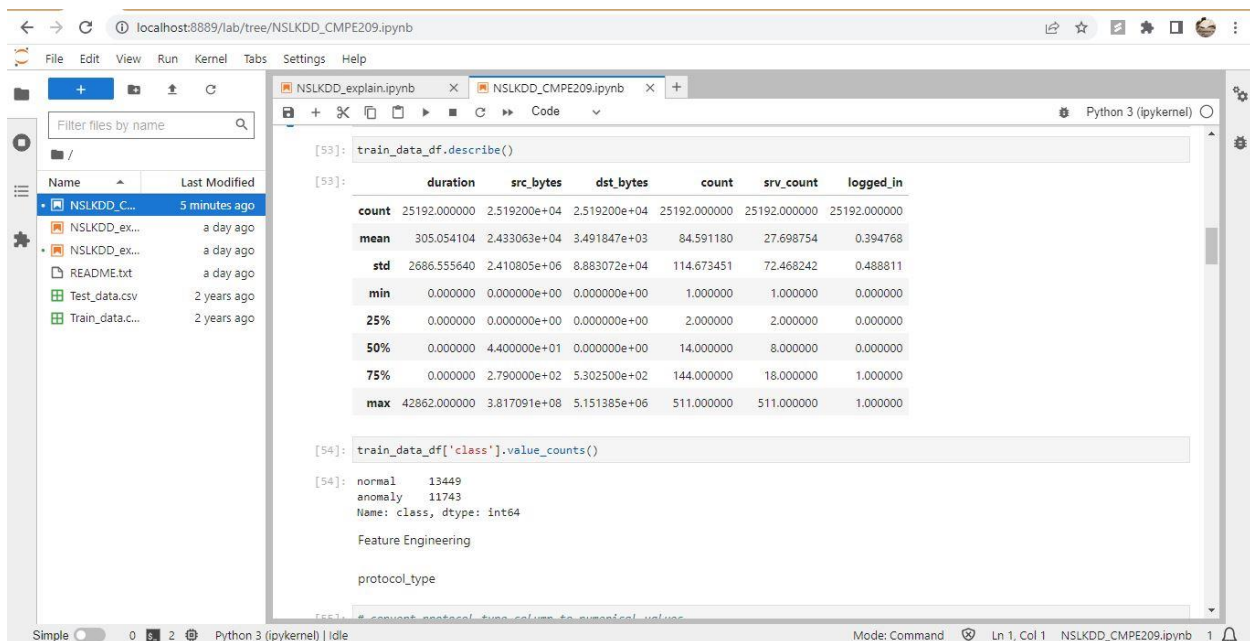
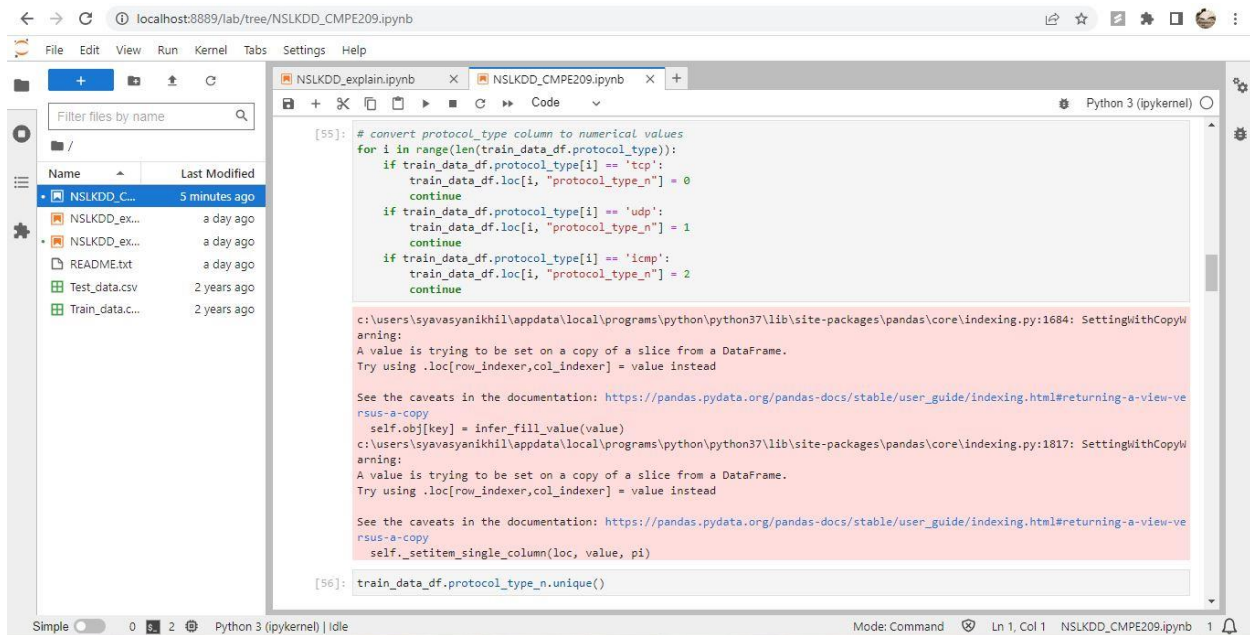


