

Capstone Project

Mobile Price Range Prediction

Supervised Machine Learning (Classification)

INDEX

- **Problem Statement**
- **Data Description**
- **Data Wrangling**
- **EDA (Exploratory Data Analysis)**
- **Feature Engineering**
- **Model Selection and Evaluation**
- **Evaluation of models**
- **Feature Importance**
- **Challenges**
- **Conclusion.**



Problem Statement

- Mobile phones have become a necessity for every individual nowadays. People want more features and best specifications in a phone and that too at cheaper prices.
- Mobile phones come in all sorts of prices, features, specifications and all. 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.
- 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 (e.g.:- 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.
- The main objective of this project is to build a model which will classify the price range of mobile phones based on the specifications of mobile phones.



❖ Data Description

Total Rows= 2000 , Total features=21

- **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.

❖ Data Description

- Pc - Primary Camera mega pixels.
- Px_height and Px_width - Pixel Resolution Height and width.
- Ram - Random Access Memory in Mega Bytes.
- Sc_h and Sc_w - Screen Height and 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.

Target variable:

- Price_range - This is the target variable with value of 0(low cost),1(medium cost),2(high cost) and3(very high cost).

❖ Data Wrangling

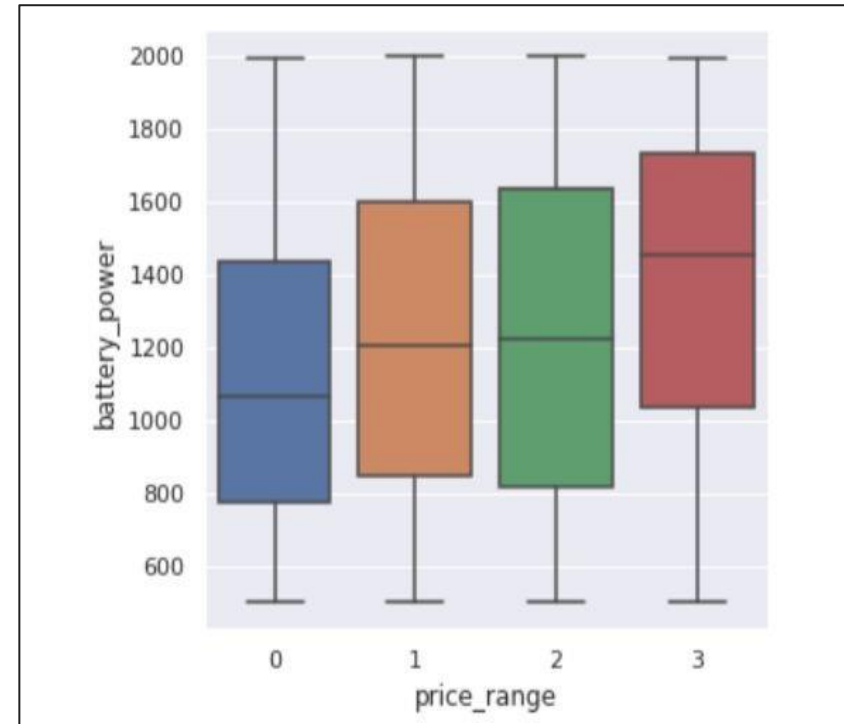
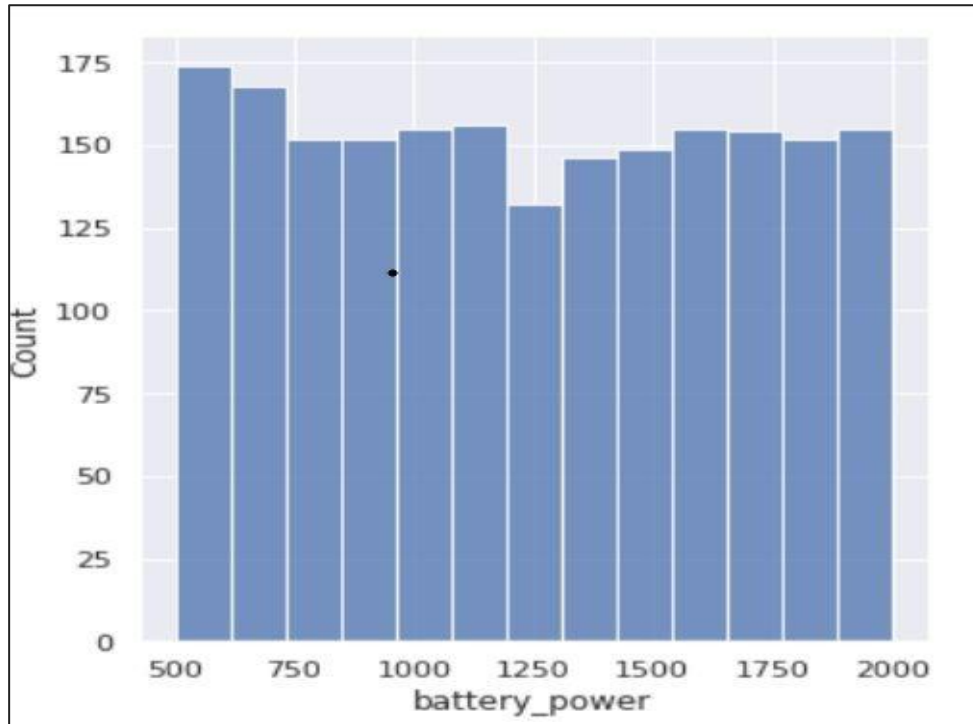
➤ Handling missing values present in data

	count	mean	std	min	25%	50%	75%	max
sc_width	2000.0	5.76700	4.356398	0.0	2.00	5.0	9.00	18.0
px_height	2000.0	645.10800	443.780811	0.0	282.75	564.0	947.25	1960.0

