CLASSIFICATION PROJECT

# Mobile Price Range Prediction

*Team Data Defenders*

SHUBHAM SARTAPE

LOKESH TOKAS

SARASWAT MUKHERJEE

CHARAN

# Content

# Problem Statement

- To predict the price range of Mobile Phones based on the available features such as RAM, camera, battery, internal memory, cores, clock speed, etc...
- The Target Variables are classified into 4 types as below.

  *0 - Low Cost Phones*

  *1 - Medium Cost Phones*

  *2 - High Cost Phones*

  *3 - Very High Cost Phones*
- This will help mobile phone market companies to understand sales data of mobile phones and factors which drive the prices.
- The objective is to find out some relation between features of a mobile phone(eg:- RAM, Internal Memory, etc) and its selling price.

# Data Summary

## Independent Variables

**Battery Power** - Total energy a battery can store in one time measured in mAh

**Blue** - Has Bluetooth or not

**Clock_speed** - speed at which microprocessor executes instructions

**Dual_sim** - Has dual sim support or not

**Fc** - Front Camera mega pixels

**Four_g** - Has 4G or not

**Int_memory** - Internal Memory in Gigabytes

**M_dep** - Mobile Depth in cm

## Independent Variables

**Mobile_wt** - Weight of mobile phone

**N_cores** - Number of cores of processor

**Pc** - Primary Camera mega pixels

**Px_height** - Pixel Resolution Height

**Px_width** - Pixel Resolution Width

**Ram** - Random Access Memory in Mega Bytes

**Sc_h** - Screen Height of mobile in cm

**Sc_w** - Screen Width of mobile in cm

**Talk_time** - longest time that a single battery charge will last over a call

**Three_g** - Has 3G or not

# Data Summary

_Independent Variables_

**Touch_screen** - Has touch screen or
not

**Wifi** - Has wifi or not

_Dependent Variables_

**Price_range -** This is the target variable
with value of
0 (low cost), 1 (medium cost), 2 (high
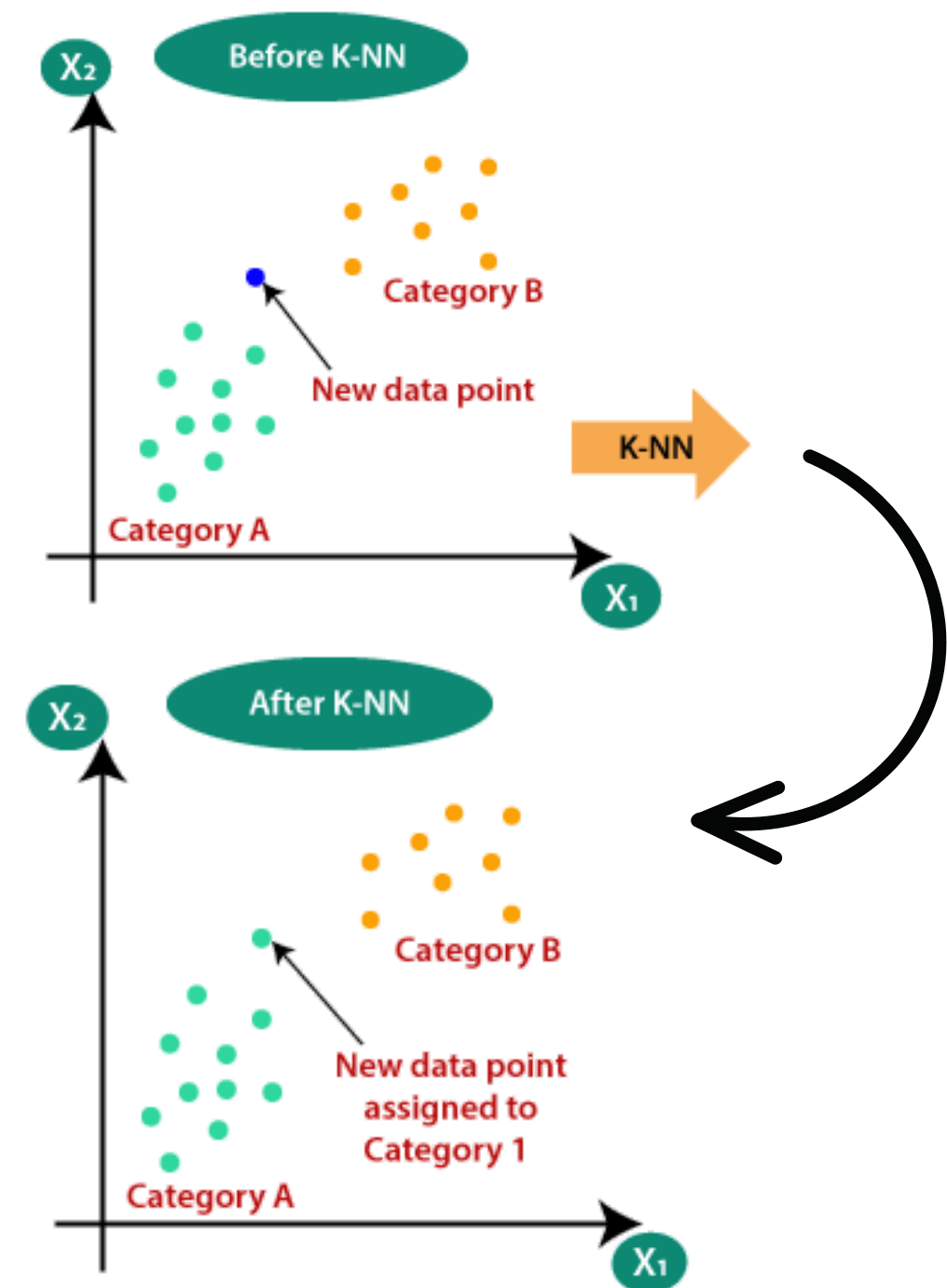cost) and 3 (very high cost).

# EDA - Data Cleaning

*Detecting Data Anomaly*

Following Anomalies were found

---

> px_height (Pixel Height) = 0
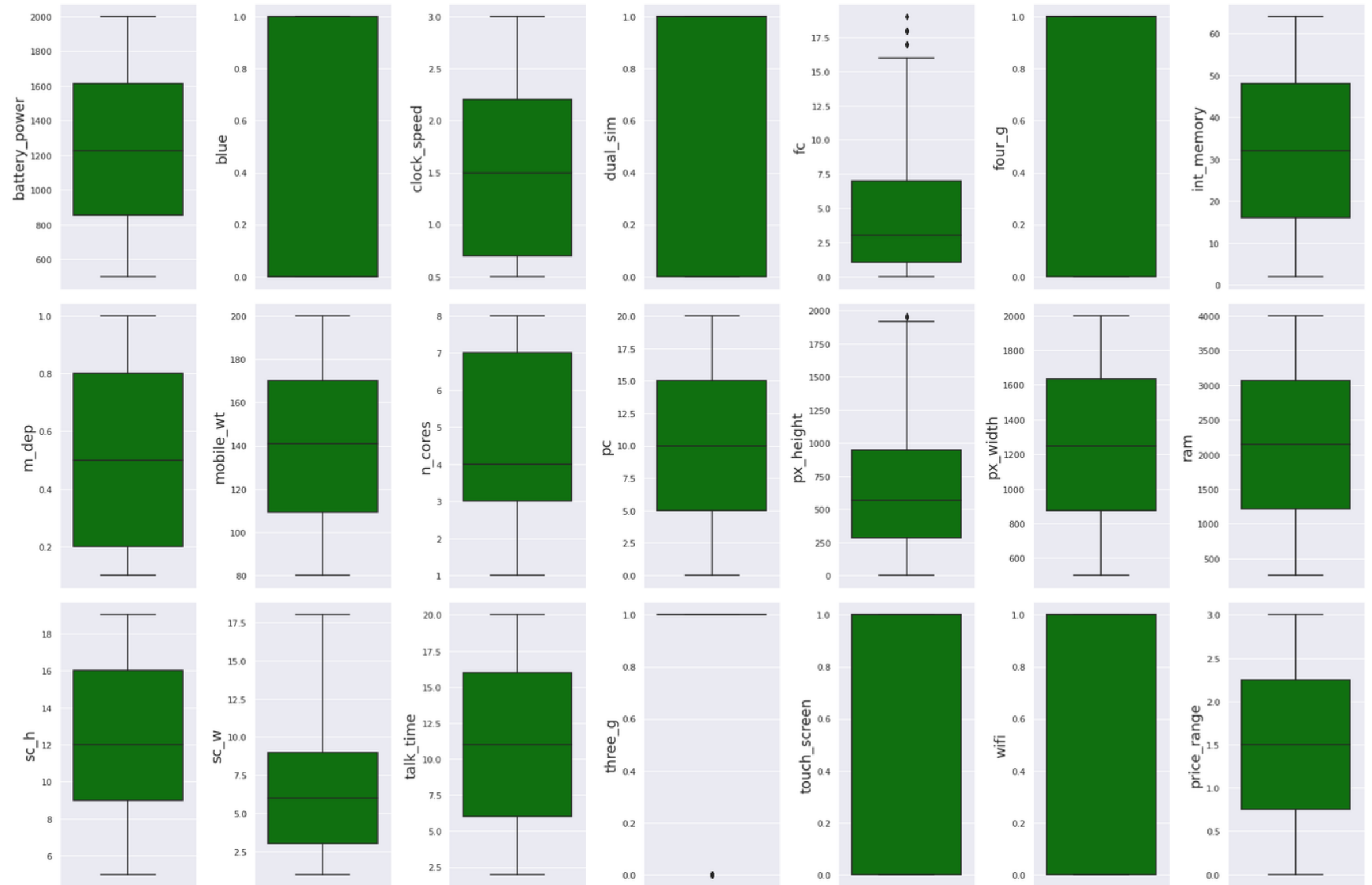
> sc_w (screen width) = 0

---

Replaced the anomalous values using KNN Imputer by assigning nearest possible value and not Mean/Avg Value.
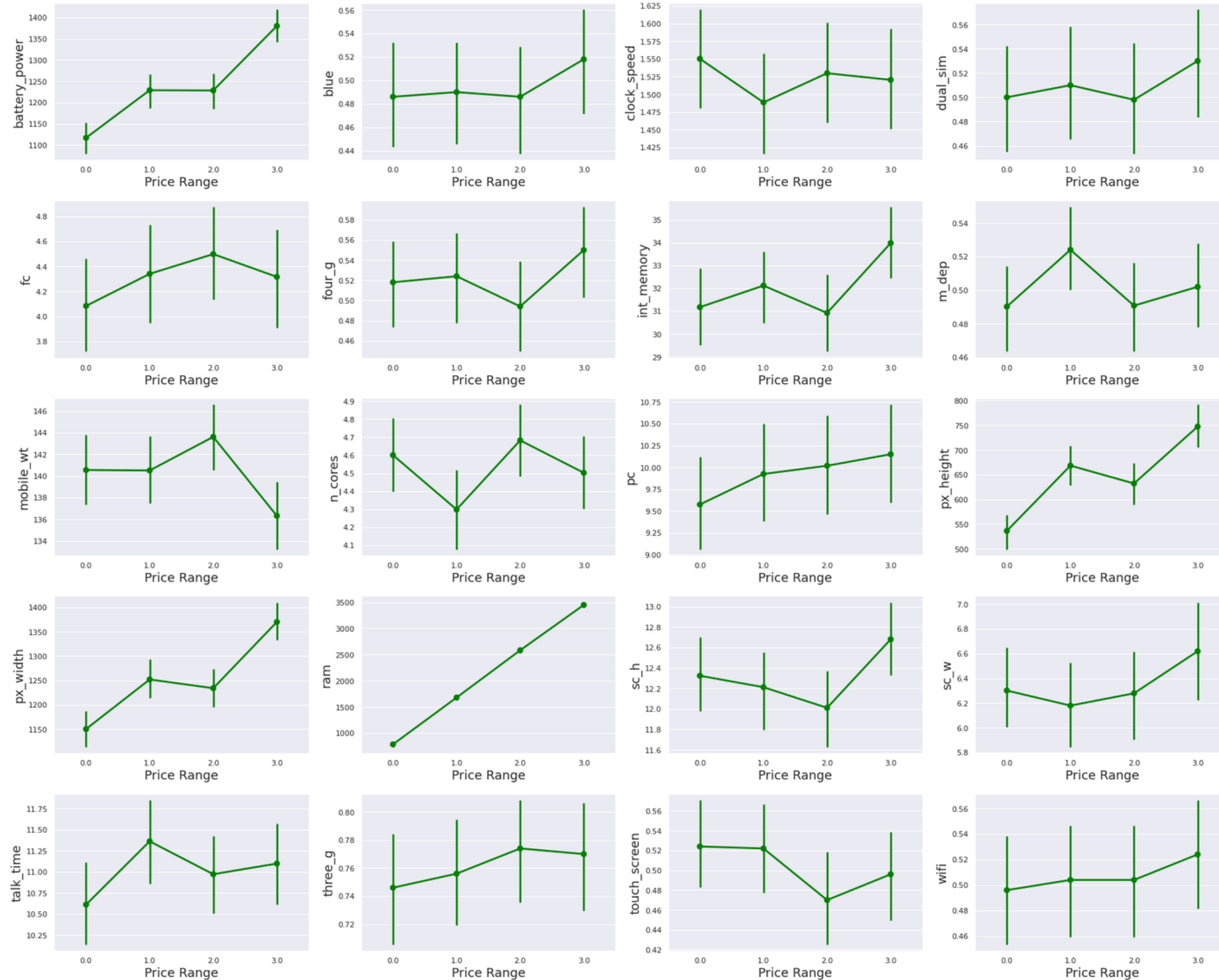
# EDA Outlier Detection