Fig. 6: Improved code 4

7. Figure 7 shows the continued code.



The screenshot shows a Jupyter Notebook with a file explorer on the left and a code editor on the right. The file explorer shows files like NSLKDD_C..., NSLKDD_ex..., README.txt, Test_data.csv, and Train_data.csv. The code editor shows the following code:

```
[55]: # convert protocol_type column to numerical values
for i in range(len(train_data_df.protocol_type)):
    if train_data_df.protocol_type[i] == 'tcp':
        train_data_df.loc[i, "protocol_type_n"] = 0
        continue
    if train_data_df.protocol_type[i] == 'udp':
        train_data_df.loc[i, "protocol_type_n"] = 1
        continue
    if train_data_df.protocol_type[i] == 'icmp':
        train_data_df.loc[i, "protocol_type_n"] = 2
        continue

c:\users\syavasyanikhil\appdata\local\programs\python\python37\lib\site-packages\pandas\core\indexing.py:1684: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

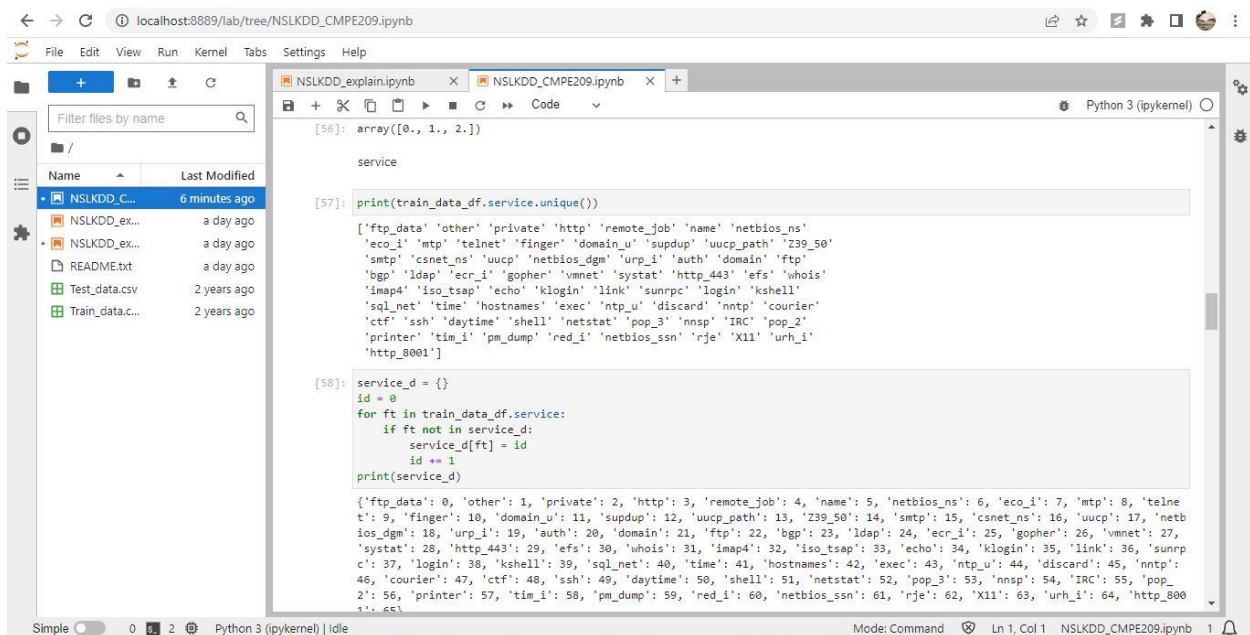
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
self.obj[key] = infer_fill_value(value)
c:\users\syavasyanikhil\appdata\local\programs\python\python37\lib\site-packages\pandas\core\indexing.py:1817: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
self.setitem_single_column(loc, value, pi)

[56]: train_data_df.protocol_type_n.unique()
```

