# COMPUTER NETWORKS LAB PROJECT

## NETWORK INTRUSION DETECTION USING MACHINE LEARNING

Group Members:

Vatsal Gupta (17104060)
Arjav Jain (17104070)
Ananya Sharma (17104038)

# Title: Network Intrusion Detection using Machine Learning

## Problem Statement

With the rapid development of information technology in the past two decades, computer networks are being widely used by industry, business and various fields of human life. Therefore, building reliable networks has become an essential task for IT administrators. However, the rapid development of the IT sector has brought to the forefront several challenges towards building reliable computer networks. There are many types of attacks threatening the availability, integrity and confidentiality of computer networks. The Denial of Service attacks (DoS), Remote to Local (R2L) and, User to Root (U2R) attacks are considered to be some of the most harmful attacks to any network.

## Introduction and Proposed Model

With the advancement in technology, millions of people are now connected through one or other form of a network where they share lots of valuable data. Hence the need for security to safeguard data integrity and confidentiality is increased rapidly. Although effort has been made to secure data transmission at the same time, the attack technique for breaching the network continued to evolve. Thus it leads to the need for such a system which can adapt with this ever-changing attack technique. In this paper, we have purposed a system which is based on machine learning. We aim to find the most suitable machine learning algorithm which can predict the type of network attack with the highest accuracy. The algorithms which we have compared are Decision Tree, Naïve Bayes, K Nearest Neighbor (KNN), Random Forest and Logistic Regression. The dataset used for training the model is NSL - KDD revised dataset. Machine learning has been used to detect the intrusion attacks due to its flexibility; for example, if any new type of attack is developed in future, the system can be trained for predicting that attack. In our project, we have made a knowledge-based intrusion detection system which is also known as the anomaly-based system. It registers the anomalies and in future predicts such malicious network to send out an alert. This way, the network can disconnect to such a connection and then have only secured connections.

## Hardware Requirements:

- Minimum memory requirement for smooth working of the code: 4 GB RAM
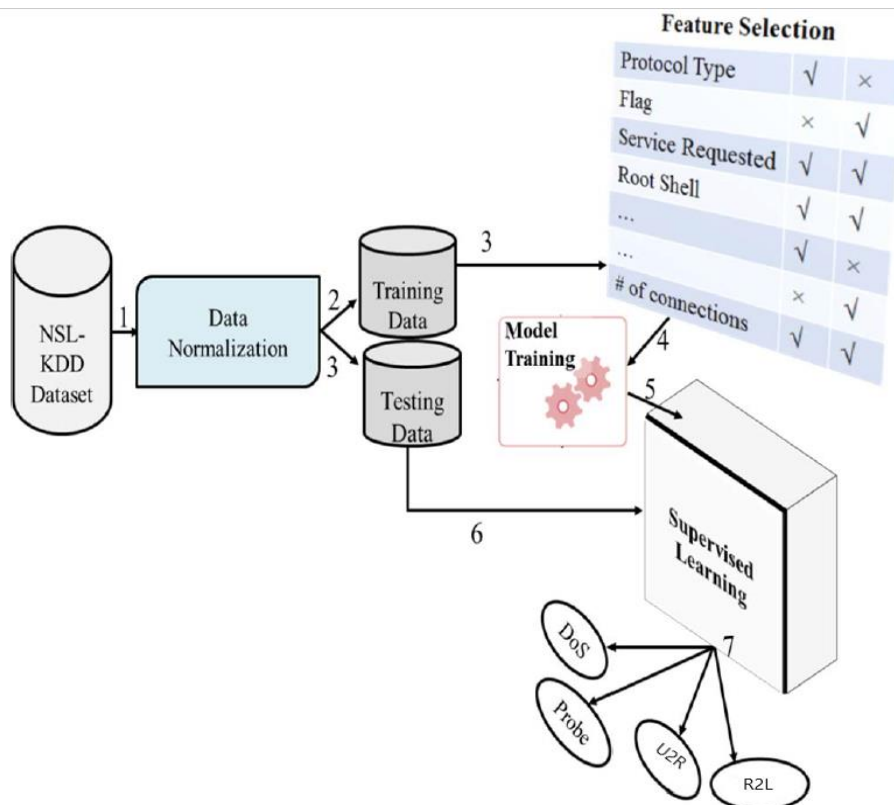- Processor: Preferably intel i5 (atleast dual core) and later