**No Extreme Outliers detected.**

- In 'px_height', 1 possible outlier was highlighted but after examination it was observed to be within reason.
- In 'fc' few observations were out of bound but normal to have high mp for experimental purpose or emerging technology.
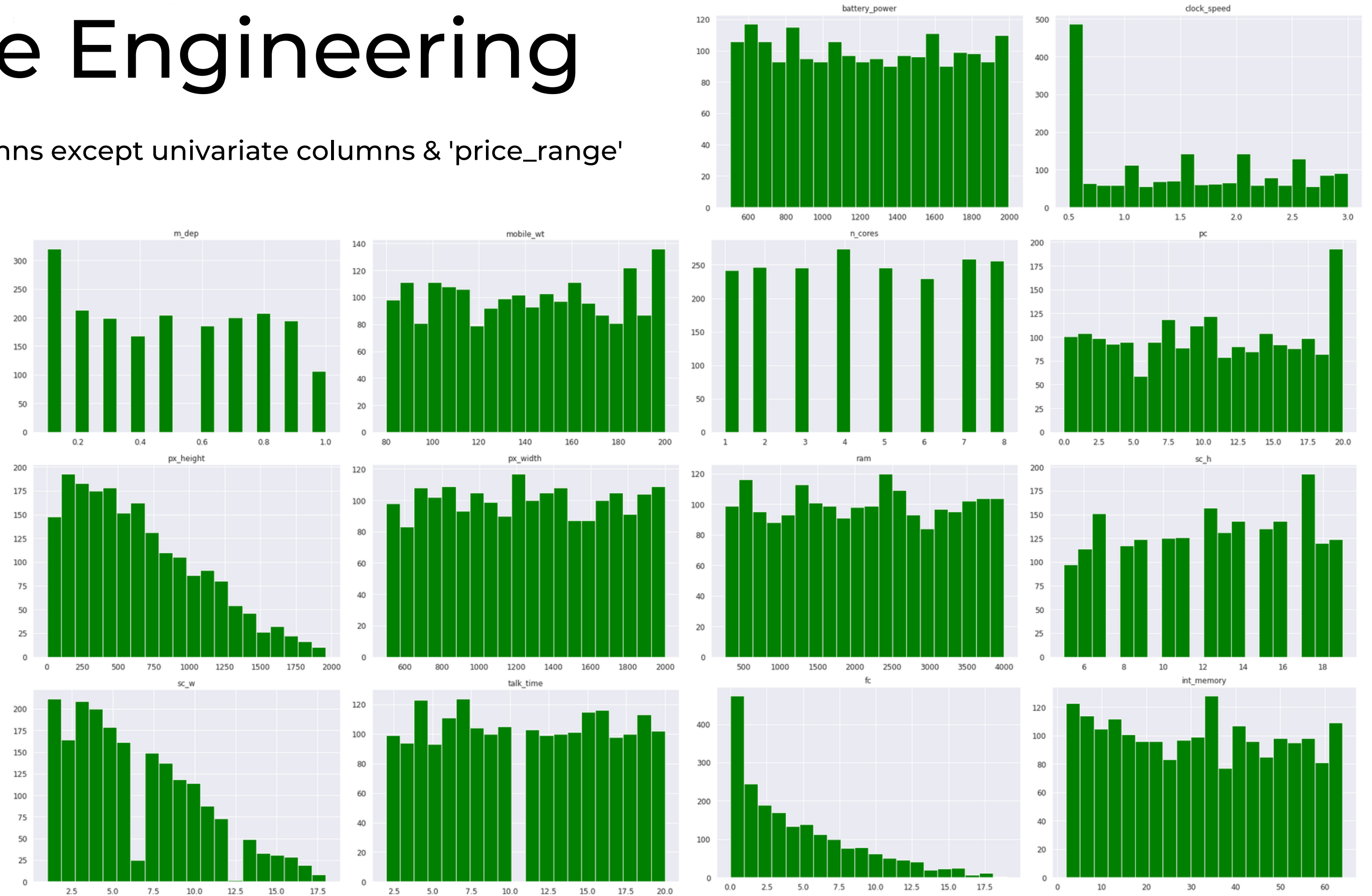
# Feature Engineering

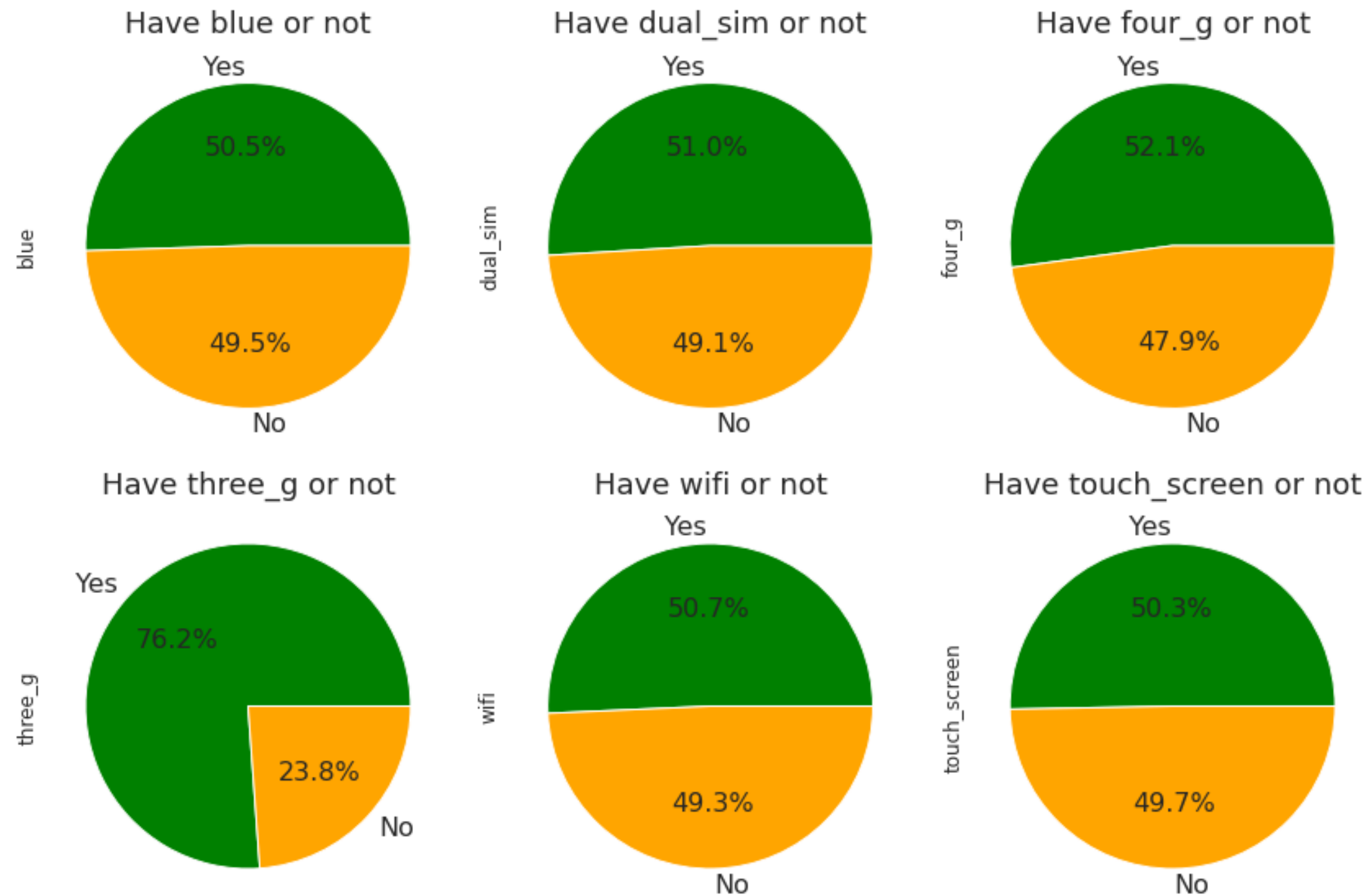Point Plot of all features on 'price_range'

# Feature Engineering

Bar Plot of all coloumns except univariate columns & 'price_range'
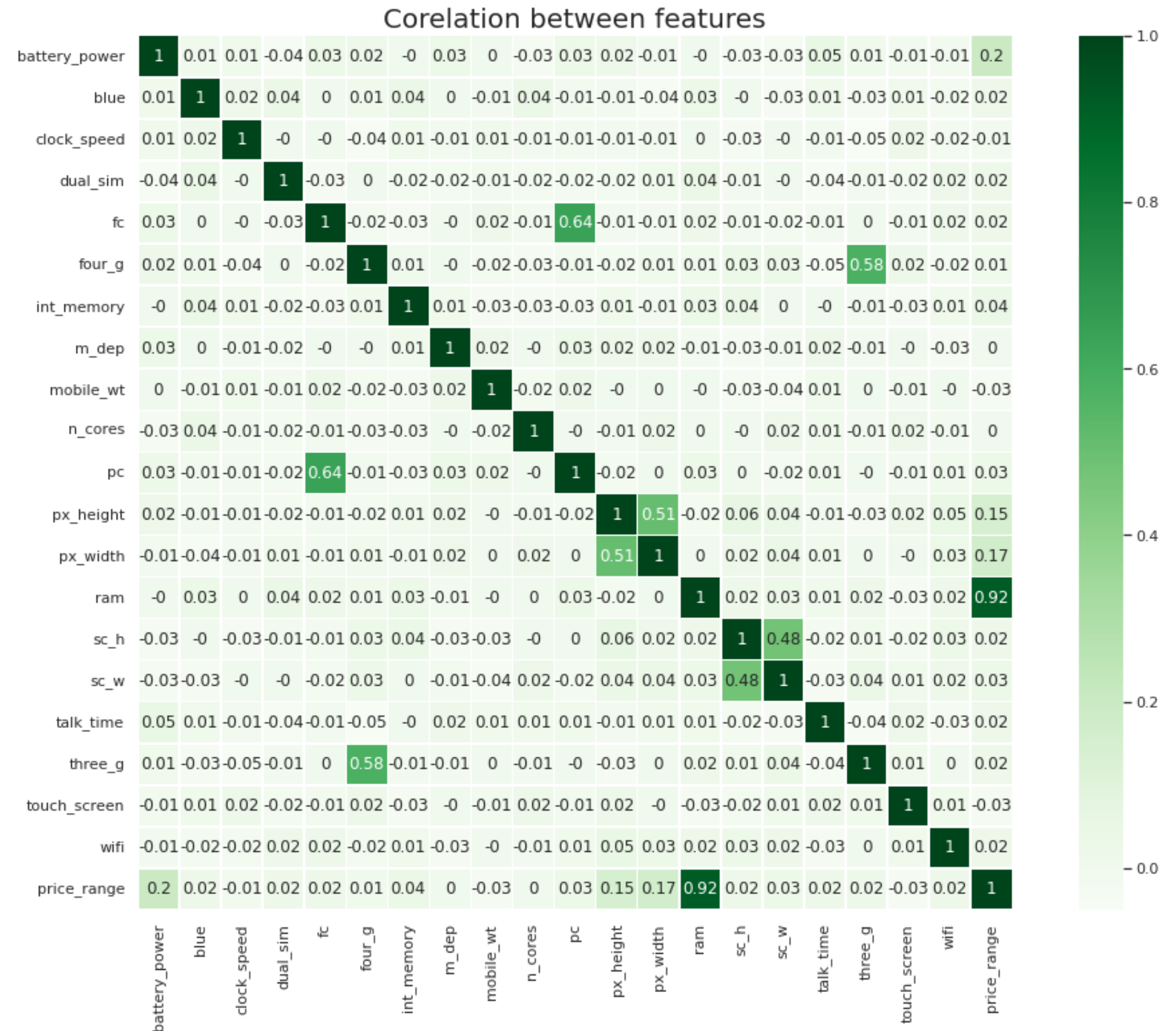
# Feature Engineering

