

```

import pandas as pd
import numpy as np
import scikitplot as skplt
import matplotlib.pyplot as plt
import matplotlib.pyplot as plt
import seaborn as sns

In [2]: train = pd.read_csv("UNSW_NB15_training-set.csv")
train.head()

Out [2]:
id      dur      proto  service  state  spkts  dpkts  sbytes  dbytes  rate ... ct_dst_sport_ltm  ct_dst_src_ltm  is_flp_login  ct_flp_cmd  ct_flw_http_mthd  ct_src_ltm  ct_srv_dst  is_sm_ips_ports  attack_cat  label
0  1  0.00011  udp      -      INT      2      0  496      0  9090.0902 ...      1      2      0      0      0      0      1      2      0  Normal  0
1  2  0.00008  udp      -      INT      2      0  1762      0  25000.0001 ...      1      2      0      0      0      0      1      2      0  Normal  0
2  3  0.00005  udp      -      INT      2      0  1068      0  20000.0001 ...      1      3      0      0      0      0      1      3      0  Normal  0
3  4  0.00005  udp      -      INT      2      0  1000      0  16666.6666 ...      1      3      0      0      0      0      2      3      0  Normal  0
4  5  0.00012  udp      -      INT      2      0  2126      0  180000.0025 ...      1      3      0      0      0      0      2      3      0  Normal  0

5 rows x 45 columns

In [3]: test = pd.read_csv("UNSW_NB15_testing-set.csv")
test.head()

Out [3]:
id      dur      proto  service  state  spkts  dpkts  sbytes  dbytes  rate ... ct_dst_sport_ltm  ct_dst_src_ltm  is_flp_login  ct_flp_cmd  ct_flw_http_mthd  ct_src_ltm  ct_srv_dst  is_sm_ips_ports  attack_cat  label
0  1  0.121478  tcp      -      FIN      6      4  258      172  74.087490 ...      1      1      0      0      0      0      1      1      0  Normal  0
1  2  0.649902  tcp      -      FIN      14  38  734  42014  78.473372 ...      1      2      0      0      0      0      1      6      0  Normal  0
2  3  1.621129  tcp      -      FIN      8  16  364  13186  14.170161 ...      1      3      0      0      0      0      2      6      0  Normal  0
3  4  1.681642  tcp      flp  FIN      12  12  628  770  13.671108 ...      1      3      1      1      1      0      2      1      0  Normal  0
4  5  0.444854  tcp      -      FIN      10  6  534  268  33.738262 ...      1      40      0      0      0      0      2      39      0  Normal  0

5 rows x 45 columns

In [4]: train.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 82332 entries, 0 to 82331
Data columns (total 45 columns):
#   Column                Non-Null Count  Dtype  ---
0  id                     82332 non-null  int64
1  dur                    82332 non-null  float64
2  proto                  82332 non-null  object
3  service                 82332 non-null  object
4  state                  82332 non-null  object
5  spkts                  82332 non-null  int64
6  dpkts                  82332 non-null  int64
7  sbytes                 82332 non-null  int64
8  dbytes                 82332 non-null  int64
9  rate                   82332 non-null  float64
10 sttl                   82332 non-null  int64
11 dttl                   82332 non-null  float64
12 sload                  82332 non-null  float64
13 dload                  82332 non-null  float64
14 sloss                  82332 non-null  float64
15 dloss                  82332 non-null  float64
16 sinpkt                 82332 non-null  float64
17 dpktpkt                82332 non-null  float64
18 sjit                   82332 non-null  float64
19 djit                   82332 non-null  float64
20 swin                   82332 non-null  int64
21 scpb                   82332 non-null  int64
22 dcpb                   82332 non-null  int64
23 dwin                   82332 non-null  int64
24 tcprtt                 82332 non-null  float64
25 synack                 82332 non-null  float64
26 ackdat                 82332 non-null  float64
27 smean                  82332 non-null  int64
28 dmean                  82332 non-null  int64
29 trans_depth            82332 non-null  int64
30 response_body_len      82332 non-null  int64
31 ct_srv_src              82332 non-null  int64
32 ct_state_ttl            82332 non-null  int64
33 ct_dst_ltm              82332 non-null  int64
34 ct_src_sport_ltm       82332 non-null  int64
35 ct_dst_sport_ltm       82332 non-null  int64
36 ct_src_ltm              82332 non-null  int64
37 is_flp_login           82332 non-null  int64
38 ct_flp_cmd              82332 non-null  int64
39 ct_flw_http_mthd        82332 non-null  int64
40 ct_src_ltm              82332 non-null  int64
41 ct_srv_dst              82332 non-null  int64
42 is_sm_ips_ports        82332 non-null  int64
43 attack_cat              82332 non-null  object
44 label                   82332 non-null  int64
dtypes: float64(13), int64(38), object(4)
memory usage: 29.3+ MB

In [5]: test.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 175341 entries, 0 to 175340
Data columns (total 45 columns):
#   Column                Non-Null Count  Dtype  ---
0  id                     175341 non-null  int64