## Software Requirements:

- Python 3.x
- Any python IDE (preferably Spyder IDE 3.3.2 or PyCharm)

## Packages/Modules Used (can be installed using pip installer)

- warnings
- matplotlib
- pandas
- numpy
- seaborn
- sklearn
- imblearn

## Data Flow Diagram

# Dataset Description

The machine learning algorithms for intrusion detection have been applied on the freely available NSL – KDD revised dataset. It is an easily accessible dataset that is being relied upon by many researchers for their works on intrusion detection. The dataset contains many attack types like the DOS, U2R, R2L, Probe and normal (no attack). There are 21 types of attacks inside the main categories mentioned above.

| Categories of Attack | Attack name | Number of instances |
| --- | --- | --- |
| DOS | SMURF | 2807886 |
| | NEPTUNE | 1072017 |
| | Back | 2203 |
| | POD | 264 |
| | Teardrop | 979 |
| U2R | Buffer overflow | 30 |
| | Load Module | 9 |
| | PERL | 3 |
| | Rootkit | 10 |
| R2L | FTP Write | 8 |
| | Guess Passwd | 53 |
| | IMAP | 12 |
| | MulitHop | 7 |
| | PHF | 4 |
| | SPY | 2 |
| | Warez client | 1020 |
| | Warez Master | 20 |
| PROBE | IPSWEEP | 12481 |
| | NMAP | 2316 |
| | PORTSWEEP | 10413 |
| | SATAN | 15892 |
| normal | | 972781 |

The dataset contains a total of 41 attributes which could be used to determine if the attack is malicious or not at all an attack.

Description of flag values in the dataset is given in the following table:

| Flag | Description |
|---|---|
| RSTOS0 | Originator sent a SYN followed by a RST, never see a SYN ACK from the responder |
| RSTR | Established, responder aborted |
| RSTO | Connection established, originator aborted (sent a RST) |
| OTH | No SYN seen, just midstream traffic (a "partial connection" that was not later closed) |
| REJ | Connection attempt rejected |
| S0 | Connection attempt seen, no reply |
| S1 | Connection established, not terminated |
| S2 | Connection established and close attempt by originator seen (but no reply from responder) |
| S3 | Connection established and close attempt by responder seen (but no reply from originator) |
| SF | Normal establishment and termination |
| SH | Originator sent a SYN followed by a FIN (finish 'flag'), never saw a SYN ACK from the responder (hence the connection was "half" open) |

## Working Model

The NSL - KDD is a huge dataset, and the one that we have used in our research is under the folder 'Dataset'. Our aim is to not only to find the best algorithm suited for the intrusion detection problem but also to implement it using the programming language python. The first process of applying machine learning is always pre-processing the data. This process involves various steps like 1) removing the redundant rows from the dataset 2) to see whether there are any missing values and then to remove those corresponding rows too.
The next process is to use this dataset and put it across various machine-learning algorithms that might give good results by correctly classifying the instances.

From our results, it can be inferred that the best algorithm that can be used for the network intrusion detection is the Random Forest. Random Forest algorithm is a classification algorithm based on ensemble learning. It works by building multiple decision trees at training, and the developed decision trees form the output function. The drawback of this algorithm is that is computation intensive and consequently has high computational time.

This program, when developed fully, will act as a filter to determine if a network is secure and will continuously learn from its own series of data making it better and stronger with each type of attack.