Pie Plot for Univariate Analysis

# Feature Engineering

Heatmap to visualise - Corelation

- 'ram' & "price range"(target variable) is highly correlated.
  More ram = Higher Price
- 'three_g' & 'four_g' is moderately correlated.
- 'pc' (primary camera) & 'fc' (front camera) is moderately correlated.
- 'px_height' & 'px_width' is moderately correlated.
- 'sc_h' & 'sc_w' (screen height & screen width) is moderately correlated.
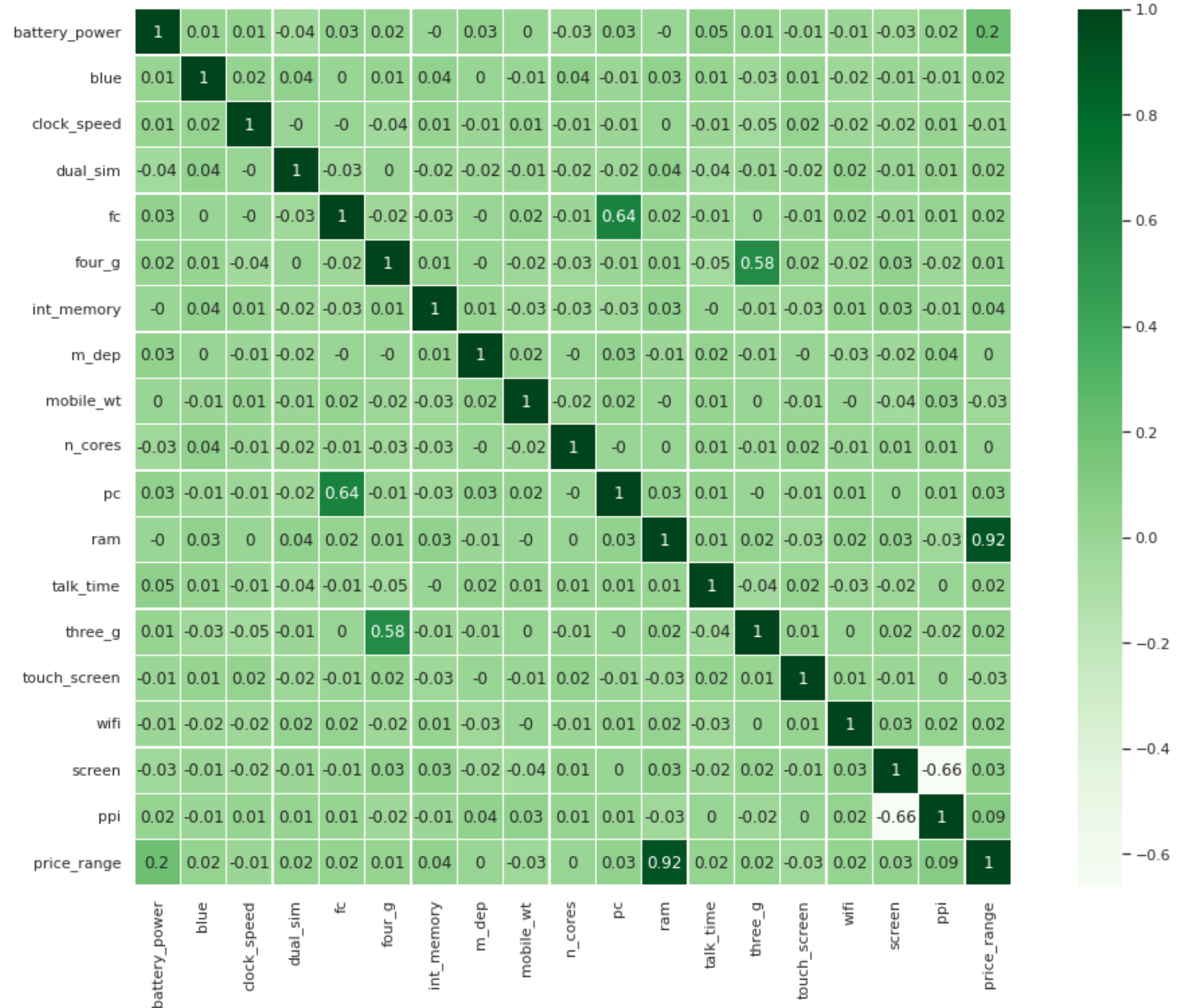


Corelation between features

# Feature Engineering

- Converting 'sc_w' and 'sc_h' to a single variable named 'screen'
- Converting 'px_width' and 'px_height' to 'ppi' (pixel per inch).

Observation
- Now we also have negative corelation between 'screen' & 'ppi'

# Feature Engineering



Dropping irrelevant columns

wifi,
blue,
m_dep
touch_screen

Feature importance obtained from coefficients

mobile_wt, clock_speed, dual_sim, pc, n_cores, wifi, blue, m_dep, touch_screen, four_g, three_g, talk_time, fc, int_memory, screen, ppi, battery_power, ram

Feature importance obtained from coefficients

mobile_wt, clock_speed, dual_sim, pc, n_cores, three_g, four_g, talk_time, fc, int_memory, screen, ppi, battery_power, ram

# Model Selection

- According to the data, we need to select suitable classification models.
- We will be comparing between 4 models
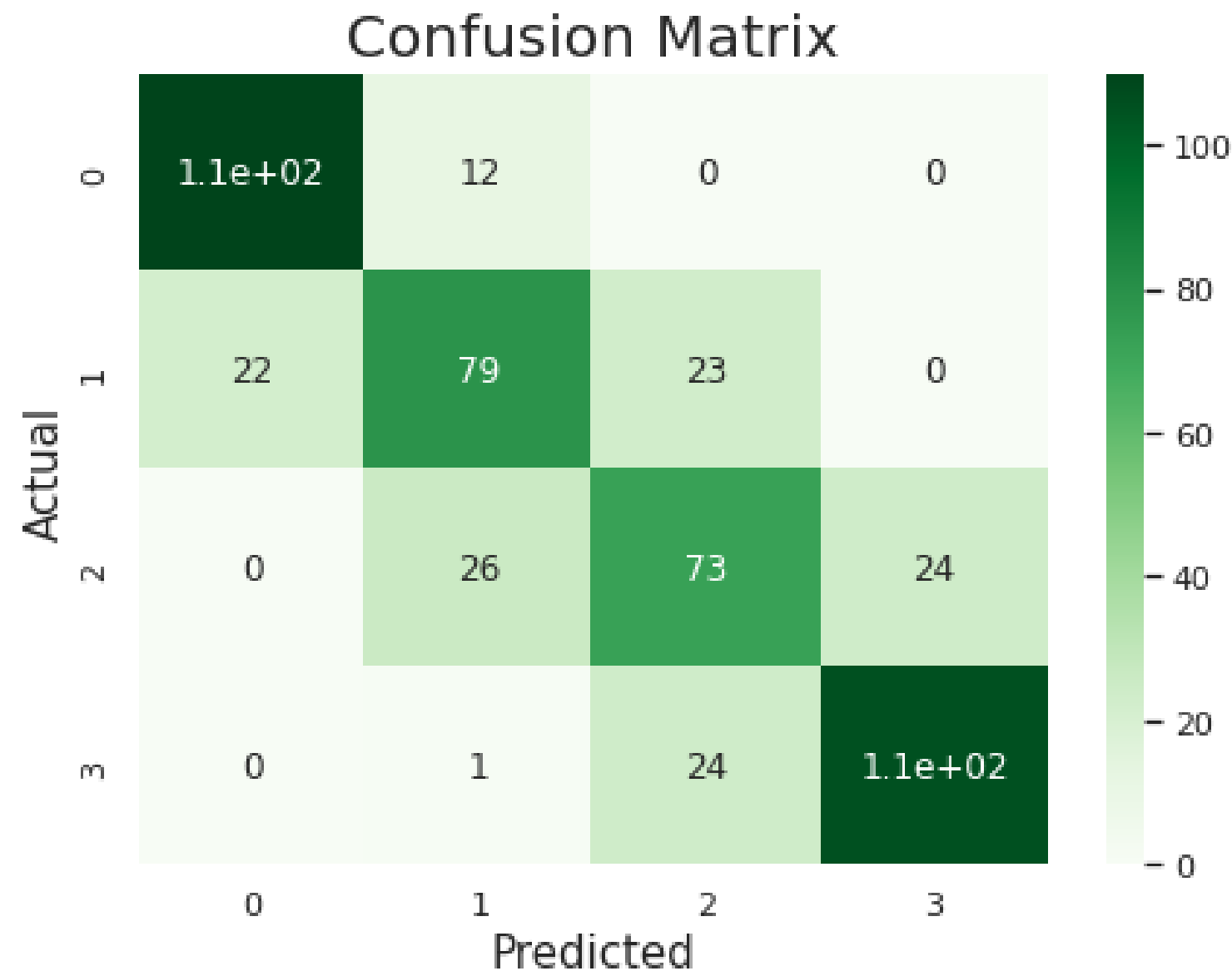
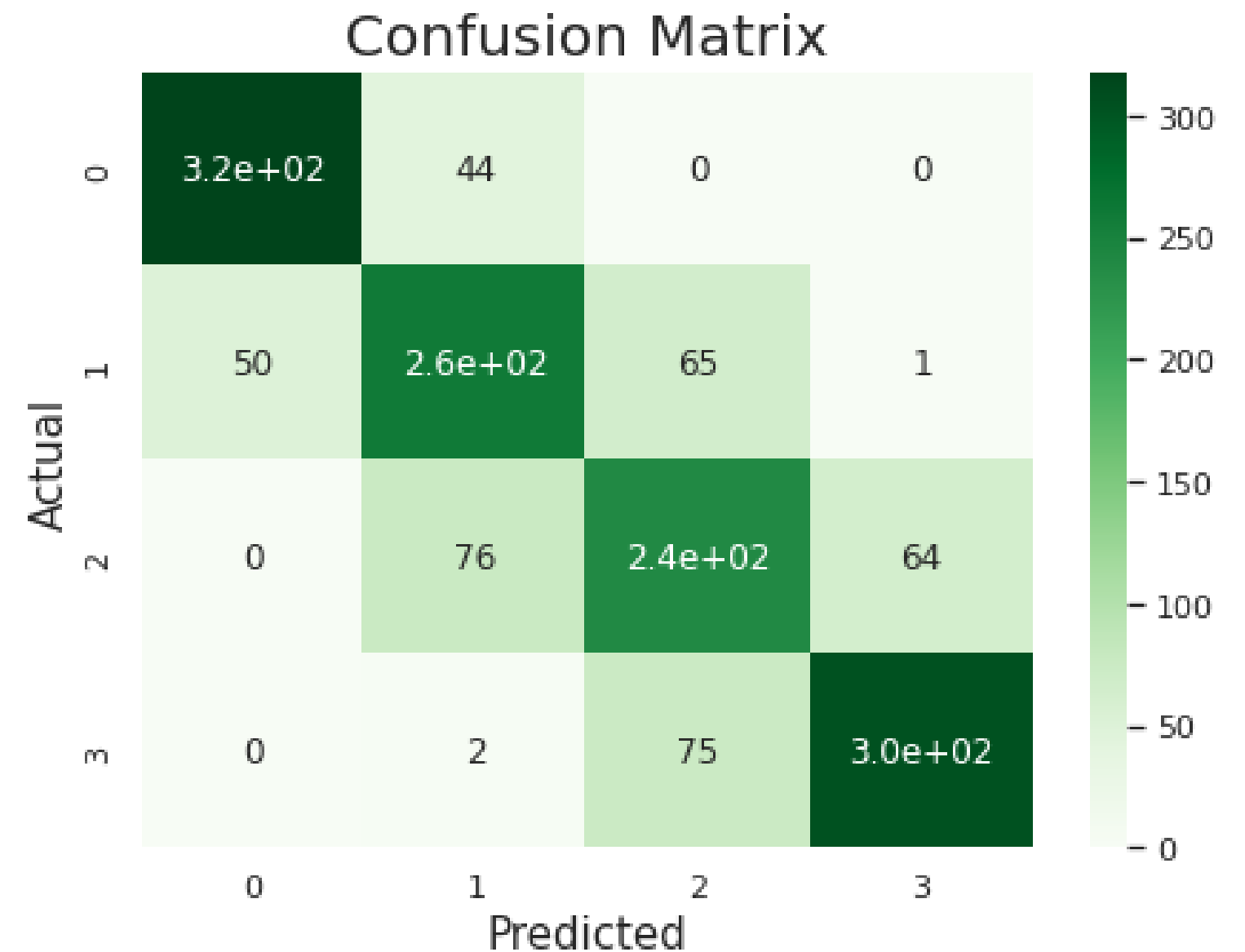Logistic Regression

Random Forest

KNN

SVM

# Logistic Regression

The accuracy on train data is : 0.7486666666666667
The accuracy on test data is : 0.736



Confusion Matrix of Test Set
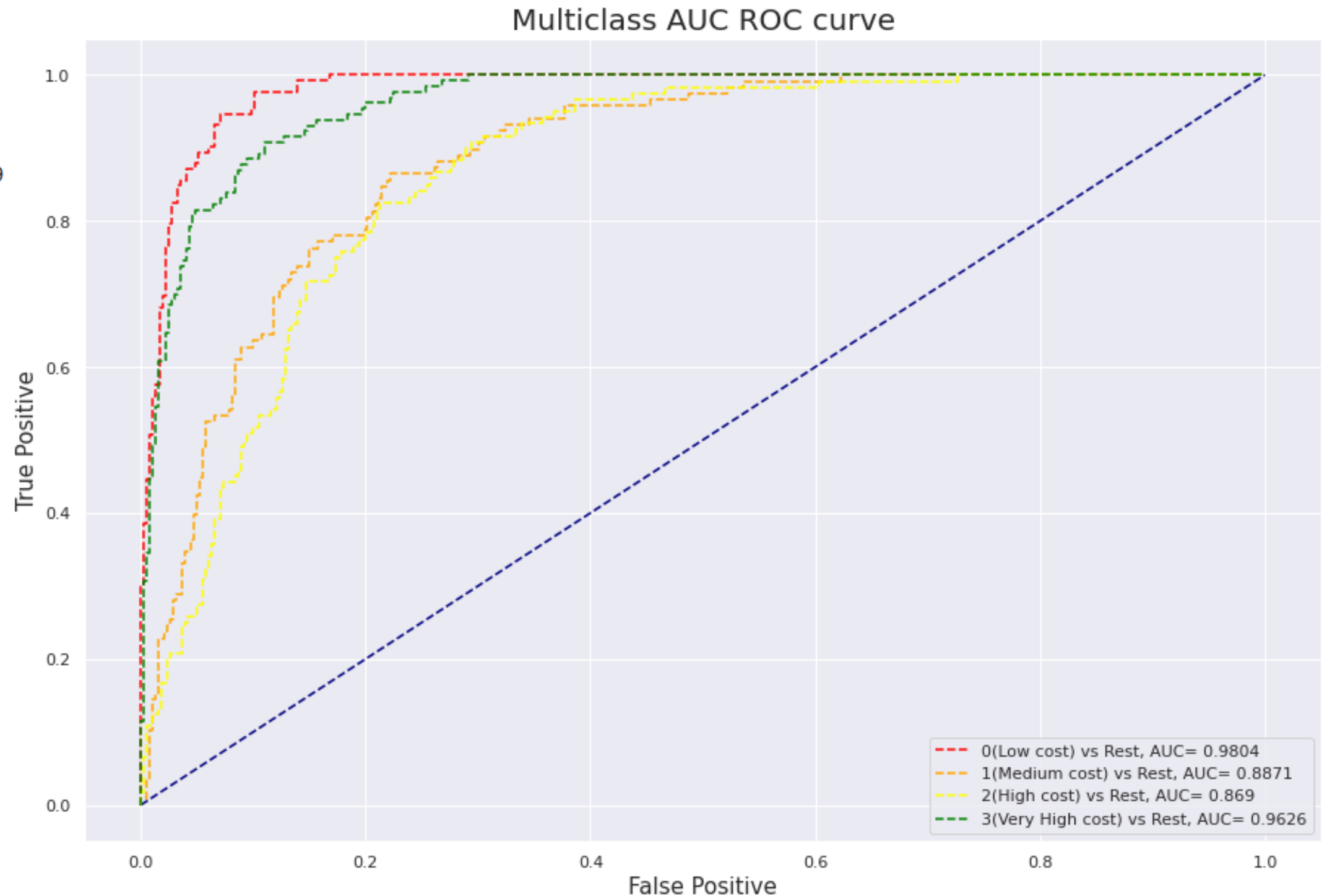


Confusion Matrix of Train Set

# Logistic Regression

The ROC AUC score on the train data is:  0.9306323173995289

The ROC AUC score on the test data is:  0.9247836188284699

Observations
- Prediction of price range of Low & Very High Cost is excellent
- Prediction of price range of Medium and High Cost is good.
- Prediction accuracy is very good but not excellent, Not to optimise. bold text



