

# Capstone Project - 3

## MOBILE PRICE RANGE PREDICTION

### SUPERVISED MACHINE LEARNING (CLASSIFICATION)



BY  
RUSHIKESH ARJUN MANE  
(Cohort Seattle)

# Contain

- 1. Problem Statement**
- 2. Work Flow**
- 3. Data Collection and Understanding**
- 4. Data Wrangling and Feature Engineering**
- 5. EDA (Exploratory Data Analysis)**
- 6. Preparation of Data for Model Building**
- 7. Model Selection and Evaluation**
- 8. Conclusion**

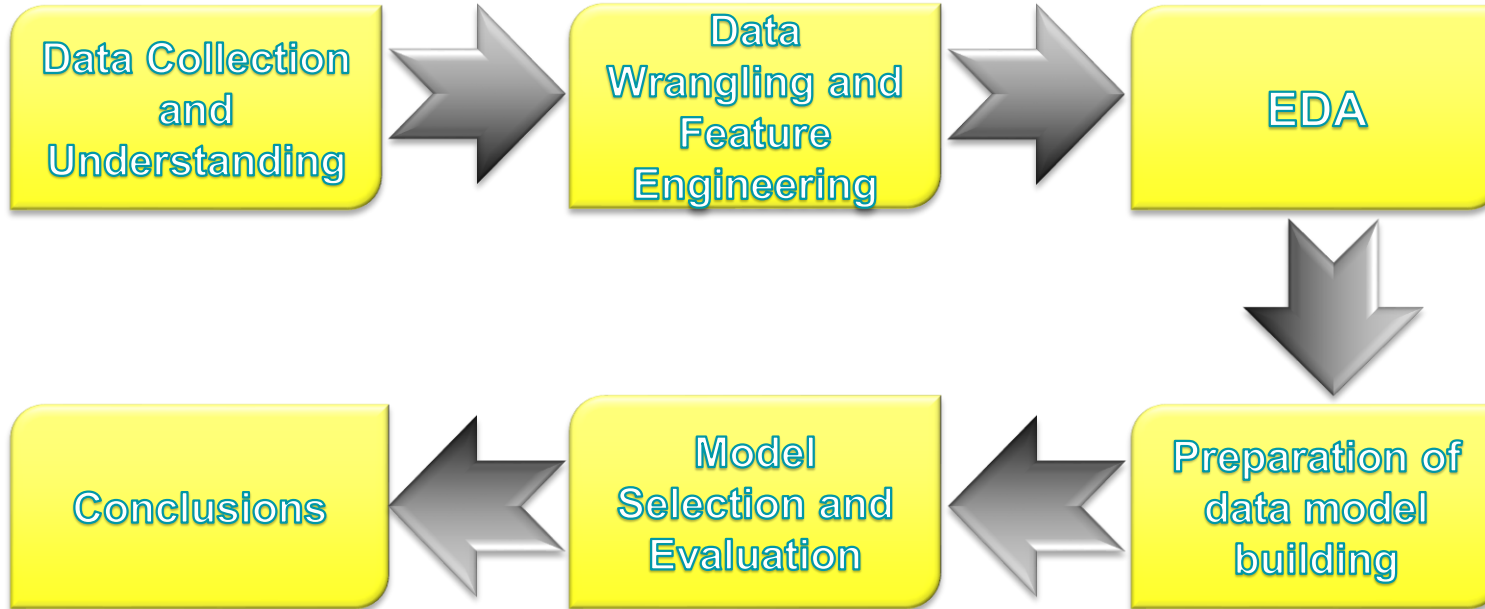
# **Introduction**

- ▶ In this Modern Era, Smartphones are an integral part of the lives of human beings.
- ▶ When a smartphone is purchased ,many factors like the Display, Processor, Memory, Camera, Thickness, Battery, Connectivity and others are taken into account .
- ▶ Price estimation and prediction is an important part of consumer strategy. Deciding on the correct price of a product is very important for the market success of a product. A new product that has to be launched must have the correct price so that consumers find it appropriate to buy the product.

# **Problem Statement**

- In the competitive mobile phone market companies want 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.
- In this problem, we do not have to predict the actual price but a price range indicating how high the price is.

# Work Flow



# Data Collection and Understanding

- ❑ For analysis and modeling we are having Mobile Price Range Data.
- ❑ The dataset contains 2000 rows and 21 columns.
  
- ❑ **DATA DESCRIPTION:**
  - **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

- **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 when you are
- **Three\_g** - Has 3G or not

- ***Touch\_screen*** - Has touch screen or not
- ***Wifi*** - Has wifi or not
- ***Price\_range*** - This is the target variable with value of
  - 0(low cost),
  - 1(medium cost),
  - 2(high cost) and
  - 3(very high cost).
- Thus our target variable has 4 categories so basically it is a Multiclass classification problem.



# Data Wrangling and Feature Engineering

## ❖ Dealing with mismatch values,

	count	mean	std	min	25%	50%	75%	max
px_height	2000.000000	645.108000	443.780811	0.000000	282.750000	564.000000	947.250000	1960.000000
sc_w	2000.000000	5.767000	4.356398	0.000000	2.000000	5.000000	9.000000	18.000000

```
# Checking How many observations having screen width value as 0.
print(df_main[df_main['sc_w']==0].shape[0])
```

180

```
# Checking How many observations having px_hieght value as 0.
print(df_main[df_main['px_height']==0].shape[0])
```

2

```
# As there are only 2 observations having px_height=0. so we will drop it.
df_main=df_main[df_main['px_height']!=0]
```

```
# Checking How many observations having sc_w value as 0.
df_main[df_main['sc_w']==0].shape[0]
```

0

- ❖ Here we used K-Nearest Neighbours approach to impute the missing values where a Euclidean distance is used to find the nearest neighbours.

```
df_main.isnull().sum()
```