Fig. 7: Improved code 5

8. Figure 8 shows the continued code.



The screenshot shows a Jupyter Notebook with a file explorer on the left and a code editor on the right. The file explorer shows files like NSLKDD_C..., NSLKDD_ex..., README.txt, Test_data.csv, and Train_data.csv. The code editor shows the following code:

```
[56]: array([0., 1., 2.])

service

[57]: print(train_data_df.service.unique())

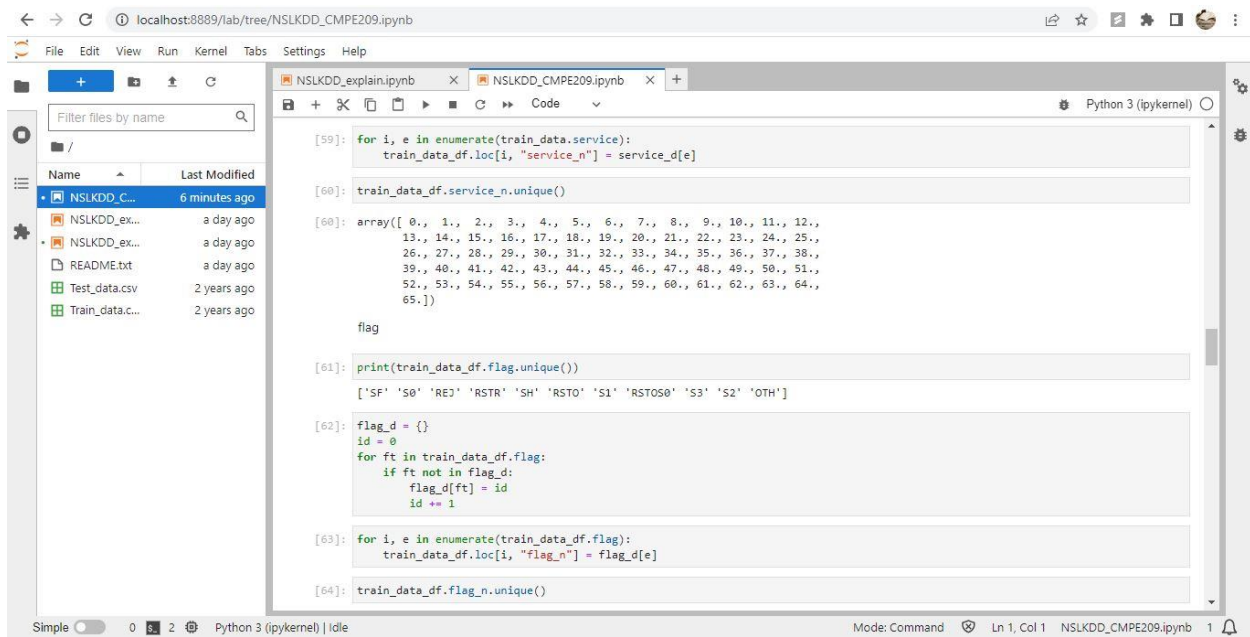
['ftp_data' 'other' 'private' 'http' 'remote_job' 'name' 'netbios_ns'
'eco_i' 'mtp' 'telnet' 'finger' 'domain_u' 'supdup' 'uucp_path' 'Z39_50'
'smtp' 'csnet_ns' 'uucp' 'netbios_dgm' 'urp_i' 'auth' 'domain' 'ftp'
'bgp' 'ldap' 'ecr_i' 'gopher' 'vmnet' 'sysstat' 'http_443' 'efs' 'whois'
'imap4' 'iso_tsap' 'echo' 'klogin' 'link' 'sunrpc' 'login' 'kshell'
'sql_net' 'time' 'hostnames' 'exec' 'ntp_u' 'discard' 'nntp' 'courier'
'ctf' 'ssh' 'daytime' 'shell' 'netstat' 'pop_3' 'nnsf' 'IRC' 'pop_2'
'printer' 'tim_i' 'pm_dump' 'red_i' 'netbios_ssn' 'rje' 'X11' 'urh_i'
'http_8001']

[58]: service_d = {}
id = 0
for ft in train_data_df.service:
    if ft not in service_d:
        service_d[ft] = id
        id += 1
print(service_d)

{'ftp_data': 0, 'other': 1, 'private': 2, 'http': 3, 'remote_job': 4, 'name': 5, 'netbios_ns': 6, 'eco_i': 7, 'mtp': 8, 'telnet': 9, 'finger': 10, 'domain_u': 11, 'supdup': 12, 'uucp_path': 13, 'Z39_50': 14, 'smtp': 15, 'csnet_ns': 16, 'uucp': 17, 'netbios_dgm': 18, 'urp_i': 19, 'auth': 20, 'domain': 21, 'ftp': 22, 'bgp': 23, 'ldap': 24, 'ecr_i': 25, 'gopher': 26, 'vmnet': 27, 'sysstat': 28, 'http_443': 29, 'efs': 30, 'whois': 31, 'imap4': 32, 'iso_tsap': 33, 'echo': 34, 'klogin': 35, 'link': 36, 'sunrpc': 37, 'login': 38, 'kshell': 39, 'sql_net': 40, 'time': 41, 'hostnames': 42, 'exec': 43, 'ntp_u': 44, 'discard': 45, 'nntp': 46, 'courier': 47, 'ctf': 48, 'ssh': 49, 'daytime': 50, 'shell': 51, 'netstat': 52, 'pop_3': 53, 'nnsf': 54, 'IRC': 55, 'pop_2': 56, 'printer': 57, 'tim_i': 58, 'pm_dump': 59, 'red_i': 60, 'netbios_ssn': 61, 'rje': 62, 'X11': 63, 'urh_i': 64, 'http_8001': 65}
```