Multiclass AUC ROC curve

- 0(Low cost) vs Rest, AUC= 0.9804
- 1(Medium cost) vs Rest, AUC= 0.8871
- 2(High cost) vs Rest, AUC= 0.869
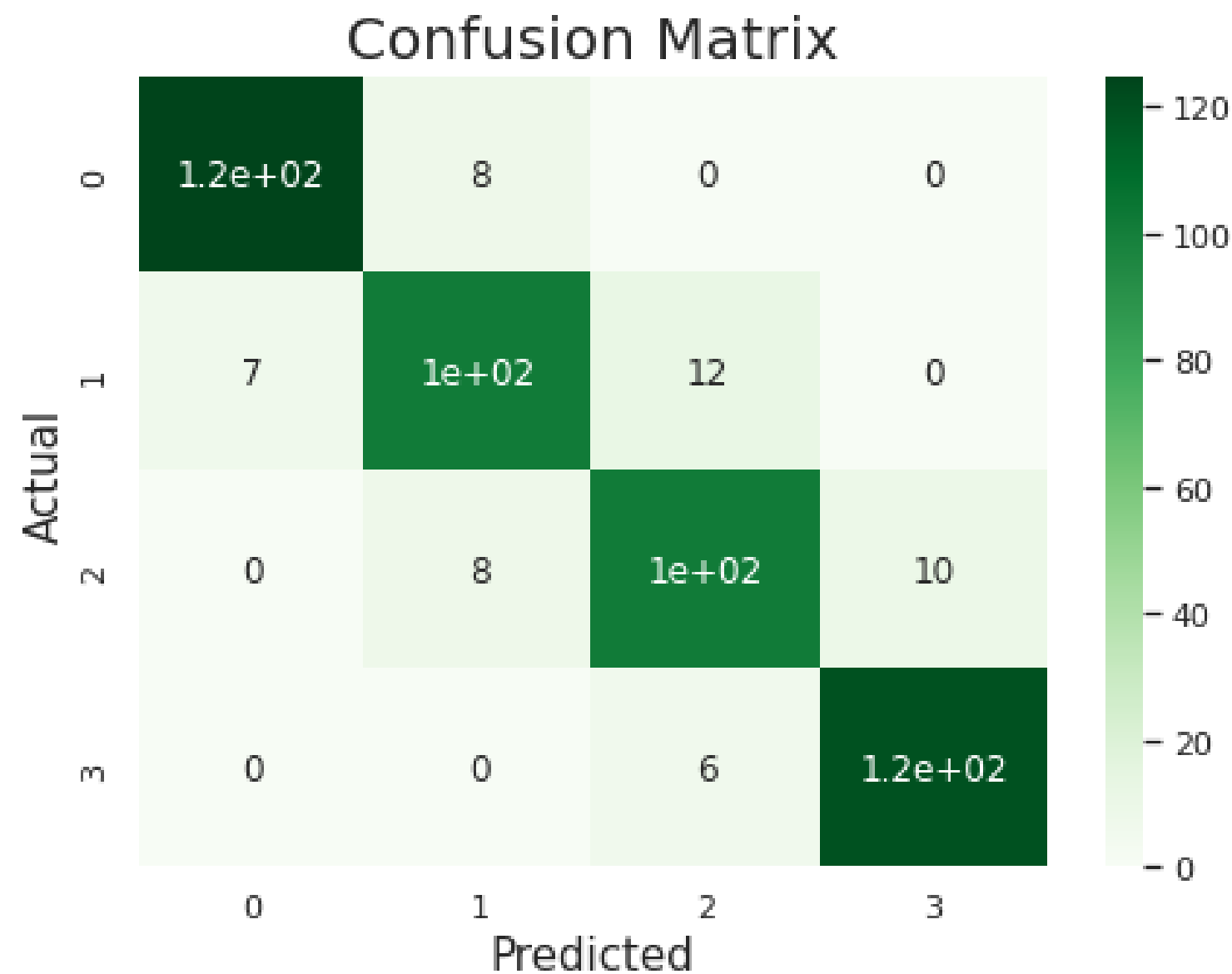- 3(Very High cost) vs Rest, AUC= 0.9626
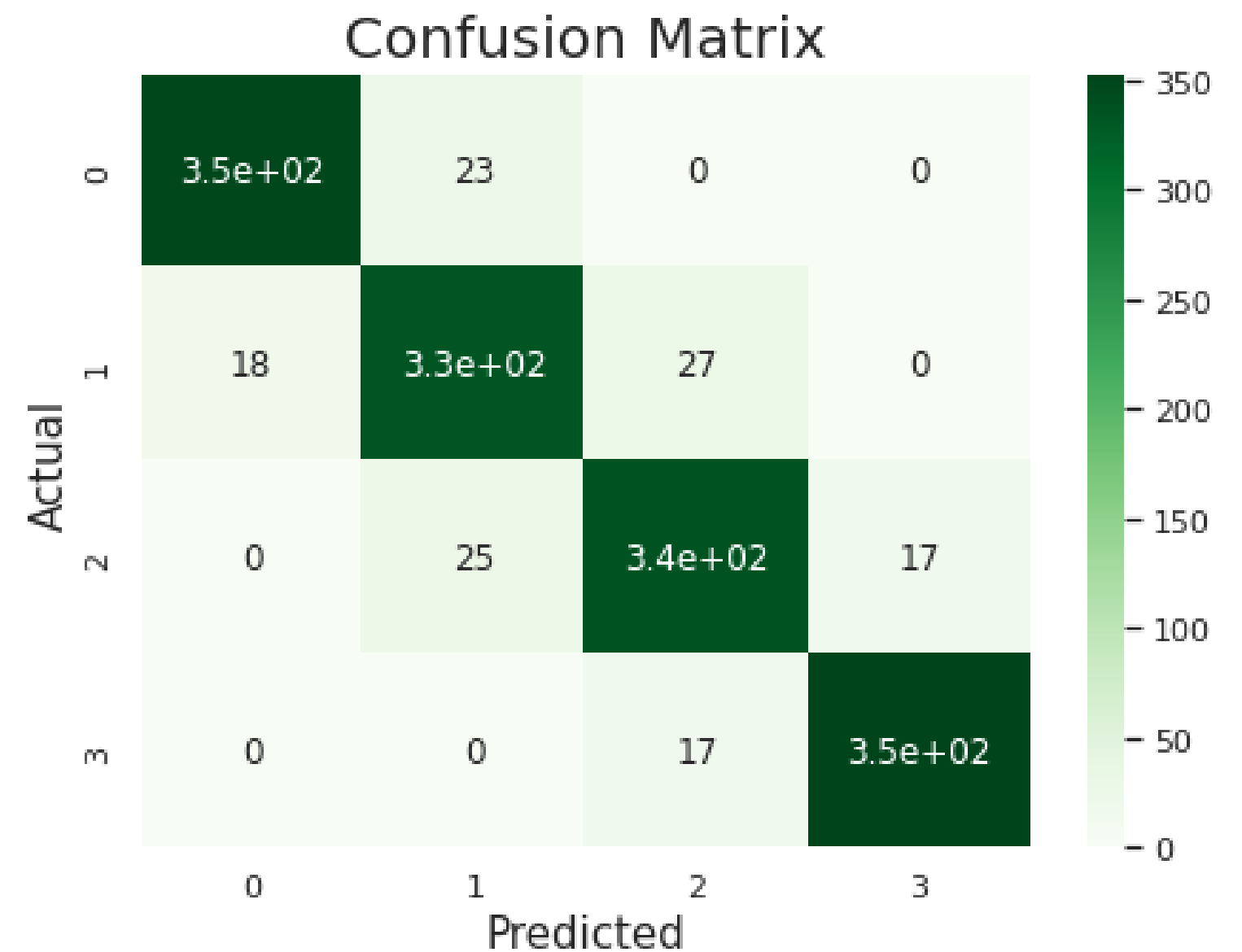
# Logistic Regression

Hyper Parameter Tuning and Cross Validation of Logistic Regression Classification

The accuracy on train data is : 0.9153333333333333
The accuracy on test data is : 0.898



Confusion Matrix of Test Set
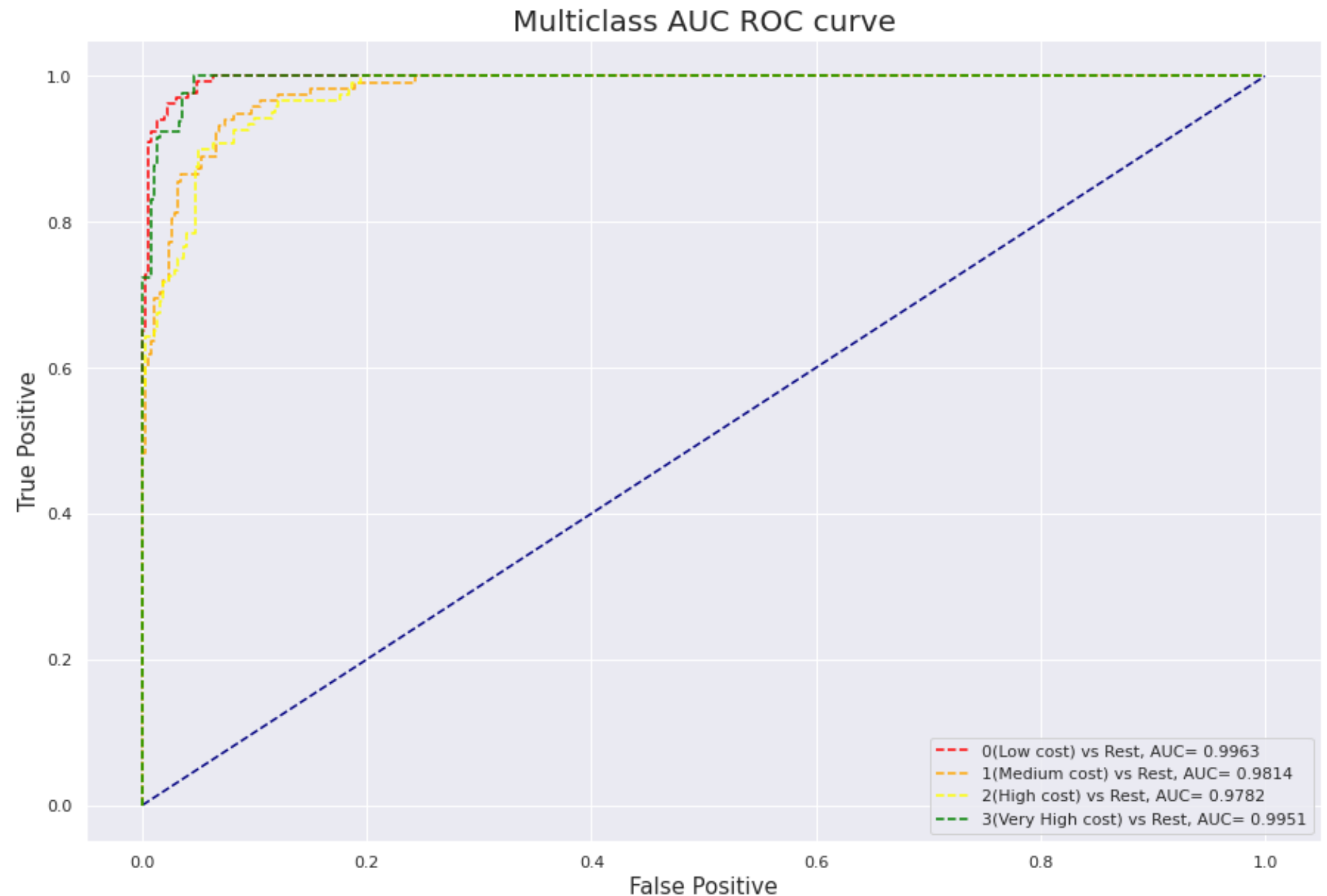
Confusion Matrix of Train Set