- From above information, min screen width & min pixel height is zero. So that means this information contains missing values so let's check there respective counts.
- Total phones with `sc_w = 0` is 180
- Total phones with `px_height = 0` is 2
- Where there is `sc_W` and `px_height` is zero ,assigning mean value to that value in order to handle missing values.
- In This data set there is not any duplicate values present.

❖ Exploratory Data Analysis

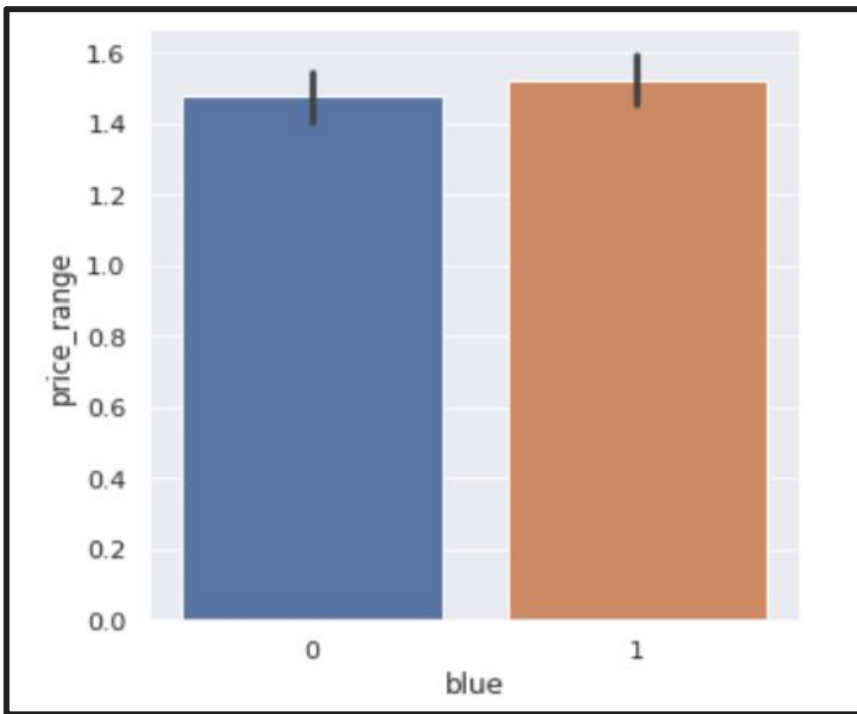
This Graph shows that as battery power increase there is gradual increase in mobile price range.