**Attack Class Distribution**



**Screenshots of the code:**

```python
1 # Ignore warnings
2 import warnings
3 warnings.filterwarnings('ignore')
4
5 # import relevant modules
6 import matplotlib
7 import matplotlib.pyplot as plt
8 import pandas as pd
9 import numpy as np
10 import seaborn as sns
11 import sklearn
12 import imblearn
13
14
15 # Settings
16 pd.set_option('display.max_columns', None)
17 np.set_printoptions(threshold=0)
18 np.set_printoptions(precision=3)
19 sns.set(style="darkgrid")
20 plt.rcParams['axes.labelsize'] = 14
21 plt.rcParams['xtick.labelsize'] = 12
22 plt.rcParams['ytick.labelsize'] = 12
23
24 # Dataset field names
25 datacols = ["duration","protocol_type","service","flag","src_bytes",
26     "dst_bytes","land","wrong_fragment","urgent","hot","num_failed_logins",
27     "logged_in","num_compromised","root_shell","su_attempted","num_root",
28     "num_file_creations","num_shells","num_access_files","num_outbound_cmds",
29     "is_host_login","is_guest_login","count","srv_count","serror_rate",
30     "srv_serror_rate","rerror_rate","srv_rerror_rate","same_srv_rate",
31     "diff_srv_rate","srv_diff_host_rate","dst_host_count","dst_host_srv_count",
32     "dst_host_same_srv_rate","dst_host_diff_srv_rate","dst_host_same_src_port_rate",
33     "dst_host_srv_diff_host_rate","dst_host_serror_rate","dst_host_srv_serror_rate",
34     "dst_host_rerror_rate","dst_host_srv_rerror_rate","attack", "last_flag"]
```

