In [2]: columns=[] columns =["index", "duration", "protocol\_type", "service", "flag", "src\_bytes", "dst\_bytes", "land" , "wrong\_fragment", "urgent", "hot", "num\_failed\_logins", "logged\_in", "num\_compromised", "root\_she 11", "su\_attempted", "num\_root", "num\_file\_creations", "num\_shells", "num\_access\_files", "num\_outb ound\_cmds", "is\_host\_login", "is\_guest\_login", "count", "srv\_count", "serror\_rate", "srv\_serror\_ra te", "rerror\_rate", "srv\_rerror\_rate", "same\_srv\_rate", "diff\_srv\_rate", "srv\_diff\_host\_rate", "ds t\_host\_count", "dst\_host\_srv\_count", "dst\_host\_same\_srv\_rate", "dst\_host\_diff\_srv\_rate", "dst\_host\_host\_srv\_rate", "dst\_host\_same\_srv\_rate", "dst\_ho st\_same\_src\_port\_rate", "dst\_host\_srv\_diff\_host\_rate", "dst\_host\_serror\_rate", "dst\_host\_srv\_se rror\_rate", "dst\_host\_rerror\_rate", "dst\_host\_srv\_rerror\_rate"] test\_columns=columns train\_columns=columns+['target'] print(len(columns)) 42 In [3]: | df = pd.read\_csv("full.csv", names = train\_columns[1:]) df.head() Out[3]: duration protocol\_type service flag src\_bytes dst\_bytes land wrong\_fragment urgent hot ... dst\_host\_srv\_count ds 0 0 http SF 215 45076 0 0 ... 0 0 ... 1 0 http SF 162 4528 0 0 0 1 tcp http SF 2 0 2 236 1228 0 ... 2032 0 ... 3 0 http SF 233 0 0 3 tcp 4 0 http SF 239 486 0 0 ... 5 rows × 42 columns finding categorical columns and then encoding them In [4]: | numeric\_cols = df.\_get\_numeric\_data().columns categorical\_cols = list(set(df.columns)-set(numeric\_cols)) categorical\_cols.remove('target') print(categorical\_cols) ['service', 'protocol\_type', 'flag'] In [5]: df["protocol\_type"]= df["protocol\_type"].astype('category') df["protocol\_type"] = df["protocol\_type"].cat.codes df["service"] = df["service"].astype('category') df["service"] = df["service"].cat.codes # df2["service"] df["flag"]= df["flag"].astype('category') df["flag"] = df["flag"].cat.codes feature selection using correlation metrics In [6]: corr = df.corr() plt.figure(figsize =(15, 12)) sns.heatmap(corr,cmap='BrBG') plt.show() protocol\_type service src bytes 0.75 dst\_bytes land wrong\_fragment urgent hot num\_failed\_logins - 0.50 logged\_in num\_compromised su attempted num root 0.25 num\_file\_creations num shells num access files num\_outbound\_cmds : is\_host\_login -0.00 is\_guest\_login <sup>-</sup> count srv\_count serror\_rate srv\_serror\_rate rerror rate --0.25 srv\_rerror\_rate same\_srv\_rate diff\_srv\_rate srv\_diff\_host\_rate dst host count -0.50dst\_host\_srv\_count dst\_host\_same\_srv\_rate dst\_host\_diff\_srv\_rate = dst\_host\_same\_src\_port\_rate = dst\_host\_srv\_diff\_host\_rate -0.75dst\_host\_serror\_rate dst host srv serror rate dst\_host\_rerror\_rate dst\_host\_srv\_rerror\_rate selecting highly correlated features and dropping the insignificant ones In [7]: list\_to\_drop=["num\_compromised", "srv\_serror\_rate", "rerror\_rate", "dst\_host\_srv\_serror\_rate", "dst\_host\_serror\_rate", "srv\_rerror\_rate", "dst\_host\_srv\_rerror\_rate", "dst\_host\_same\_srv\_rate" df.drop(list\_to\_drop, axis = 1, inplace = True) normalising data In [8]: y = df[["target"]] X = df.drop(['target'], axis = 1) sc = MinMaxScaler()  $X = sc.fit_transform(X)$ splitting data In [9]: X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size = 0.33, random\_state = 4 