

## Import modules

```
In [1]: 1 %matplotlib inline
2 import matplotlib
3 import matplotlib.pyplot as plt
4 import pandas as pd
5 import numpy as np
6 import seaborn as sns
7 import sklearn
8 import imblearn
9 import sys
10
11 # Ignore warnings
12 import warnings
13 warnings.filterwarnings('ignore')
14
15 # Config
16 pd.set_option('display.max_columns', None)
17 np.set_printoptions(threshold=sys.maxsize)
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 print("pandas : {}".format(pd.__version__))
25 print("numpy : {}".format(np.__version__))
26 print("matplotlib : {}".format(matplotlib.__version__))
27 print("seaborn : {}".format(sns.__version__))
28 print("sklearn : {}".format(sklearn.__version__))
29 print("imblearn : {}".format(imblearn.__version__))
```

```
pandas : 1.2.4
numpy : 1.20.1
matplotlib : 3.3.4
seaborn : 0.11.1
sklearn : 0.24.1
imblearn : 0.8.0
```

## Load DataSet

```
In [2]: 1 # Dataset field names
2 datacols = ["duration", "protocol_type", "service", "flag", "src_bytes",
3             "dst_bytes", "land", "wrong_fragment", "urgent", "hot", "num_failed_logins",
4             "logged_in", "num_compromised", "root_shell", "su_attempted", "num_root",
5             "num_file_creations", "num_shells", "num_access_files", "num_outbound_cmds",
6             "is_host_login", "is_guest_login", "count", "srv_count", "error_rate",
7             "srv_error_rate", "error_rate", "srv_error_rate", "same_srv_rate",
8             "diff_srv_rate", "srv_diff_host_rate", "dst_host_count", "dst_host_srv_coun
9             "dst_host_same_srv_rate", "dst_host_diff_srv_rate", "dst_host_same_src_por
10            "dst_host_srv_diff_host_rate", "dst_host_error_rate", "dst_host_srv_serro
11            "dst_host_error_rate", "dst_host_srv_error_rate", "attack", "last_flag"]
12
13 # Load train dataset
14 dfkdd_train = pd.read_table("./NSL_KDD_dataset/KDDTrain.txt", sep=",", names=
15 dfkdd_train = dfkdd_train.iloc[:, :-1] # removes an unwanted extra field
16
17 # Load test dataset
18 dfkdd_test = pd.read_table("./NSL_KDD_dataset/KDDTest.txt", sep=",", names=d
19 dfkdd_test = dfkdd_test.iloc[:, :-1]
```

## Train dataset

```
In [3]: 1 # View train data
2 dfkdd_train.head(3)
3
4 # train set dimension
5 print('Train set dimension: {} rows, {} columns'.format(dfkdd_train.shape[0]
```

Train set dimension: 125973 rows, 42 columns

## Test dataset

```
In [4]: 1 # View test data
2 dfkdd_test.head(3)
3
4 # test set dimension
5 print('Test set dimension: {} rows, {} columns'.format(dfkdd_test.shape[0],
```

Test set dimension: 22481 rows, 42 columns

## Data Preprocessing

```
In [5]: 1 # Map attack type to attack group
2 mapping = {'ipsweep': 'Probe', 'satan': 'Probe', 'nmap': 'Probe', 'portsweep':
3           'teardrop', 'DoS', 'pod': 'DoS', 'land': 'DoS', 'back': 'DoS', 'neptune'
4           'udpstorm': 'DoS', 'apache2': 'DoS', 'processtable': 'DoS',
5           'perl': 'U2R', 'loadmodule': 'U2R', 'rootkit': 'U2R', 'buffer_overflow'
6           'sqlattack': 'U2R', 'httptunnel': 'U2R',
7           'ftp_write': 'R2L', 'phf': 'R2L', 'guess_passwd': 'R2L', 'warezmaster':
8           'spy': 'R2L', 'multihop': 'R2L', 'named': 'R2L', 'snmpguess': 'R2L', 'wo
9           'xsnoop': 'R2L', 'xlock': 'R2L', 'sendmail': 'R2L',
10          'normal': 'Normal'
11          }
```

```
In [6]: 1 # Apply attack type mappings to the dataset
2 dfkdd_train['attack_class'] = dfkdd_train['attack'].apply(lambda v: mapping[
3 dfkdd_test['attack_class'] = dfkdd_test['attack'].apply(lambda v: mapping[v])
```

```
In [7]: 1 # Drop attack field from both train and test data
2 dfkdd_train.drop(['attack'], axis=1, inplace=True)
3 dfkdd_test.drop(['attack'], axis=1, inplace=True)
```

```
In [8]: 1 # check new train data
2 dfkdd_train.head()
```

```
Out[8]:
```

|   | duration | protocol_type | service  | flag | src_bytes | dst_bytes | land | wrong_fragment | urgent | hot |
|---|----------|---------------|----------|------|-----------|-----------|------|----------------|--------|-----|
| 0 | 0        | tcp           | ftp_data | SF   | 491       | 0         | 0    | 0              | 0      | 0   |
| 1 | 0        | udp           | other    | SF   | 146       | 0         | 0    | 0              | 0      | 0   |
| 2 | 0        | tcp           | private  | S0   | 0         | 0         | 0    | 0              | 0      | 0   |
| 3 | 0        | tcp           | http     | SF   | 232       | 8153      | 0    | 0              | 0      | 0   |
| 4 | 0        | tcp           | http     | SF   | 199       | 420       | 0    | 0              | 0      | 0   |

```
In [9]: 1 # Descriptive statistics
2 dfkdd_train.describe()
```

```
Out[9]:
```

|       | duration     | src_bytes    | dst_bytes    | land          | wrong_fragment | urgent        |
|-------|--------------|--------------|--------------|---------------|----------------|---------------|
| count | 125973.00000 | 1.259730e+05 | 1.259730e+05 | 125973.000000 | 125973.000000  | 125973.000000 |
| mean  | 287.14465    | 4.556674e+04 | 1.977911e+04 | 0.000198      | 0.022687       | 0.000111      |
| std   | 2604.51531   | 5.870331e+06 | 4.021269e+06 | 0.014086      | 0.253530       | 0.014366      |
| min   | 0.00000      | 0.000000e+00 | 0.000000e+00 | 0.000000      | 0.000000       | 0.000000      |
| 25%   | 0.00000      | 0.000000e+00 | 0.000000e+00 | 0.000000      | 0.000000       | 0.000000      |
| 50%   | 0.00000      | 4.400000e+01 | 0.000000e+00 | 0.000000      | 0.000000       | 0.000000      |
| 75%   | 0.00000      | 2.760000e+02 | 5.160000e+02 | 0.000000      | 0.000000       | 0.000000      |
| max   | 42908.00000  | 1.379964e+09 | 1.309937e+09 | 1.000000      | 3.000000       | 3.000000      |

```
In [10]: 1 dfkdd_train['num_outbound_cmds'].value_counts()
        2 dfkdd_test['num_outbound_cmds'].value_counts()
```

```
Out[10]: 0    22481
        Name: num_outbound_cmds, dtype: int64
```