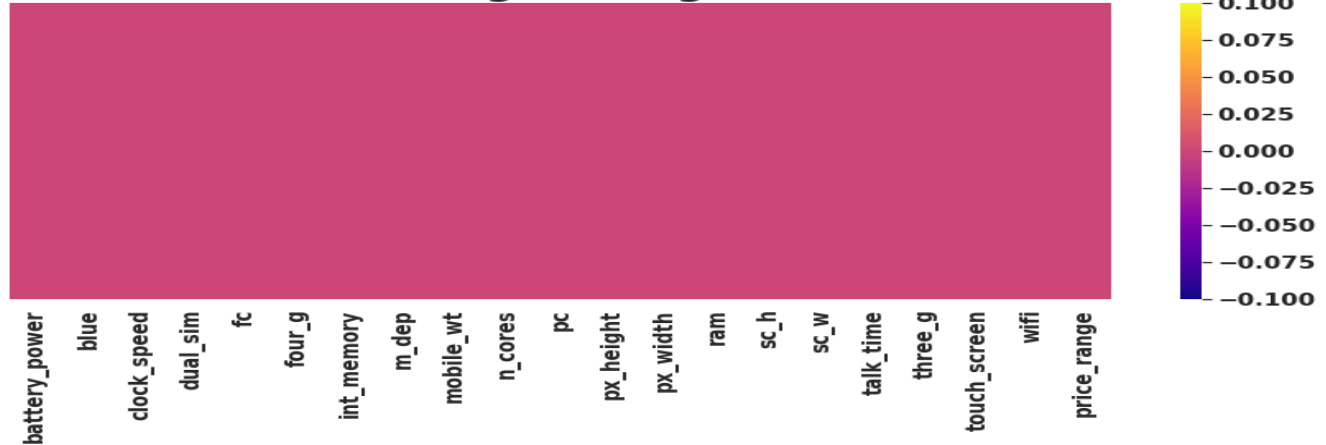
```
battery_power    0
blue             0
clock_speed      0
dual_sim         0
fc              0
four_g           0
int_memory       0
m_dep           0
mobile_wt        0
n_cores          0
pc              0
px_height        0
px_width         0
ram             0
sc_h            0
sc_w            0
talk_time        0
three_g         0
touch_screen     0
wifi            0
price_range      0
dtype: int64
```

```
## Checking Duplicate rows in our Dataset.
duplicates=df_main.duplicated().sum()
print(f"We are haveing {duplicates} rows in our Dataframe.")
```

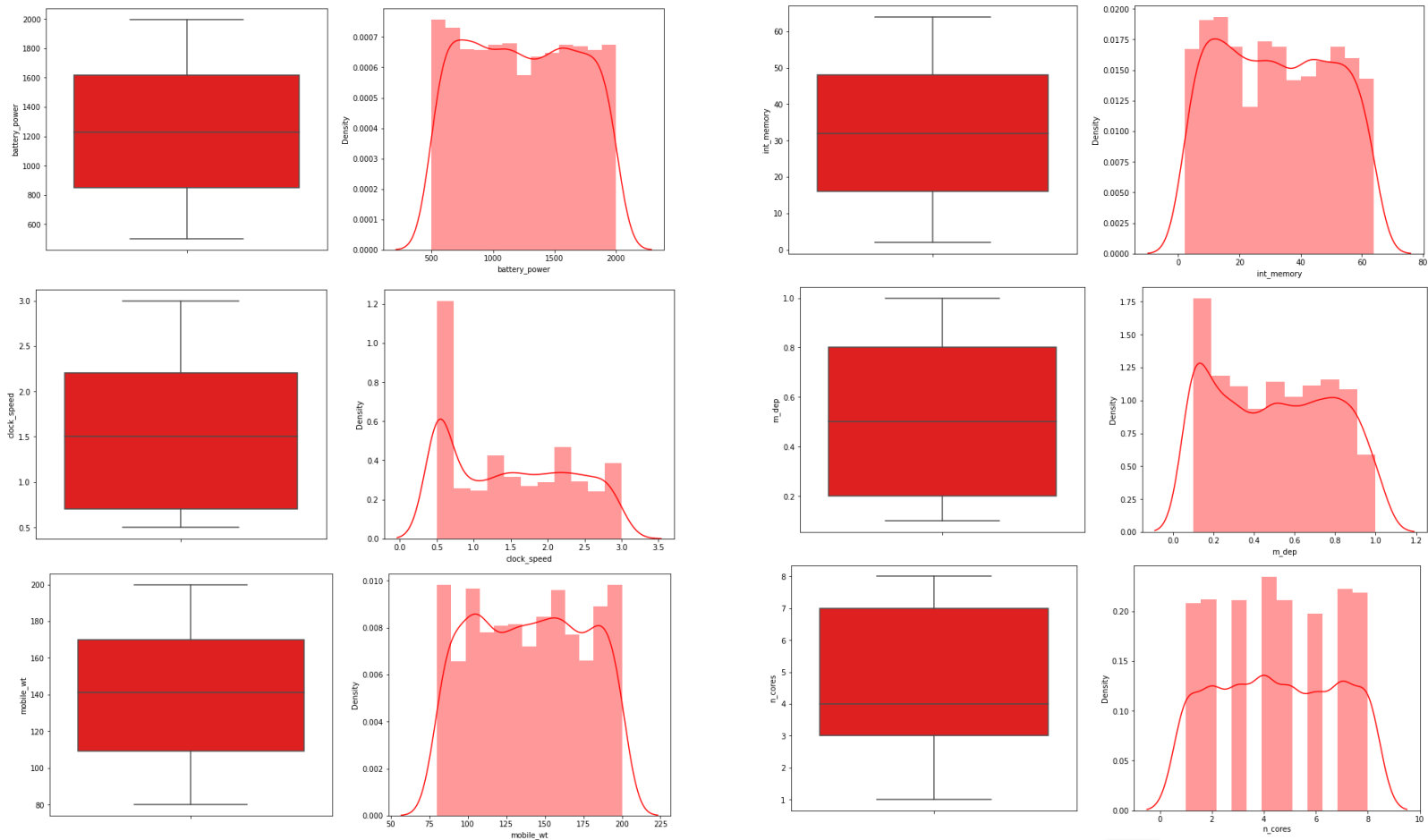
We are haveing 0 rows in our Dataframe.