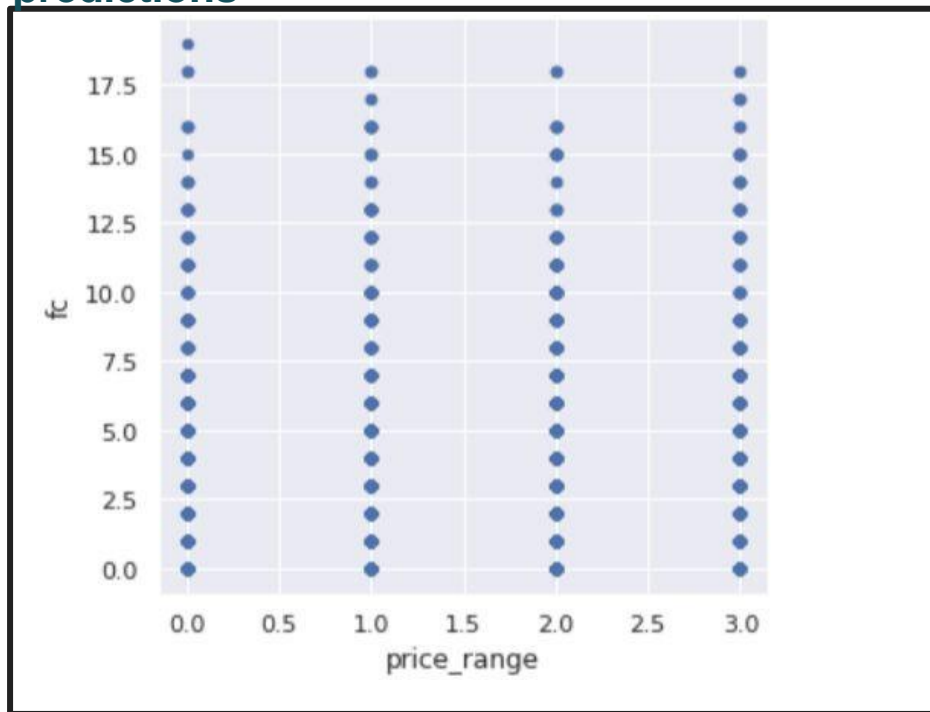


Exploratory Data Analysis

Having bluetooth or not do not affect the price of mobile

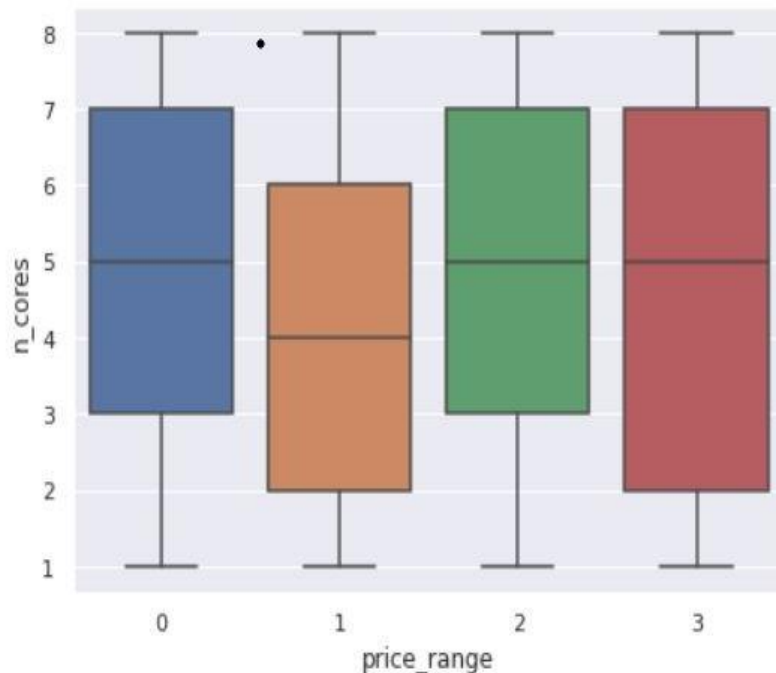
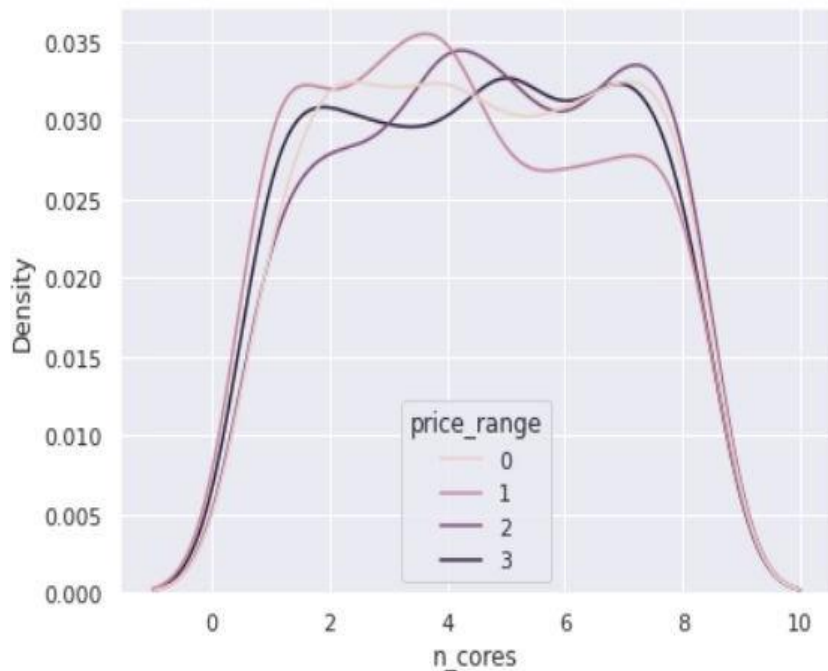


Front camera megapixel distribution is almost similar along all the price ranges variable, it may not be helpful in making predictions