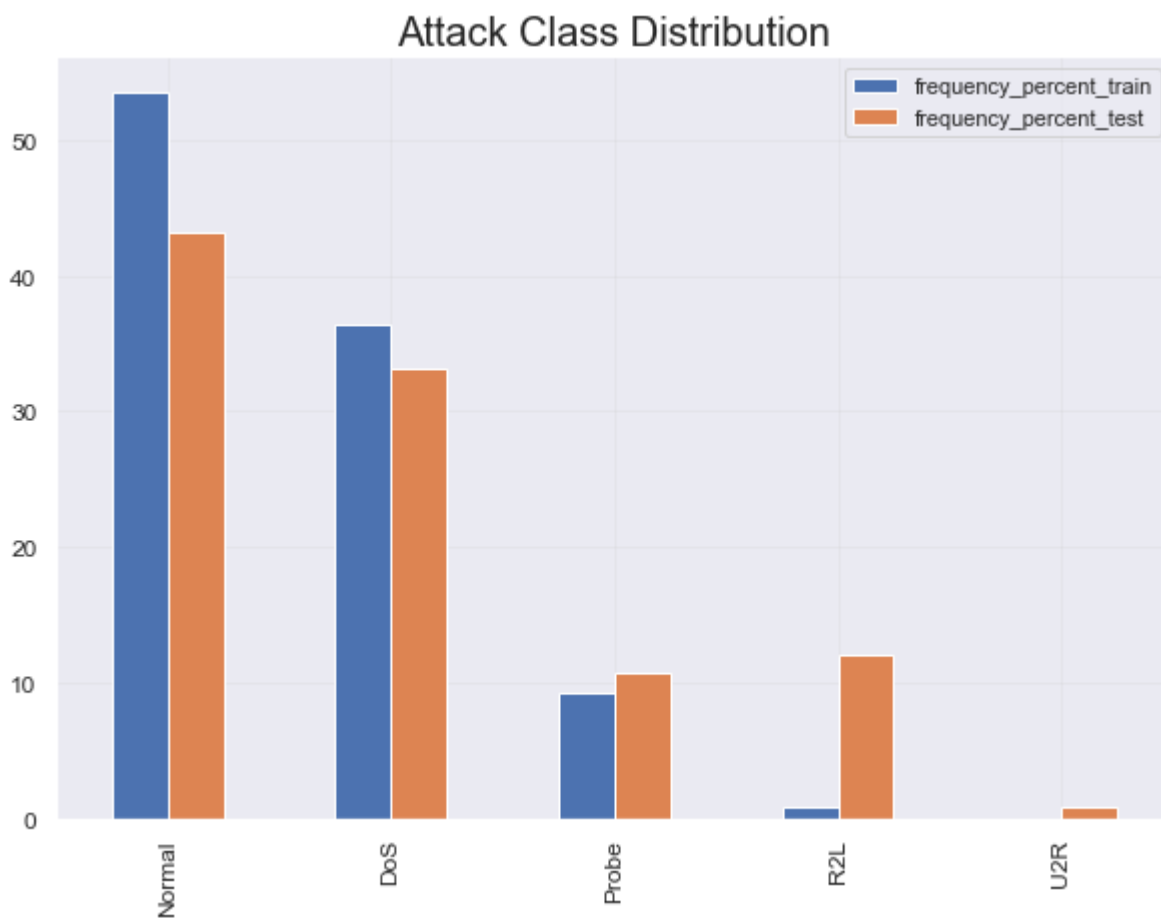
```
In [11]: 1 # 'num_outbound_cmds' field has value 0 for all rows. Hence, it is removed f
        2 dfkdd_train.drop(['num_outbound_cmds'], axis=1, inplace=True)
        3 dfkdd_test.drop(['num_outbound_cmds'], axis=1, inplace=True)
```

```
In [12]: 1 # checking attack group Distribution
        2 attack_class_freq_train = dfkdd_train[['attack_class']].apply(lambda x: x.va
        3 attack_class_freq_test = dfkdd_test[['attack_class']].apply(lambda x: x.valu
        4 attack_class_freq_train['frequency_percent_train'] = round((100 * attack_cla
        5 attack_class_freq_test['frequency_percent_test'] = round((100 * attack_class
        6
        7 attack_class_dist = pd.concat([attack_class_freq_train, attack_class_freq_tes
        8 attack_class_dist
```

```
Out[12]:
```

|               | attack_class | frequency_percent_train | attack_class | frequency_percent_test |
|---------------|--------------|-------------------------|--------------|------------------------|
| <b>Normal</b> | 67343        | 53.46                   | 9711         | 43.20                  |
| <b>DoS</b>    | 45927        | 36.46                   | 7458         | 33.17                  |
| <b>Probe</b>  | 11656        | 9.25                    | 2421         | 10.77                  |
| <b>R2L</b>    | 995          | 0.79                    | 2754         | 12.22                  |
| <b>U2R</b>    | 52           | 0.04                    | 200          | 0.89                   |

```
In [13]: 1 # Attack group bar plot
2 plot = attack_class_dist[['frequency_percent_train', 'frequency_percent_test']
3 plot.set_title("Attack Class Distribution", fontsize=20);
4 plot.grid('on', color='lightgray', alpha=0.25);
```



In [14]:

```
1 dfkdd_train.head()
```

Out[14]:

|   | duration | protocol_type | service  | flag | src_bytes | dst_bytes | land | wrong_fragment | urgent | hot |
|---|----------|---------------|----------|------|-----------|-----------|------|----------------|--------|-----|
| 0 | 0        | tcp           | ftp_data | SF   | 491       | 0         | 0    | 0              | 0      | 0   |
| 1 | 0        | udp           | other    | SF   | 146       | 0         | 0    | 0              | 0      | 0   |
| 2 | 0        | tcp           | private  | S0   | 0         | 0         | 0    | 0              | 0      | 0   |
| 3 | 0        | tcp           | http     | SF   | 232       | 8153      | 0    | 0              | 0      | 0   |
| 4 | 0        | tcp           | http     | SF   | 199       | 420       | 0    | 0              | 0      | 0   |

## Scaling Numerical Attributes

In [15]:

```
1 from sklearn.preprocessing import StandardScaler
2 scaler = StandardScaler()
3
4 # extract numerical features and scale it to have zero mean and unit variance
5 cols = dfkdd_train.select_dtypes(include=['float64', 'int64']).columns
6 sc_train = scaler.fit_transform(dfkdd_train.select_dtypes(include=['float64',
7 sc_test = scaler.fit_transform(dfkdd_test.select_dtypes(include=['float64', '
8
9 # turn the result back to a dataframe
10 sc_traindf = pd.DataFrame(sc_train, columns = cols)
11 sc_testdf = pd.DataFrame(sc_test, columns = cols)
```

## Encoding of Categorical Attributes

In [16]:

```
1 from sklearn.preprocessing import LabelEncoder
2 encoder = LabelEncoder()
3
4 # extract categorical features from both training and test sets
5 cattrain = dfkdd_train.select_dtypes(include=['object']).copy()
6 cattest = dfkdd_test.select_dtypes(include=['object']).copy()
7
8 # encoding
9 traincat = cattrain.apply(encoder.fit_transform)
10 testcat = cattest.apply(encoder.fit_transform)
11
12 # separate target column from encoded data
13 enctrain = traincat.drop(['attack_class'], axis=1)
14 encctest = testcat.drop(['attack_class'], axis=1)
15
16 cat_Ytrain = traincat[['attack_class']].copy()
17 cat_Ytest = testcat[['attack_class']].copy()
```

## Data Sampling

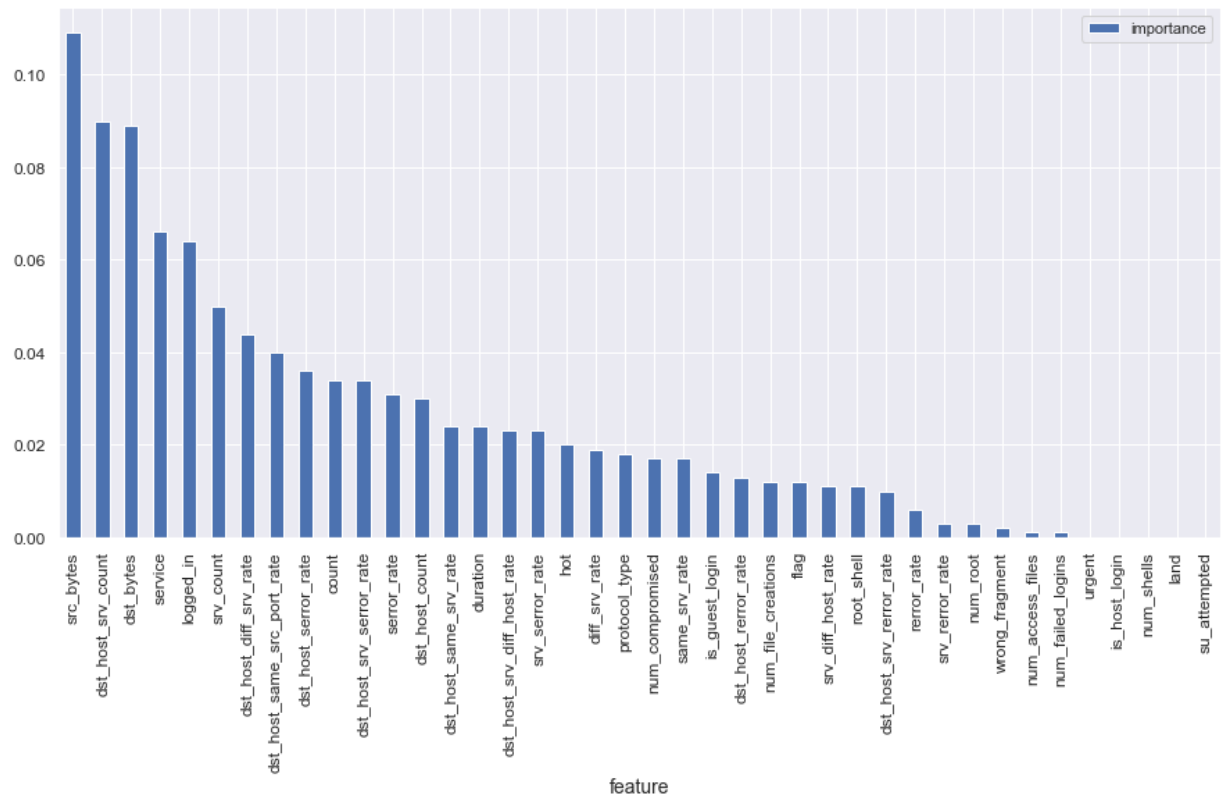
```
In [17]: 1 from imblearn.over_sampling import RandomOverSampler
2 from collections import Counter
3
4 # define columns and extract encoded train set for sampling
5 sc_traindf = dfkdd_train.select_dtypes(include=['float64','int64'])
6 refclasscol = pd.concat([sc_traindf, enctrain], axis=1).columns
7 refclass = np.concatenate((sc_train, enctrain.values), axis=1)
8 X = refclass
9
10 # reshape target column to 1D array shape
11 c, r = cat_Ytest.values.shape
12 y_test = cat_Ytest.values.reshape(c,)
13
14 c, r = cat_Ytrain.values.shape
15 y = cat_Ytrain.values.reshape(c,)
16
17 # apply the random over sampling
18 ros = RandomOverSampler(random_state=42)
19 X_res, y_res = ros.fit_resample(X, y)
20 print('Original dataset shape {}'.format(Counter(y)))
21 print('Resampled dataset shape {}'.format(Counter(y_res)))
```

