



Informatics Institute of Technology

Department of Computing

Bsc(Hons) Artificial Intelligence and Data Science

Module: CM2604 Machine Learning

Module Coordinator: Mr. Prasan Yapa

Coursework Report

Vinuwara Ronath Jayasuriya

RGU ID – 2119942

IIT Student No. – 20210167

Spam dataset is done using the Decision Trees Classification and K Nearest Neighbors Classification (KNN)

Preprocessing techniques

Corpus Preparation

Data Cleaning: This was done by removing duplicates and null values from the dataset.

Removing duplicate in dataset.

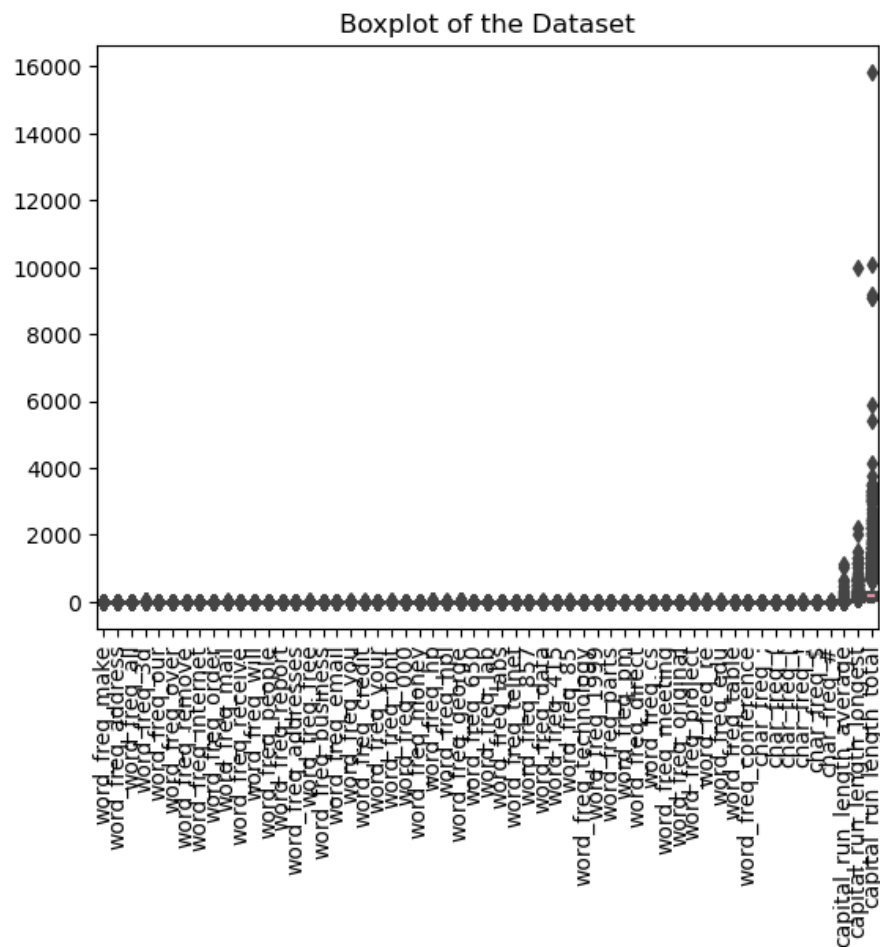
Getting the duplicates in the dataset	
Data set.duplicated()	
0	False
1	False
2	False
3	False
4	False
...	
4596	False
4597	False
4598	False
4599	False
4600	False
Length: 4601, dtype: bool	
removing the duplicate values	
Data_set.drop_duplicates(inplace=True)	