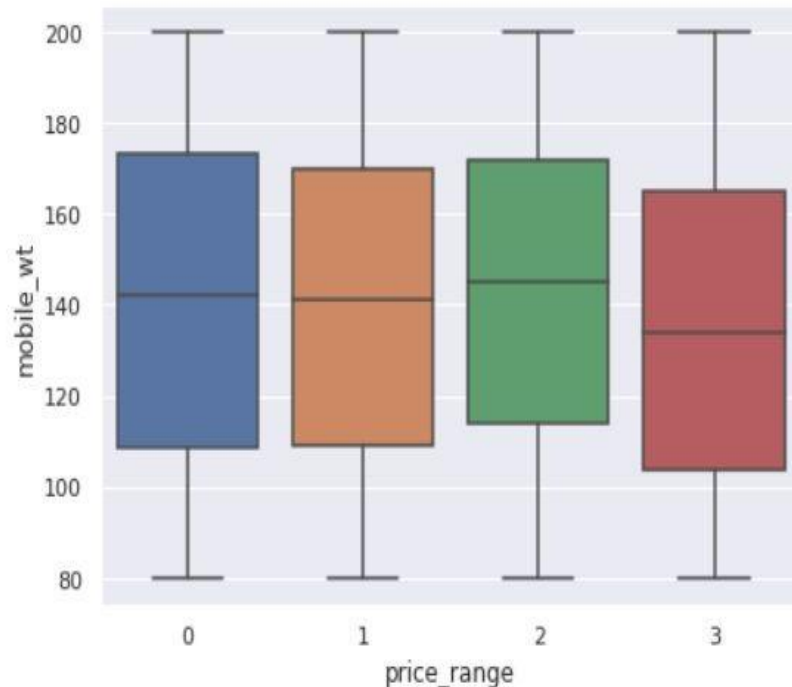
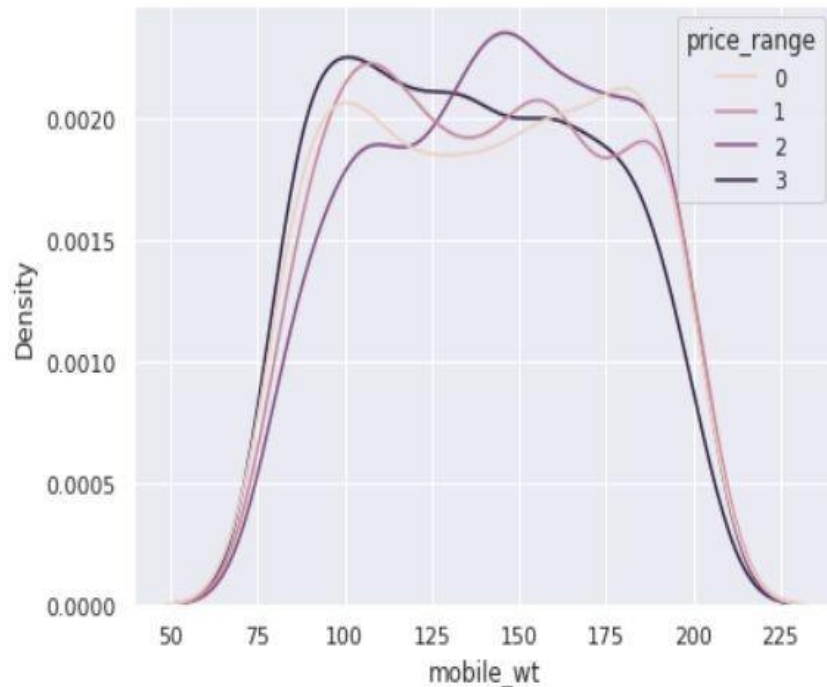
❖ Exploratory Data Analysis

No of core Processors are showing a little variation along the target categories that helps in price prediction.



✦ Exploratory Data Analysis

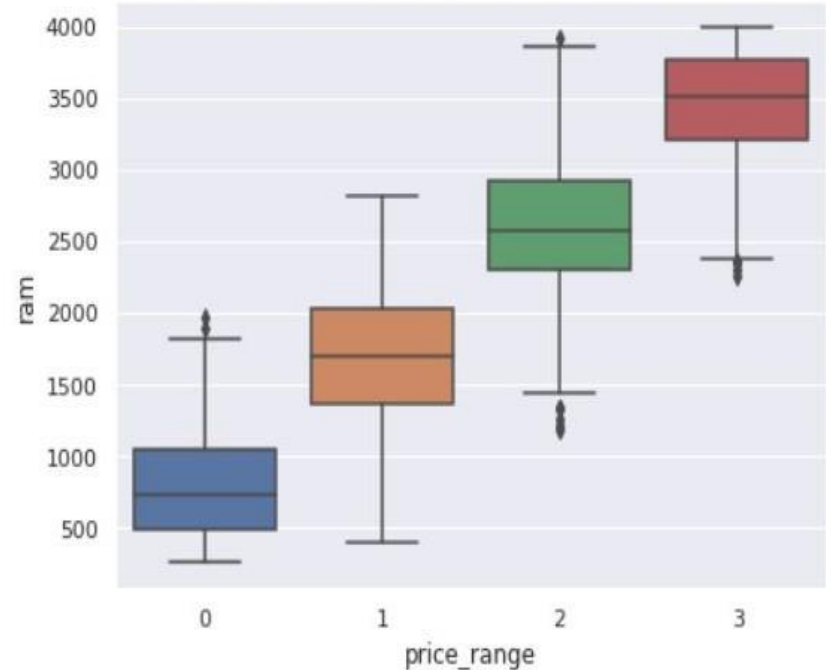
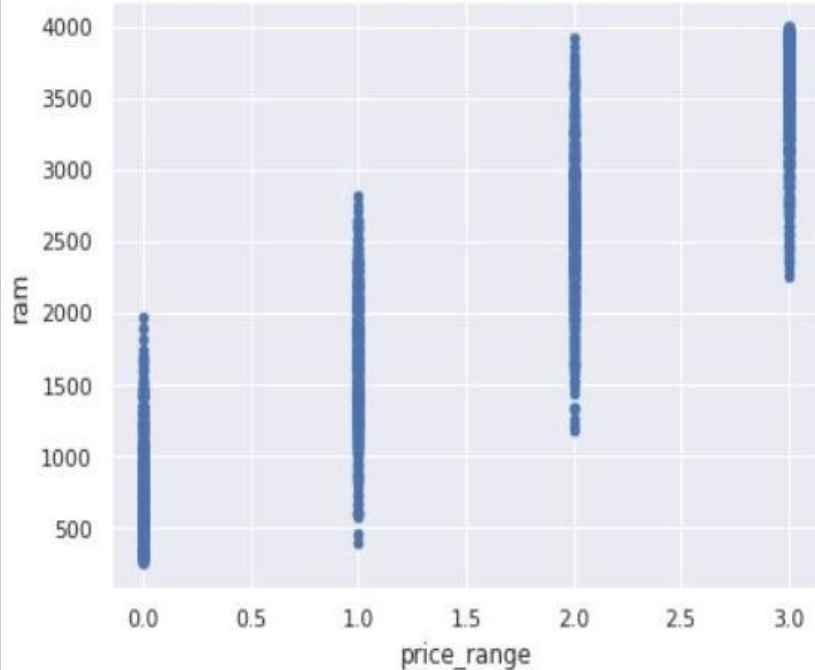
Costly phones are lighter. mobile weight affect the pricing of mobile