---

```python
31     "diff_srv_rate","srv_diff_host_rate","dst_host_count","dst_host_srv_count",
32     "dst_host_same_srv_rate","dst_host_diff_srv_rate","dst_host_same_src_port_rate",
33     "dst_host_srv_diff_host_rate","dst_host_serror_rate","dst_host_srv_serror_rate",
34     "dst_host_rerror_rate","dst_host_srv_rerror_rate","attack", "last_flag"]
35
36 # Load NSL_KDD train dataset
37 dfkdd_train = pd.read_table("KDDTrain.txt", sep=",", names=datacols) # change path to where the dataset is located.
38 dfkdd_train = dfkdd_train.iloc[:,:-1] # removes an unwanted extra field
39
40 # Load NSL_KDD test dataset
41 dfkdd_test = pd.read_table("KDDTest.txt", sep=",", names=datacols)
42 dfkdd_test = dfkdd_test.iloc[:,:-1]
43
44 mapping = {'ipsweep': 'Probe','satan': 'Probe','nmap': 'Probe','portsweep': 'Probe','saint': 'Probe','mscan': 'Probe',
45         'teardrop': 'DoS','pod': 'DoS','land': 'DoS','back': 'DoS','neptune': 'DoS','smurf': 'DoS','mailbomb': 'DoS',
46         'udpstorm': 'DoS','apache2': 'DoS','processtable': 'DoS',
47         'perl': 'U2R','loadmodule': 'U2R','rootkit': 'U2R','buffer_overflow': 'U2R','xterm': 'U2R','ps': 'U2R',
48         'sqlattack': 'U2R','httptunnel': 'U2R',
49         'ftp_write': 'R2L','phf': 'R2L','guess_passwd': 'R2L','warezmaster': 'R2L','warezclient': 'R2L','imap': 'R2L',
50         'spy': 'R2L','multihop': 'R2L','named': 'R2L','snmpguess': 'R2L','worm': 'R2L','snmpgetattack': 'R2L',
51         'xsnoop': 'R2L','xlock': 'R2L','sendmail': 'R2L',
52         'normal': 'Normal'
53         }
54
55 # Apply attack class mappings to the dataset
56 dfkdd_train['attack_class'] = dfkdd_train['attack'].apply(lambda v: mapping[v])
57 dfkdd_test['attack_class'] = dfkdd_test['attack'].apply(lambda v: mapping[v])
58
59 # Drop attack field and 'num_output_cmds' field from both training and testing data
60 dfkdd_train.drop(['attack'], axis=1, inplace=True)
61 dfkdd_test.drop(['attack'], axis=1, inplace=True)
62 dfkdd_train.drop(['num_outbound_cmds'], axis=1, inplace=True)
63 dfkdd_test.drop(['num_outbound_cmds'], axis=1, inplace=True)
64
```

```python
62 dfkdd_train.drop(['num_outbound_cmds'], axis=1, inplace=True)
63 dfkdd_test.drop(['num_outbound_cmds'], axis=1, inplace=True)
64
65 # Attack Class Distribution
66 attack_class_freq_train = dfkdd_train[['attack_class']].apply(lambda x: x.value_counts())
67 attack_class_freq_test = dfkdd_test[['attack_class']].apply(lambda x: x.value_counts())
68 attack_class_freq_train['frequency_percent_train'] = round((100 * attack_class_freq_train / attack_class_freq_train.sum()),2)
69 attack_class_freq_test['frequency_percent_test'] = round((100 * attack_class_freq_test / attack_class_freq_test.sum()),2)
70
71 attack_class_dist = pd.concat([attack_class_freq_train,attack_class_freq_test], axis=1)
72 attack_class_dist
73
74 # Attack class bar plot
75 plot = attack_class_dist[['frequency_percent_train', 'frequency_percent_test']].plot(kind="bar");
76 plot.set_title("Attack Class Distribution", fontsize=20);
77 plot.grid(color='lightgray', alpha=0.5);
78
79 from sklearn.preprocessing import StandardScaler
80 scaler = StandardScaler()
81
82 # Standardisation - extract numerical attributes and scale it to have zero mean and unit variance
83 cols = dfkdd_train.select_dtypes(include=['float64','int64']).columns
84 sc_train = scaler.fit_transform(dfkdd_train.select_dtypes(include=['float64','int64']))
85 sc_test = scaler.fit_transform(dfkdd_test.select_dtypes(include=['float64','int64']))
86
87 # turn the result back to a dataframe
88 sc_traindf = pd.DataFrame(sc_train, columns = cols)
89 sc_testdf = pd.DataFrame(sc_test, columns = cols)
90
91 from sklearn.preprocessing import LabelEncoder
92 encoder = LabelEncoder()
93
94 # extract categorical attributes from both training and test sets
95 cattrain = dfkdd_train.select_dtypes(include=['object']).copy()
```

```python
93
94 # extract categorical attributes from both training and test sets
95 cattrain = dfkdd_train.select_dtypes(include=['object']).copy()
96 cattest = dfkdd_test.select_dtypes(include=['object']).copy()
97
98 # encode the categorical attributes
99 traincat = cattrain.apply(encoder.fit_transform)
100 testcat = cattest.apply(encoder.fit_transform)
101
102 # separate target column from encoded data
103 enctrain = traincat.drop(['attack_class'], axis=1)
104 enctest = testcat.drop(['attack_class'], axis=1)
105
106 cat_Ytrain = traincat[['attack_class']].copy()
107 cat_Ytest = testcat[['attack_class']].copy()
108