# Logistic Regression

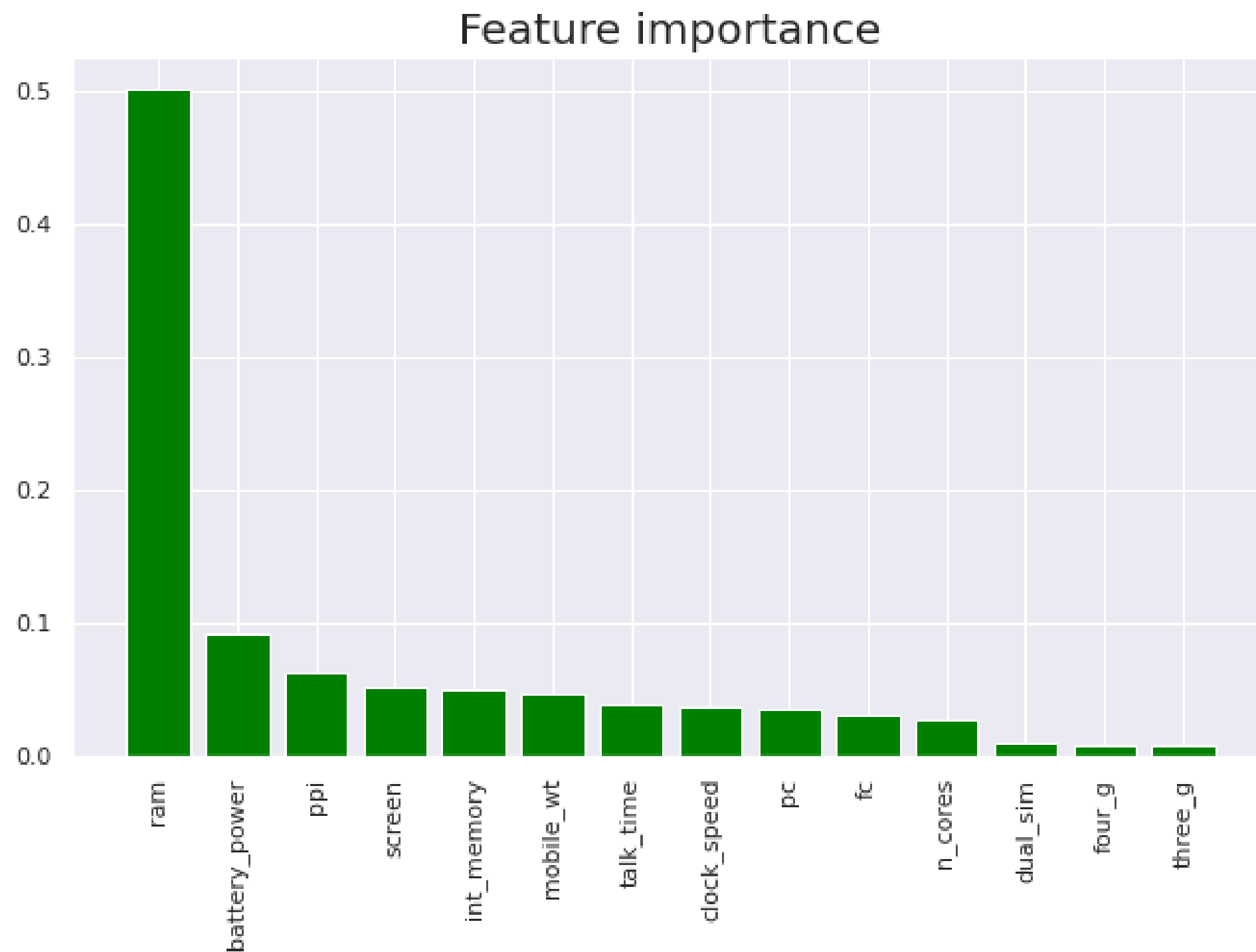The ROC AUC score on the train data is:   0.9906551817725504

The ROC AUC score on the test data is:   0.9877287010315632

Observations
- Overall accuracy has improved to excellent levels
- Prediction accuracy for price_range 1 & 2 has also increased to excellent levels
- Accuracy can be further improved if max_iterations are increased, but we have capped it to 100 to keep run time under control for GOOGLE COLAB(Free Version)



Multiclass AUC ROC curve

- - - 0(Low cost) vs Rest, AUC= 0.9963
- - - 1(Medium cost) vs Rest, AUC= 0.9814
- - - 2(High cost) vs Rest, AUC= 0.9782
- - - 3(Very High cost) vs Rest, AUC= 0.9951

# Random Forest Classification

### Feature importance



```
The accuracy on train data is : 1.0
The accuracy on test data is : 0.832


The confusion matrix on the train data is :
[[368   0   0   0]
 [  0 382   0   0]
 [  0   0 380   0]
 [  0   0   0 370]]


The confusion matrix on the test data is :
[[127  13   0   0]
 [  5  95  18   0]
 [  0  10  84  20]
 [  0   0  18 110]]
```

# Random Forest Classification

```
The classification report on the train data is :
              precision    recall  f1-score   support

         0.0       1.00      1.00      1.00       368
         1.0       1.00      1.00      1.00       382
         2.0       1.00      1.00      1.00       380
         3.0       1.00      1.00      1.00       370

    accuracy                           1.00      1500
   macro avg       1.00      1.00      1.00      1500
weighted avg       1.00      1.00      1.00      1500
```

```
The classification report on the test data is :
              precision    recall  f1-score   support

         0.0       0.96      0.91      0.93       140
         1.0       0.81      0.81      0.81       118
         2.0       0.70      0.74      0.72       114
         3.0       0.85      0.86      0.85       128

    accuracy                           0.83       500
   macro avg       0.83      0.83      0.83       500
weighted avg       0.84      0.83      0.83       500
```

The ROC AUC score on the train data is:  1.0
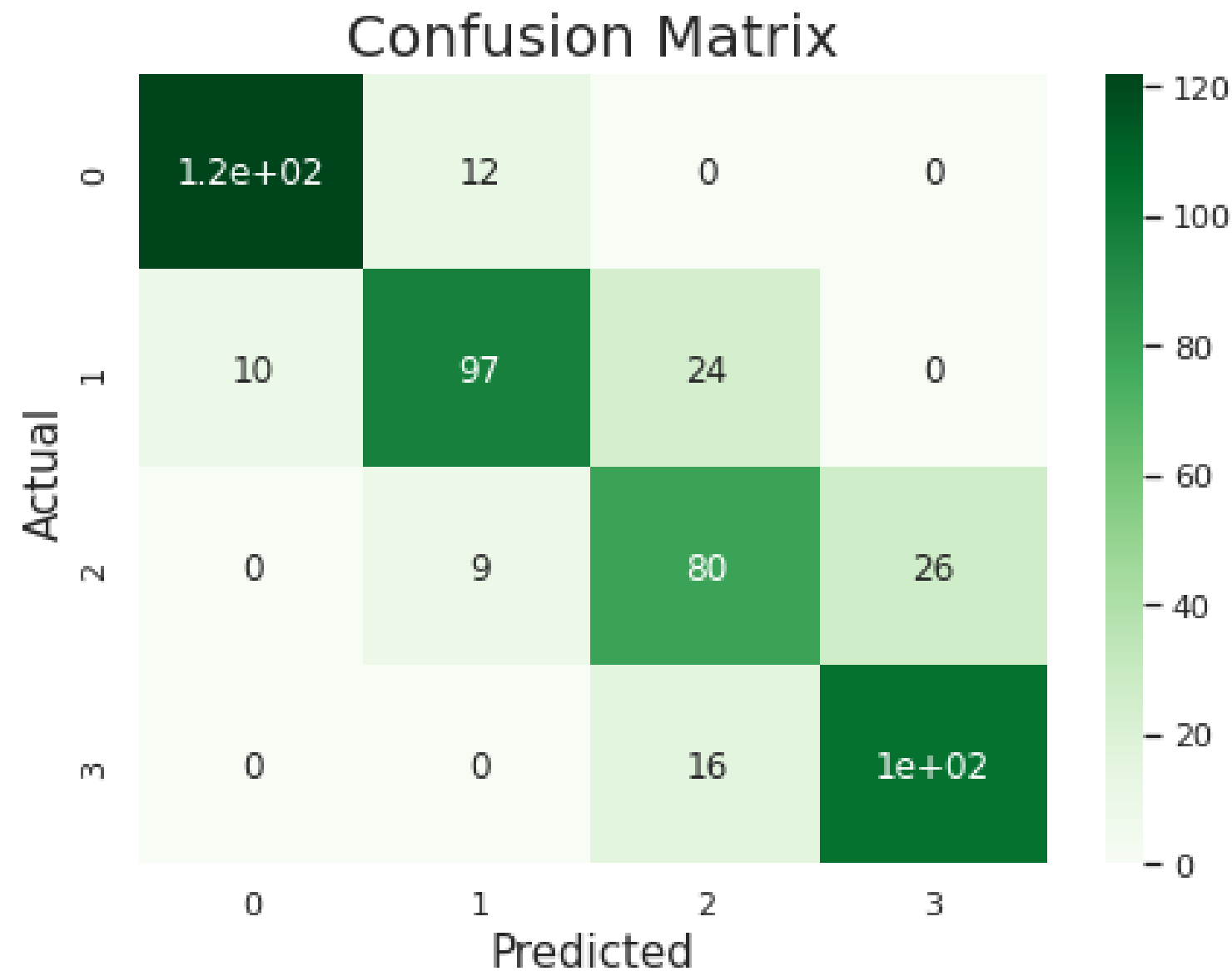The ROC AUC score on the test data is:  0.9658171969858055

Observations
- Training accuracy is 1, Random Forest is overfitting
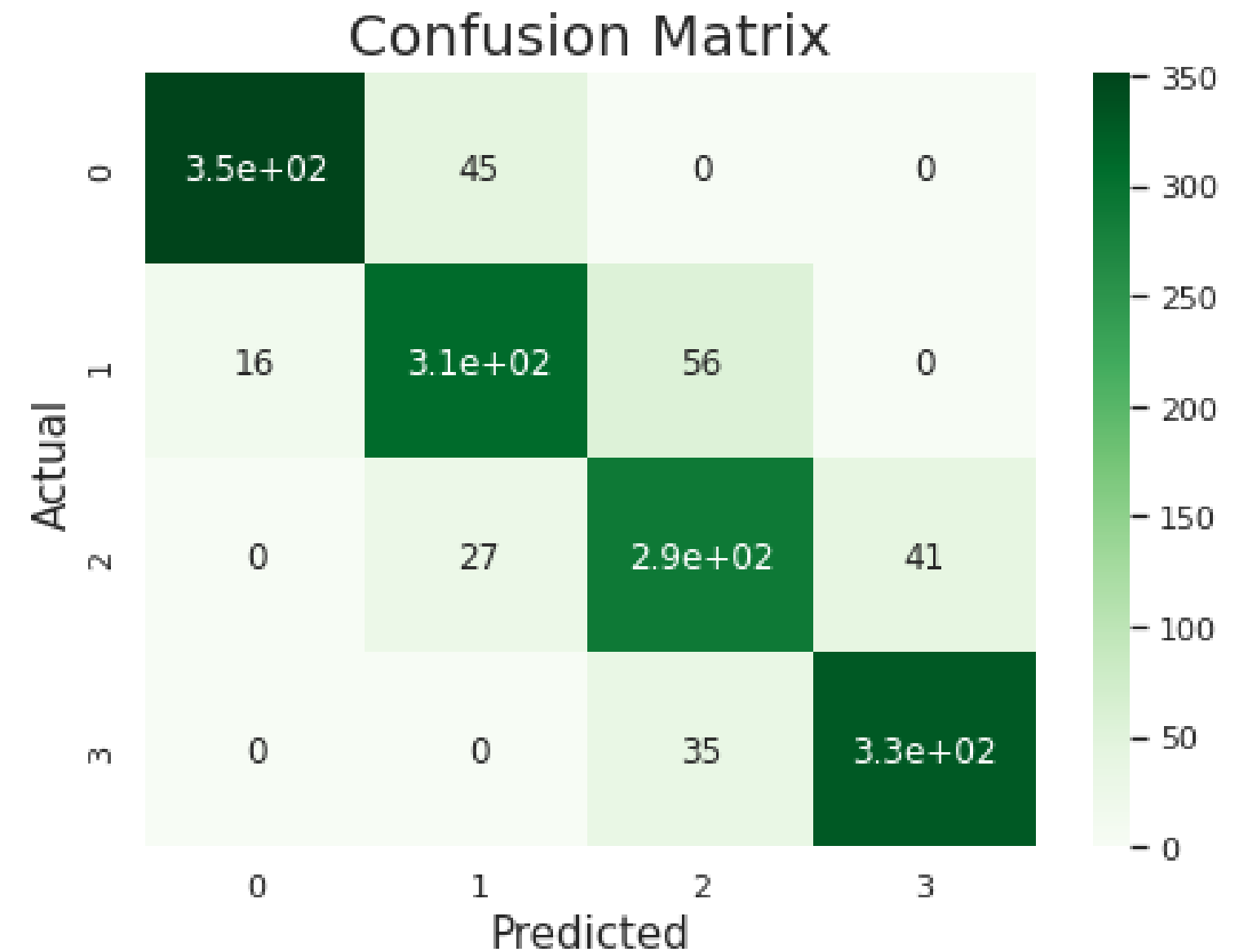- Dropping random Forest and moving to next model

# KNN Classification

The accuracy on train data is : 0.8533333333333334
The accuracy on test data is : 0.806