Exploratory Data Analysis

- Mobiles having RAM more than 3000MB falls under Very high cost category. As RAM increases price range also increases.
- Mobiles having RAM less than 1000 MB falls under low cost category.

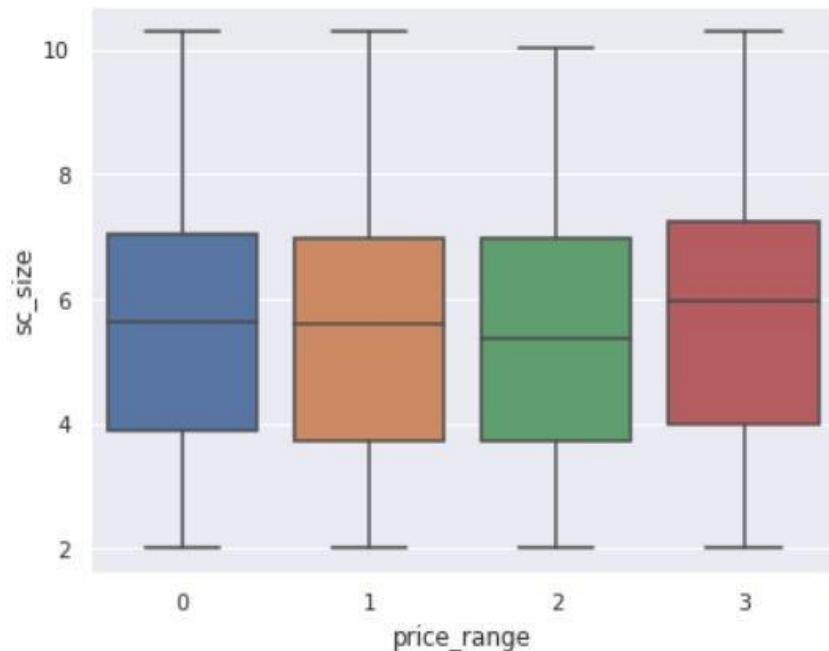
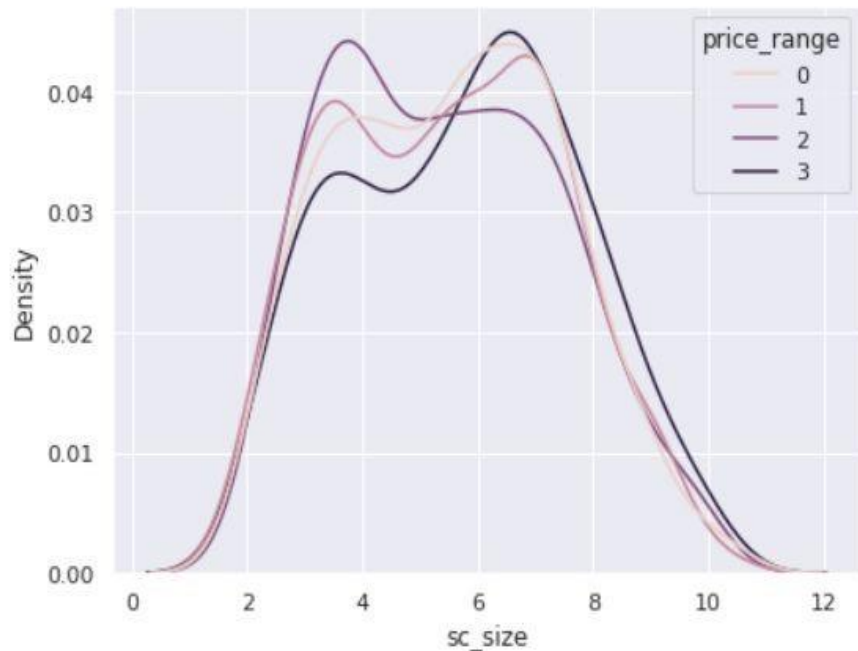




Exploratory Data Analysis

AI

Screen Size shows little variation along the target variables. This can be helpful in predicting the target variable that is mobile price.

