

Software Network Intrusion Detection

A close-up, low-angle shot of a person's hand typing on a laptop keyboard. The keyboard is backlit with a warm, yellowish light. The background is dark and out of focus. The title text is overlaid on the left side of the image.

Rutuja Kute, Ruthwik

Overview

Intrusion Detection Systems and Prevention Systems are crucial tools for network users to defend against various online threats. Our project aims to develop a system that can protect SDNs against various types of online threats.

Goal is to develop an accurate and efficient network intrusion detection system for Software Defined Networks (SDNs) using machine learning algorithms. We are using machine learning algorithms to classify four types of online threats in our project: DDoS attacks, Web Attack Brute Force, Web Attack XSS, and Web Attack SQL Injection.

Flow of the Code

- 01 Importing essential libraries
- 02 Exploring the SDN dataset
- 03 Preprocessing dataset
- 04 Feature Selection
- 05 Visualization
- 06 Algorithms Used
- 07 Learning Curve
- 08 Optimising Hyperparameters
- 09 Confusion Matrix

Importing Libraries

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import MinMaxScaler

from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score, confusion_matrix
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import learning_curve
from joblib import parallel_backend
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression
from xgboost import XGBClassifier
```

Exploring Dataset

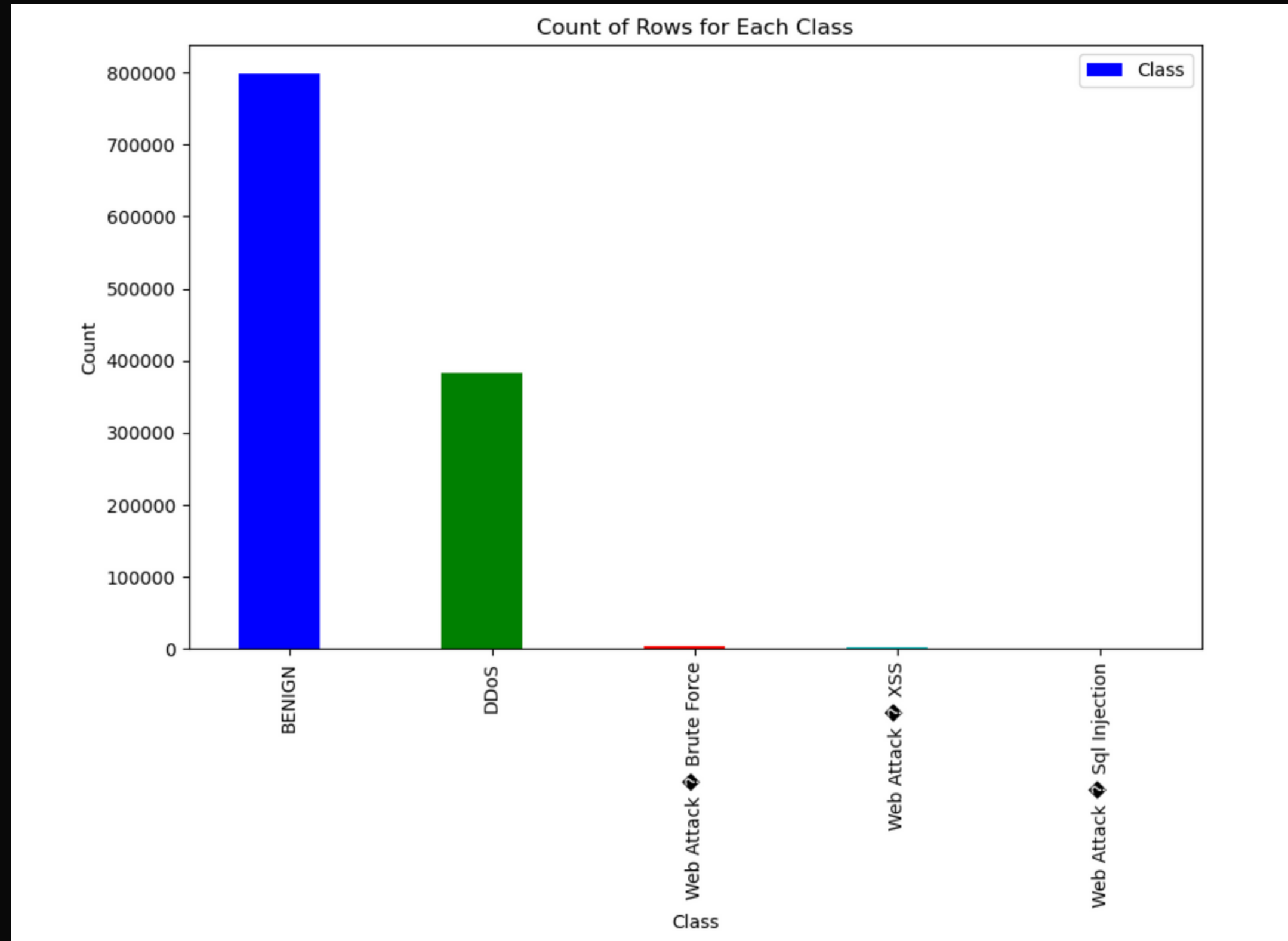
`df.info()`

`df.describe()`