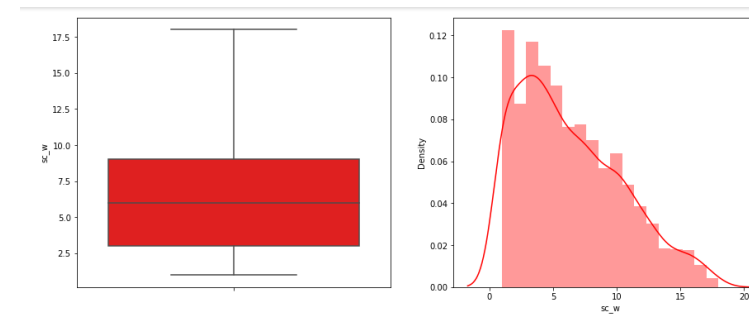
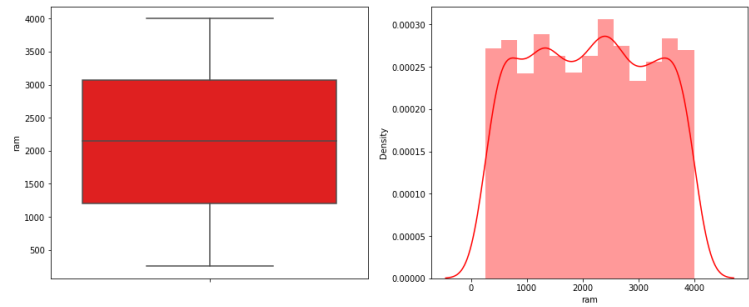
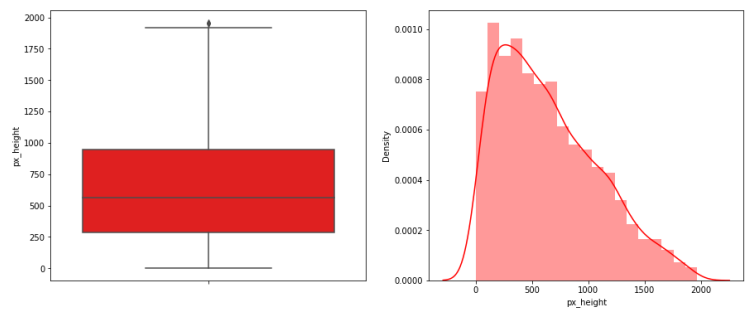
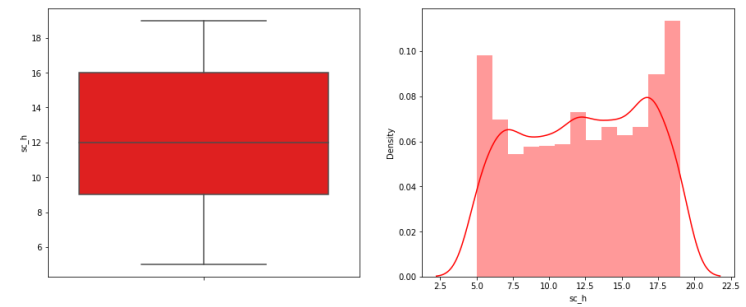
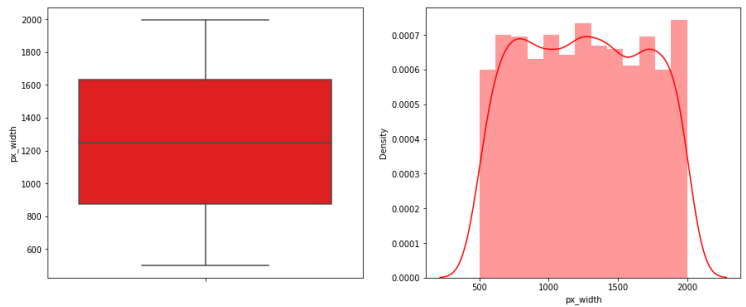
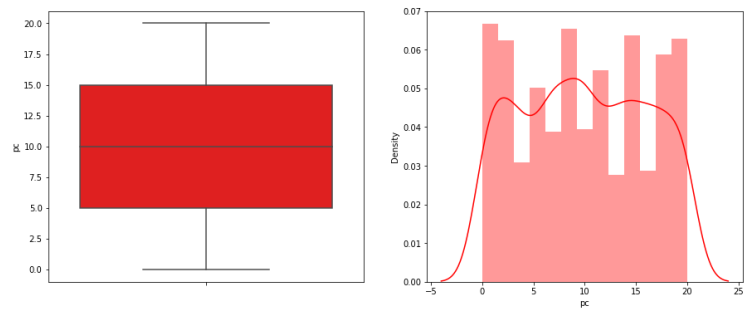
- From the output we can see that we no having any missing value in our dataset after handling mismatching form our data.
- And also we can see that there in no duplicate value.

### Visualising Missing Values

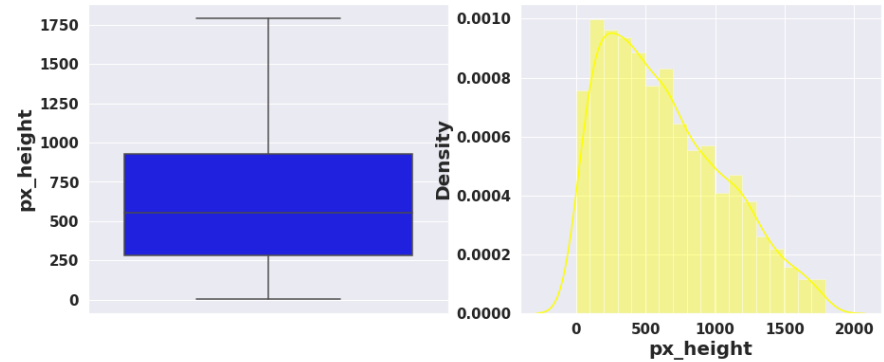
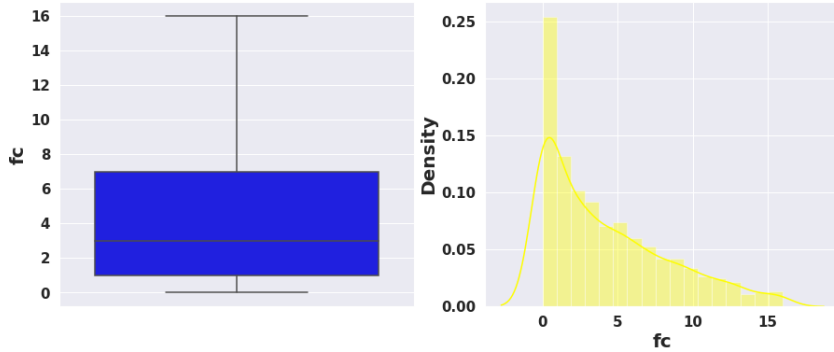