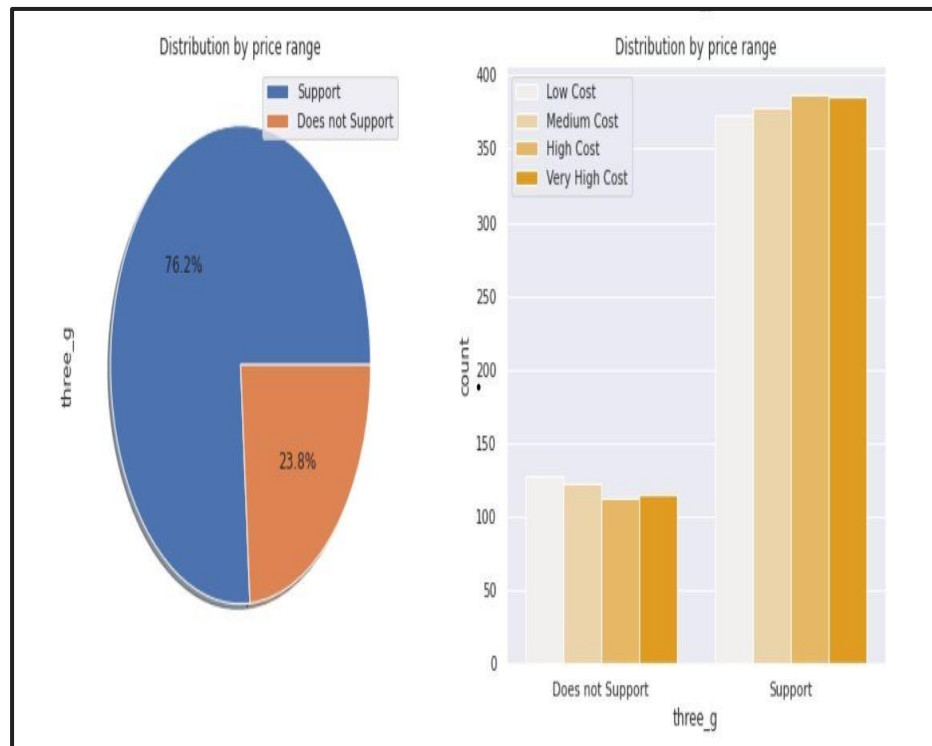
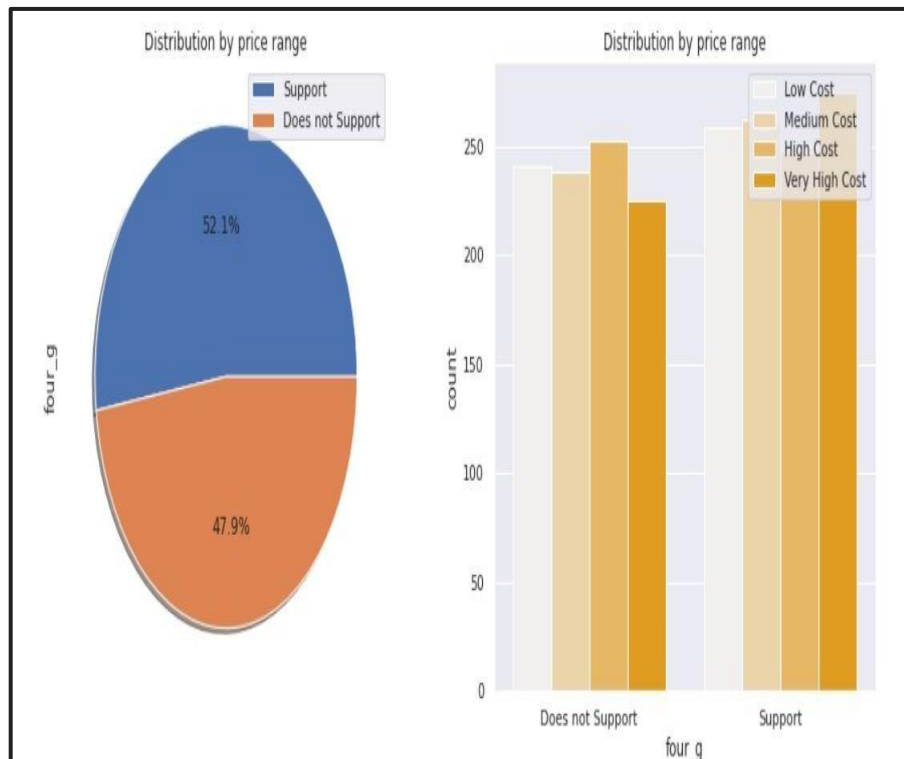




Exploratory Data Analysis



- 50% of the phones support 4_g and 76% of phones support 3_g
- feature 'three_g' play an important feature in prediction of mobile price.