Original dataset shape Counter({1: 67343, 0: 45927, 2: 11656, 3: 995, 4: 52})  
Resampled dataset shape Counter({1: 67343, 0: 67343, 3: 67343, 2: 67343, 4: 67343})

## Feature Selection

```
In [18]: 1 from sklearn.ensemble import RandomForestClassifier
2 rfc = RandomForestClassifier();
3
4 # fit random forest classifier on the training set
5 rfc.fit(X_res, y_res);
6 # extract important features
7 score = np.round(rfc.feature_importances_,3)
8 importances = pd.DataFrame({'feature':refclasscol,'importance':score})
9 importances = importances.sort_values('importance',ascending=False).set_inde
```

```
In [19]: 1 # plot importances
2 plt.rcParams['figure.figsize'] = (15, 7)
3 importances.plot(kind='bar');
```



```
In [20]: 1 from sklearn.feature_selection import RFE
2 import itertools
3 rfc = RandomForestClassifier()
4
5 # create the RFE model and select 10 attributes
6 rfe = RFE(rfc, n_features_to_select=10)
7 rfe = rfe.fit(X_res, y_res)
8
9 # summarize the selection of the attributes
10 feature_map = [(i, v) for i, v in itertools.zip_longest(rfe.get_support(), range(35))]
11 selected_features = [v for i, v in feature_map if i==True]
```



```
In [21]: 1 selected_features = ['src_bytes', 'dst_bytes', 'logged_in', 'count', 'srv_count',
```

## Dataset Partition

```
In [22]: 1 # define columns to new dataframe
2 newcol = list(refclasscol)
3 newcol.append('attack_class')
4
5 # add a dimension to target
6 new_y_res = y_res[:, np.newaxis]
7
8 # create a dataframe from sampled data
9 res_arr = np.concatenate((X_res, new_y_res), axis=1)
10 res_df = pd.DataFrame(res_arr, columns = newcol)
11
12 # create test dataframe
13 reftest = pd.concat([sc_testdf, testcat], axis=1)
14 reftest['attack_class'] = reftest['attack_class'].astype(np.float64)
15 reftest['protocol_type'] = reftest['protocol_type'].astype(np.float64)
16 reftest['flag'] = reftest['flag'].astype(np.float64)
17 reftest['service'] = reftest['service'].astype(np.float64)
18
19 res_df.shape
20 reftest.shape
```

Out[22]: (22481, 41)

```
In [23]: 1 from collections import defaultdict
2 classdict = defaultdict(list)
3
4 # create two-target classes (normal class and an attack class)
5 attacklist = [('DoS', 0.0), ('Probe', 2.0), ('R2L', 3.0), ('U2R', 4.0)]
6 normalclass = [('Normal', 1.0)]
7
8 def create_classdict():
9     '''This function subdivides train and test dataset into two-class attack
10     for j, k in normalclass:
11         for i, v in attacklist:
12             restrain_set = res_df.loc[(res_df['attack_class'] == k) | (res_d
13             classdict[j + '_' + i].append(restrain_set)
14             # test labels
15             reftest_set = reftest.loc[(reftest['attack_class'] == k) | (reft
16             classdict[j + '_' + i].append(reftest_set)
17
18 create_classdict()
```

```
In [24]: 1 for k, v in classdict.items():
2         print(k)
```

Normal\_DoS  
Normal\_Probe  
Normal\_R2L  
Normal\_U2R

```
In [25]: 1 pretrain = classdict['Normal_Probe'][0]
2 pretest = classdict['Normal_DoS'][1]
3 grpclass = 'Normal_DoS'
```

## Finalize data preprocessing for training

```
In [26]: 1 from sklearn.preprocessing import OneHotEncoder
2 enc = OneHotEncoder(handle_unknown='ignore')
3
4 #Xresdf = pretrain
5 #newtest = pretest
6
7 Xresdf = res_df
8 newtest = reftest
9
10 Xresdfnew = Xresdf[selected_features]
11 Xresdfnum = Xresdfnew.drop(['service'], axis=1)
12 Xresdfcat = Xresdfnew[['service']].copy()
13
14 Xtest_features = newtest[selected_features]
15 Xtestdfnum = Xtest_features.drop(['service'], axis=1)
16 Xtestcat = Xtest_features[['service']].copy()
17
18
19 # Fit train data
20 enc.fit(Xresdfcat)
21
22 # Transform train data
23 X_train_1hotenc = enc.transform(Xresdfcat).toarray()
24
25 # Transform test data
26 X_test_1hotenc = enc.transform(Xtestcat).toarray()
27
28 X_train = np.concatenate((Xresdfnum.values, X_train_1hotenc), axis=1)
29 X_test = np.concatenate((Xtestdfnum.values, X_test_1hotenc), axis=1)
30
31 y_train = Xresdf[['attack_class']].copy()
32 c, r = y_train.values.shape
33 Y_train = y_train.values.reshape(c,)
34
35 y_test = newtest[['attack_class']].copy()
36 c, r = y_test.values.shape
37 Y_test = y_test.values.reshape(c,)
```

## Train Models

In [27]:

```
1 from sklearn.svm import SVC
2 from sklearn.naive_bayes import BernoulliNB
3 from sklearn import tree
4 from sklearn.model_selection import cross_val_score
5 from sklearn.neighbors import KNeighborsClassifier
6 from sklearn.linear_model import LogisticRegression
7 from sklearn.ensemble import VotingClassifier
8
9 # Train KNeighborsClassifier Model
10 KNN_Classifier = KNeighborsClassifier(n_jobs=-1)
11 KNN_Classifier.fit(X_train, Y_train);
12
13 # Train Gaussian Naive Baye Model
14 BNB_Classifier = BernoulliNB()
15 BNB_Classifier.fit(X_train, Y_train)
16
17 # Train Decision Tree Model
18 DTC_Classifier = tree.DecisionTreeClassifier(criterion='entropy', max_depth
19 DTC_Classifier.fit(X_train, Y_train);
20
21 # Train RandomForestClassifier Model
22 RF_Classifier = RandomForestClassifier(criterion='entropy', n_jobs=-1, n_est
23 RF_Classifier.fit(X_train, Y_train);
24
25 # Train Ensemble Model (This method combines all the individual models above
26 combined_model = [('Naive Baye Classifier', BNB_Classifier),
27                   ('RF_Classifier', RF_Classifier),
28                   ('Decision Tree Classifier', DTC_Classifier),
29                   ('KNeighborsClassifier', KNN_Classifier),
30                   ]
31 EnsembleClassifier = VotingClassifier(estimators = combined_model, voting =
32 EnsembleClassifier.fit(X_train, Y_train);
```

## Evaluate Models

```

In [28]: 1 from sklearn import metrics
2 from sklearn.metrics import plot_confusion_matrix
3
4 def plotHeatmap(cm):
5     fig, ax = plt.subplots(figsize=(10,10))
6     heatmap = sns.heatmap(cm, annot=True, linewidths=0.5, linecolor="lightgr
7     heatmap.set_xticklabels(['DoS', 'Normal', 'Probe', 'R2L', 'U2R'])
8     heatmap.set_yticklabels(['DoS', 'Normal', 'Probe', 'R2L', 'U2R'])
9     plt.xlabel("Predicted Group", labelpad=20)
10    plt.ylabel("True Group", labelpad=20)
11    plt.show()
12
13
14 models = []
15 models.append(('Naive Baye Classifier', BNB_Classifier))
16 models.append(('Decision Tree Classifier', DTC_Classifier))
17 models.append(('RandomForest Classifier', RF_Classifier))
18 #models.append(('KNeighborsClassifier', KNN_Classifier))
19 models.append(('EnsembleClassifier', EnsembleClassifier))
20
21 for i, v in models:
22     Y_pred = v.predict(X_train)
23     accuracy = metrics.accuracy_score(Y_train, Y_pred)
24     confusion_matrix = metrics.confusion_matrix(Y_train, Y_pred, labels = v.
25     classification = metrics.classification_report(Y_train, Y_pred)
26     print()
27     print('===== {} Model Evaluation =====')
28     print()
29     print("Model Accuracy:" "\n", accuracy)
30     print()
31     print("Confusion matrix:" "\n", confusion_matrix)
32     plotHeatmap(confusion_matrix)
33     print()
34     print("Classification report:" "\n", classification)
35     print()