1  dur                    175341 non-null  float64
2  proto                  175341 non-null  object
3  service                 175341 non-null  object
4  state                  175341 non-null  object
5  spkts                  175341 non-null  int64
6  dpkts                  175341 non-null  int64
7  sbytes                 175341 non-null  int64
8  dbytes                 175341 non-null  int64
9  rate                   175341 non-null  float64
10 sttl                   175341 non-null  int64
11 dttl                   175341 non-null  float64
12 sload                  175341 non-null  float64
13 dload                  175341 non-null  float64
14 sloss                  175341 non-null  float64
15 dloss                  175341 non-null  float64
16 sinpkt                 175341 non-null  float64
17 dpktpkt                175341 non-null  float64
18 sjit                   175341 non-null  float64
19 djit                   175341 non-null  float64
20 swin                   175341 non-null  int64
21 scpb                   175341 non-null  int64
22 dcpb                   175341 non-null  int64
23 dwin                   175341 non-null  int64
24 tcprtt                 175341 non-null  float64
25 synack                 175341 non-null  float64
26 ackdat                 175341 non-null  float64
27 smean                  175341 non-null  int64
28 dmean                  175341 non-null  int64
29 trans_depth            175341 non-null  int64
30 response_body_len      175341 non-null  int64
31 ct_srv_src              175341 non-null  int64
32 ct_state_ttl            175341 non-null  int64
33 ct_dst_ltm              175341 non-null  int64
34 ct_src_sport_ltm       175341 non-null  int64
35 ct_dst_sport_ltm       175341 non-null  int64
36 ct_src_ltm              175341 non-null  int64
37 is_flp_login           175341 non-null  int64
38 ct_flp_cmd              175341 non-null  int64
39 ct_flw_http_mthd        175341 non-null  int64
40 ct_src_ltm              175341 non-null  int64
41 ct_srv_dst              175341 non-null  int64
42 is_sm_ips_ports        175341 non-null  int64
43 attack_cat              175341 non-null  object
44 label                   175341 non-null  int64
dtypes: float64(13), int64(38), object(4)
memory usage: 60.2+ MB

In [6]: #delete proto
train = train.drop('proto', axis=1)
test = test.drop('proto', axis=1)
train.state = train.state.map({'FIN':0, 'INT':1, 'CON':2, 'REQ':3, 'ACC':4, 'RST':5, 'CLO':6})
train.service = train.service.map({'':0, 'dns':1, 'http':2, 'smtp':3, 'ftp':4, 'ftp-data':5, 'pop3':6, 'ssh':7, 'ssl':8, 'snmp':9, 'dhcp':10, 'radius':11, 'irc':12})
train.state = test.state.map({'FIN':0, 'INT':1, 'CON':2, 'REQ':3, 'ACC':4, 'RST':5, 'CLO':6})
test.service = test.service.map({'':0, 'dns':1, 'http':2, 'smtp':3, 'ftp':4, 'ftp-data':5, 'pop3':6, 'ssh':7, 'ssl':8, 'snmp':9, 'dhcp':10, 'radius':11, 'irc':12})
train.attack_cat = train.attack_cat.map({'Normal':0, 'Generic':1, 'Exploits':2, 'Fuzzers':3, 'DoS':4, 'Reconnaissance':5, 'Analysis':6, 'Backdoor':7, 'Shellcode':8, 'Worms':9})
test.attack_cat = test.attack_cat.map({'Normal':0, 'Generic':1, 'Exploits':2, 'Fuzzers':3, 'DoS':4, 'Reconnaissance':5, 'Analysis':6, 'Backdoor':7, 'Shellcode':8, 'Worms':9})

In [7]: train = train[np.isfinite(train).all()]
test = test[np.isfinite(test).all()]

In [8]: percent_missing = train.isnull().sum() * 100 / len(train)
missing_value_df = pd.DataFrame({'column_name': train.columns, 'percent_missing': percent_missing})
missing_value_df.sort_values('percent_missing', inplace=True)

Out [8]:
column_name  percent_missing
id            0.0
synack        0.0
ackdat        0.0
smean         0.0
dmean         0.0
trans_depth   0.0
response_body_len  0.0
ct_srv_src    0.0
ct_state_ttl  0.0
tcprtt        0.0
ct_dst_ltm    0.0
ct_dst_sport_ltm  0.0
ct_dst_src_ltm  0.0
is_flp_login  0.0
ct_flp_cmd    0.0
ct_flw_http_mthd  0.0
ct_src_ltm    0.0
ct_srv_dst    0.0
is_sm_ips_ports  0.0
ct_src_sport_ltm  0.0
dwin          0.0
dcpb          0.0
scpb          0.0
dur           0.0
service       0.0
state         0.0
spkts         0.0
dpkts         0.0
sbytes        0.0
dbytes        0.0
rate          0.0
sttl          0.0
dttl          0.0
sload         0.0
dload         0.0
sloss         0.0
dloss         0.0
sinpkt        0.0
dpktpkt       0.0
sjit          0.0
djit          0.0
swin          0.0
scpb          0.0
attack_cat    0.0
label         0.0

In [9]: percent_missing = test.isnull().sum() * 100 / len(test)
missing_value_df = pd.DataFrame({'column_name': test.columns, 'percent_missing': percent_missing})
missing_value_df.sort_values('percent_missing', inplace=True)

Out [9]:
column_name  percent_missing
id            0.0
synack        0.0
ackdat        0.0
smean         0.0
dmean         0.0
trans_depth   0.0
response_body_len  0.0
ct_srv_src    0.0
ct_state_ttl  0.0
tcprtt        0.0
ct_dst_ltm    0.0
ct_dst_sport_ltm  0.0
ct_dst_src_ltm  0.0
is_flp
```