```
[4]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1188333 entries, 0 to 1188332
Data columns (total 80 columns):
 #   Column                                          Non-Null Count  Dtype  
---  --
 0   Unnamed: 0                                    1188333 non-null  int64  
 1   Destination Port                             1188333 non-null  int64  
 2   Flow Duration                                1188333 non-null  int64  
 3   Total Fwd Packets                            1188333 non-null  int64  
 4   Total Backward Packets                      1188333 non-null  int64  
 5   Total Length of Fwd Packets                  1188333 non-null  int64  
 6   Total Length of Bwd Packets                  1188333 non-null  int64  
 7   Fwd Packet Length Max                       1188333 non-null  int64  
 8   Fwd Packet Length Min                       1188333 non-null  int64  
 9   Fwd Packet Length Mean                      1188333 non-null  float64 
10   Fwd Packet Length Std                       1188333 non-null  float64 
11   Bwd Packet Length Max                       1188333 non-null  int64  
12   Bwd Packet Length Min                       1188333 non-null  int64  
13   Bwd Packet Length Mean                      1188333 non-null  float64 
14   Bwd Packet Length Std                       1188333 non-null  float64 
15   Flow Bytes/s                                 1188262 non-null  float64 
16   Flow Packets/s                              1188333 non-null  float64 
17   Flow IAT Mean                               1188333 non-null  float64 
18   Flow IAT Std                                1188333 non-null  float64 
19   Flow IAT Max                                1188333 non-null  int64  
20   Flow IAT Min                                1188333 non-null  int64  
21   Fwd IAT Total                               1188333 non-null  int64  
22   Fwd IAT Mean                               1188333 non-null  float64 
23   Fwd IAT Std                                1188333 non-null  float64 
24   Fwd IAT Max                                1188333 non-null  int64  
25   Fwd IAT Min                                1188333 non-null  int64  
26   Bwd IAT Total                               1188333 non-null  int64  
27   Bwd IAT Mean                               1188333 non-null  float64 
28   Bwd IAT Std                                1188333 non-null  float64 
29   Bwd IAT Max                                1188333 non-null  int64  
30   Bwd IAT Min                                1188333 non-null  int64  
31   Fwd PSH Flags                               1188333 non-null  int64  
32   Bwd PSH Flags                               1188333 non-null  int64  
33   Fwd URG Flags                               1188333 non-null  int64  
34   Bwd URG Flags                               1188333 non-null  int64  
35   Fwd Header Length                           1188333 non-null  int64  
36   Bwd Header Length                           1188333 non-null  int64  
37   Fwd Packets/s                              1188333 non-null  float64 
38   Bwd Packets/s                              1188333 non-null  float64 
39   Min Packet Length                           1188333 non-null  int64  
40   Max Packet Length                           1188333 non-null  int64  
41   Packet Length Mean                          1188333 non-null  float64 
42   Packet Length Std                          1188333 non-null  float64 
43   Packet Length Variance                      1188333 non-null  float64
```

Visualization



Preprocessing Dataset

Determine columns containing nulls