```

```

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

```

```

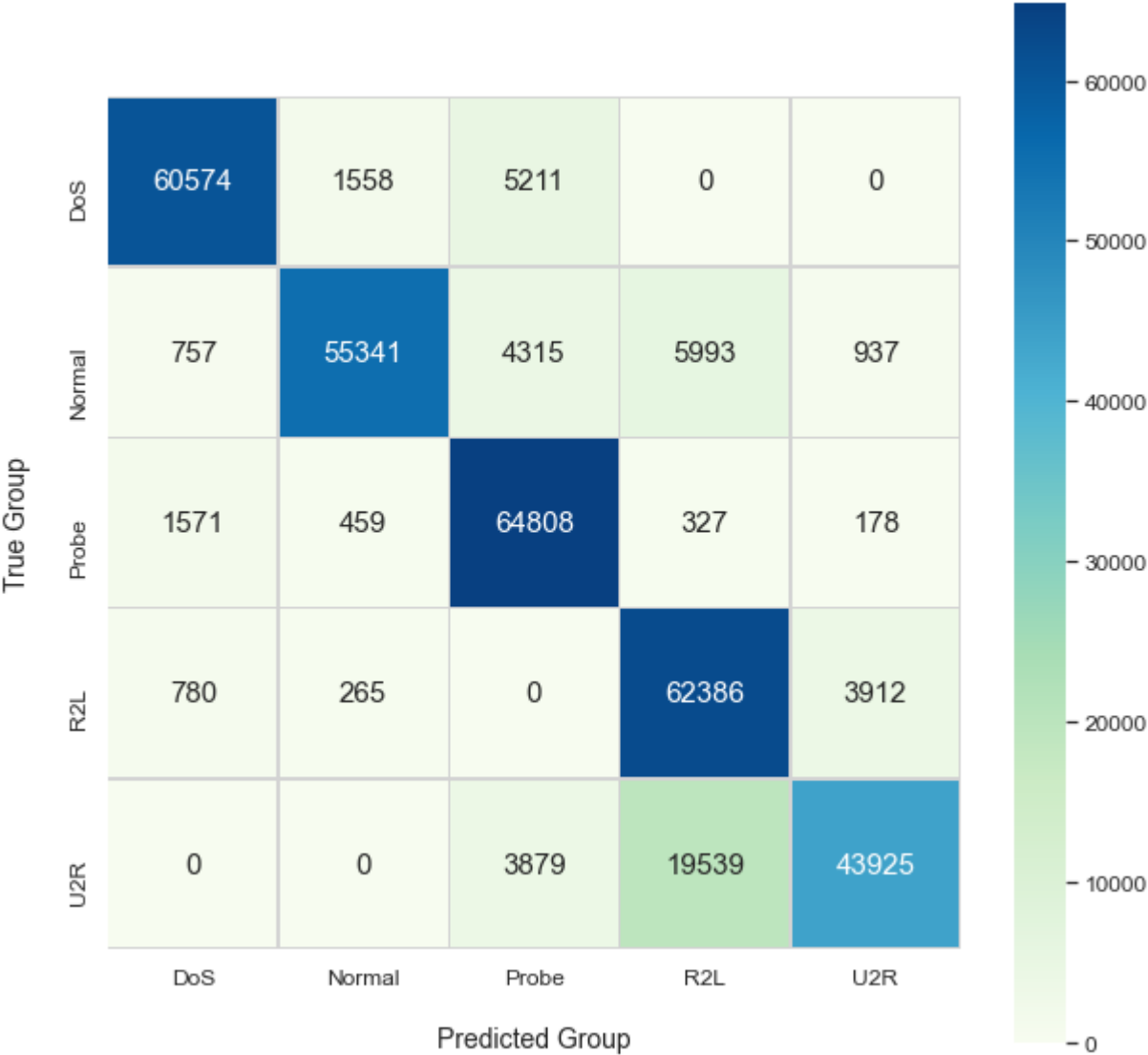
Model Accuracy:
0.8524538556345871

```

```

Confusion matrix:
[[60574 1558 5211      0      0]
 [ 757 55341 4315 5993  937]
 [ 1571  459 64808  327  178]
 [  780  265      0 62386 3912]
 [      0      0 3879 19539 43925]]

```



Classification report:

|          | precision | recall | f1-score | support |
|----------|-----------|--------|----------|---------|
| 0.0      | 0.95      | 0.90   | 0.92     | 67343   |
| 1.0      | 0.96      | 0.82   | 0.89     | 67343   |
| 2.0      | 0.83      | 0.96   | 0.89     | 67343   |
| 3.0      | 0.71      | 0.93   | 0.80     | 67343   |
| 4.0      | 0.90      | 0.65   | 0.76     | 67343   |
| accuracy |           |        | 0.85     | 336715  |

```

macro avg      0.87      0.85      0.85      336715
weighted avg    0.87      0.85      0.85      336715

```

```

===== Decision Tree Classifier Model Evaluation =====
=====

```

```

Model Accuracy:
0.8474615030515421

```

```

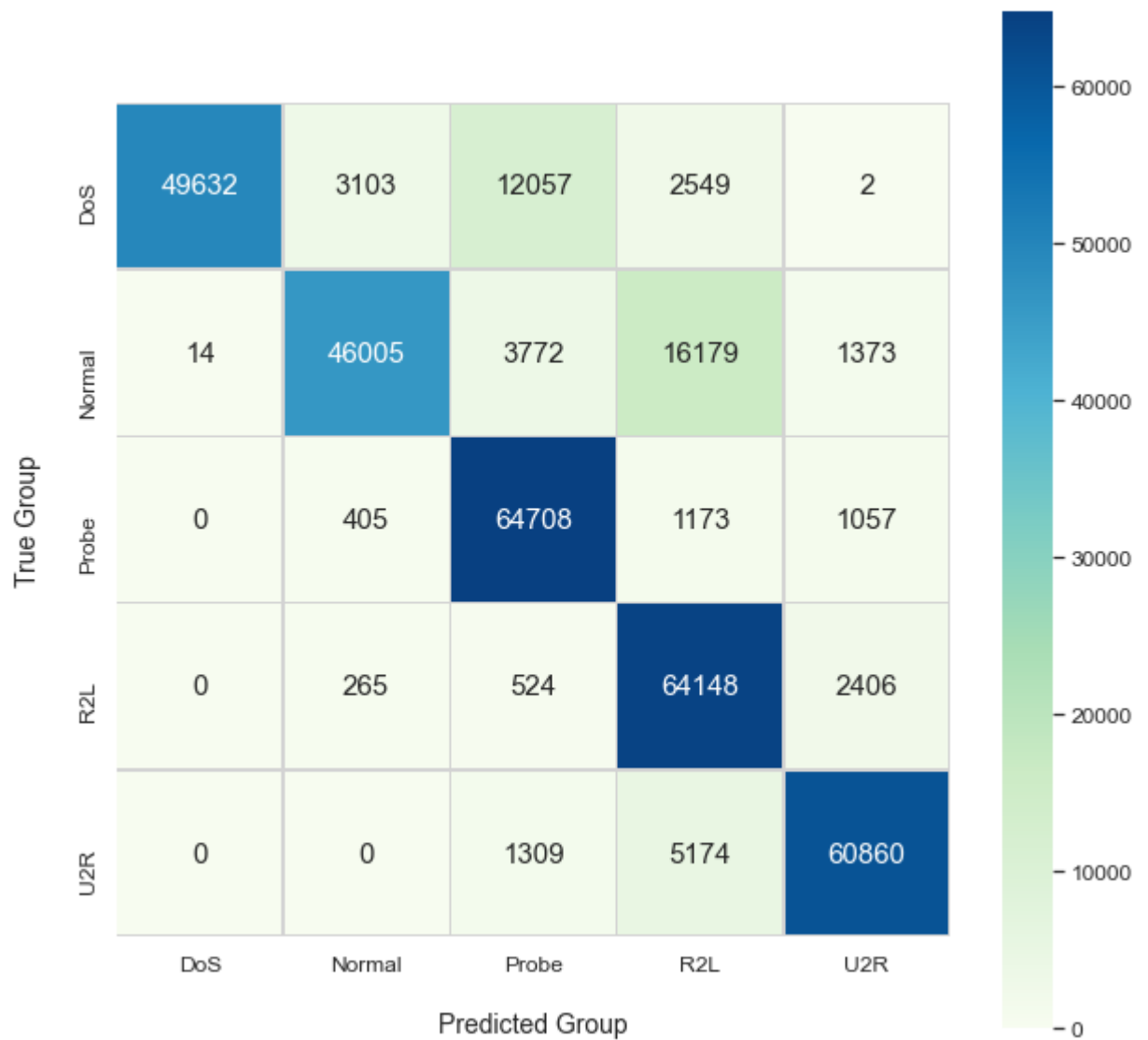
Confusion matrix:

```

```

[[49632  3103 12057  2549    2]
 [   14 46005  3772 16179 1373]
 [    0   405 64708  1173 1057]
 [    0   265   524 64148 2406]
 [    0    0  1309  5174 60860]]

```



## Classification report:

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0.0          | 1.00      | 0.74   | 0.85     | 67343   |
| 1.0          | 0.92      | 0.68   | 0.79     | 67343   |
| 2.0          | 0.79      | 0.96   | 0.86     | 67343   |
| 3.0          | 0.72      | 0.95   | 0.82     | 67343   |
| 4.0          | 0.93      | 0.90   | 0.91     | 67343   |
| accuracy     |           |        | 0.85     | 336715  |
| macro avg    | 0.87      | 0.85   | 0.85     | 336715  |
| weighted avg | 0.87      | 0.85   | 0.85     | 336715  |