109
110 from imblearn.over_sampling import RandomOverSampler # Over Sampling is done for better prediciton
111 from collections import Counter
112
113 # define columns and extract encoded train set for sampling
114 sc_traindf = dfkdd_train.select_dtypes(include=['float64','int64'])
115 refclasscol = pd.concat([sc_traindf, enctrain], axis=1).columns
116 refclass = np.concatenate((sc_train, enctrain.values), axis=1)
117 X = refclass
118
119 # reshape target column to 1D array shape
120 c, r = cat_Ytest.values.shape
121 y_test = cat_Ytest.values.reshape(c,)
122
123 c, r = cat_Ytrain.values.shape
124 y = cat_Ytrain.values.reshape(c,)
125
126 # apply the random over-sampling
```

```python
124  y = cat_Ytrain.values.reshape(c,)
125
126  # apply the random over-sampling
127  ros = RandomOverSampler(random_state=42)
128  X_res, y_res = ros.fit_sample(X, y)
129  print('Original dataset shape {}'.format(Counter(y)))
130  print('Resampled dataset shape {}'.format(Counter(y_res)))
131
132
133  from sklearn.ensemble import RandomForestClassifier
134  rfc = RandomForestClassifier();
135
136  # fit random forest classifier on the training set
137  rfc.fit(X_res, y_res);
138  # extract important features
139  score = np.round(rfc.feature_importances_,3)
140  importances = pd.DataFrame({'feature':refclasscol,'importance':score})
141  importances = importances.sort_values('importance',ascending=False).set_index('feature')
142  # plot importances
143  plt.rcParams['figure.figsize'] = (11, 4)
144  importances.plot.bar();
145
146  from sklearn.feature_selection import RFE
147  import itertools
148  rfc = RandomForestClassifier()
149
150  # create the RFE model and select 10 attributes
151  rfe = RFE(rfc, n_features_to_select=10)
152  rfe = rfe.fit(X_res, y_res)
153
154  # summarize the selection of the attributes
155  feature_map = [(i, v) for i, v in itertools.zip_longest(rfe.get_support(), refclasscol)]
156  selected_features = [v for i, v in feature_map if i==True]
157  # define columns to new dataframe
```

```python
153
154  # summarize the selection of the attributes
155  feature_map = [(i, v) for i, v in itertools.zip_longest(rfe.get_support(), refclasscol)]
156  selected_features = [v for i, v in feature_map if i==True]
157  # define columns to new dataframe
158  newcol = list(refclasscol)
159  newcol.append('attack_class')
160
161  # add a dimension to target
162  new_y_res = y_res[:, np.newaxis]
163
164  # create a dataframe from sampled data
165  res_arr = np.concatenate((X_res, new_y_res), axis=1)
166  res_df = pd.DataFrame(res_arr, columns = newcol)
167
168  # create test dataframe
169  reftest = pd.concat([sc_testdf, testcat], axis=1)
170  reftest['attack_class'] = reftest['attack_class'].astype(np.float64)
171  reftest['protocol_type'] = reftest['protocol_type'].astype(np.float64)
172  reftest['flag'] = reftest['flag'].astype(np.float64)
173  reftest['service'] = reftest['service'].astype(np.float64)
174
175  res_df.shape
176  reftest.shape
177
178  from collections import defaultdict
179  classdict = defaultdict(list)
180
181  # create two-target classes (normal class and an attack class)
182  attacklist = [('DoS', 0.0), ('Probe', 2.0), ('R2L', 3.0), ('U2R', 4.0)]
183  normalclass = [('Normal', 1.0)]
184
185  def create_classdict():
186      '''This function subdivides train and test dataset into two-class attack labels'''
```

```python
187     for j, k in normalclass:
188         for i, v in attacklist:
189             restrain_set = res_df.loc[(res_df['attack_class'] == k) | (res_df['attack_class'] == v)]
190             classdict[j +'_' + i].append(restrain_set)
191             # test Labels
192             reftest_set = reftest.loc[(reftest['attack_class'] == k) | (reftest['attack_class'] == v)]
193             classdict[j +'_' + i].append(reftest_set)
194
195 create_classdict()
196
197 for k, v in classdict.items():
198     k
199
200 pretrain = classdict['Normal_DoS'][0]
201 pretest = classdict['Normal_DoS'][1]
202 grpclass = 'Normal_DoS'
203
204 from sklearn.preprocessing import OneHotEncoder
205 enc = OneHotEncoder()
206
207 Xresdf = pretrain
208 newtest = pretest
209
210 Xresdfnew = Xresdf[selected_features]
211 Xresdfnum = Xresdfnew.drop(['service'], axis=1)
212 Xresdfcat = Xresdfnew[['service']].copy()
213
214 Xtest_features = newtest[selected_features]
215 Xtestdfnum = Xtest_features.drop(['service'], axis=1)
216 Xtestcat = Xtest_features[['service']].copy()
217
218
219 # Fit train data
220 enc.fit(Xresdfcat)
```