# ➤ Checking Distribution of Numerical Variables and outliers:

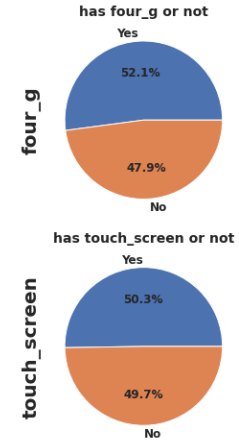
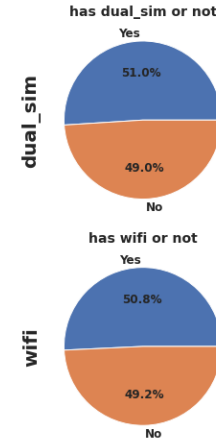
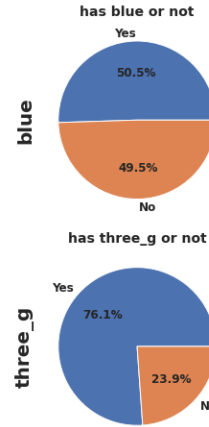
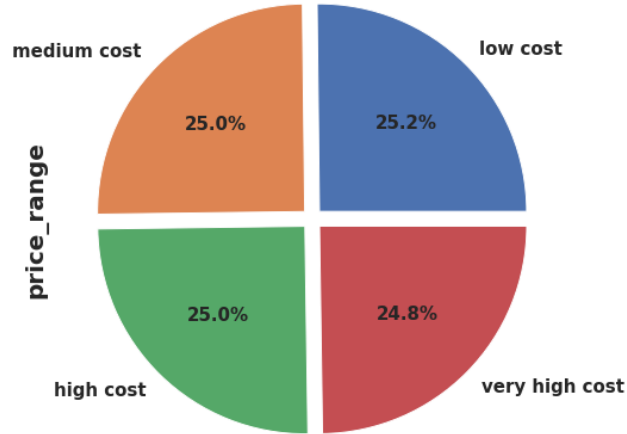




- Form the above chart we can see that our data is well distributed.
- And form boxplot we can see that fc and px\_height has some outliers. So we removed those outlier and we can conform form below output that outlier are removed.

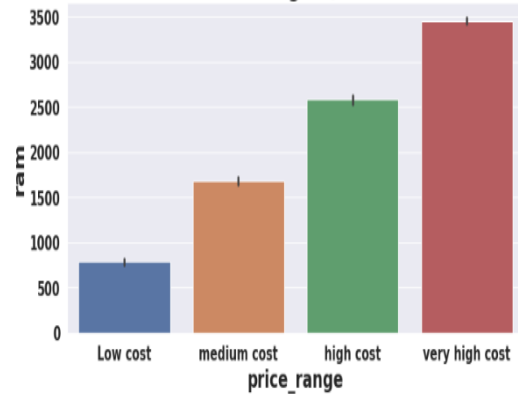


# EDA :

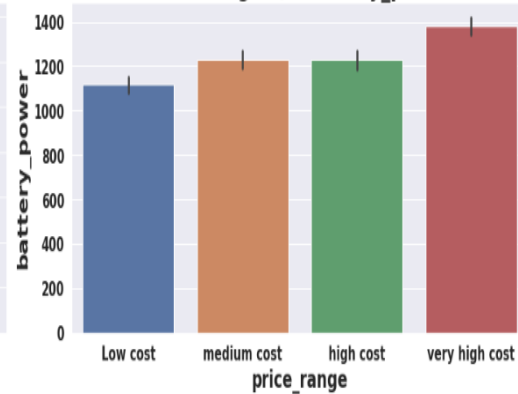


- Target variable has equal number of observation in each category and target variable is nearly equally distributed.
- Percentage Distribution of Mobiles having bluetooth, dual sim, 4G, wifi and touch screen are almost 50 % and very few mobiles (23.8%) do not have 3G.