Confusion Matrix of Test Set
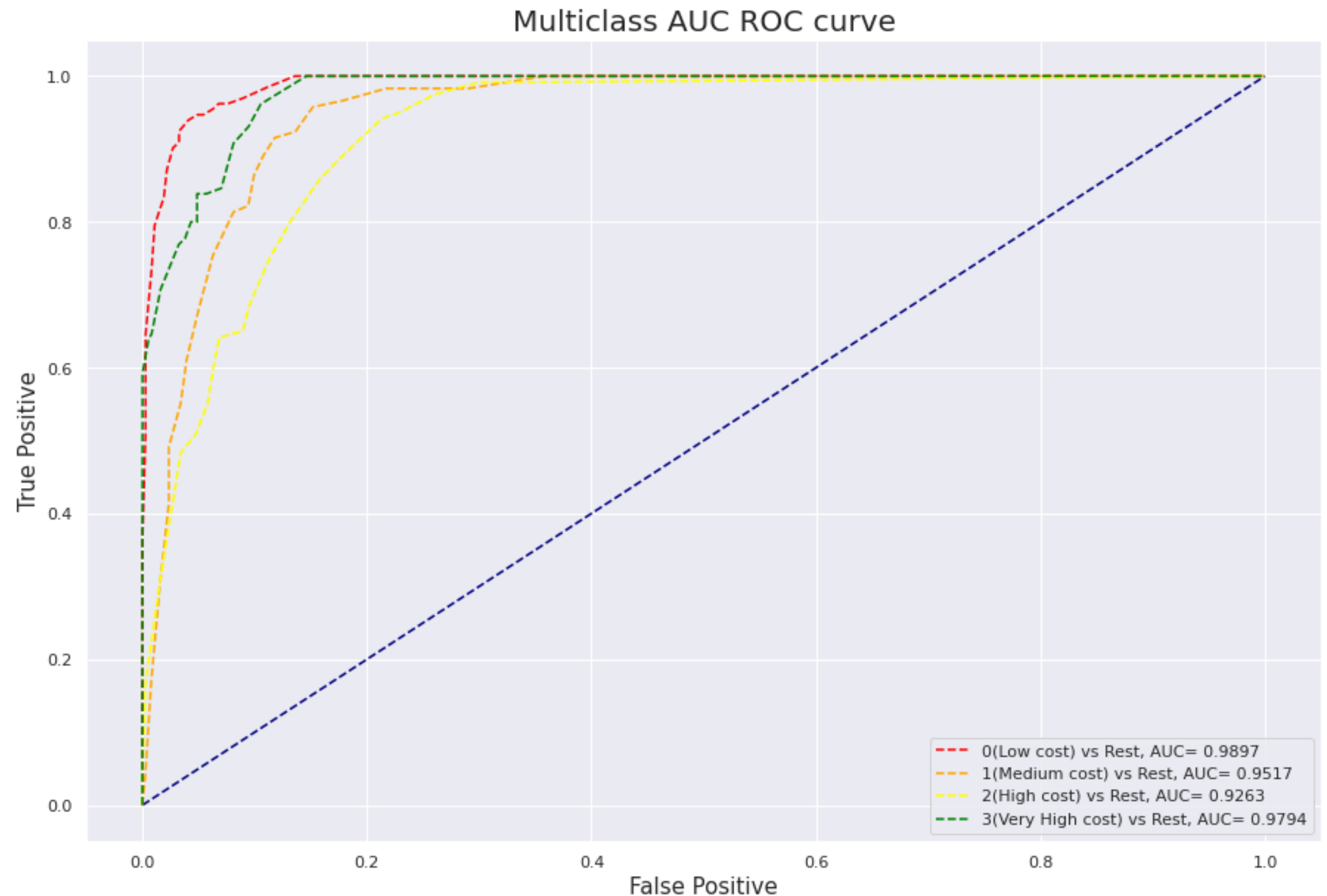
Confusion Matrix of Train Set

# KNN Classification

The ROC AUC score on the train data is:  0.9746939230833361

The ROC AUC score on the test data is:  0.9617833679300767

Observations
- Prediction accuracy is less than optimised Logistic classification
- AUC-ROC score is better than logistc classification on average for all 4 price_range
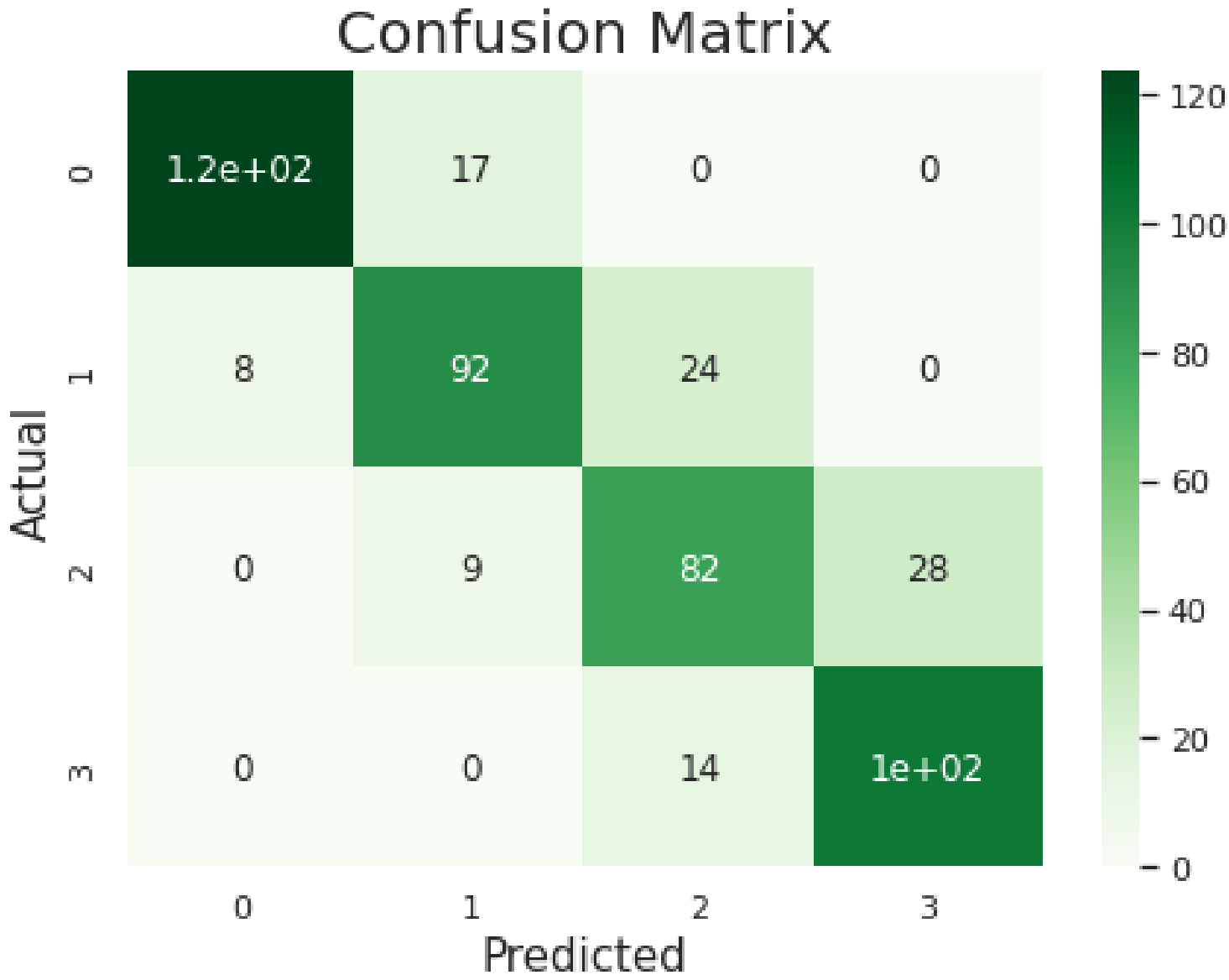- Need to evaluate further with parameter tuning



Multiclass AUC ROC curve
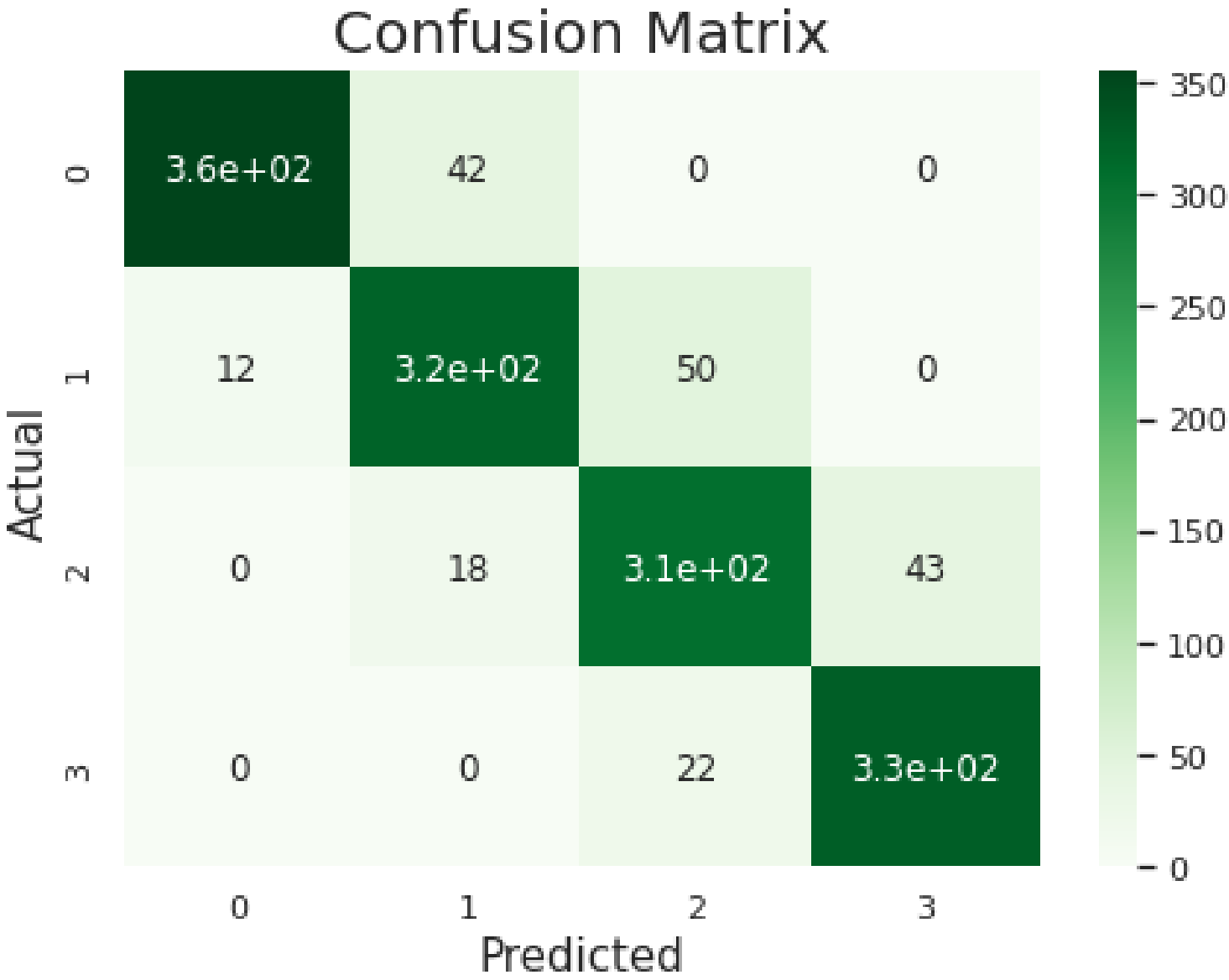
# KNN Classification

Hyper Parameter Tuning and Cross Validation of KNN Classification

```
The accuracy on train data is : 0.8753333333333333
The accuracy on test data is : 0.8
```