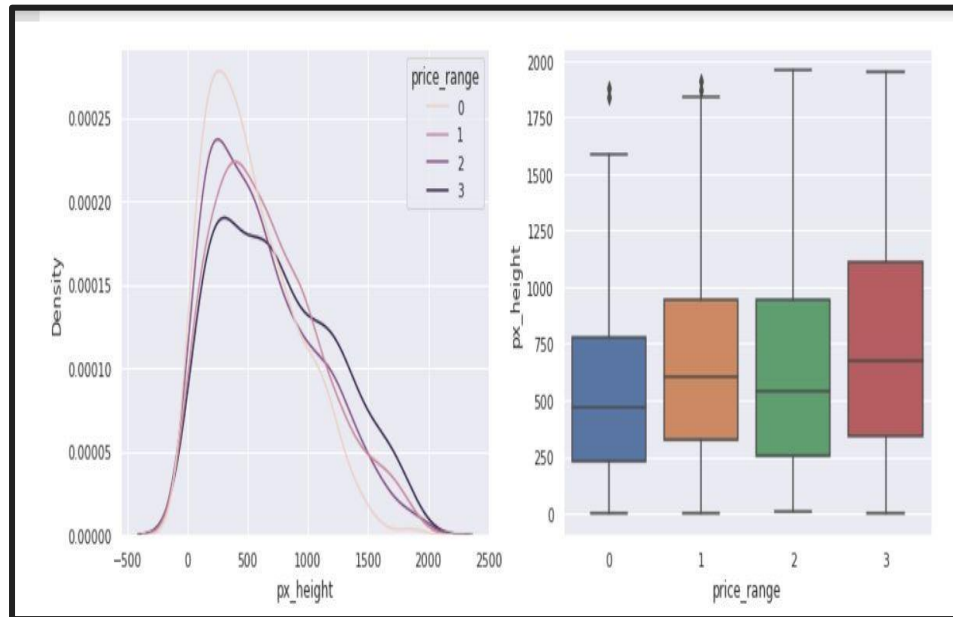
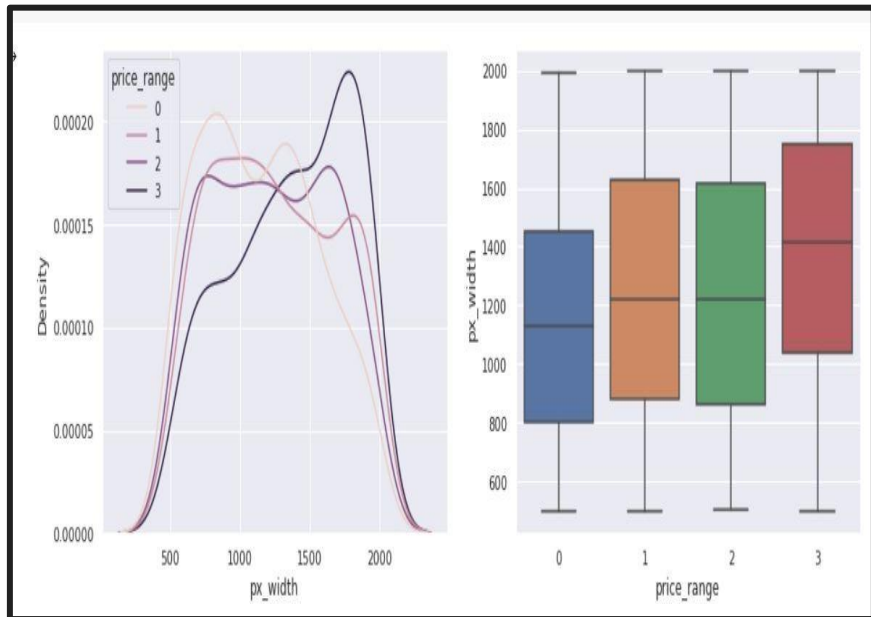




Exploratory Data Analysis

AI

- Pixel height and pixel width shows little variation in mobile price range.
- There is not a continuous increase in pixel width & height as we move from Low cost to Very high cost. Mobiles with 'Medium cost' and 'High cost' has almost equal pixel width & height





Feature Engineering

- RAM and price_range shows high correlation which is a good sign, it signifies that RAM will play major deciding factor in estimating the price range.
- There is some collinearity in features ('pc', 'fc') and ('px_width', 'px_height') and (sc_w, sc_h). Both correlations are justified since there are good chances that if front camera of a phone is good, the back camera would also be good.
- Also, if px_height increases, pixel width also increases, that means the overall pixels in the screen. We can replace these two features with one feature. Front Camera megapixels and Primary camera megapixels are different entities despite of showing colinearity. So we'll be keeping them as they are.
- Also, if sc_h increases, sc w also increases, that means the overall screen size increase. We can replace these two features with one feature i.e screen_size