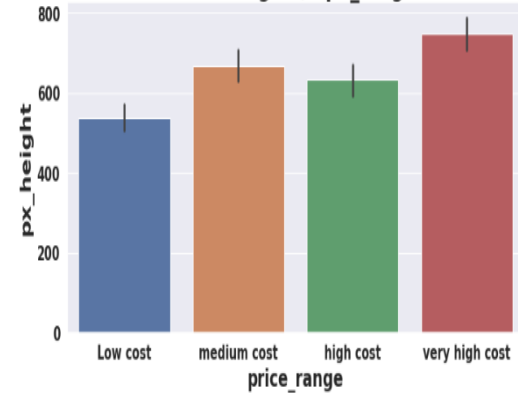
Price range v/s ram



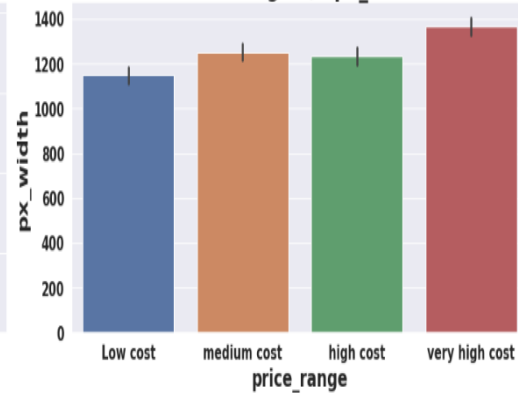
Price range v/s battery\_power



Price range v/s px\_height

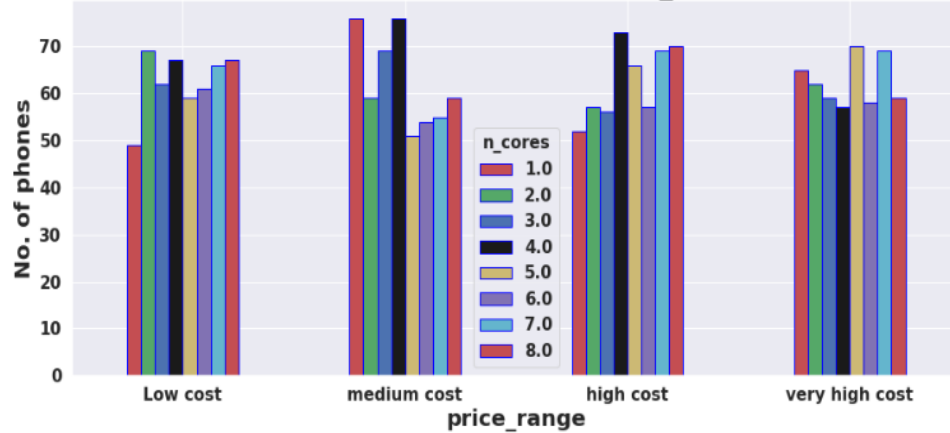


Price range v/s px\_width



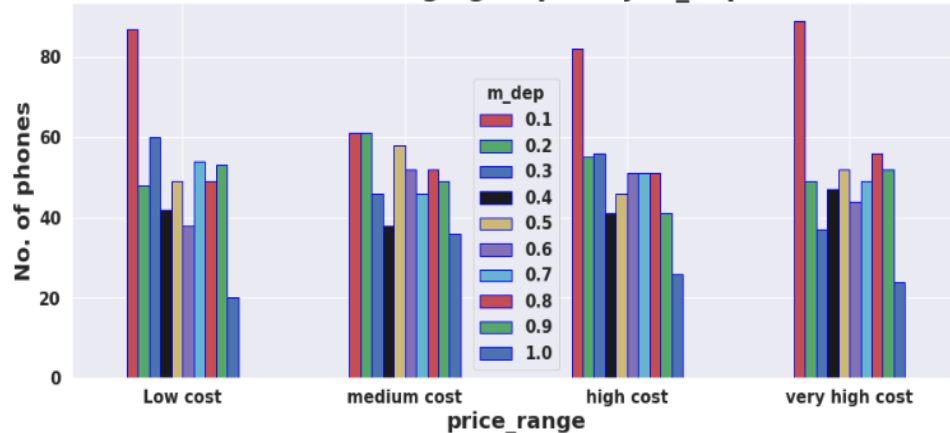
- I. Mobiles having RAM less than 1GB are comes under low cost category.
- II. Mobiles having more than 3GB RAM comes under very high cost category. So we can see that as RAM increases price of mobiles also increases.
- III. Mobiles more than 700 pixel high and width more than 1300 has very high cost.
- IV. Mobiles with battery power more than 1300 mAh has very high cost. And Mobiles with battery power between 1200 and 1300 mAh falls under medium and high cost category.