Removing null values

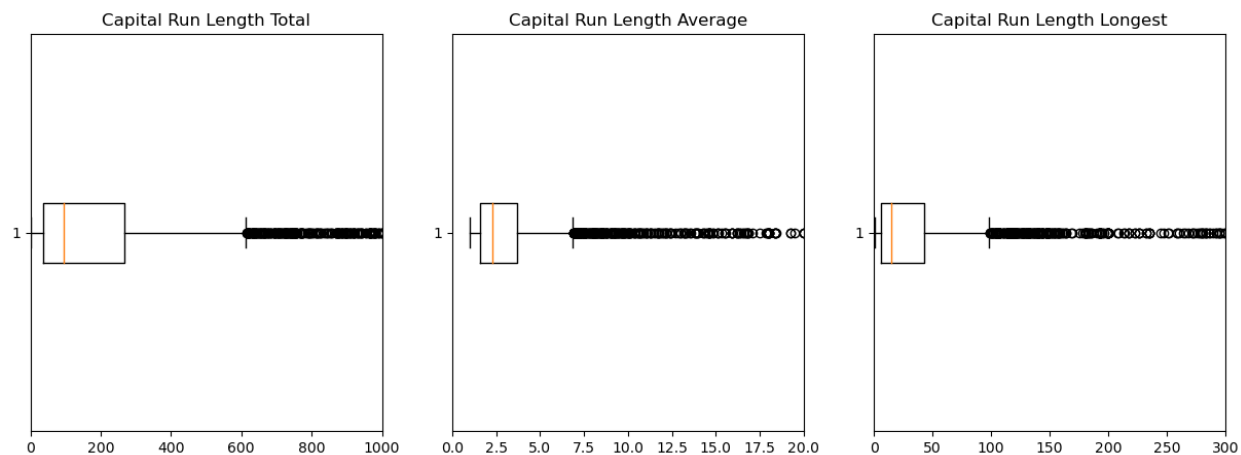
Data set.isna().sum()	
word_freq_make	0
word_freq_address	0
word_freq_all	0
word_freq_3d	0
word_freq_our	0
word_freq_over	0
word_freq_remove	0
word_freq_internet	0
word_freq_order	0
word_freq_mail	0
word_freq_receive	0
word_freq_will	0
word_freq_people	0
word_freq_report	0
word_freq_addresses	0
word_freq_free	0
word_freq_business	0
word_freq_email	0
word_freq_you	0

word_freq_credit	0
word_freq_your	0
word_freq_font	0
word_freq_000	0
word_freq_money	0
word_freq_hp	0
word_freq_hpl	0
word_freq_george	0
word_freq_650	0
word_freq_lab	0
word_freq_labs	0
word_freq_telnet	0
word_freq_857	0
word_freq_data	0
word_freq_415	0
word_freq_85	0
word_freq_technology	0
word_freq_1999	0
word_freq_parts	0
word_freq_pm	0
word_freq_direct	0
word_freq_cs	0
word_freq_meeting	0
word_freq_original	0
word_freq_project	0
word_freq_re	0
word_freq_edu	0
word_freq_table	0
word_freq_conference	0
char_freq_;	0
char_freq_(0
char_freq_[0
char_freq_!	0
char_freq_\$	0
char_freq_#	0
capital_run_length_average	0
capital_run_length_longest	0
capital_run_length_total	0
spam	0
dtype: int64	
Removing the target column	
data=Data_set.drop(labels=['spam'], axis=1)	
data.head()	

Data Transformation: outliers were removed using standard scaling.



Since capital run length avarage, capital run length longest and Capital run length total are outliers