```
: #Dropping missing rows  
df.dropna(inplace = True)  
df.drop('Unnamed: 0', axis = 1, inplace = True)  
df
```

```
df.shape
```

```
(1188262, 79)
```

Encode the target column i.e Class

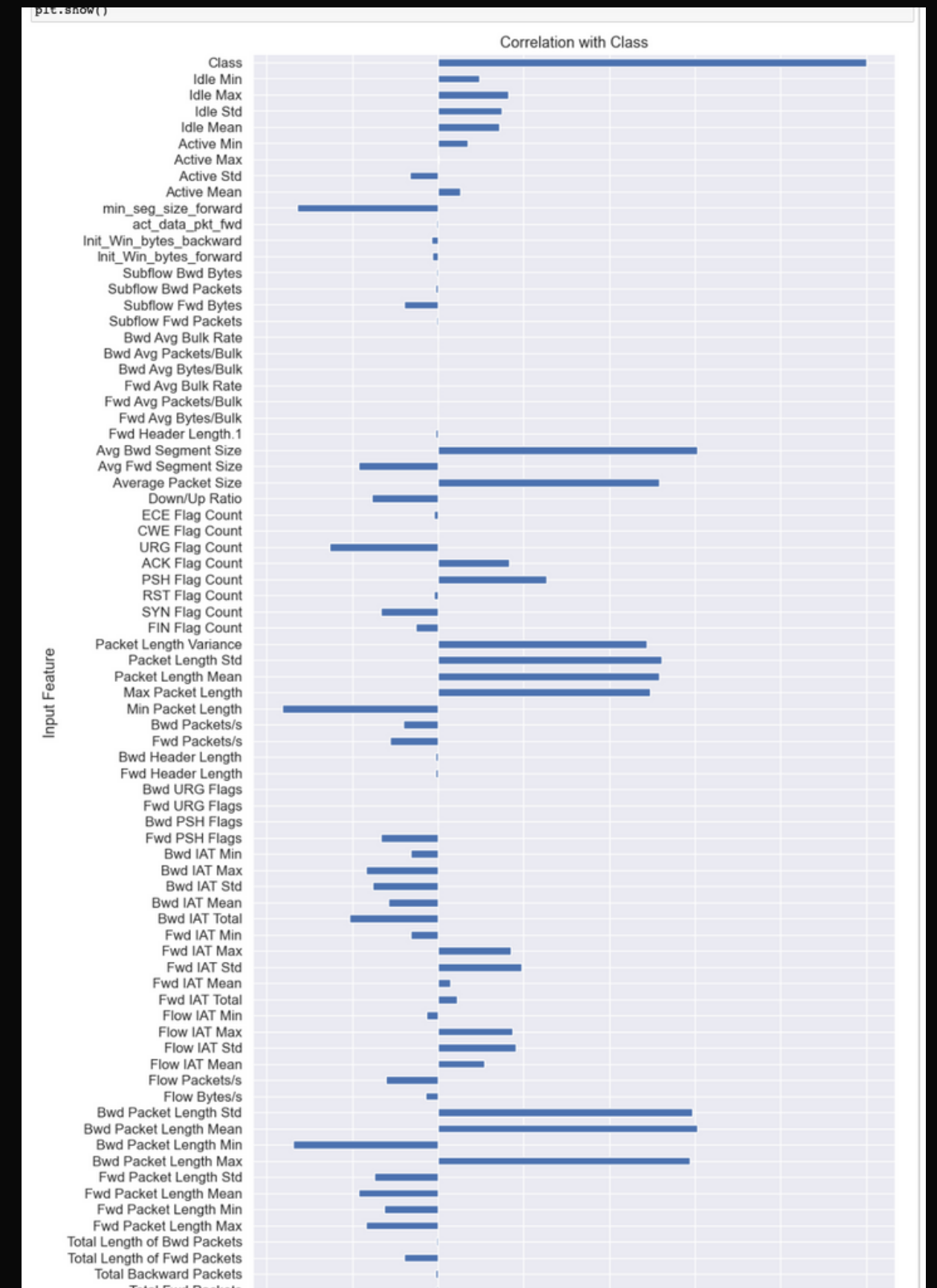
```
7]: #Encoding the Class labels to numerical values to calculate correlation  
le = LabelEncoder()  
df_1['Class'] = le.fit_transform(df['Class'])  
df_1
```

Feature Selection