print(X\_train.shape, X\_test.shape) print(y\_train.shape, y\_test.shape) (3281948, 33) (1616483, 33) (3281948, 1) (1616483, 1) In [10]: def get\_score(yactual,ypred): recall=(recall\_score(yactual, ypred,average='micro')) accuracy=(accuracy\_score(yactual, ypred)) f1=(f1\_score(yactual, ypred, average='micro')) precision=(precision\_score(yactual, ypred, average='micro')) return recall, accuracy, f1, precision **Training with different ML models** In [11]: # DecisionTreeClassifier from sklearn.tree import DecisionTreeClassifier clfd = DecisionTreeClassifier(criterion ="entropy", max\_depth = 4) start\_time = time.time() clfd.fit(X\_train, y\_train.values.ravel()) end\_time = time.time() print("Training time: ", end\_time-start\_time) Training time: 18.238346099853516 In [12]: start\_time = time.time() y\_test\_pred = clfd.predict(X\_test) end\_time = time.time() print("Testing time: ", end\_time-start\_time) Testing time: 0.4967329502105713 In [13]: recall, accuracy, f1, precision=get\_score(y\_test, y\_test\_pred) print("Performance in model: DecisionTree") print("Recall:", recall) print("Accuracy:", accuracy) print("f1 score:",f1) print("precision:", precision) Performance in model: DecisionTree Recall: 0.9949414871668926 Accuracy: 0.9949414871668926 f1 score: 0.9949414871668926 precision: 0.9949414871668926 In [14]: # RandomForestClassifier from sklearn.ensemble import RandomForestClassifier clfr = RandomForestClassifier(n\_estimators = 30) start\_time = time.time() clfr.fit(X\_train, y\_train.values.ravel()) end\_time = time.time() print("Training time: ", end\_time-start\_time) Training time: 162.05097031593323 In [15]: start\_time = time.time() y\_test\_pred = clfr.predict(X\_test) end\_time = time.time() print("Testing time: ", end\_time-start\_time) Testing time: 10.296815156936646 In [16]: recall, accuracy, f1, precision=get\_score(y\_test, y\_test\_pred) print("Performance in model: RandomForestClassifier") print("Recall:", recall) print("Accuracy:", accuracy) print("f1 score:",f1) print("precision:", precision) Performance in model: RandomForestClassifier Recall: 0.9998979265479438 Accuracy: 0.9998979265479438 f1 score: 0.9998979265479438 precision: 0.9998979265479438 In [17]: # LogisticRegression from sklearn.linear\_model import LogisticRegression clfl = LogisticRegression(max\_iter = 1000000) start\_time = time.time() clfl.fit(X\_train, y\_train.values.ravel()) end\_time = time.time() print("Training time: ", end\_time-start\_time) Training time: 2332.4507892131805 In [18]: start\_time = time.time() y\_test\_pred = clfl.predict(X\_test) end\_time = time.time() print("Testing time: ", end\_time-start\_time) Testing time: 1.1773362159729004 In [19]: recall, accuracy, f1, precision=get\_score(y\_test, y\_test\_pred) print("Performance in model: LogisticRegression") print("Recall:", recall) print("Accuracy:", accuracy) print("f1 score:",f1) print("precision:", precision) Performance in model: LogisticRegression Recall: 0.9986575794487168 Accuracy: 0.9986575794487168 f1 score: 0.9986575794487168 precision: 0.9986575794487168 Selecting the best model of the above based on train-test scores - Random Forest classifier Training RFC on the complete data In [21]: # RandomForestClassifier from sklearn.ensemble import RandomForestClassifier clfr = RandomForestClassifier(n\_estimators = 50) start\_time = time.time() clfr.fit(X, y.values.ravel()) # clfr.fit(X\_train, y\_train.values.ravel()) end\_time = time.time() print("Training time: ", end\_time-start\_time) Training time: 181.85842990875244 Reading test data, data cleaning, encoding and feature selection In [22]: | df2=pd.read\_csv("test.csv", names = test\_columns) df2.drop('index', axis = 1, inplace = True) df2.head() /opt/anaconda3/envs/n/lib/python3.7/site-packages/IPython/core/interactiveshell.py:3147: Dtyp eWarning: Columns (26,27,28,30,31,34,36,37,38,39,40,41) have mixed types. Specify dtype option on import or set low\_memory=False. interactivity=interactivity, compiler=compiler, result=result) Out[22]: duration protocol\_type service flag src\_bytes dst\_bytes land wrong\_fragment urgent hot ... dst\_host\_count dst\_hc 0 0.1 0.3 0.4 ... udp private SF 105 146 0.2 255 0.0 0.0 ... 1 0 udp private SF 105 146 0.0 0.0 255 0 SF 105 146 0.0 0.0 0.0 0.0 ... 255 udp private 3 0 private SF 105 146 0.0 0.0 0.0 0.0 ... 255 private SF 105 146 0.0 0.0 0.0 0.0 ... 255 udp 5 rows × 41 columns In [23]: df2=df2.iloc[1:] df2.drop(list\_to\_drop, axis = 1, inplace = True) df2["flag"]= df2["flag"].astype('category') df2["flag"] = df2["flag"].cat.codes df2["service"]= df2["service"].astype('category') df2["service"] = df2["service"].cat.codes df2["protocol\_type"]= df2["protocol\_type"].astype('category') df2["protocol\_type"] = df2["protocol\_type"].cat.codes In [24]: df2\_x=sc.fit\_transform(df2) **Predicting on the test data** In [25]: start\_time = time.time() y\_test\_pred = clfr.predict(df2\_x) end\_time = time.time() print("Testing time: ", end\_time-start\_time) Testing time: 3.469721794128418 In [26]: df=pd.DataFrame(y\_test\_pred) Out[26]: **0** normal. **2** normal. 3 normal. 4 normal. **311023** normal. **311024** normal. **311025** normal. **311026** normal. **311027** normal. 311028 rows × 1 columns In [27]: df5=pd.DataFrame([["normal."]]) df6=pd.concat([df5,df]) df6 # df5 Out[27]: 0 **0** normal. 0 normal. **1** normal. 2 normal. **3** normal. **311023** normal. **311024** normal. **311025** normal. **311026** normal. **311027** normal. 311029 rows × 1 columns In [28]: df6.insert(0, "index", range(311029), True) Out[28]: index 0 normal. 1 normal. 2 normal. 3 normal. 4 normal. **311023** 311024 normal. **311024** 311025 normal. **311025** 311026 normal. **311026** 311027 normal. **311027** 311028 normal. 311029 rows × 2 columns In [29]: | df7=df6.rename(columns={0:"target"},inplace=False) df7 Out[29]: index target 0 normal. 1 normal. 2 normal. 3 normal. 4 normal. **311023** 311024 normal. **311024** 311025 normal. **311025** 311026 normal. **311026** 311027 normal. **311027** 311028 normal. 311029 rows × 2 columns In [ ]: In [30]: # df7=pd.read\_csv("testLabel2.csv") df7.to\_csv("testLabel6.csv", index=False) In [31]: set(df7["target"]) Out[31]: {'back.', 'ipsweep.', 'land.', 'neptune.',

'nmap.',
'normal.',
'pod.',

In [ ]:

'portsweep.',
'satan.',
'smurf.',
'teardrop.',
'warezclient.',
'warezmaster.'}

In [1]: import csv

import os

import time

import pandas as pd
import numpy as np

import seaborn as sns

defining column names

import matplotlib.pyplot as plt

from sklearn.model\_selection import train\_test\_split

from sklearn.metrics import accuracy\_score,f1\_score,recall\_score,precision\_score

from sklearn.preprocessing import MinMaxScaler