Remove the outliers using the IQR method

code

```
import numpy as np

# Calculate the interquartile range for each column
Q1 = Data_set.quantile(0.25)
Q3 = Data_set.quantile(0.75)
IQR = Q3 - Q1

# Remove the outliers using the IQR method
Data_set_outliers_removed = Data_set[~((Data_set < (Q1 - 1.5 * IQR)) | (Data_set > (Q3 + 1.5 * IQR))).any(axis=1)]

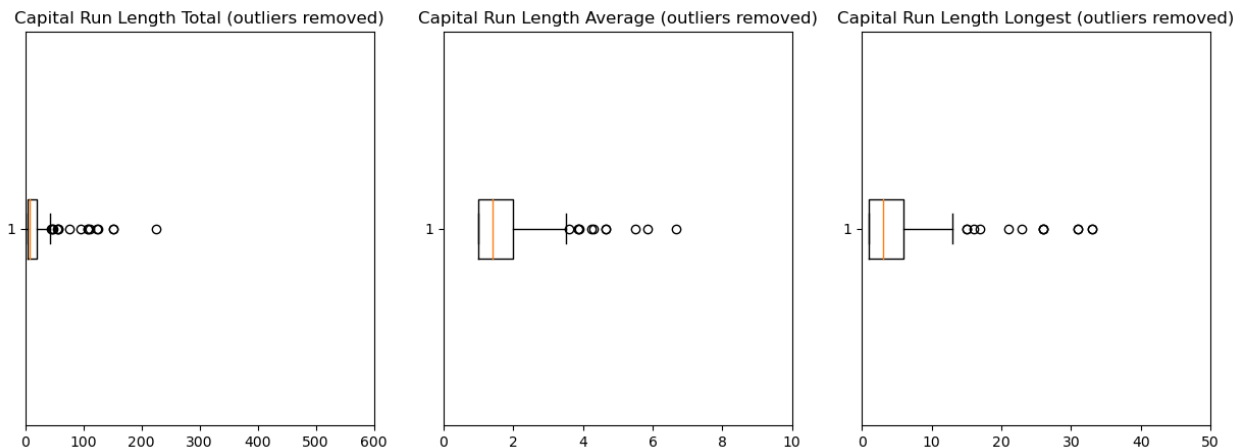
# Create individual box plots for the columns with outliers removed
fig, axs = plt.subplots(1, 3, figsize=(15, 5))

axs[0].boxplot(Data_set_outliers_removed['capital_run_length_total'], vert=False)
axs[0].set_title('Capital Run Length Total (outliers removed)')
axs[0].set_xlabel("")
axs[0].set_xlim([0, 600])

axs[1].boxplot(Data_set_outliers_removed['capital_run_length_average'], vert=False)
axs[1].set_title('Capital Run Length Average (outliers removed)')
axs[1].set_xlabel("")
axs[1].set_xlim([0, 10])

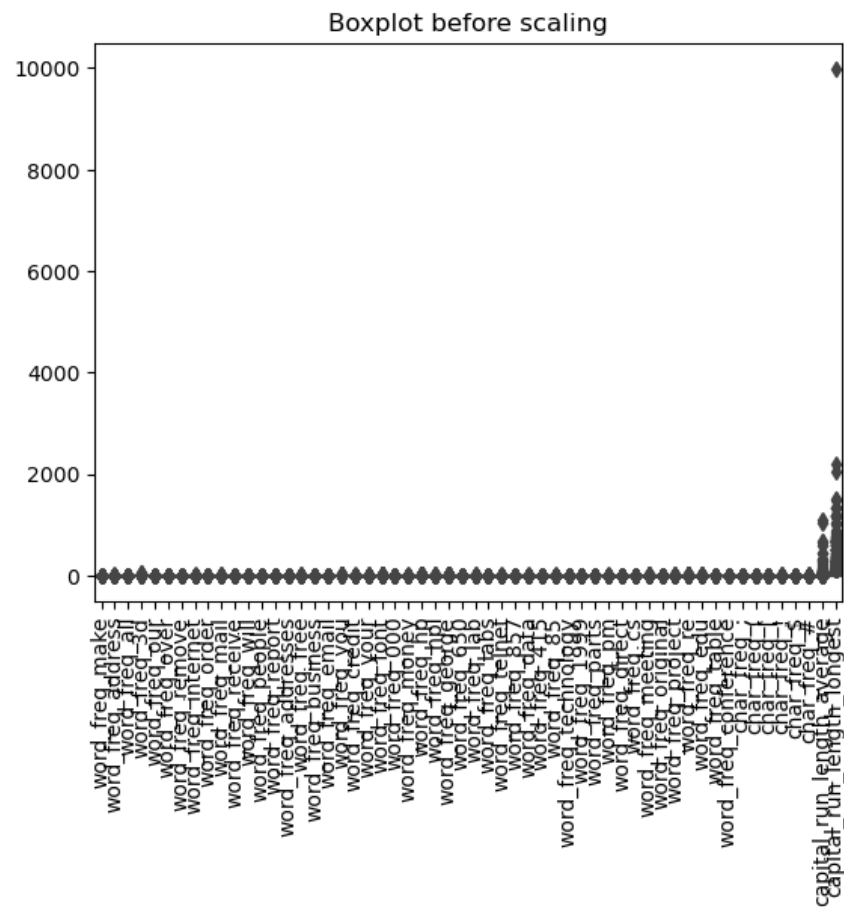
axs[2].boxplot(Data_set_outliers_removed['capital_run_length_longest'], vert=False)
axs[2].set_title('Capital Run Length Longest (outliers removed)')
axs[2].set_xlabel("")
axs[2].set_xlim([0, 50])

plt.show()
```