confusion    1 of 8

Permissions: RW    End-of-lines: CRLF    Encoding: ASCII    Line: 23    Column: 1    Memory: 68 %

```python
218
219 # Fit train data
220 enc.fit(Xresdfcat)
221
222 # Transform train data
223 X_train_1hotenc = enc.transform(Xresdfcat).toarray()
224
225 # Transform test data
226 X_test_1hotenc = enc.transform(Xtestcat).toarray()
227
228 X_train = np.concatenate((Xresdfnum.values, X_train_1hotenc), axis=1)
229 X_test = np.concatenate((Xtestdfnum.values, X_test_1hotenc), axis=1)
230
231 y_train = Xresdf[['attack_class']].copy()
232 c, r = y_train.values.shape
233 Y_train = y_train.values.reshape(c,)
234
235 y_test = newtest[['attack_class']].copy()
236 c, r = y_test.values.shape
237 Y_test = y_test.values.reshape(c,)
238
239 from sklearn.naive_bayes import BernoulliNB
240 from sklearn import tree
241 from sklearn.model_selection import cross_val_score
242 from sklearn.neighbors import KNeighborsClassifier
243 from sklearn.linear_model import LogisticRegression
244 from sklearn.ensemble import VotingClassifier
245
246 # Train KNeighborsClassifier Model
247 KNN_Classifier = KNeighborsClassifier(n_jobs=-1)
248 KNN_Classifier.fit(X_train, Y_train);
249
250 # Train LogisticRegression Model
251 LGR_Classifier = LogisticRegression(n_jobs=-1, random_state=0)
```

confusion    1 of 8

Permissions: RW    End-of-lines: CRLF    Encoding: ASCII    Line: 23    Column: 1    Memory: 68 %

```python
249
250  # Train LogisticRegression Model
251  LGR_Classifier = LogisticRegression(n_jobs=-1, random_state=0)
252  LGR_Classifier.fit(X_train, Y_train);
253
254  # Train Gaussian Naive Baye Model
255  BNB_Classifier = BernoulliNB()
256  BNB_Classifier.fit(X_train, Y_train);
257
258  # Train Decision Tree Model
259  DTC_Classifier = tree.DecisionTreeClassifier(criterion='entropy', random_state=0)
260  DTC_Classifier.fit(X_train, Y_train);
261
262  from sklearn import metrics
263
264  models = []
265
266  models.append(('Naive Bayes Classifier', BNB_Classifier))
267  models.append(('Decision Tree Classifier', DTC_Classifier))
268  models.append(('KNeighborsClassifier', KNN_Classifier))
269  models.append(('LogisticRegression', LGR_Classifier))
270
271
272  for i, v in models:
273      scores = cross_val_score(v, X_train, Y_train, cv=10)
274      accuracy = metrics.accuracy_score(Y_train, v.predict(X_train))
275      confusion_matrix = metrics.confusion_matrix(Y_train, v.predict(X_train))
276      classification = metrics.classification_report(Y_train, v.predict(X_train))
277      print()
278      print('============================== {} {} Model Evaluation =============================='.format(grpclass, i))
279      print()
280      print ("Cross Validation Mean Score:" "\n", scores.mean())
281      print()
282      print ("Model Accuracy:" "\n", accuracy)
```