Fig. 8: Improved code 6

9. Figure 9 shows the continued code.



```
[59]: for i, e in enumerate(train_data.service):
      train_data_df.loc[i, "service_n"] = service_d[e]

[60]: train_data_df.service_n.unique()

[60]: array([ 0.,  1.,  2.,  3.,  4.,  5.,  6.,  7.,  8.,  9., 10., 11., 12.,
        13., 14., 15., 16., 17., 18., 19., 20., 21., 22., 23., 24., 25.,
        26., 27., 28., 29., 30., 31., 32., 33., 34., 35., 36., 37., 38.,
        39., 40., 41., 42., 43., 44., 45., 46., 47., 48., 49., 50., 51.,
        52., 53., 54., 55., 56., 57., 58., 59., 60., 61., 62., 63., 64.,
        65.])

      flag

[61]: print(train_data_df.flag.unique())

['SF' 'S0' 'REJ' 'RSTR' 'SH' 'RSTO' 'S1' 'RSTO50' 'S3' 'S2' 'OTH']

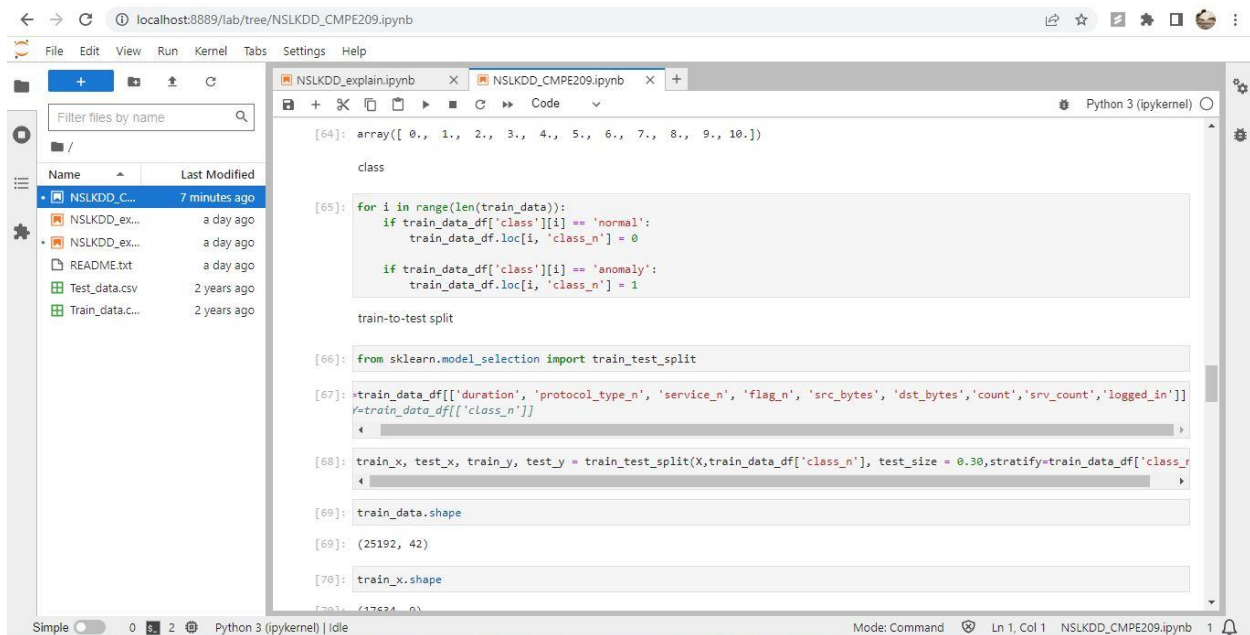
[62]: flag_d = {}
      id = 0
      for ft in train_data_df.flag:
          if ft not in flag_d:
              flag_d[ft] = id
              id += 1

[63]: for i, e in enumerate(train_data_df.flag):
      train_data_df.loc[i, "flag_n"] = flag_d[e]

[64]: train_data_df.flag_n.unique()
```