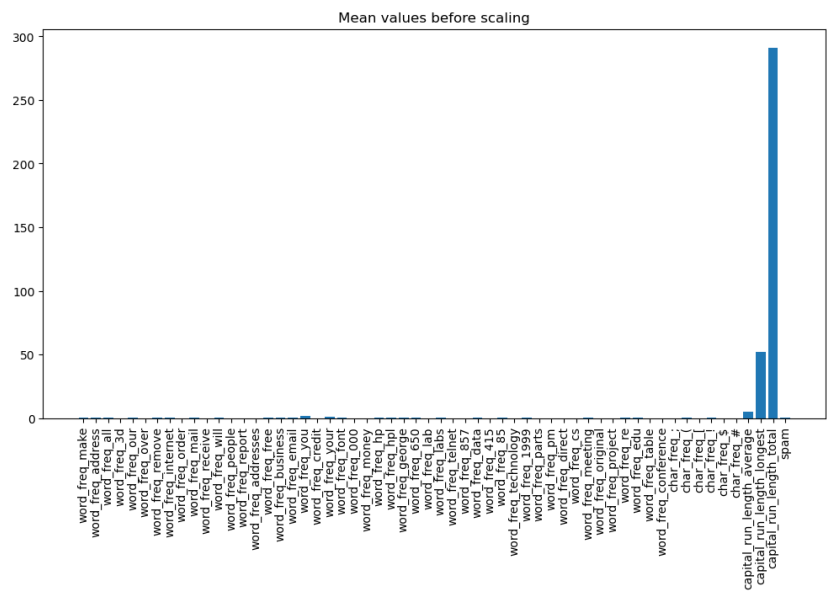


Used Standard Scaler on the dataset

Before standard scaler



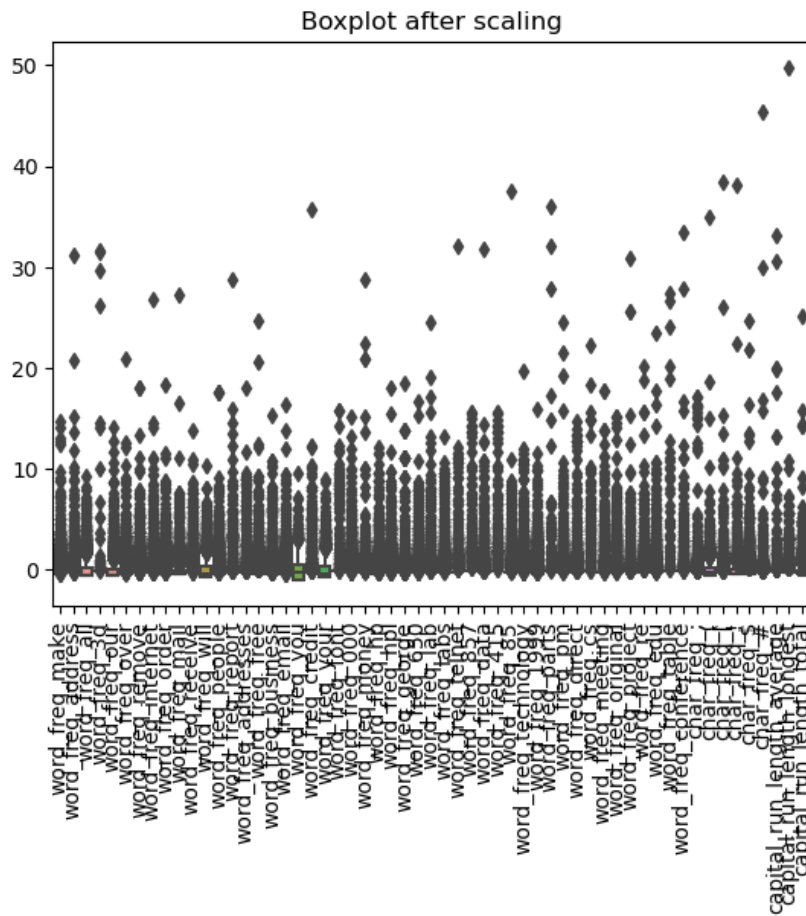
Mean



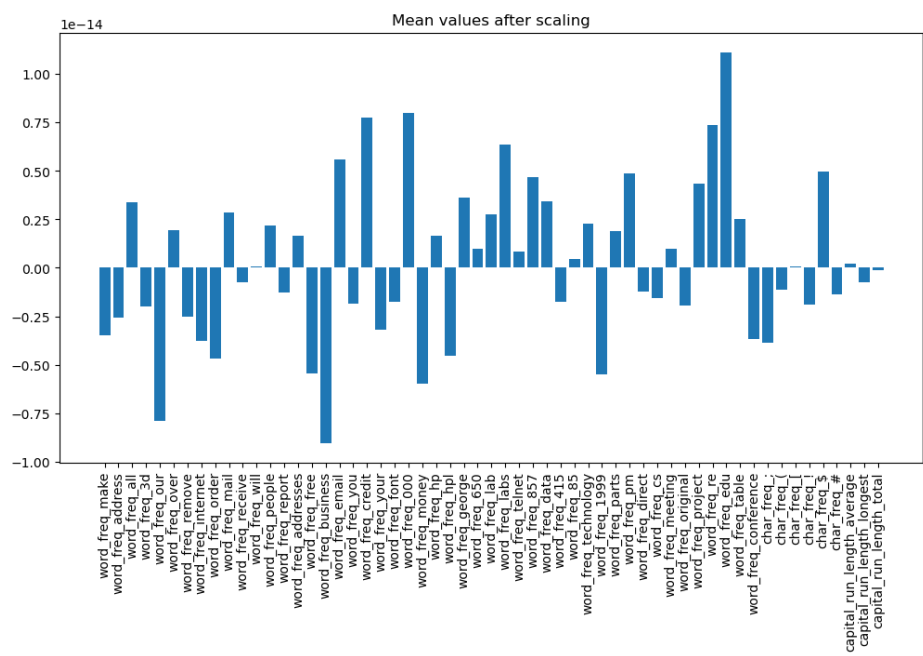
Code – performing standard scaler on the dataset

```