Price range grouped by n\_cores



1. There are very few mobiles in price range low cost and medium cost are having lesser number of cores.

Price range grouped by m\_dep



2. Most of the mobiles in price range between high cost and very high cost are with more number of cores.



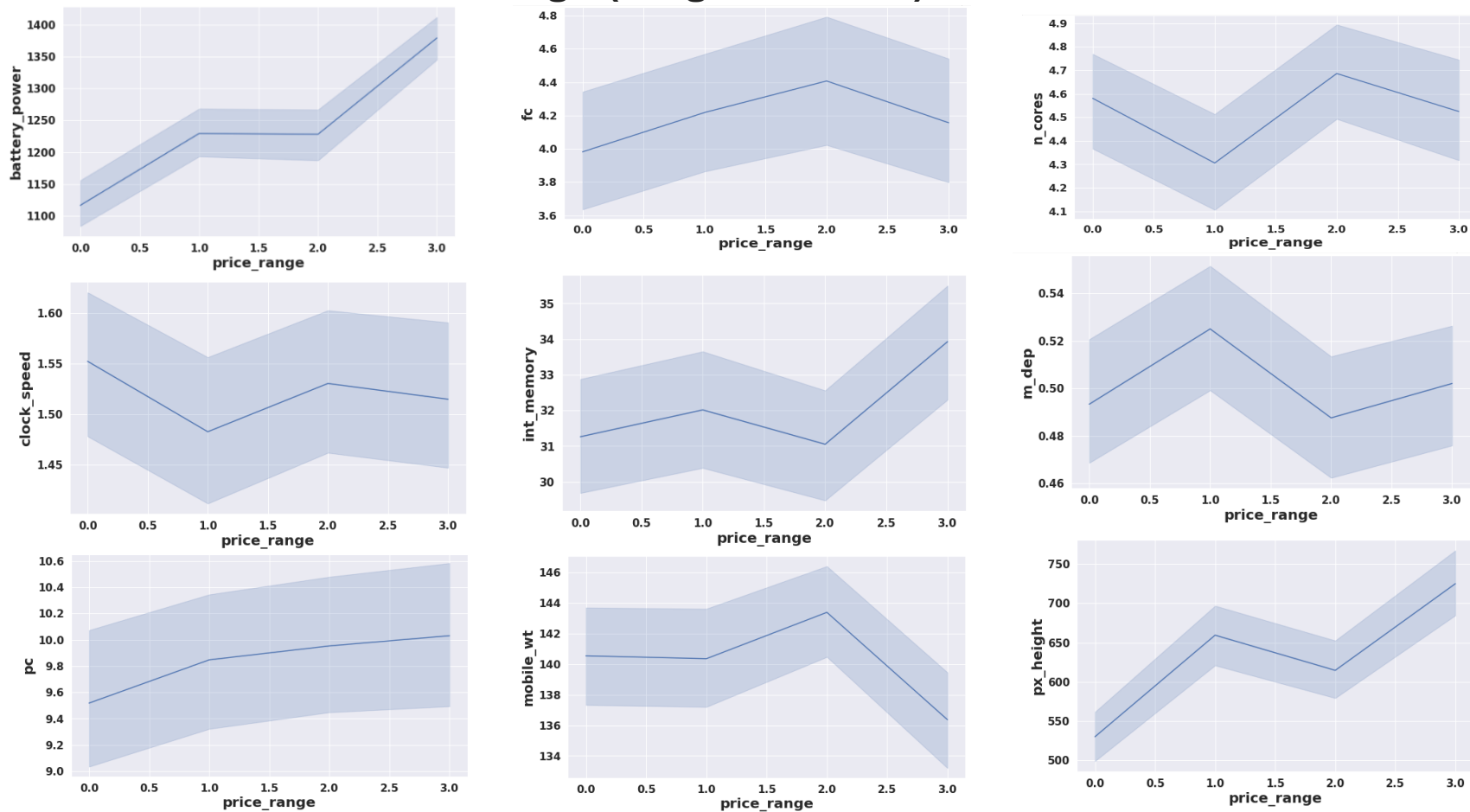
battery_power	1	0.0064	0.0073	-0.039	0.015	0.019	-0.0086	0.031	0.0043	-0.026	0.021	0.014	-0.011	-0.0056	-0.022	-0.01	0.048	0.012	-0.011	-0.011	0.2
blue	0.0064	1	0.022	0.035	0.0014	0.016	0.036	0.0027	-0.01	0.038	-0.01	-0.0077	-0.04	0.024	0.0015	-0.024	0.012	-0.03	0.01	-0.02	0.018
clock_speed	0.0073	0.022	1	0.0026	-0.0052	-0.041	0.0066	-0.01	0.0086	-0.0081	-0.011	-0.012	-0.01	0.00044	-0.03	0.004	-0.013	-0.043	0.017	-0.024	-0.0089
dual_sim	-0.039	0.035	0.0026	1	-0.034	0.0013	-0.021	-0.02	-0.0044	-0.027	-0.018	-0.015	0.019	0.046	-0.01	0.0021	-0.044	-0.014	-0.011	0.025	0.023
fc	0.015	0.0014	-0.0052	-0.034	1	-0.017	-0.028	0.001	0.014	-0.0016	0.64	-0.027	-0.013	0.018	0.0029	0.0019	-0.0061	-0.0012	-0.022	0.012	0.019
four_g	0.019	0.016	-0.041	0.0013	-0.017	1	0.008	-0.0011	-0.016	-0.032	-0.0047	-0.031	0.0016	0.0067	0.027	0.039	-0.048	0.58	0.018	-0.021	0.014
int_memory	0.0086	0.036	0.0066	-0.021	-0.028	0.008	1	0.0076	-0.032	-0.027	-0.029	0.008	-0.013	0.033	0.044	0.0016	-0.0072	-0.0097	-0.028	0.011	0.043
m_dep	0.031	0.0027	-0.01	-0.02	0.001	-0.0011	0.0076	1	0.019	-0.0033	0.028	0.021	0.022	-0.012	-0.024	-0.0077	0.017	-0.013	-0.0041	-0.03	-0.0043
mobile_wt	0.0043	-0.01	0.0086	-0.0044	0.014	-0.016	-0.032	0.019	1	-0.021	0.014	-0.0014	0.00053	-0.0029	-0.034	-0.041	0.013	0.0054	-0.015	-0.0032	-0.03
n_cores	-0.026	0.038	-0.0081	-0.027	-0.0016	-0.032	-0.027	-0.0033	-0.021	1	0.007	-0.0048	0.026	0.0094	-0.0029	0.018	0.012	-0.015	0.027	-0.0063	0.01
pc	0.021	-0.01	-0.011	-0.018	0.64	-0.0047	-0.029	0.028	0.014	0.007	1	-0.026	0.0016	0.028	0.012	-0.0098	0.016	-0.0034	-0.012	-0.002	0.031
px_height	0.014	-0.0077	-0.012	-0.015	-0.027	-0.031	0.008	0.021	-0.0014	-0.0048	-0.026	1	0.49	-0.023	0.055	0.049	-0.015	-0.041	0.009	0.041	-0.14
px_width	-0.011	-0.04	-0.01	0.019	-0.013	0.0016	-0.013	0.022	0.00053	0.026	0.0016	0.49	1	0.0033	0.017	0.051	0.0018	-0.0052	-0.0097	0.027	0.16
ram	-0.0056	0.024	0.00044	0.046	0.018	0.0067	0.033	-0.012	-0.0029	0.0094	0.028	-0.023	0.0033	1	0.019	0.025	0.011	0.018	-0.036	0.022	0.92
sc_h	-0.022	0.0015	-0.03	-0.01	0.0029	0.027	0.044	-0.024	-0.034	-0.0029	0.012	0.055	0.017	0.019	1	0.46	-0.02	0.013	-0.018	0.023	0.027
sc_w	-0.01	-0.024	0.004	0.0021	0.0019	0.039	0.0016	-0.0077	-0.041	0.018	-0.0098	0.049	0.051	0.025	0.46	1	-0.014	0.043	0.0053	0.016	0.037
talk_time	0.048	0.012	-0.013	-0.044	-0.0061	-0.048	-0.0072	0.017	0.013	0.012	0.016	-0.015	0.0018	0.011	-0.02	-0.014	1	-0.046	0.015	-0.028	0.02
three_g	0.012	-0.03	-0.043	-0.014	-0.0012	0.58	-0.0097	-0.013	0.0054	-0.015	-0.0034	-0.041	-0.0052	0.018	0.013	0.043	-0.046	1	0.014	0.0028	0.025
touch_screen	-0.011	0.01	0.017	-0.011	-0.022	0.018	-0.028	-0.0041	-0.015	0.027	-0.012	0.009	-0.0097	-0.036	-0.018	0.0053	0.015	0.014	1	0.011	-0.038
wifi	-0.011	-0.02	-0.024	0.025	0.012	-0.021	0.011	-0.03	-0.0032	-0.0063	-0.002	0.041	0.027	0.022	0.023	0.016	-0.028	0.0028	0.011	1	0.016
price_range	0.2	0.018	-0.0089	0.023	0.019	0.014	0.043	-0.0043	-0.03	0.01	0.031	0.14	0.16	0.92	0.027	0.037	0.02	0.025	-0.038	0.016	1