Encoding the Class column

Examining Correlation of Features

Dropping uncorrelated and highly correlated columns (to avoid overfitting)

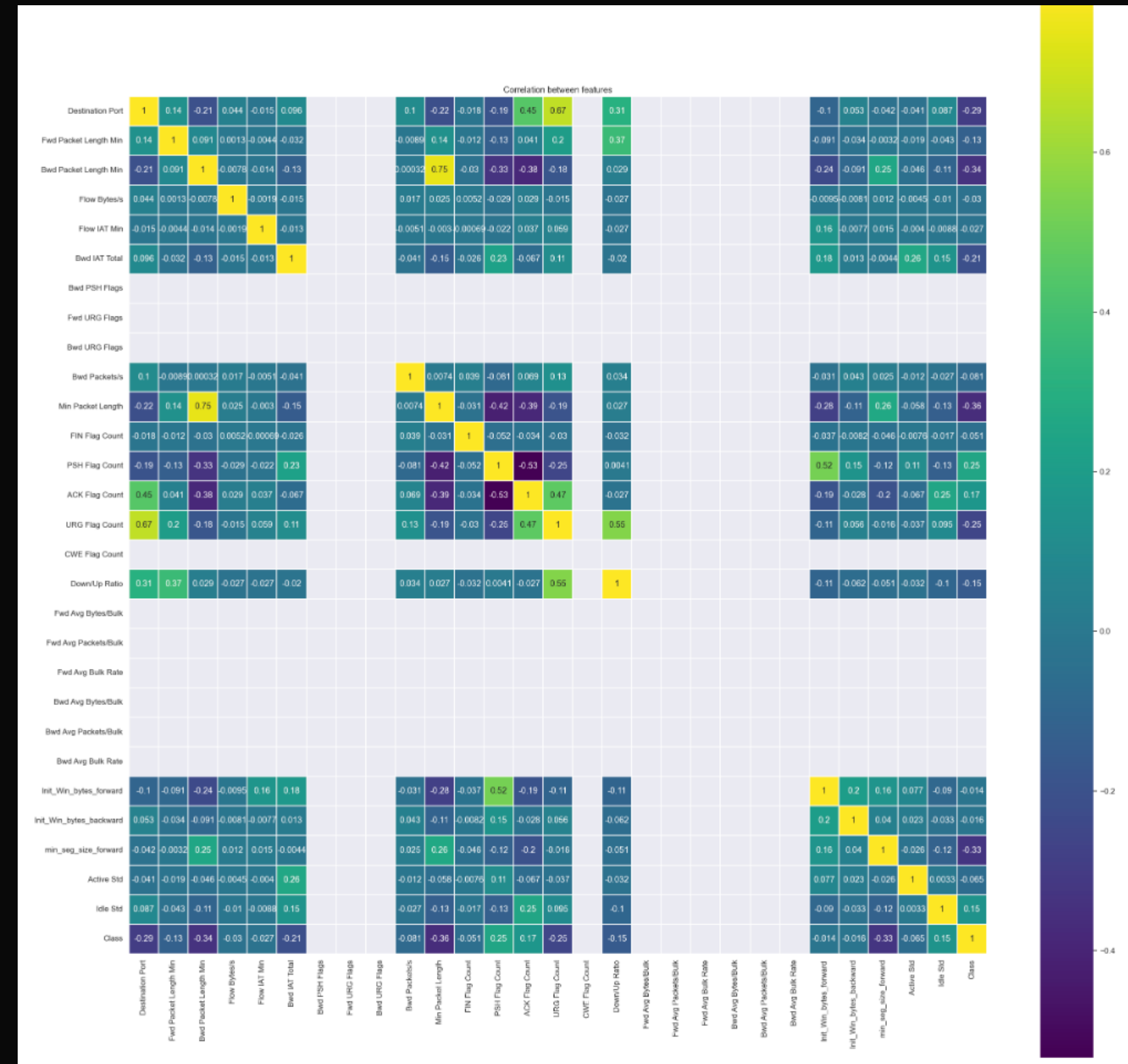


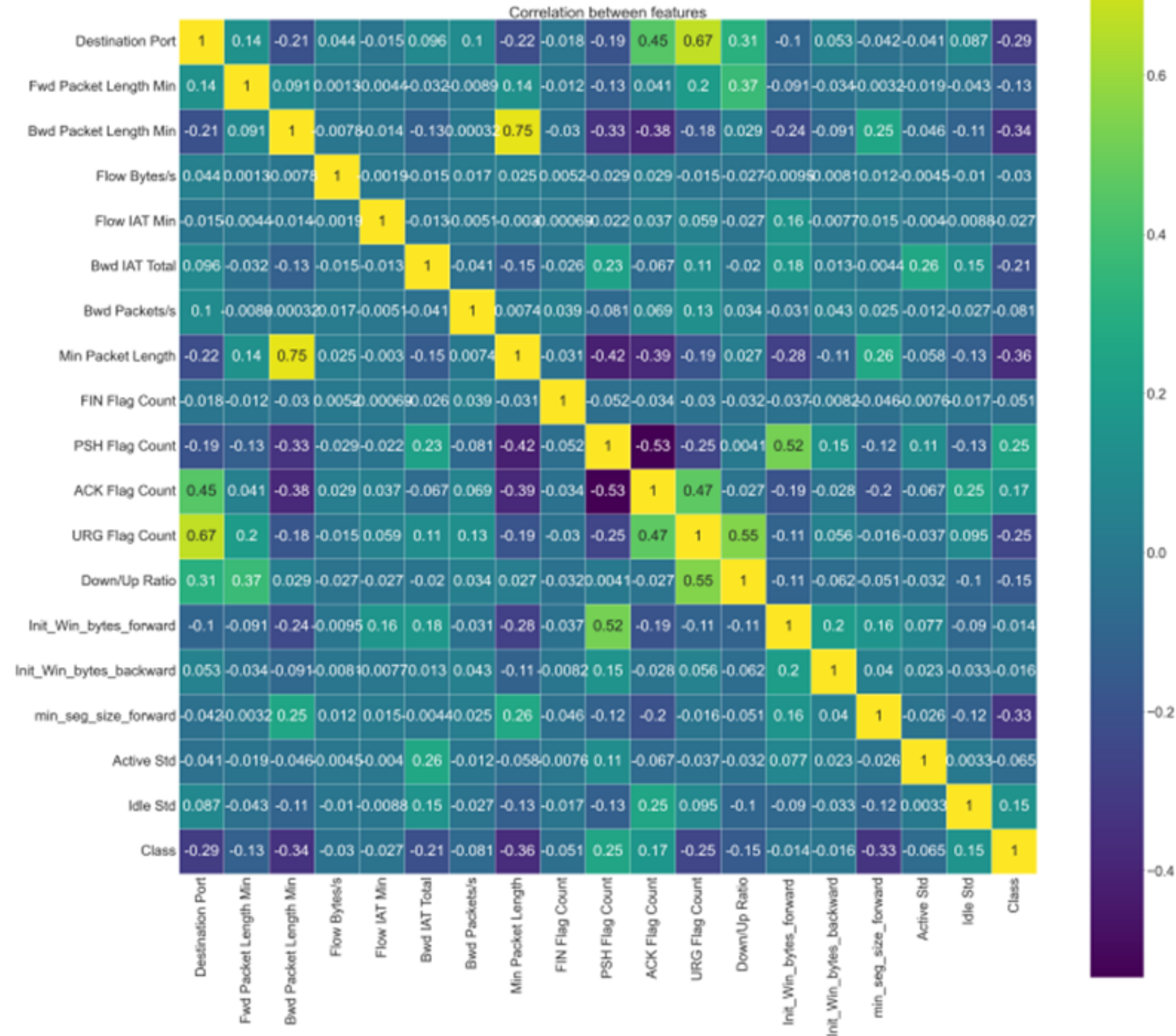
Visualizations

The first step is to compute the correlation matrix using the `corr()` function from pandas, which returns a square matrix containing the Pearson correlation coefficients between each pair of features.

Seaborn library is used to create a heatmap of the correlation matrix

Dropping uncorrelated and highly correlated columns (to avoid overfitting)





Algorithms Used

Supervised Classification Algorithms

Decision Tree Classifier

Random Forest Classifier

Logistic Regression Classifier

XGBoost Classifier

Naive Bayes Classifier

Normalization:

Normalization of Data