===== RandomForest Classifier Model Evaluation =====  
=====

## Model Accuracy:

0.8708967524464309

## Confusion matrix:

```
[[65136  858   761    19   569]
 [ 1726 48498  5165  3461  8493]
 [ 8805    97 58057   101   283]
 [    0   265   884 60693  5501]
 [    0    0  2645  3838 60860]]
```



Classification report:

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0.0          | 0.86      | 0.97   | 0.91     | 67343   |
| 1.0          | 0.98      | 0.72   | 0.83     | 67343   |
| 2.0          | 0.86      | 0.86   | 0.86     | 67343   |
| 3.0          | 0.89      | 0.90   | 0.90     | 67343   |
| 4.0          | 0.80      | 0.90   | 0.85     | 67343   |
| accuracy     |           |        | 0.87     | 336715  |
| macro avg    | 0.88      | 0.87   | 0.87     | 336715  |
| weighted avg | 0.88      | 0.87   | 0.87     | 336715  |

===== KNeighborsClassifier Model Evaluation =====  
 =====

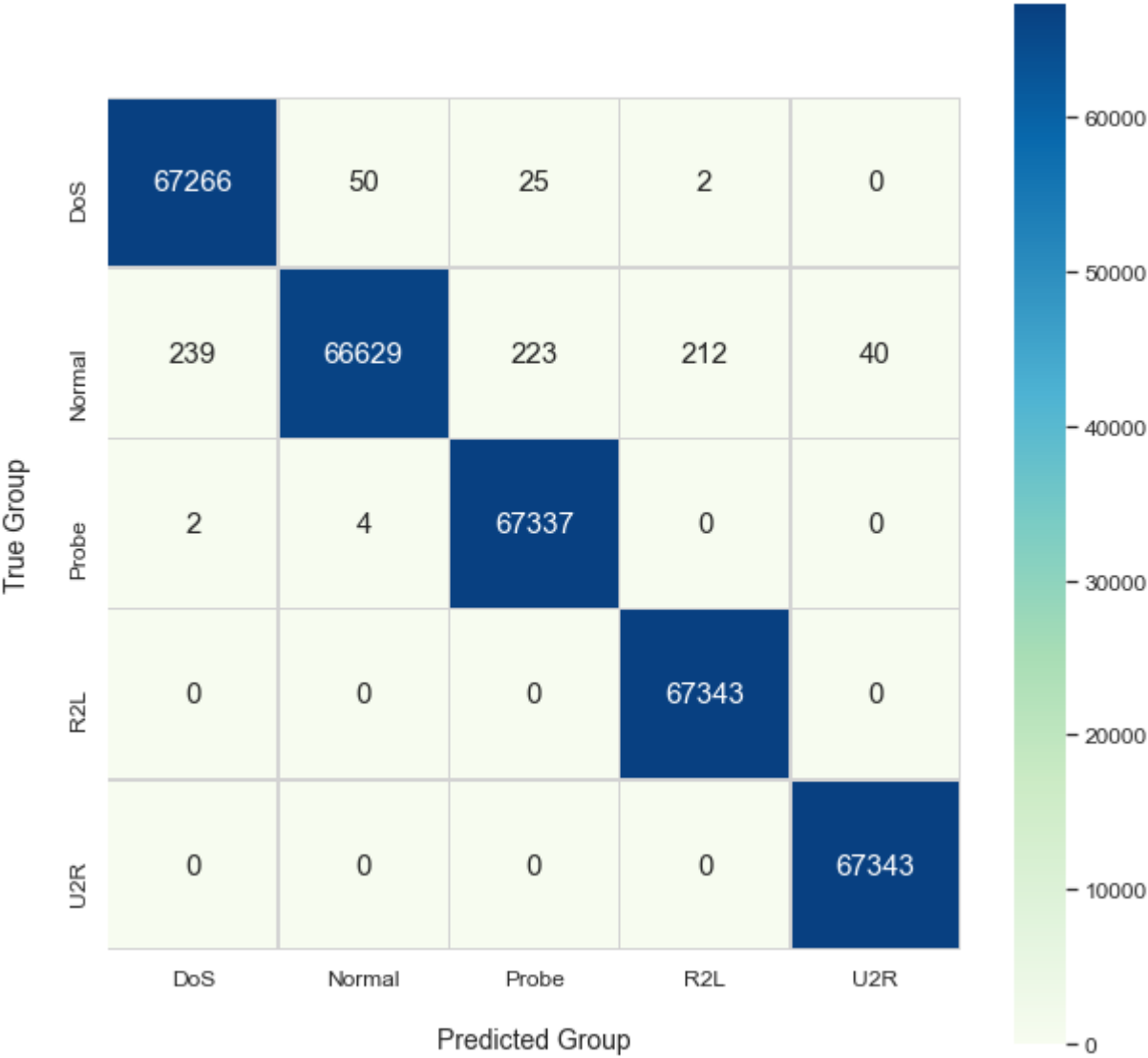
Model Accuracy:

0.9576330130822803



Confusion matrix:

```
[[67266  50  25  2  0]
 [ 239 66629 223 212 40]
 [  2  4 67337  0  0]
 [  0  0  0 67343  0]
 [  0  0  0  0 67343]]
```



## Classification report:

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0.0          | 1.00      | 1.00   | 1.00     | 67343   |
| 1.0          | 1.00      | 0.99   | 0.99     | 67343   |
| 2.0          | 1.00      | 1.00   | 1.00     | 67343   |
| 3.0          | 1.00      | 1.00   | 1.00     | 67343   |
| 4.0          | 1.00      | 1.00   | 1.00     | 67343   |
| accuracy     |           |        | 1.00     | 336715  |
| macro avg    | 1.00      | 1.00   | 1.00     | 336715  |
| weighted avg | 1.00      | 1.00   | 1.00     | 336715  |

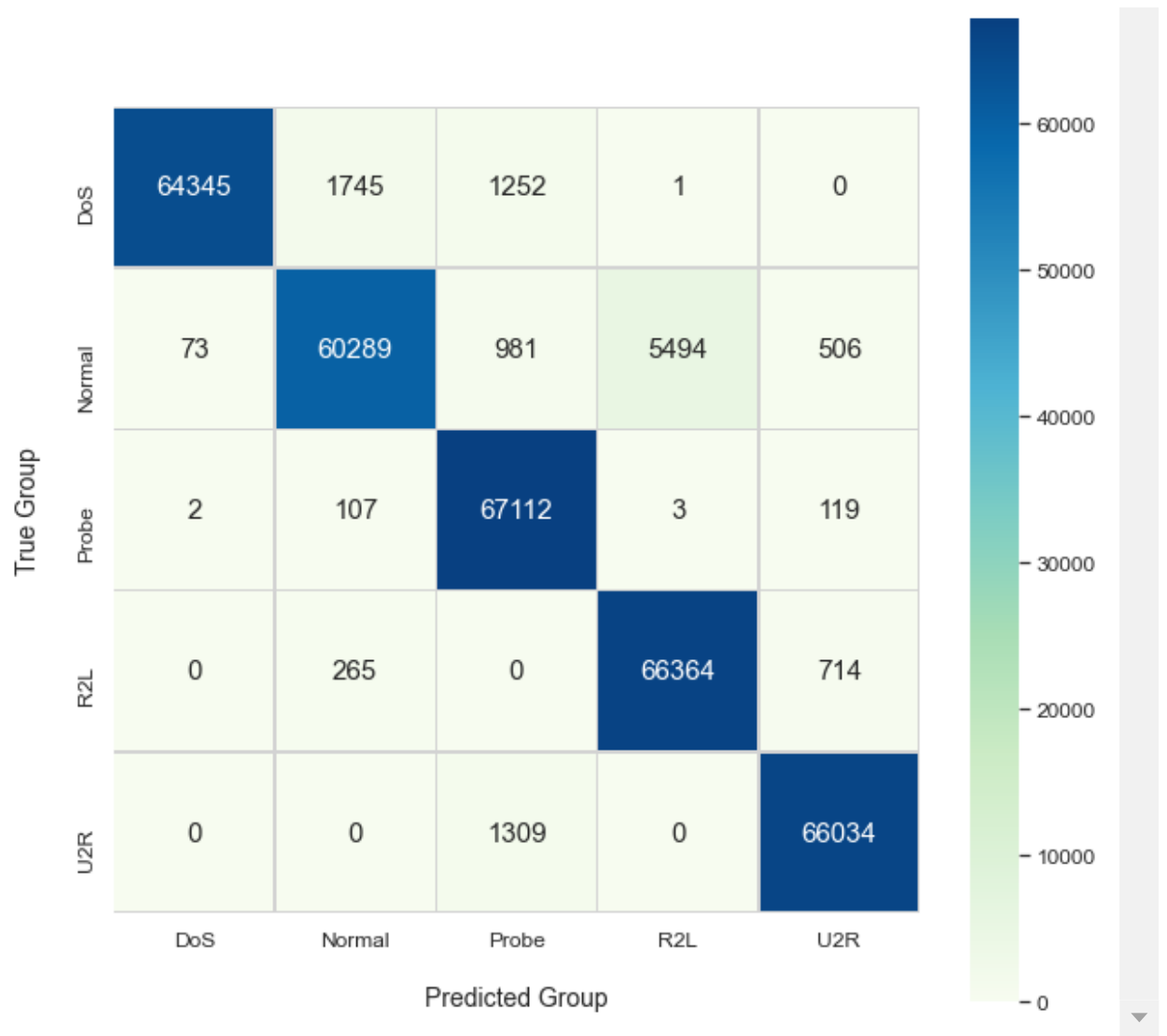
=====  
 ===== EnsembleClassifier Model Evaluation =====  
 =====