# Perform Standard Scaling on the data
scaler = StandardScaler()
scaled_data = scaler.fit_transform(Data_set.iloc[:, :-1])
Data_scaled = pd.DataFrame(data=scaled_data, columns=Data_set.columns[:-1])
```

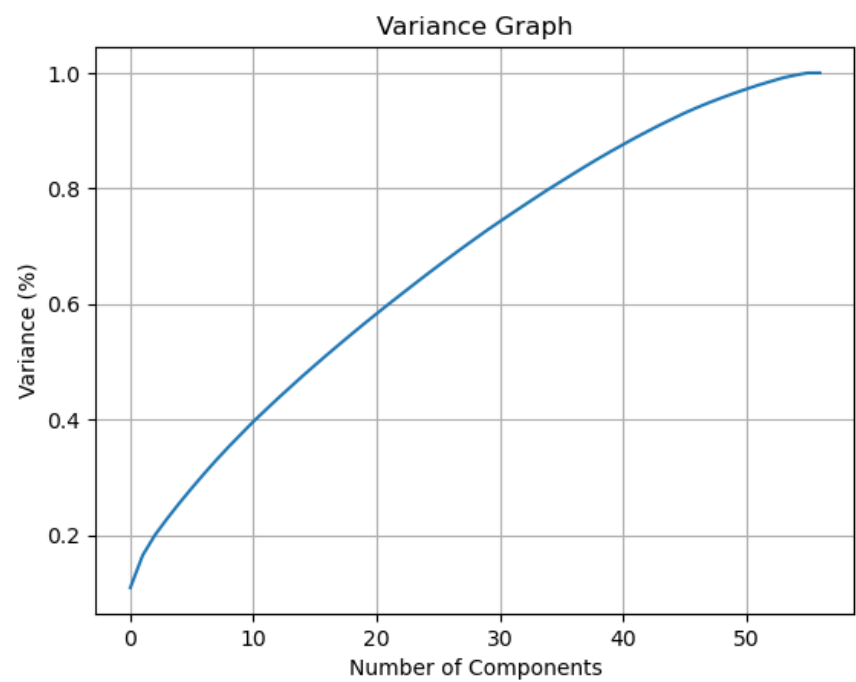
Boxplot after performing the standard scaler on the dataset