Confusion Matrix of Test Set
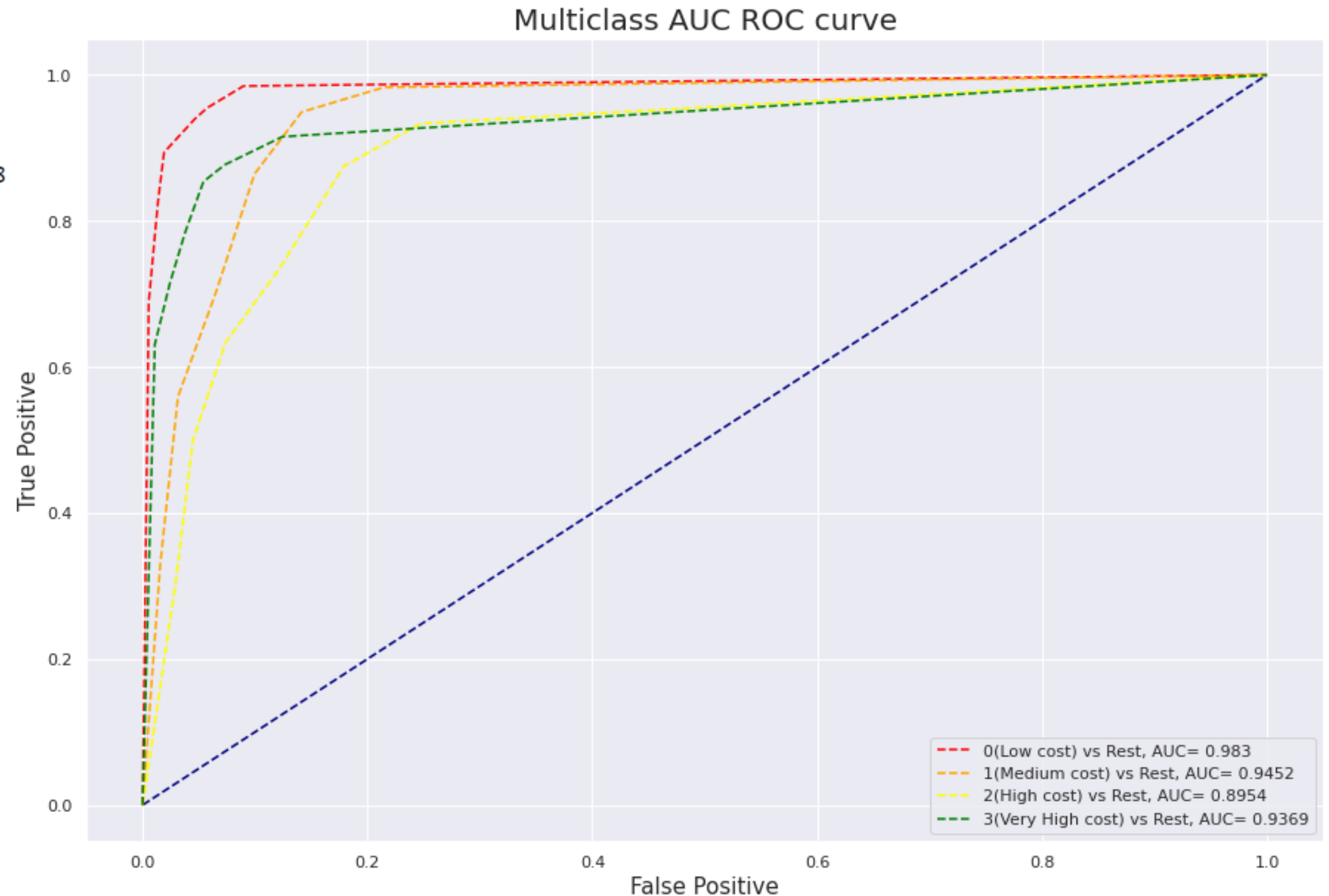


Confusion Matrix of Train Set

# KNN Classification

The ROC AUC score on the train data is: 0.9831038848505038

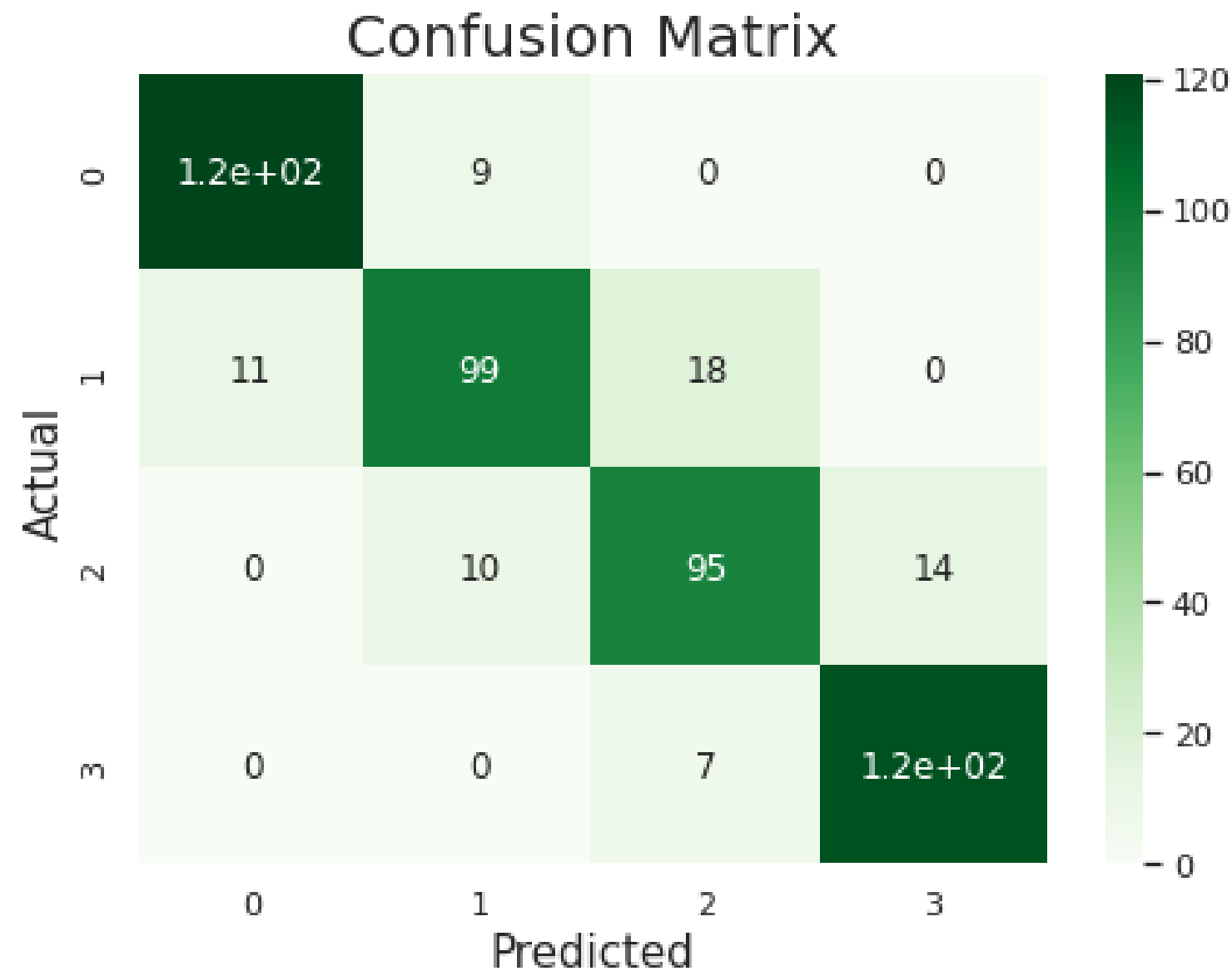The ROC AUC score on the test data is: 0.9401348940851196

Observations
- After optimisation KNN improved very well.
- But for multiclass 'price_range' prediction of price_range = 2 is slightly less in comparisson with optimised logistic regression.



Multiclass AUC ROC curve

True Positive / False Positive

- - - 0(Low cost) vs Rest, AUC= 0.983
- - - 1(Medium cost) vs Rest, AUC= 0.9452
- - - 2(High cost) vs Rest, AUC= 0.8954
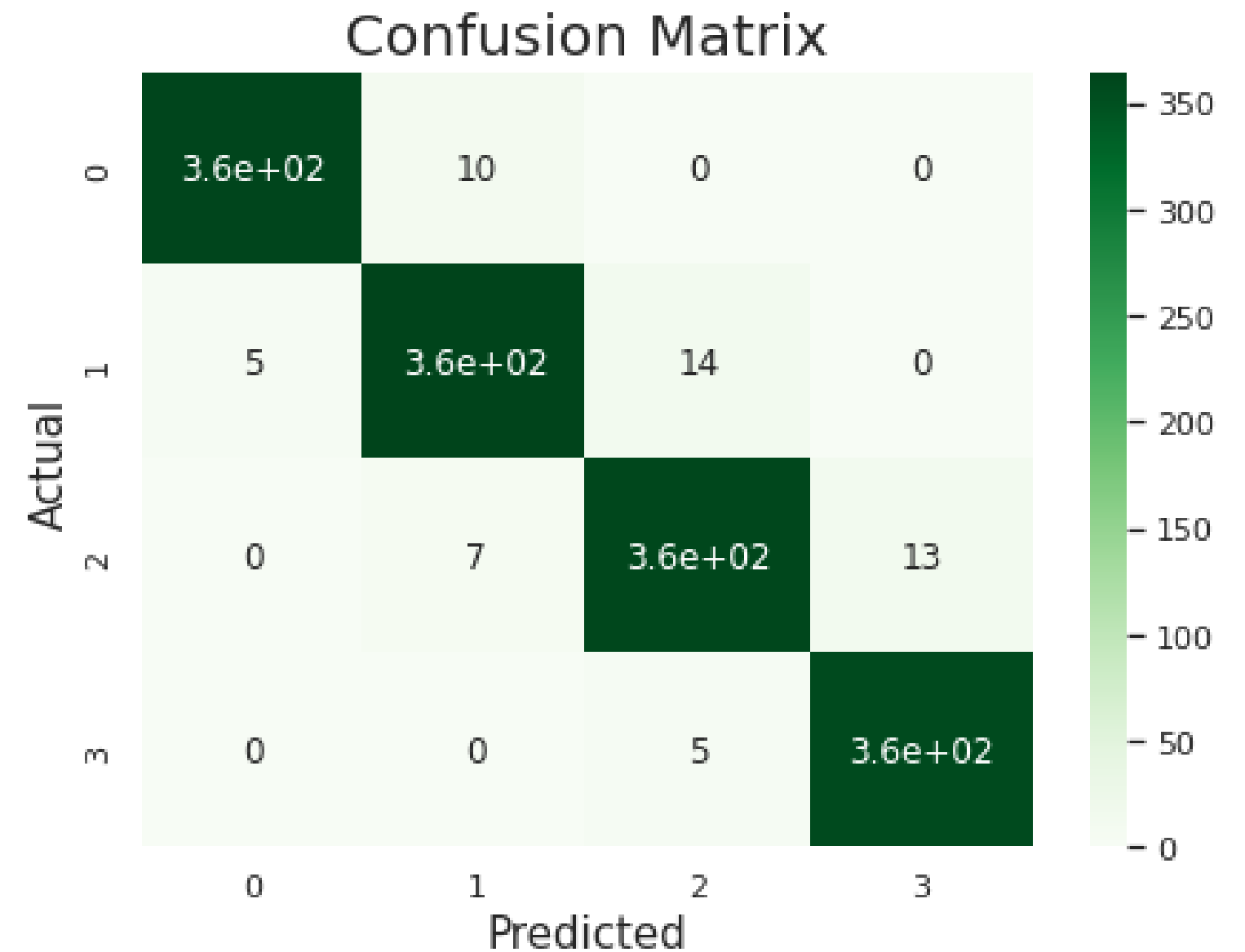- - - 3(Very High cost) vs Rest, AUC= 0.9369

# SVM Classification

The accuracy on train data is : 0.964
The accuracy on test data is : 0.862



Confusion Matrix of Test Set
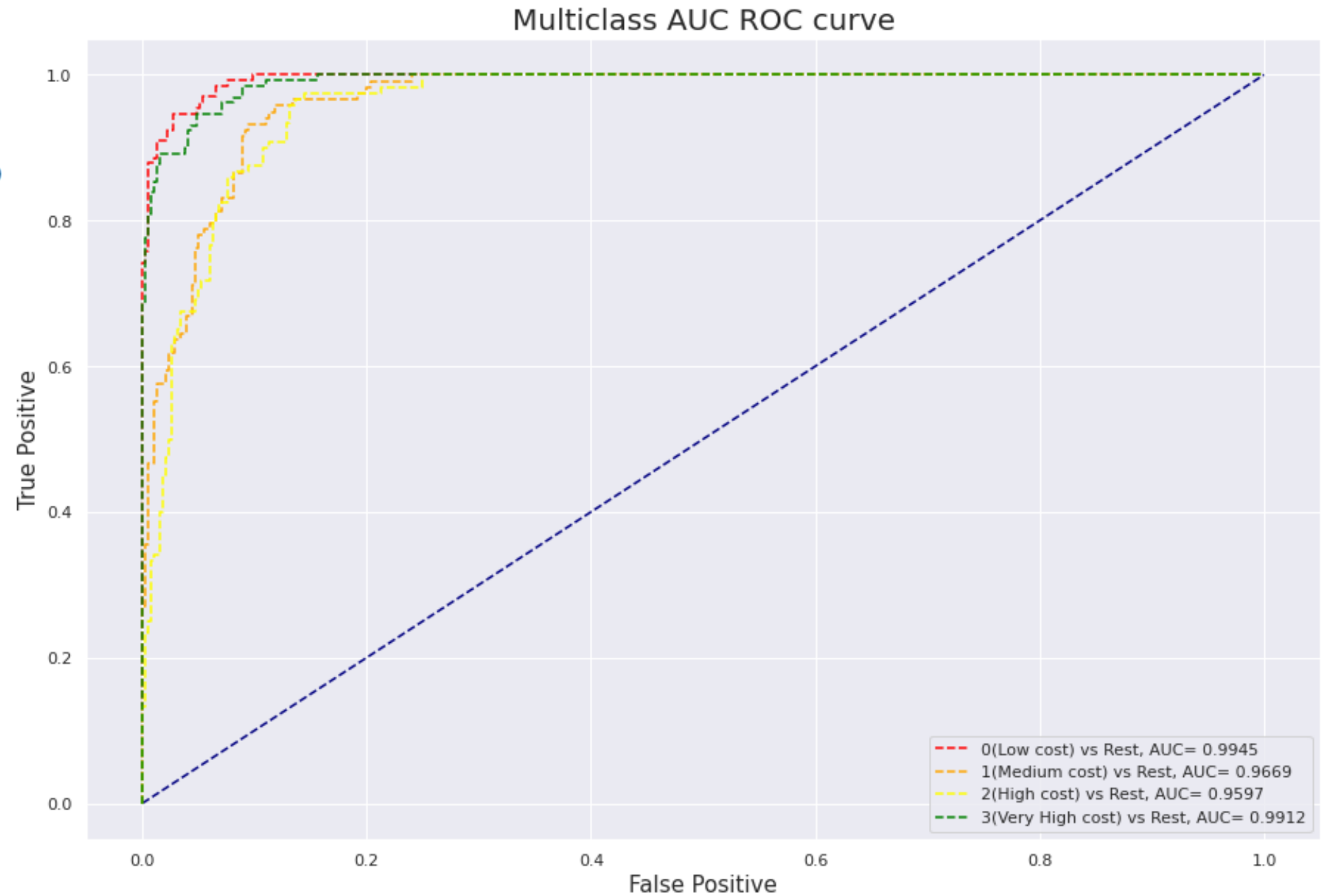
Confusion Matrix of Train Set