Model Accuracy:

0.9626657559063303

Confusion matrix:

```
[[64345 1745 1252    1    0]
 [  73 60289  981 5494  506]
 [   2  107 67112    3  119]
 [   0  265    0 66364  714]
 [   0    0 1309    0 66034]]
```



Classification report:

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0.0          | 1.00      | 0.96   | 0.98     | 67343   |
| 1.0          | 0.97      | 0.90   | 0.93     | 67343   |
| 2.0          | 0.95      | 1.00   | 0.97     | 67343   |
| 3.0          | 0.92      | 0.99   | 0.95     | 67343   |
| 4.0          | 0.98      | 0.98   | 0.98     | 67343   |
| accuracy     |           |        | 0.96     | 336715  |
| macro avg    | 0.96      | 0.96   | 0.96     | 336715  |
| weighted avg | 0.96      | 0.96   | 0.96     | 336715  |

```
In [25]: 1 import pickle
2 for i, v in models:
3     pickle.dump(v, open(i, 'wb'))
4
5 KNN_Classifier = pickle.load(open("./models/multiclass/KNeighborsClassifier"
6 BNB_Classifier = pickle.load(open("./models/multiclass/Naive Baye Classifier
7 DTC_Classifier = pickle.load(open("./models/multiclass/Decision Tree Classif
8 RF_Classifier = pickle.load(open("./models/multiclass/RandomForest Classifie
9 EnsembleClassifier = pickle.load(open("./models/multiclass/EnsembleClassifie
```

```
In [28]: 1 from sklearn import metrics
2 from sklearn.metrics import plot_confusion_matrix
3 models = []
4 models.append(('Naive Baye Classifier', BNB_Classifier))
5 models.append(('Decision Tree Classifier', DTC_Classifier))
6 models.append(('RandomForest Classifier', RF_Classifier))
7 models.append(('KNeighborsClassifier', KNN_Classifier))
8 models.append(('EnsembleClassifier', EnsembleClassifier))
```

## Test Models

```

In [29]: 1 def plotHeatmap(cm):
2         fig, ax = plt.subplots(figsize=(10,10))
3         heatmap = sns.heatmap(cm, annot=True, linewidths=0.5, linecolor="lightgr
4         heatmap.set_xticklabels(['DoS', 'Normal', 'Probe', 'R2L', 'U2R'])
5         heatmap.set_yticklabels(['DoS', 'Normal', 'Probe', 'R2L', 'U2R'])
6         plt.xlabel("Predicted Group", labelpad=20)
7         plt.ylabel("True Group", labelpad=20)
8         plt.show()
9
10        for i, v in models:
11            accuracy = metrics.accuracy_score(Y_test, v.predict(X_test))
12            confusion_matrix = metrics.confusion_matrix(Y_test, v.predict(X_test))
13            classification = metrics.classification_report(Y_test, v.predict(X_test)
14            print()
15            print('===== {} Model Test Results =====')
16            print()
17            print("Model Accuracy:" "\n", accuracy)
18            print()
19            print("Confusion matrix:" "\n", confusion_matrix)
20            plotHeatmap(confusion_matrix)
21            print()
22            print("Classification report:" "\n", classification)
23            print()
24

```

```

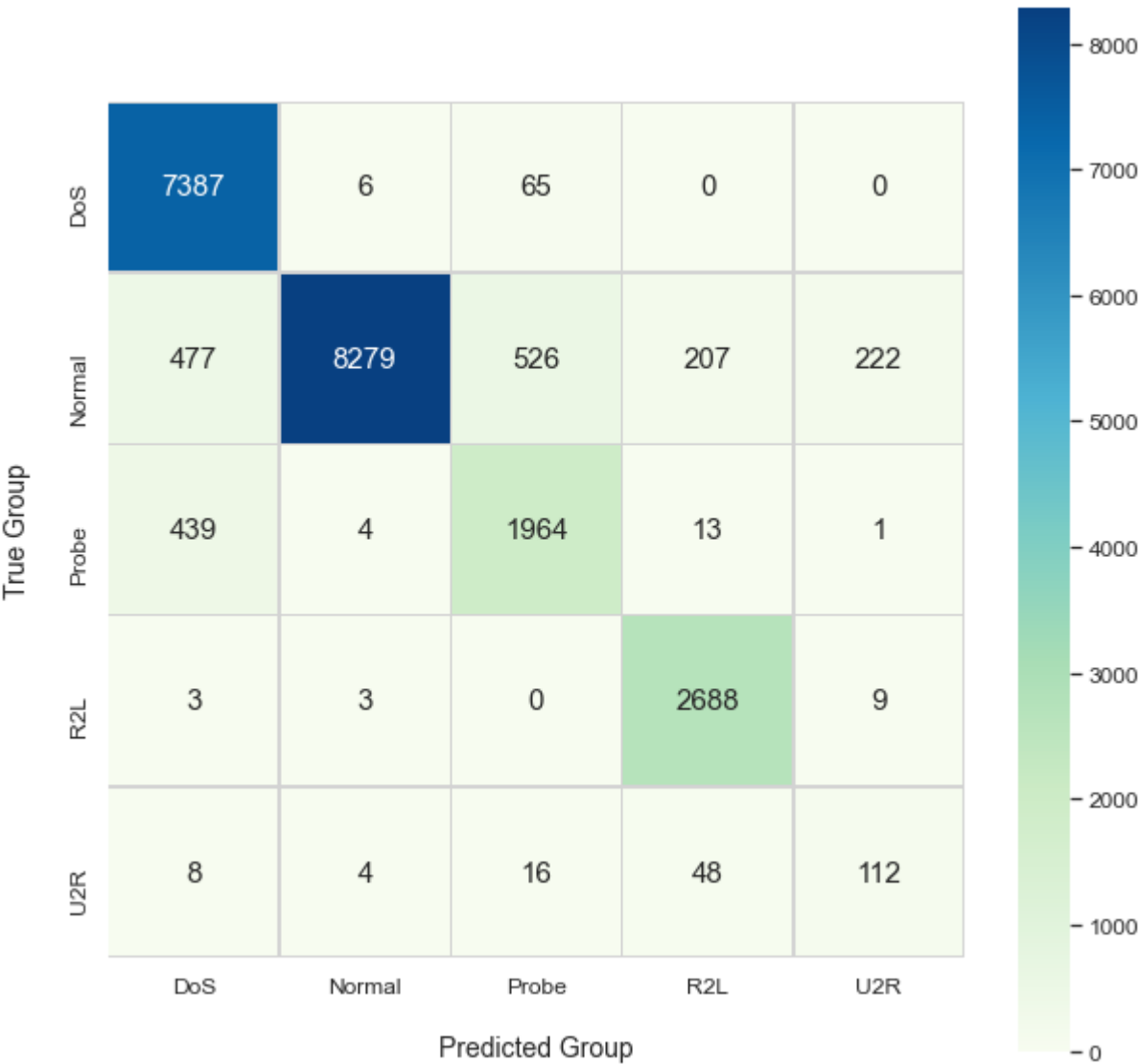
===== Naive Baye Classifier Model Test Results =====
=====

```

Model Accuracy:  
0.9087674035852498

Confusion matrix:

|       |      |      |      |       |
|-------|------|------|------|-------|
| [7387 | 6    | 65   | 0    | 0]    |
| [ 477 | 8279 | 526  | 207  | 222]  |
| [ 439 | 4    | 1964 | 13   | 1]    |
| [ 3   | 3    | 0    | 2688 | 9]    |
| [ 8   | 4    | 16   | 48   | 112]] |



Classification report:

|     | precision | recall | f1-score | support |
|-----|-----------|--------|----------|---------|
| 0.0 | 0.89      | 0.99   | 0.94     | 7458    |
| 1.0 | 1.00      | 0.85   | 0.92     | 9711    |
| 2.0 | 0.76      | 0.81   | 0.79     | 2421    |
| 3.0 | 0.91      | 0.99   | 0.95     | 2703    |