---

```python
270
271
272  for i, v in models:
273      scores = cross_val_score(v, X_train, Y_train, cv=10)
274      accuracy = metrics.accuracy_score(Y_train, v.predict(X_train))
275      confusion_matrix = metrics.confusion_matrix(Y_train, v.predict(X_train))
276      classification = metrics.classification_report(Y_train, v.predict(X_train))
277      print()
278      print('============================== {} {} Model Evaluation =============================='.format(grpclass, i))
279      print()
280      print ("Cross Validation Mean Score:" "\n", scores.mean())
281      print()
282      print ("Model Accuracy:" "\n", accuracy)
283      print()
284      print("Confusion matrix:" "\n", confusion_matrix) # measures false negatives and false positives
285      print()
286      print("Classification report:" "\n", classification)
287      print()
288
289  for i, v in models:
290      accuracy = metrics.accuracy_score(Y_test, v.predict(X_test))
291      confusion_matrix = metrics.confusion_matrix(Y_test, v.predict(X_test))
292      classification = metrics.classification_report(Y_test, v.predict(X_test))
293      print()
294      print('============================== {} {} Model Test Results =============================='.format(grpclass, i))
295      print()
296      print ("Model Accuracy:" "\n", accuracy)
297      print()
298      print("Confusion matrix:" "\n", confusion_matrix)
299      print()
300      print("Classification report:" "\n", classification)
301      print()
302
303
```

# Screenshots of the results:

```
============================== Normal_DoS LogisticRegression Model Test Results
==============================

Model Accuracy:
 0.8421573766672491

Confusion matrix:
 [[5963 1495]
 [1215 8496]]

Classification report:
              precision    recall  f1-score   support

         0.0       0.83      0.80      0.81      7458
         1.0       0.85      0.87      0.86      9711

    accuracy                           0.84     17169
   macro avg       0.84      0.84      0.84     17169
weighted avg       0.84      0.84      0.84     17169


============================== Normal_DoS Naive Baye Classifier Model Evaluation
==============================

Cross Validation Mean Score:
 0.9737760721391011

Model Accuracy:
 0.9737686173767133

Confusion matrix:
 [[65346  1997]
 [ 1536 65807]]

Classification report:
              precision    recall  f1-score   support

         0.0       0.98      0.97      0.97     67343
         1.0       0.97      0.98      0.97     67343

    accuracy                           0.97    134686
   macro avg       0.97      0.97      0.97    134686
weighted avg       0.97      0.97      0.97    134686
```

```
============================== Normal_DoS Decision Tree Classifier Model Evaluation
==============================

Cross Validation Mean Score:
 0.9997698368967857

Model Accuracy:
 0.9999480272634127

Confusion matrix:
 [[67343     0]
 [    7 67336]]

Classification report:
              precision    recall  f1-score   support

         0.0       1.00      1.00      1.00     67343
         1.0       1.00      1.00      1.00     67343

    accuracy                           1.00    134686
   macro avg       1.00      1.00      1.00    134686
weighted avg       1.00      1.00      1.00    134686
```

```
============================== Normal_DoS KNeighborsClassifier Model Evaluation
==============================

Cross Validation Mean Score:
 0.9965698163470991

Model Accuracy:
 0.9977577476500898

Confusion matrix:
 [[67287    56]
 [  246 67097]]

Classification report:
              precision    recall  f1-score   support

         0.0       1.00      1.00      1.00     67343
         1.0       1.00      1.00      1.00     67343

    accuracy                           1.00    134686
   macro avg       1.00      1.00      1.00    134686
weighted avg       1.00      1.00      1.00    134686
```

# References:

**Research Papers Referred**

- Kevric, J., Jukic, S., & Subasi, A. (2017). An effective combining classifier approach using tree algorithms for network intrusion detection. *Neural Computing and Applications*, *28*(1), 1051-1058.

- Lee, C. H., Su, Y. Y., Lin, Y. C., & Lee, S. J. (2017, September). Machine learning based network intrusion detection. In *2017 2nd IEEE International Conference on Computational Intelligence and Applications (ICCIA)* (pp. 79-83). IEEE.

- Taher, K. A., Jisan, B. M. Y., & Rahman, M. M. (2019, January). Network Intrusion Detection using Supervised Machine Learning Technique with Feature Selection. In *2019 International Conference on Robotics, Electrical and Signal Processing Techniques (ICREST)* (pp. 643-646). IEEE.

## For Dataset

https://www.unb.ca/cic/datasets/nsl.html