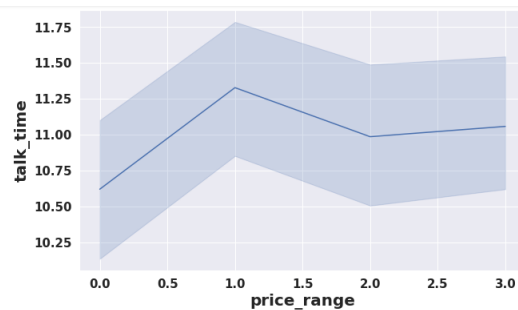
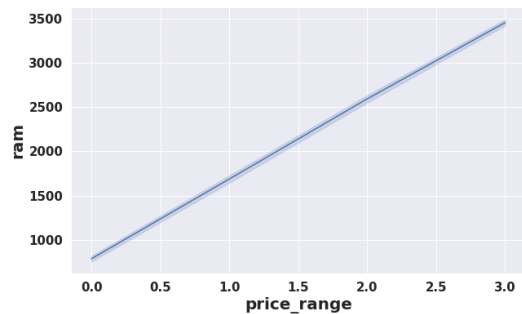
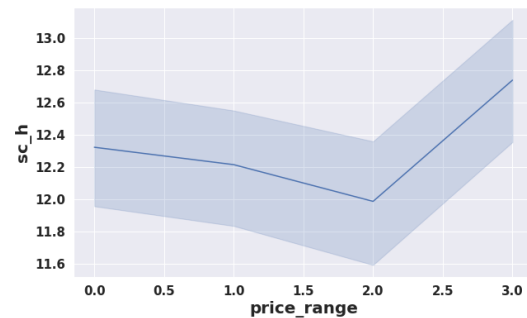
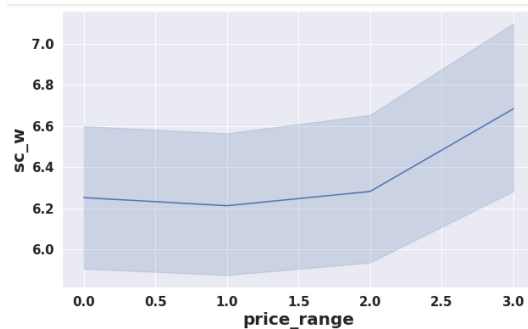
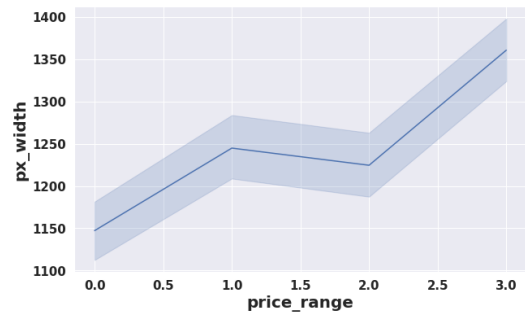


Form the correlation matrix,

1. Price range is having strong positive correlation with RAM. As we know mobiles are costly with high RAM.
2. Price range has positive correlation with Battery power. Mostly high price mobiles are having great battery backup.
3. Also px\_height and px\_width (Pixel Resolution Height and width) are positively correlated. Generally High price range mobiles have good resolutions.
4. 3G and 4G are positively correlated with each other i.e. because nowadays all mobiles is rather 3G or 4G .

## Different trends of price range (Target variable) with other attributes:





# **Model Selection and Evaluation:**

Before doing analysis with models we performed the train test split. We kept 25% of the data for test and remaining 75% of the data for training the model.

We compared 6 Classification algorithms and evaluated them based on the overall accuracy score and the recall of the individual classes.

- Accuracy is the ratio of the total number of correct predictions and the total number of predictions.
- The recall is the measure of our model correctly identifying True Positives.

1. **Decision Tree**
2. **Random Forest Classifier**
3. **Gradient Boosting Classifier**
4. **K-Nearest Neighbor Classifier**
5. **XG Boost Classifier**
6. **Support Vector Machine (SMV)**

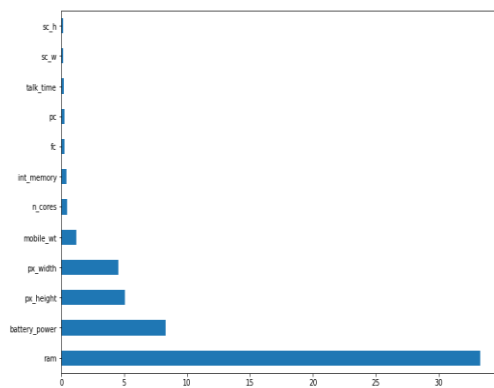
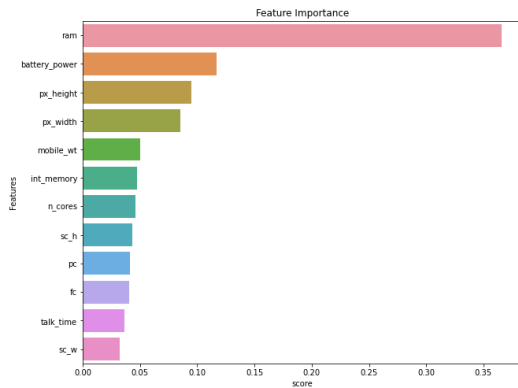
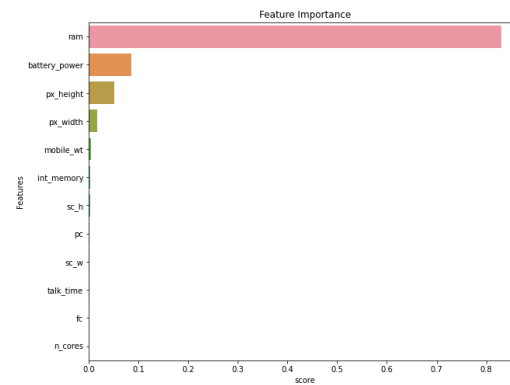
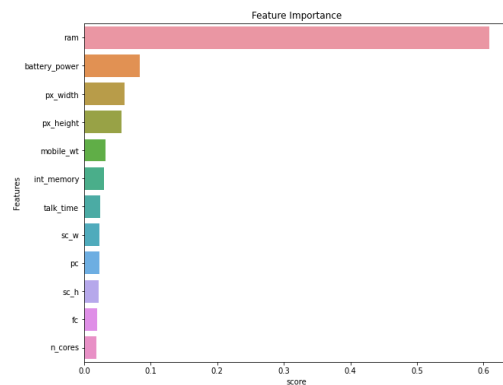
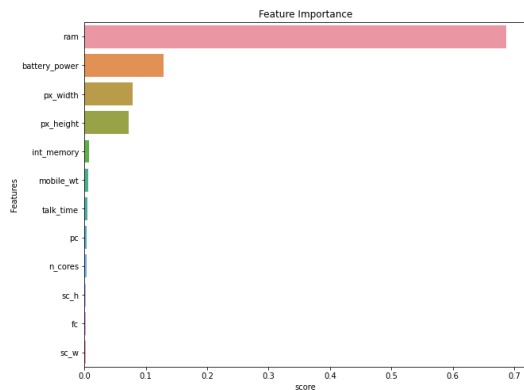
Decision Tree 2) Random Forest classifier 3) Gradient Boosting Classifier 4) K-nearest Neighbor classifier 5) XG Boost Classifier 6) Support Vector Machine(SVM)