6/14/2021

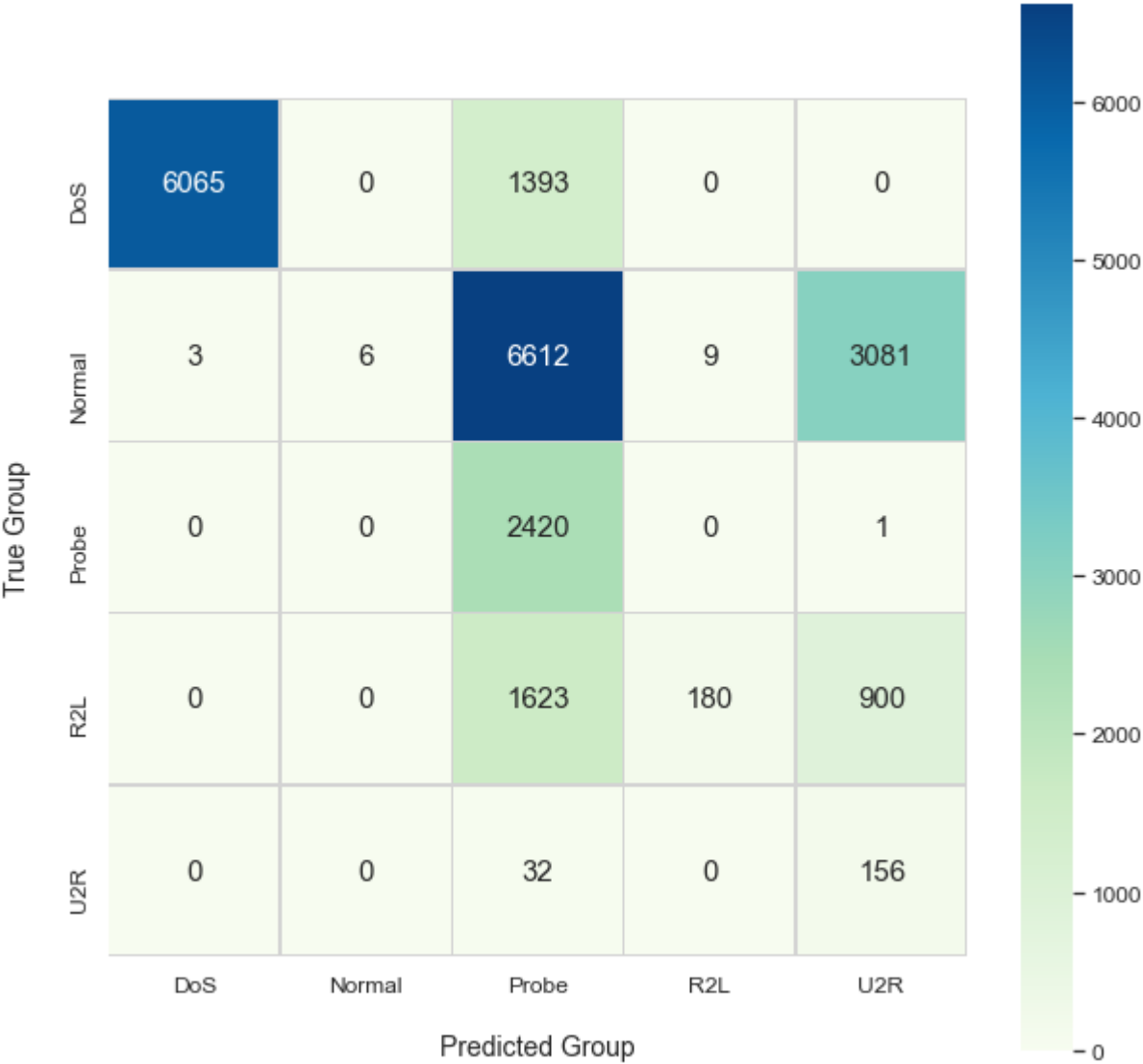
Intrusion Detection Multiclass - Jupyter Notebook

|              |      |      |      |      |       |
|--------------|------|------|------|------|-------|
|              | 4.0  | 0.33 | 0.60 | 0.42 | 188   |
| accuracy     |      |      |      | 0.91 | 22481 |
| macro avg    | 0.78 | 0.85 | 0.80 |      | 22481 |
| weighted avg | 0.92 | 0.91 | 0.91 |      | 22481 |

===== Decision Tree Classifier Model Test Results ==  
=====

Model Accuracy:  
0.3926426760375428

Confusion matrix:  
[[6065 0 1393 0 0]  
[ 3 6 6612 9 3081]  
[ 0 0 2420 0 1]  
[ 0 0 1623 180 900]  
[ 0 0 32 0 156]]



Classification report:

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0.0          | 1.00      | 0.81   | 0.90     | 7458    |
| 1.0          | 1.00      | 0.00   | 0.00     | 9711    |
| 2.0          | 0.20      | 1.00   | 0.33     | 2421    |
| 3.0          | 0.95      | 0.07   | 0.12     | 2703    |
| 4.0          | 0.04      | 0.83   | 0.07     | 188     |
| accuracy     |           |        | 0.39     | 22481   |
| macro avg    | 0.64      | 0.54   | 0.29     | 22481   |
| weighted avg | 0.90      | 0.39   | 0.35     | 22481   |

==== RandomForest Classifier Model Test Results =====

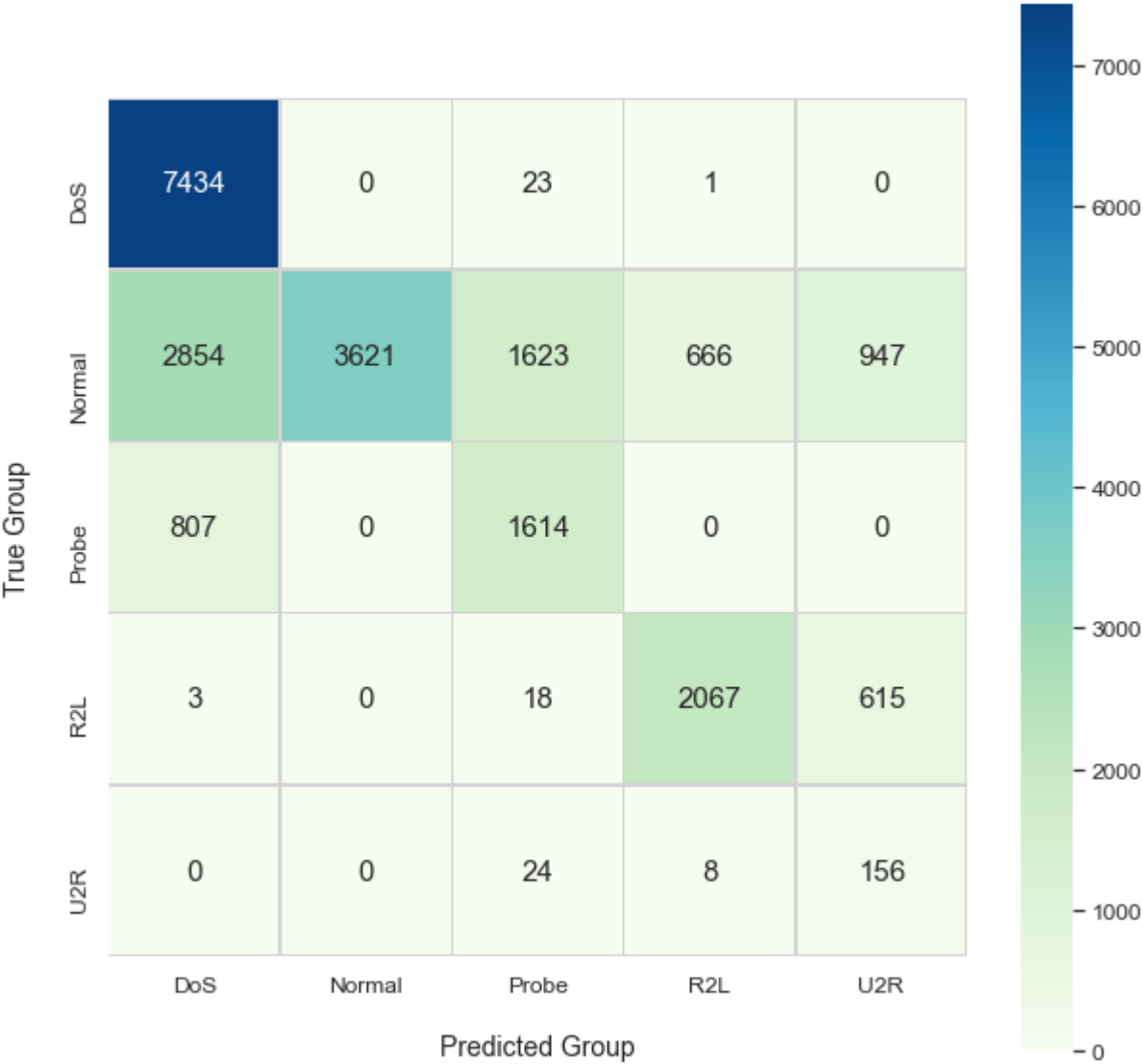
Model Accuracy:

0.6624260486633157

Confusion matrix:

```
[[7434   0   23   1   0]
 [2854 3621 1623  666  947]
 [ 807   0 1614   0   0]
 [   3   0   18 2067  615]
 [   0   0   24   8  156]]
```





## Classification report:

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0.0          | 0.67      | 1.00   | 0.80     | 7458    |
| 1.0          | 1.00      | 0.37   | 0.54     | 9711    |
| 2.0          | 0.49      | 0.67   | 0.56     | 2421    |
| 3.0          | 0.75      | 0.76   | 0.76     | 2703    |
| 4.0          | 0.09      | 0.83   | 0.16     | 188     |
| accuracy     |           |        | 0.66     | 22481   |
| macro avg    | 0.60      | 0.73   | 0.57     | 22481   |
| weighted avg | 0.80      | 0.66   | 0.65     | 22481   |