Mean



Variance graph



Performing PCA to the Dataset

Code

```
pca = PCA(n_components=44)
new_data = pca.fit_transform(Data_set)

# The new data fed to the algorithm.
principal_Df = pd.DataFrame(data = new_data, columns = [
    'PC1', 'PC2', 'PC3', 'PC4', 'PC5', 'PC6', 'PC7', 'PC8', 'PC9', 'PC10',
    'PC11', 'PC12', 'PC13', 'PC14', 'PC15', 'PC16', 'PC17', 'PC18', 'PC19', 'PC20',
    'PC21', 'PC22', 'PC23', 'PC24', 'PC25', 'PC26', 'PC27', 'PC28', 'PC29', 'PC30',
    'PC31', 'PC32', 'PC33', 'PC34', 'PC35', 'PC36', 'PC37', 'PC38', 'PC39', 'PC40',
    'PC41', 'PC42', 'PC43', 'PC44'
])
```

Splitting the dataset into the Training set and Test set

- 20 percent of the dataset was used to test the dataset.
- 80 percent of the dataset was used to train the dataset

Code

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.20, random_state=0)
```

Decision Tree Classification

Accuracy of testing dataset : 0.9560570071258907

Summary of the dataset

code

```
from sklearn.metrics import classification_report
predicted_labels = clf.predict(X_test)
classification_report = classification_report(y_test, predicted_labels)