Fig. 9: Improved code 7

10. Figure 10 shows the continued code.



```
[64]: array([ 0.,  1.,  2.,  3.,  4.,  5.,  6.,  7.,  8.,  9., 10.])

      class

[65]: for i in range(len(train_data)):
      if train_data_df['class'][i] == 'normal':
          train_data_df.loc[i, 'class_n'] = 0

      if train_data_df['class'][i] == 'anomaly':
          train_data_df.loc[i, 'class_n'] = 1

      train-to-test split

[66]: from sklearn.model_selection import train_test_split

[67]: train_data_df[['duration', 'protocol_type_n', 'service_n', 'flag_n', 'src_bytes', 'dst_bytes', 'count', 'srv_count', 'logged_in']]
      ~~~~~train_data_df[['class_n']]

[68]: train_x, test_x, train_y, test_y = train_test_split(X=train_data_df[['class_n']], test_size = 0.30, stratify=train_data_df[['class_n']])

[69]: train_data.shape

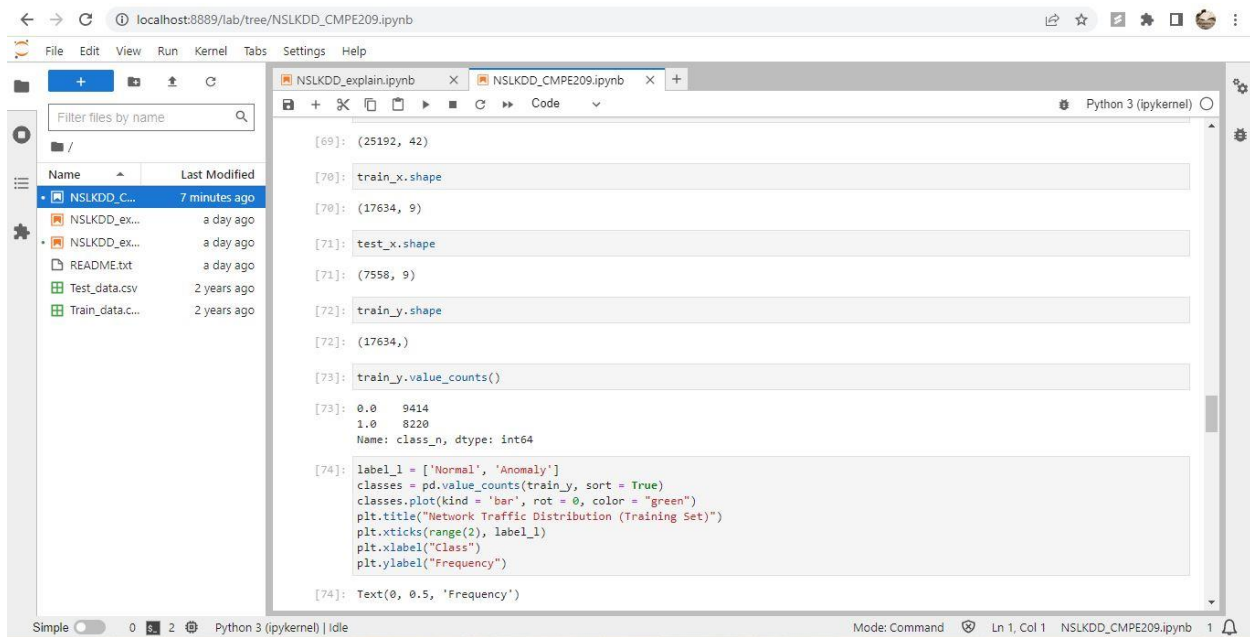
[69]: (25192, 42)

[70]: train_x.shape

[70]: (17634, 42)
```

Fig. 10: Improved code 8

11. Figure 11 shows the continued code.



The Jupyter Notebook interface shows the following code cells:

```
[69]: (25192, 42)

[70]: train_x.shape

[70]: (17634, 9)

[71]: test_x.shape

[71]: (7558, 9)

[72]: train_y.shape

[72]: (17634,)

[73]: train_y.value_counts()

[73]: 0.0    9414
      1.0    8220
      Name: class_n, dtype: int64

[74]: label_1 = ['Normal', 'Anomaly']
      classes = pd.value_counts(train_y, sort = True)
      classes.plot(kind = 'bar', rot = 0, color = "green")
      plt.title("Network Traffic Distribution (Training Set)")
      plt.xticks(range(2), label_1)
      plt.xlabel("Class")
      plt.ylabel("Frequency")

[74]: Text(0, 0.5, 'Frequency')
```