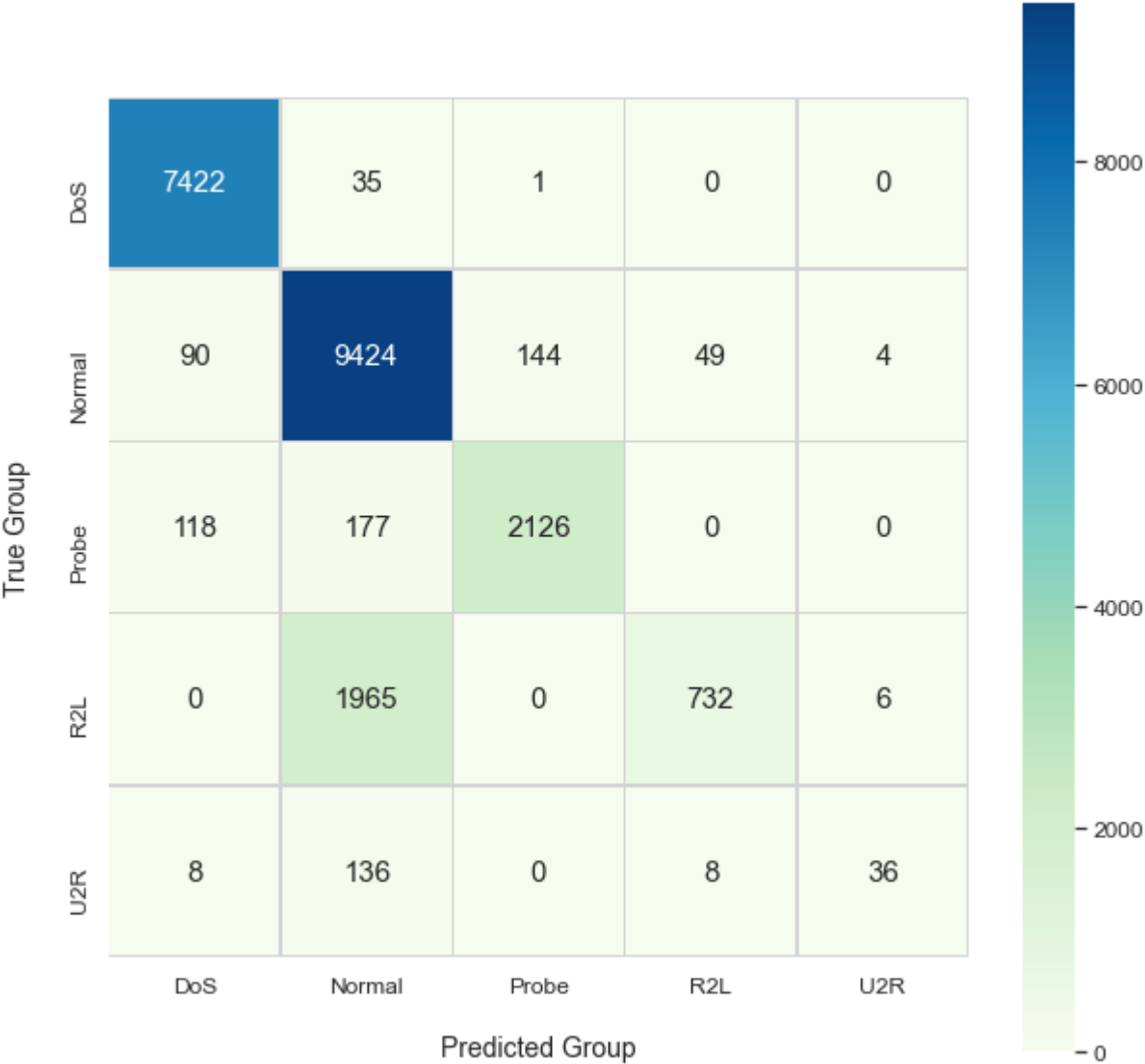
===== KNeighborsClassifier Model Test Results =====  
 =====

## Model Accuracy:

0.8780748187358214

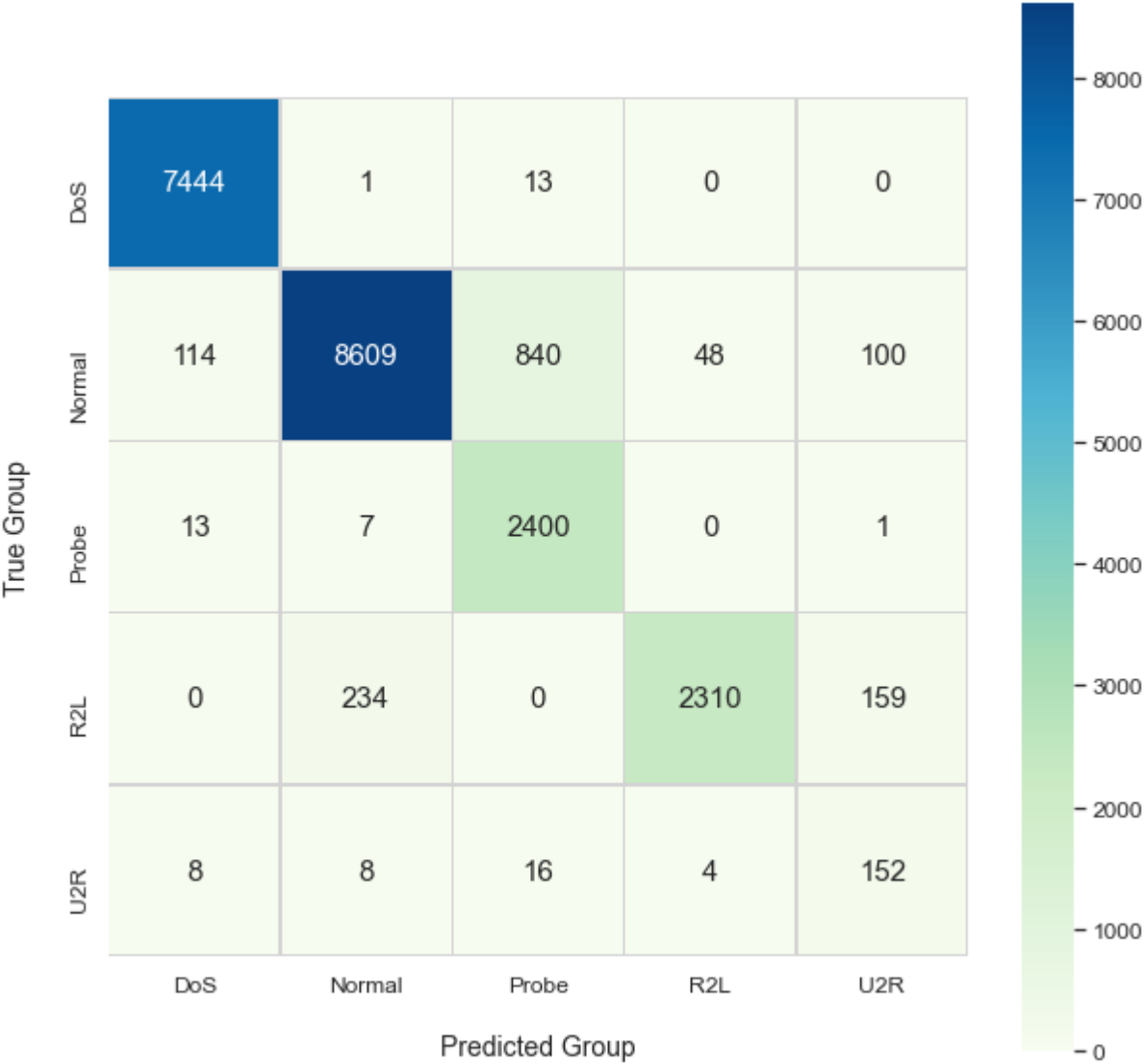
## Confusion matrix:

```
[[7422  35   1   0   0]
 [  90 9424  144  49   4]
 [ 118  177 2126   0   0]
 [   0 1965   0  732   6]
 [   8  136   0   8  36]]
```



Classification report:

|     | precision | recall | f1-score | support |
|-----|-----------|--------|----------|---------|
| 0.0 | 0.97      | 1.00   | 0.98     | 7458    |
| 1.0 | 0.80      | 0.97   | 0.88     | 9711    |
| 2.0 | 0.94      | 0.88   | 0.91     | 2421    |
| 3.0 | 0.93      | 0.27   | 0.42     | 2703    |
| 4.0 | 0.78      | 0.19   | 0.31     | 188     |



Classification report:

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0.0          | 0.98      | 1.00   | 0.99     | 7458    |
| 1.0          | 0.97      | 0.89   | 0.93     | 9711    |
| 2.0          | 0.73      | 0.99   | 0.84     | 2421    |
| 3.0          | 0.98      | 0.85   | 0.91     | 2703    |
| 4.0          | 0.37      | 0.81   | 0.51     | 188     |
| accuracy     |           |        | 0.93     | 22481   |
| macro avg    | 0.81      | 0.91   | 0.84     | 22481   |
| weighted avg | 0.95      | 0.93   | 0.93     | 22481   |