```
5]: X = pd.DataFrame(X, columns=X.columns)
    threshold = 1000
    X = X.clip(lower=-threshold, upper=threshold)
```

```
# Scale the data using MinMaxScaler
scaler = MinMaxScaler()
X = scaler.fit_transform(X)
```

```
6]: #Train Test Split
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```


Random Forest Classifier:

```
# Train a random forest classifier on the training set
clf = RandomForestClassifier(n_estimators=100, random_state=42, max_depth=2)
clf.fit(X_train, y_train)

# Evaluate the model on the testing set
y_pred = clf.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy: {:.2f}%".format(accuracy*100))
```

Accuracy: 94.53%

Optimising Hyperparameters

RandomizedSearchCV to find optimal hyperparameters

```
[55]: from sklearn.tree import DecisionTreeClassifier
      from sklearn.model_selection import RandomizedSearchCV
      from scipy.stats import randint

      # Define the decision tree model
      model = DecisionTreeClassifier()

      # Define the hyperparameters to tune
      params = {'max_depth': randint(2, 10),
                'min_samples_split': randint(2, 10),
                'min_samples_leaf': randint(1, 5)}

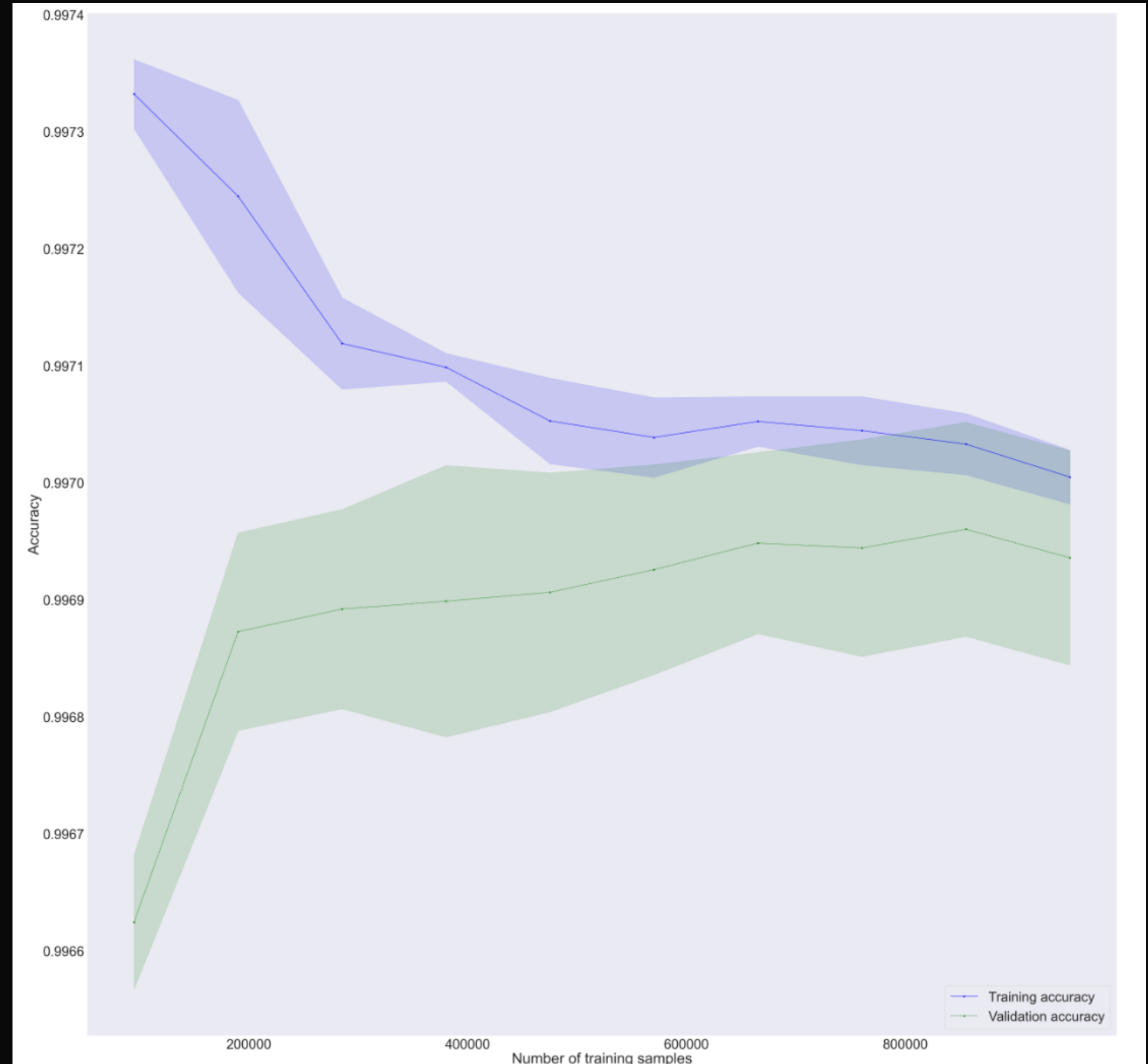
      # Use RandomizedSearchCV to find the best hyperparameters
      with parallel_backend('multiprocessing', n_jobs=4):
          random_search = RandomizedSearchCV(model, params, cv=5, n_iter=50, random_state=42)
          random_search.fit(X_train, y_train)

      # Train the model with the best hyperparameters
      best_model = random_search.best_estimator_
      best_model.fit(X_train, y_train)

[55]: DecisionTreeClassifier(max_depth=9, min_samples_split=3)
```

Learning Curve

The learning curve shows the relationship between the size of the training set and the performance of the model, such as accuracy. It helps to identify if the model is underfitting (high bias) or overfitting (high variance). The learning curve is important because it helps to determine the optimal amount of data needed to train a model effectively and to identify any issues with the model's performance.