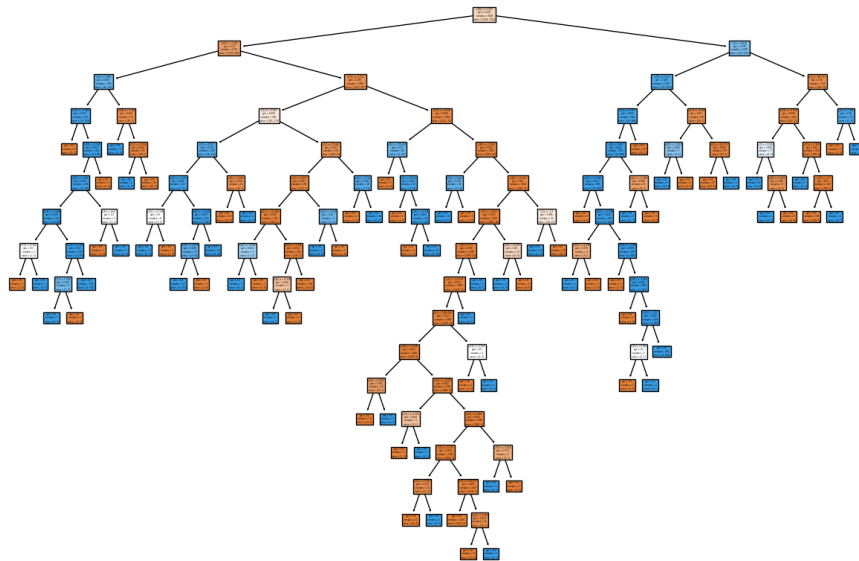
print(classification_report)
```

Classification report

	precision	recall	f1-score	support
0	0.96	0.96	0.96	495
1	0.94	0.94	0.94	347
accuracy			0.95	842
macro avg	0.95	0.95	0.95	842
weighted avg	0.95	0.95	0.95	842

Accuracy of the training dataset of the decision tree: 1.0

Decision tree



Confusion Matrix before pruning

Equations

True positive rate: measures the proportion of positive instances correctly identified as positive.

- FN = False Negative : Number of positive instances that are incorrectly classified as positive.

$$TPR = TP / (TP + FN)$$

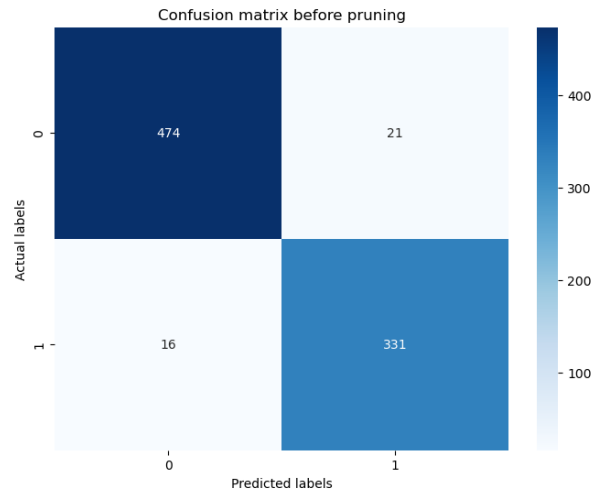
False positive rate : measures the proportion of negative instances incorrectly identified as positive.

- TN = true negative : The number of negative instances that are correctly classified as negative

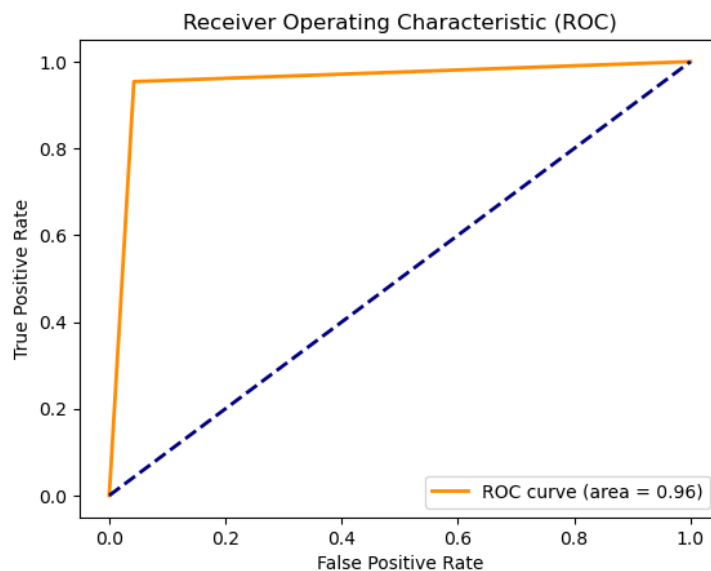
$$FPR = FP / (FP + TN)$$

Experiments

- TP :474 mails are correctly identified as not spam.
- FP : 21 mails are incorrectly identified as spam as not spam.
- TN :331 mails are correctly identified as spam.
- FN : 16 mails are incorrectly identified as not spam as spam.



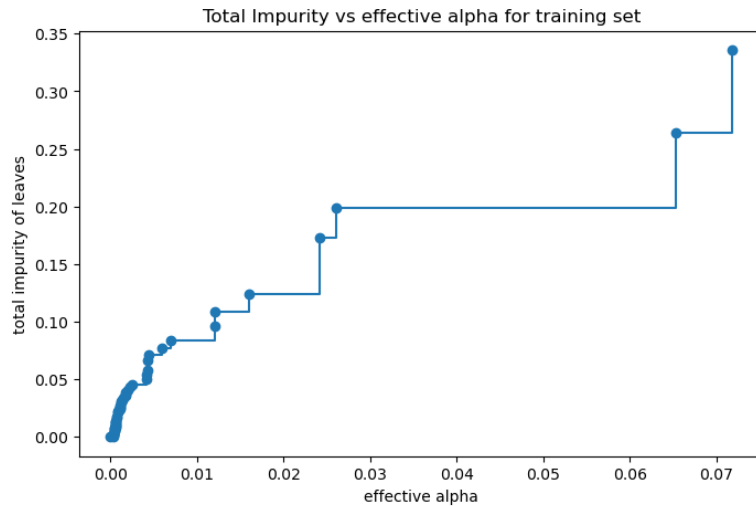
Receiver Operating Characteristics (ROC) : A plot to test sensitivity as the y coordinate versus its 1-specificity or false positive rate (FPR) as the x coordinate, is an effective method of evaluating the performance of diagnostic tests.



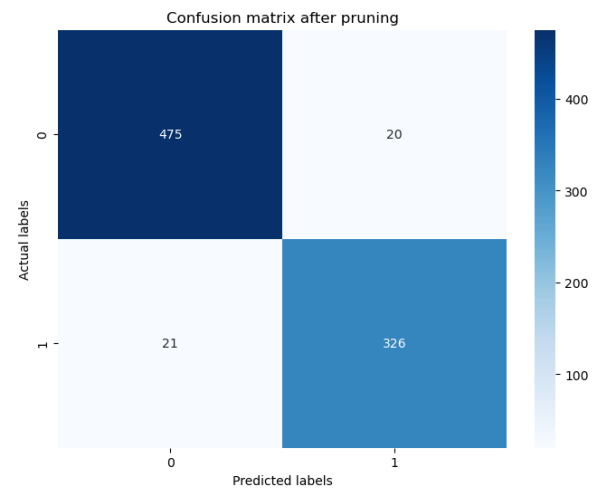
AUC score for the code: 0.96

Pruning the Decision tree : Pruning is a technique that removes the parts of the Decision Tree which prevent it from growing to its full depth. The parts that it removes from the tree are the parts that do not provide the power to classify instances.