# SVM Classification

The ROC AUC score on the train data is:  0.9980200219416429

The ROC AUC score on the test data is:  0.9780710753879851

Observations
- Model is overfitting but can be fixed by optmisation.
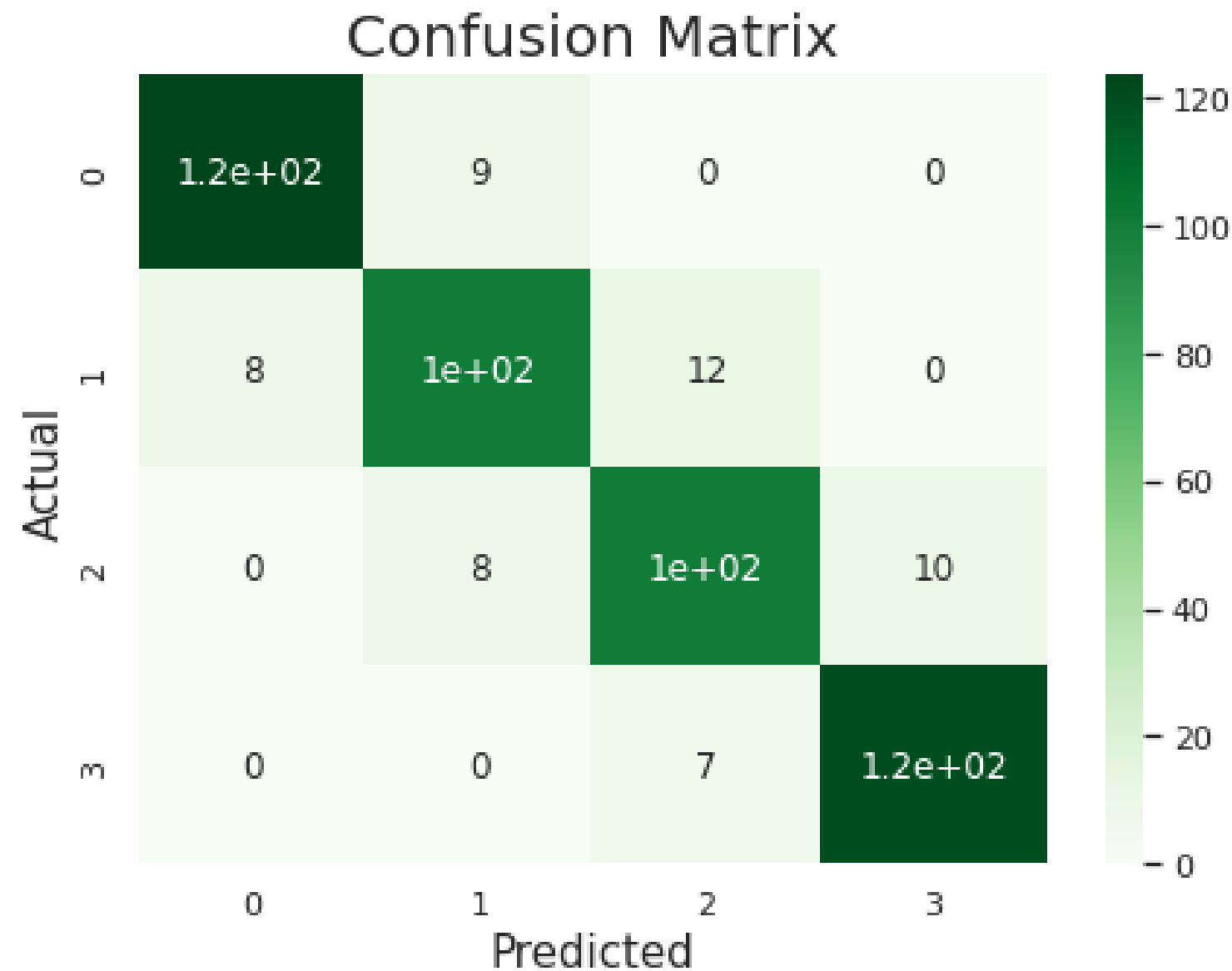


Multiclass AUC ROC curve

True Positive / False Positive

- - - 0(Low cost) vs Rest, AUC= 0.9945
- - - 1(Medium cost) vs Rest, AUC= 0.9669
- - - 2(High cost) vs Rest, AUC= 0.9597
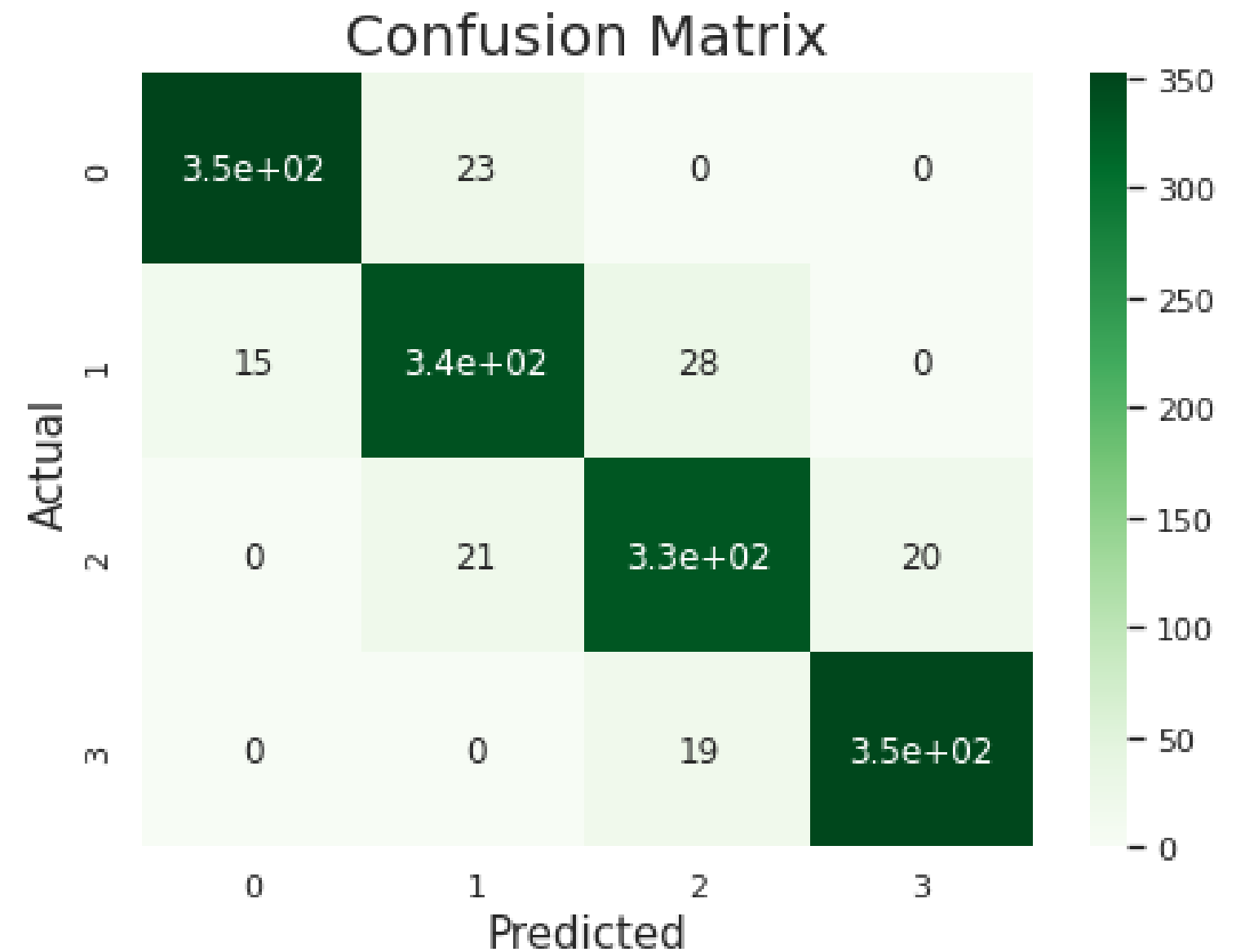- - - 3(Very High cost) vs Rest, AUC= 0.9912

# SVM Classification

Hyper Parameter Tuning and Cross Validation of SVM Classification

```
The accuracy on train data is : 0.916
The accuracy on test data is : 0.892
```



Confusion Matrix of Test Set
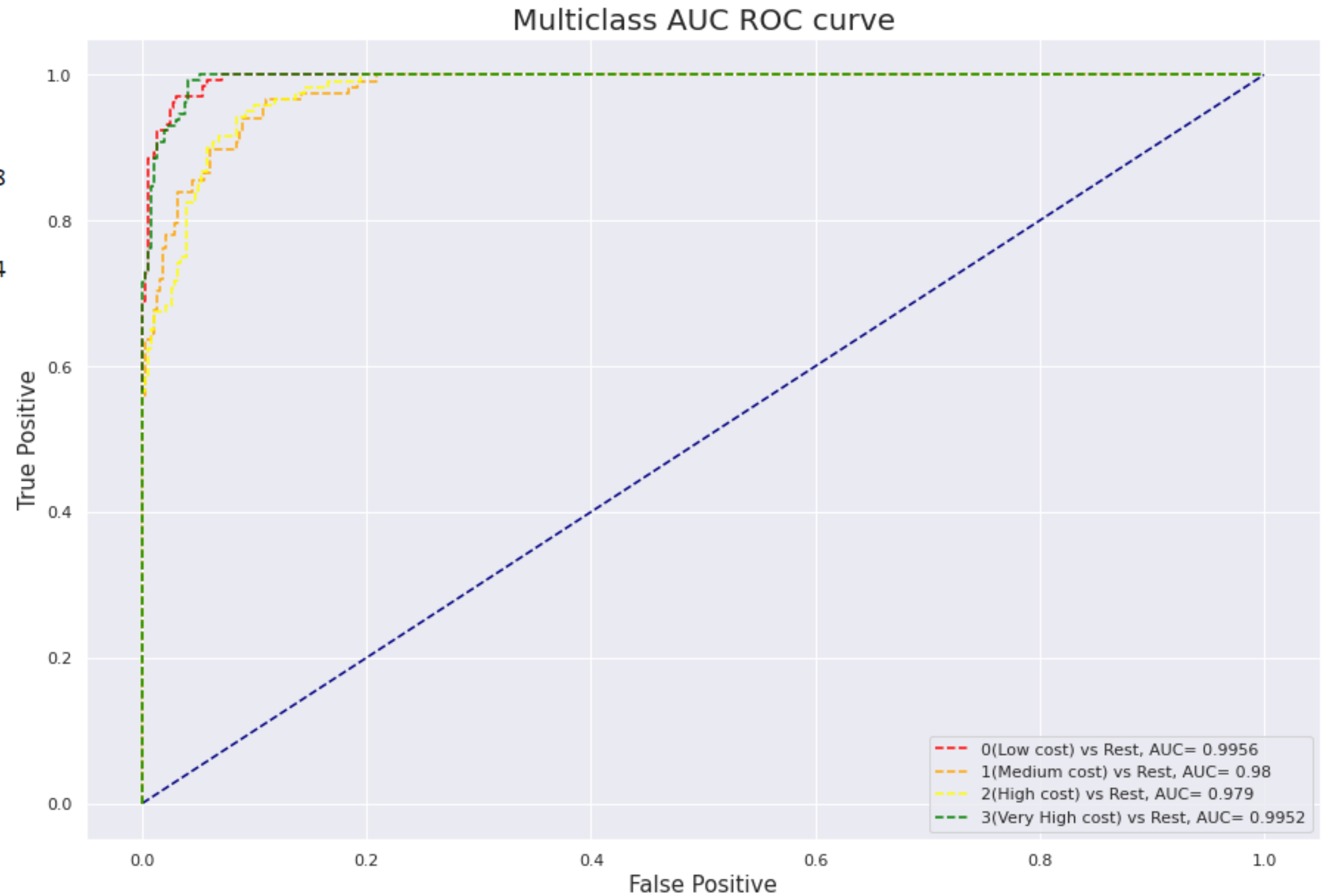
Confusion Matrix of Train Set

# SVM Classification

The ROC AUC score on the train data is:  0.990249880703548

The ROC AUC score on the test data is:  0.9874538341775304

Observations
- Overfitting is reduced
- Prediction for all 4 price_range is good



Multiclass AUC ROC curve

Legend:
- 0(Low cost) vs Rest, AUC= 0.9956
- 1(Medium cost) vs Rest, AUC= 0.98
- 2(High cost) vs Rest, AUC= 0.979
- 3(Very High cost) vs Rest, AUC= 0.9952

# Conclusion

- In EDA there were columns/features that were inter-related, we converted to new features using them.
- ram and batter_power has the highest impact on price_range.
- Using logistic regeression feature importance we observed that some of the columns were not relevant or had no impact negative/positive. Hence, they were dropped.
- Implemented various classification algorithms, Logistics and SVM accuracy was similar.
- Logistic regression classification model gave best results after hyper-parameter tuning with 91.5% train accuracy and 89.2% test accuracy score.
- SVM (Support vector machine) algorithm also gave equally best accuracy after hyper-parameter tuning with 91.6% train accuracy and 89.2 % test accuracy.
- Random Forest was Over-fitting
- KNN after optimization performed very well but for mutliclass price_range the prediction of price_range = 2 was lowered than Logistic and SVM in comparission.

*We will go forward with Logistic regression classification model as using it increases the explainability of price_rage as per business requirement.*

# Thank You!