Fig. 11: Improved code 9

12. Figure 12 shows the continued code.

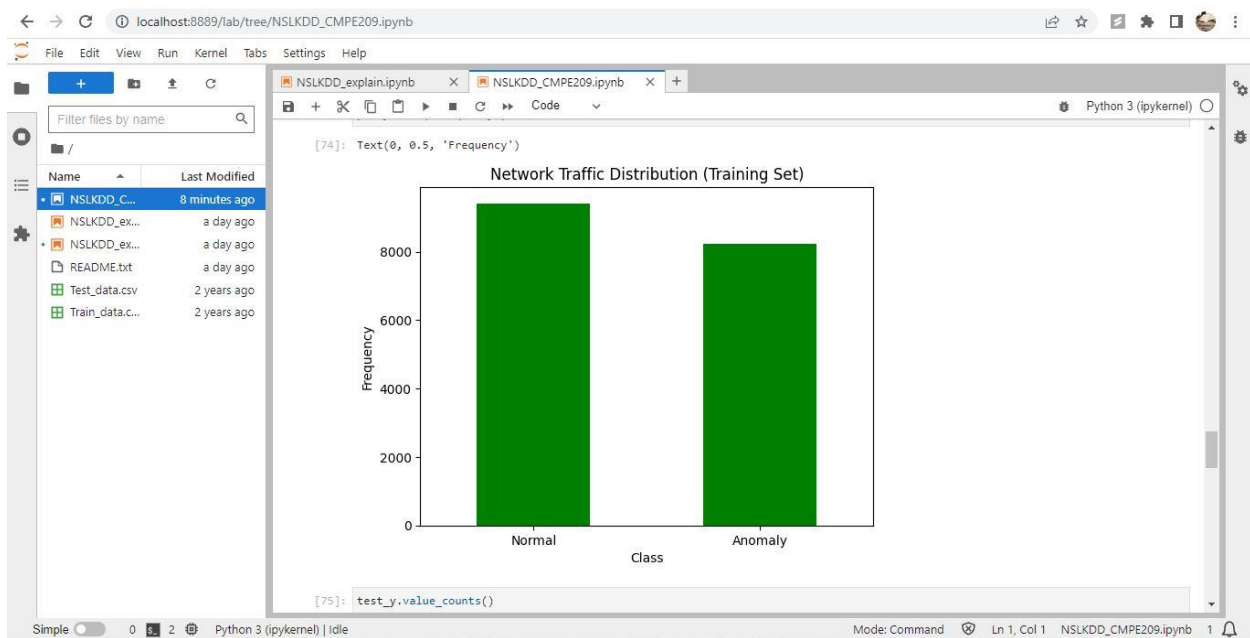


Fig. 12: Improved code 10

13. Figure 13 shows the continued code.

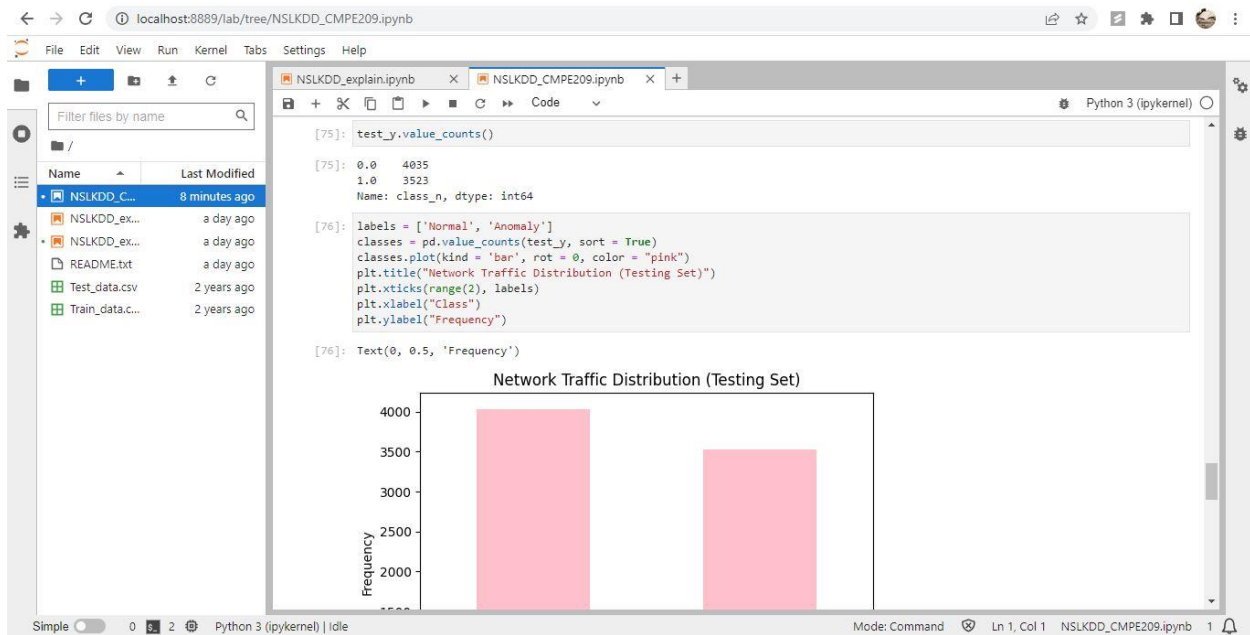


Fig. 13: Improved code 11

14. Figure 14 shows the continued code.

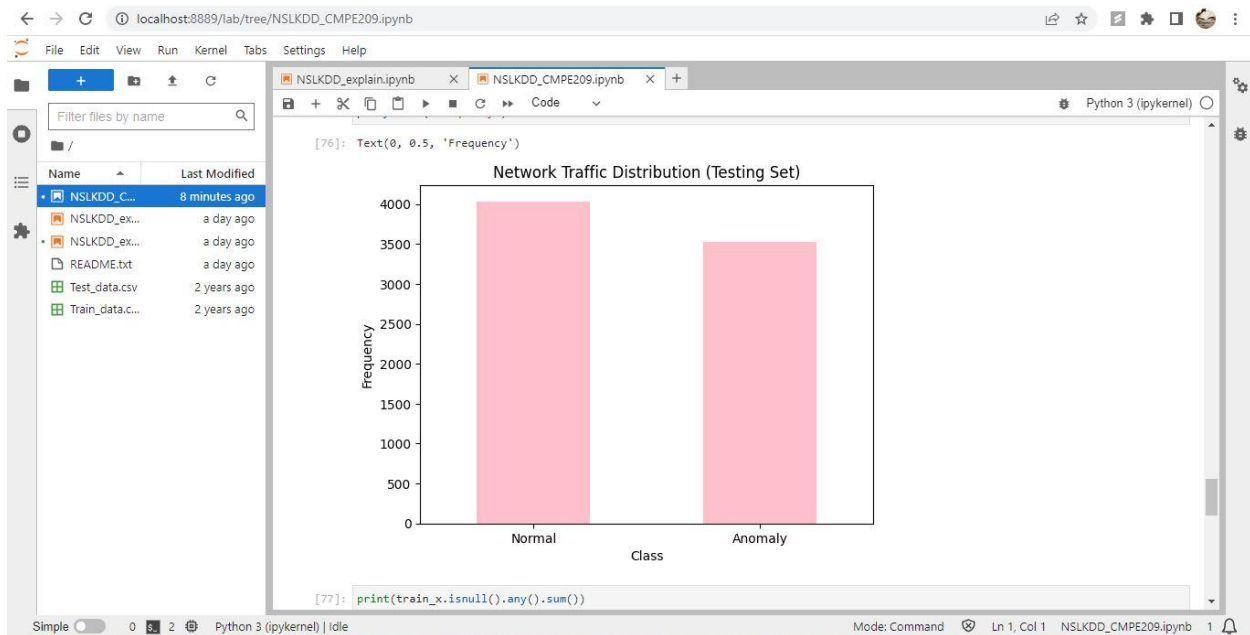


Fig. 14: Improved code 12

15. Figure 15 shows the achieved 99% output.

```
[77]: print(train_x.isnull().any().sum())
0
[78]: print(train_y.isnull().any().sum())
0
[79]: print(test_x.isnull().any().sum())
0
[80]: print(test_y.isnull().any().sum())
0

Modeling

[81]: clf = RandomForestClassifier(max_depth = 7).fit(train_x, train_y)
[82]: pred_test_y = clf.predict(test_x)
print(pred_test_y)
[1. 0. 0. ... 1. 1. 1.]
[83]: accuracy_score(test_y, pred_test_y)
[83]: 0.9886213283937549
[84]: from sklearn.metrics import confusion_matrix
print("Confusion Matrix: ")
```

Fig. 15: Output

16. Figure 16 shows the confusion matrix.

```
[79]: print(test_x.isnull().any().sum())
0
[80]: print(test_y.isnull().any().sum())
0

Modeling

[81]: clf = RandomForestClassifier(max_depth = 7).fit(train_x, train_y)
[82]: pred_test_y = clf.predict(test_x)
print(pred_test_y)
[1. 0. 0. ... 1. 1. 1.]
[83]: accuracy_score(test_y, pred_test_y)
[83]: 0.9886213283937549
[84]: from sklearn.metrics import confusion_matrix
print("Confusion Matrix: ")
confusion_matrix(test_y, pred_test_y)

Confusion Matrix:
[84]: array([[4017, 18],
[ 68, 3455]], dtype=int64)
[ ]:
```

Fig. 16: Confusion matrix

CONCLUSION:

From this project, I learned about data sets and how to improve the accuracy of Random Forest by changing a few parameters. Further, I learned how Machine Learning is used for network intrusion detection.