A graph is plotted between 'Total impurity of leaves' as Y axis and 'Effective alpha' as the X axis. Using the graph, we can find the optimal 'ccp alpha' value required for pruning. The regularization parameter ccp alpha balances the accuracy and model complexity in decision trees.

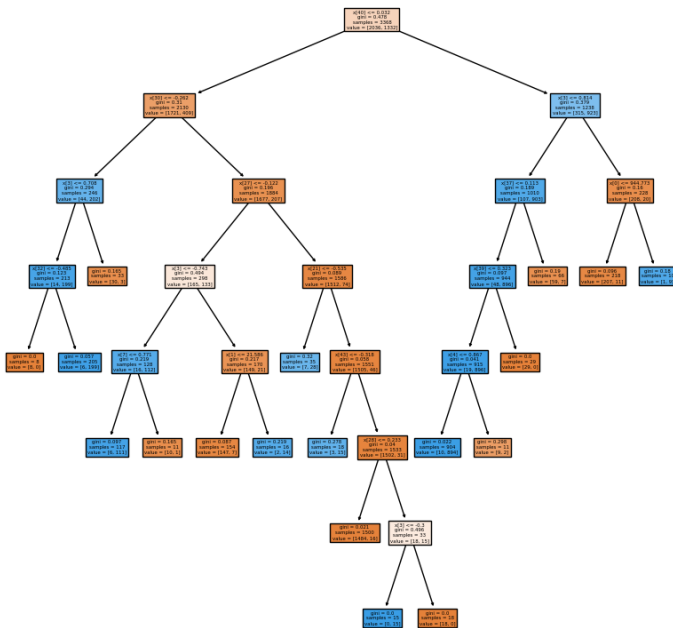


- TP :475 mails are correctly identified as not spam.
- FP : 20 mails are incorrectly identified as spam as not spam.
- TN :326 mails are correctly identified as spam.
- FN : 21 mails are incorrectly identified as not spam as spam.



	precision	recall	f1-score	support
0	0.96	0.96	0.96	495
1	0.94	0.94	0.94	347
accuracy			0.95	842
macro avg	0.95	0.95	0.95	842
weighted avg	0.95	0.95	0.95	842

Decision Tree obtaining after pruning.

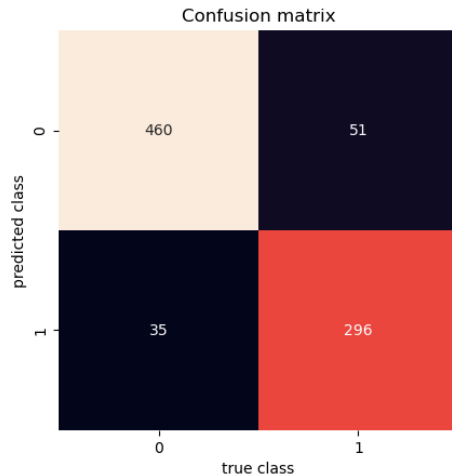


K Nearest Neighbors (KNN) Classification

Classification Report of KNN

	precision	recall	f1-score	support
0	0.90	0.93	0.91	495
1	0.89	0.85	0.87	347
accuracy			0.90	842
macro avg	0.90	0.89	0.89	842
weighted avg	0.90	0.90	0.90	842

Confusion matrix of KNN



- TP: 480 mails are correctly identified as not spam.
- FP: 51 mails are incorrectly identified as spam as not spam.
- TN: 296 mails are correctly identified as spam.
- FN: 35 mails are incorrectly identified as spam as spam.

KNN Accuracy: 0.8729216152019003

Limitations & Enhancements

Limitations

- Due to the dataset being small the accuracy of the model would be effected.
- Feature engineering step may not have included all the relevant information, making the model performance suboptimal.

Future Enhancements

- Training the model with a larger dataset would increase the accuracy of it
- Utilizing other machine learning algorithms
 - Support vector machines - SVM
 - Random Forest
- Experiment with different feature engineering techniques