# Evaluation of Models:

Algorithms	Training Set		Test Set	
	Accuracy score (%)	Recall (%)	Accuracy score (%)	Recall (Average of all 4 Classes)
Decision Tree	100	100	84	83.75
Decision Tree (Hyper-parameter Tuning)	97.62	97.5	85.13	84.75
Random Forest	100	100	88.6	88.5
Random Forest(Hyper-parameter Tuning)	100	100	89.81	89.5
Gradient Boosting	100	100	90.02	90
Gradient Boosting (Hyper-parameter Tuning)	100	100	90.42	90.5
KNN	75.86	76	59.47	59.25
KNN (Hyper-parameter Tuning)	76.61	76.75	70.26	69.75
XG-Boost	98.98	98.75	90.22	90
XG-Boost (Hyper-parameter Tuning)	100	100	92.46	92.25
SVM	98.57	98.5	89.81	89.75
SVM (Hyper-parameter Tuning)	98.3	98.5	97.96	98

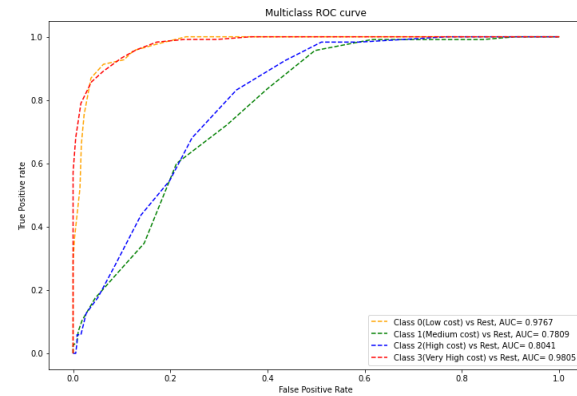
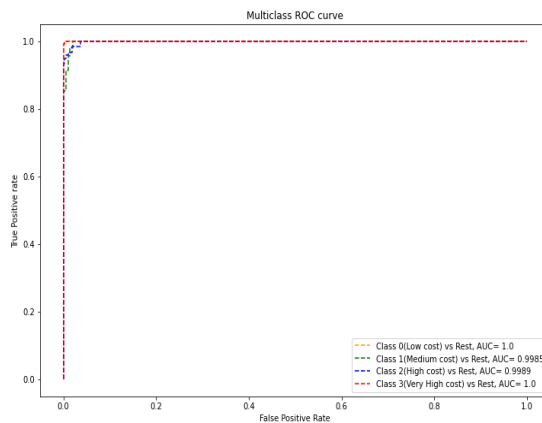
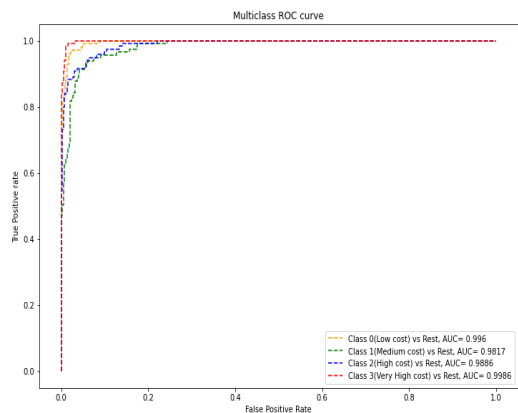
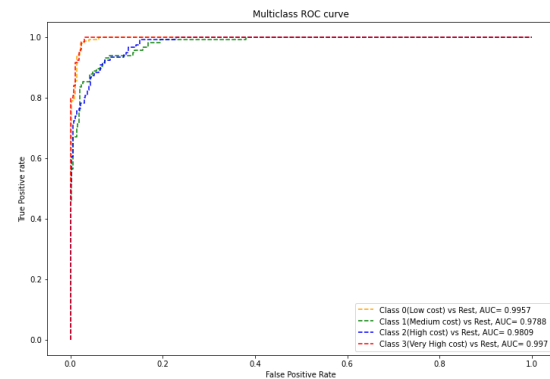
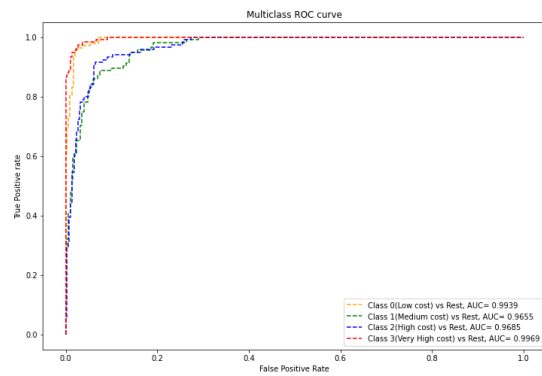
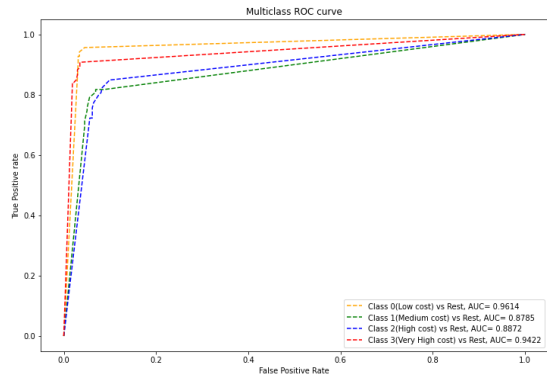
- After hyper-parameter tuning the Best Fitted model came out to be SMV.
- XG boost (Hyper-parameter Tuned) can be considered as the second most good model.
- The worst performed model is KNN model.

# Feature Importance:



**RAM, Battery Power, Pixel height and weight contributed the most in predicting the price range.**

# ● AOC ROC Curve:



## **Conclusion:**

- Target variable has equal number of observation in each category and target variable is nearly equally distributed.
- Percentage Distribution of Mobiles having bluetooth, dual sim, 4G,wifi and touch screen are almost 50 % and very few mobiles (23.8%) do not have 3G.
- Price range is having strong positive correlation with RAM. As we know mobiles are costly with high RAM.
- Price range has positive correlation with Battery power. Mostly high price mobiles are having great battery backup.



- **Support Vector Machine (SMV) algorithm gave best performance after hyper-parameter tuning with 98.3% train accuracy and 97% is test accuracy.**
- **XG boost is the second best good model which gave good performance after hyper-parameter tuning with 100% train accuracy and 92.25% test accuracy score.**
- **KNN gave very worst model performance.**

**THANK  
YOU!**