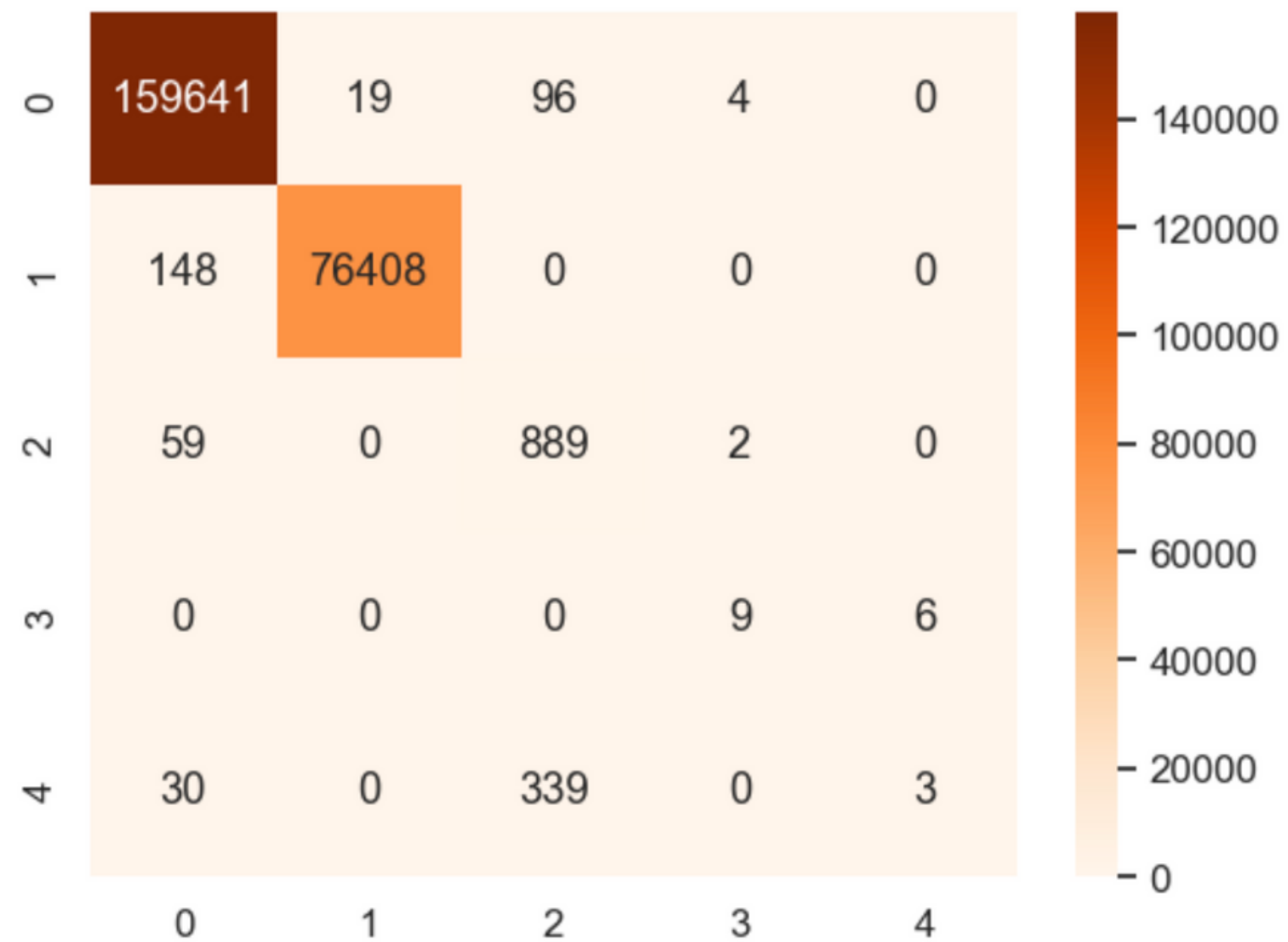


Confusion Matrix

Confusion Matrix

```
In [63]: import seaborn as sns
sns.set(font_scale=1.10)
sns.heatmap(cm, annot=True, cmap='Oranges', fmt = 'd')
```

Out[63]: <AxesSubplot:>



Cross Validation

Cross Validation

```
In [33]: from sklearn.model_selection import cross_val_score
cv_scores = cross_val_score(best_model, X_train, y_train, cv=5)

# Calculate the average cross-validation score
avg_cv_score = np.mean(cv_scores)

cv_scores
```

```
Out[33]: array([0.99684413, 0.9970177 , 0.99694407, 0.99697037, 0.99688619])
```

Conclusion

It is clear that the XGBoost Classifier outperformed the other models in terms of accuracy with an impressive score of 99.74%. However, it was observed that the XGBoost model overfit the data, as it performed significantly better on the training set than on the testing set. The Logistic Regression Classifier also performed well with an accuracy score of 97.78%. The Random Forest Classifier and Naive Bayes Classifier performed slightly lower, but still achieved respectable accuracy scores of 94.53% and 95.58% respectively.

Resources

Data Sources: <https://www.kaggle.com/datasets/subhajournal/sdn-intrusion-detection>