APPENDIX

Code:

Name: NSLKDD_CMPE209.ipynb

```
import math
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

import sklearn.preprocessing
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
```

```
train_data = pd.read_csv('Train_data.csv')
test_data = pd.read_csv('Test_data.csv')
```

```
train_data_df = pd.DataFrame(train_data)
print(train_data_df.shape)
```

```
test_data_df = pd.DataFrame(test_data)
print(test_data_df.shape)
```

```
set(train_data_df).difference(set(test_data_df))
```

```
train_data_df.columns
```

```
train_data_df
```

```
train_data_df = train_data_df[['duration', 'protocol_type', 'service', 'flag', 'src_bytes',
'dst_bytes', 'count', 'srv_count', 'logged_in', 'class']]
```

```
train_data_df.describe()
```

```
train_data_df['class'].value_counts()
```

Feature Engineering

protocol_type

```
# convert protocol_type column to numerical values
for i in range(len(train_data_df.protocol_type)):
    if train_data_df.protocol_type[i] == 'tcp':
        train_data_df.loc[i, "protocol_type_n"] = 0
        continue
    if train_data_df.protocol_type[i] == 'udp':
        train_data_df.loc[i, "protocol_type_n"] = 1
        continue
```

```
if train_data_df.protocol_type[i] == 'icmp':  
    train_data_df.loc[i, "protocol_type_n"] = 2  
    continue
```

```
train_data_df.protocol_type_n.unique()
```

service

```
print(train_data_df.service.unique())
```

```
service_d = { }  
id = 0  
for ft in train_data_df.service:  
    if ft not in service_d:  
        service_d[ft] = id  
        id += 1  
print(service_d)
```

```
for i, e in enumerate(train_data.service):  
    train_data_df.loc[i, "service_n"] = service_d[e]
```

```
train_data_df.service_n.unique()
```

flag

```
print(train_data_df.flag.unique())
```

```
flag_d = { }  
id = 0  
for ft in train_data_df.flag:  
    if ft not in flag_d:  
        flag_d[ft] = id  
        id += 1
```

```
for i, e in enumerate(train_data_df.flag):  
    train_data_df.loc[i, "flag_n"] = flag_d[e]
```

```
train_data_df.flag_n.unique()
```

class

```
for i in range(len(train_data)):   
    if train_data_df['class'][i] == 'normal':  
        train_data_df.loc[i, 'class_n'] = 0  
  
    if train_data_df['class'][i] == 'anomaly':  
        train_data_df.loc[i, 'class_n'] = 1
```

train-to-test split

```
from sklearn.model_selection import train_test_split
```

```
X=train_data_df[['duration', 'protocol_type_n', 'service_n', 'flag_n', 'src_bytes',  
'dst_bytes','count','srv_count','logged_in']]  
#Y=train_data_df[['class_n']]
```

```
train_x, test_x, train_y, test_y = train_test_split(X,train_data_df['class_n'], test_size =  
0.30,stratify=train_data_df['class_n'])
```

```
train_data.shape
```

```
train_x.shape
```

```
test_x.shape
```

```
train_y.shape
```

```
train_y.value_counts()
```

```
label_l = ['Normal', 'Anomaly']  
classes = pd.value_counts(train_y, sort = True)  
classes.plot(kind = 'bar', rot = 0, color = "green")  
plt.title("Network Traffic Distribution (Training Set)")  
plt.xticks(range(2), label_l)  
plt.xlabel("Class")  
plt.ylabel("Frequency")
```

```
test_y.value_counts()
```

```
labels = ['Normal', 'Anomaly']  
classes = pd.value_counts(test_y, sort = True)  
classes.plot(kind = 'bar', rot = 0, color = "pink")  
plt.title("Network Traffic Distribution (Testing Set)")  
plt.xticks(range(2), labels)  
plt.xlabel("Class")  
plt.ylabel("Frequency")
```

```
print(train_x.isnull().any().sum())
```

```
print(train_y.isnull().any().sum())
```

```
print(test_x.isnull().any().sum())
```

```
print(test_y.isnull().any().sum())
```


Modeling

```
clf = RandomForestClassifier(max_depth = 7).fit(train_x, train_y)
```

```
pred_test_y = clf.predict(test_x)  
print(pred_test_y)
```

```
accuracy_score(test_y, pred_test_y)
```

```
from sklearn.metrics import confusion_matrix  
print("Confusion Matrix: ")  
confusion_matrix(test_y, pred_test_y)
```