	battery_power	blue	clock_speed	dual_sim	fc	four_g	int_memory	m_dep	mobile_wt	n_cores	pc	px_height	px_width	ram	sc_h	sc_w	talk_time	three_g	touch_screen	wifi	price_range
battery_power	1	0.011	0.011	0.042	0.033	0.016	0.004	0.034	0.0018	0.03	0.031	0.015	0.0084	0.0065	0.03	0.021	0.053	0.012	0.011	0.0083	0.2
blue	0.011	1	0.021	0.035	0.0036	0.013	0.041	0.004	0.0086	0.036	0.01	0.0069	0.042	0.026	0.003	0.0006	0.014	0.03	0.01	0.022	0.021
clock_speed	0.011	0.021	1	0.0013	0.0043	0.043	0.0065	0.014	0.012	0.0057	0.0052	0.015	0.0095	0.0034	0.029	0.0074	0.011	0.046	0.02	0.024	0.0066
dual_sim	0.042	0.035	0.0013	1	0.029	0.0032	0.016	0.022	0.009	0.025	0.017	0.021	0.014	0.041	0.012	0.017	0.039	0.014	0.017	0.023	0.017
fc	0.033	0.0036	0.0043	0.029	1	0.017	0.029	0.0018	0.024	0.013	0.64	0.01	0.0052	0.015	0.011	0.012	0.0068	0.0018	0.015	0.02	0.022
four_g	0.016	0.013	0.043	0.0032	0.017	1	0.0087	0.0018	0.017	0.03	0.0056	0.019	0.0074	0.0073	0.027	0.037	0.047	0.58	0.017	0.018	0.015
int_memory	0.004	0.041	0.0065	0.016	0.029	0.0087	1	0.0069	0.034	0.028	0.033	0.01	0.0083	0.033	0.038	0.012	0.0028	0.0094	0.027	0.007	0.044
m_dep	0.034	0.004	0.014	0.022	0.0018	0.0018	0.0069	1	0.022	0.0035	0.026	0.025	0.024	0.0094	0.025	0.018	0.017	0.012	0.0026	0.028	0.0008
mobile_wt	0.0018	0.0086	0.012	0.009	0.024	0.017	0.034	0.022	1	0.019	0.019	0.0094	9e-05	0.0026	0.034	0.021	0.0062	0.0016	0.014	0.0004	0.03
n_cores	0.03	0.036	0.0057	0.025	0.013	0.03	0.028	0.0035	0.019	1	0.0012	0.0069	0.024	0.0049	0.0003	0.026	0.013	0.015	0.024	0.01	0.0044
pc	0.031	0.01	0.0052	0.017	0.64	0.0056	0.033	0.026	0.019	0.0012	1	0.018	0.0042	0.029	0.0049	0.024	0.015	0.0013	0.0087	0.0054	0.034
px_height	0.015	0.0069	0.015	0.021	0.01	0.019	0.01	0.025	0.0094	0.0069	0.018	1	0.51	0.02	0.06	0.043	0.011	0.031	0.022	0.052	0.15
px_width	0.0084	0.042	0.0095	0.014	0.0052	0.0074	0.0083	0.024	9e-05	0.024	0.0042	0.51	1	0.0041	0.022	0.035	0.0067	0.00035	0.0016	0.03	0.17
ram	0.0065	0.026	0.0034	0.041	0.015	0.0073	0.033	0.0094	0.0026	0.0049	0.029	0.02	0.0041	1	0.016	0.036	0.011	0.016	0.03	0.023	0.92
sc_h	0.03	0.003	0.029	0.012	0.011	0.027	0.038	0.025	0.034	0.0003	0.0049	0.06	0.022	0.016	1	0.51	0.017	0.012	0.02	0.026	0.023
sc_w	0.021	0.0006	0.0074	0.017	0.012	0.037	0.012	0.018	0.021	0.026	0.024	0.043	0.035	0.036	0.51	1	0.023	0.031	0.013	0.035	0.039
talk_time	0.053	0.014	0.011	0.039	0.0068	0.047	0.0028	0.017	0.0062	0.013	0.015	0.011	0.0067	0.011	0.017	0.023	1	0.043	0.017	0.03	0.022
three_g	0.012	0.03	0.046	0.014	0.0018	0.58	0.0094	0.012	0.0016	0.015	0.0013	0.031	0.0003	0.016	0.012	0.031	0.043	1	0.014	0.0043	0.024
touch_screen	0.011	0.01	0.02	0.017	0.015	0.017	0.027	0.0026	0.014	0.024	0.0087	0.022	0.0016	0.03	0.02	0.013	0.017	0.014	1	0.012	0.03
wifi	0.0083	0.022	0.024	0.023	0.02	0.018	0.007	0.028	0.0004	0.01	0.0054	0.052	0.03	0.023	0.026	0.035	0.03	0.0043	0.012	1	0.019
price_range	0.2	0.021	0.0066	0.017	0.022	0.015	0.044	0.0008	0.03	0.0044	0.034	0.15	0.17	0.92	0.023	0.039	0.022	0.024	0.03	0.019	1



Feature Engineering



- Similarly the screen size of the phone is expressed in Inches.
- We have columns '**sc_h**' and '**sc_w**' out of which we have created a new feature '**Screen_size**' which is diagonal length of the screen.

Model Selection and Evaluation :

Before building a 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 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(SVM).



Evaluation of models:

Algorithms	Training Set		Test set	
	Accracy score (%)	Recall score(%)	Accracy score (%)	Recall score(%)
Logistic Regression	92%	92%	89%	89%
Random forest	100%	100%	87%	87%
Random forest (hyper tuning)	95%	95%	84%	84%
K Nearst Neighbour	100%	100%	44%	44%
XGBoost	98%	98%	89%	89%
Support Vector Machine	94%	94%	90%	90%
Gradient Boosting	99%	99%	90%	90%
Decision Tree	88%	88%	84%	84%

- Best model came out to be Support VM.
- XG boost can be considered as the second most good model.
- KNN performed very worst.



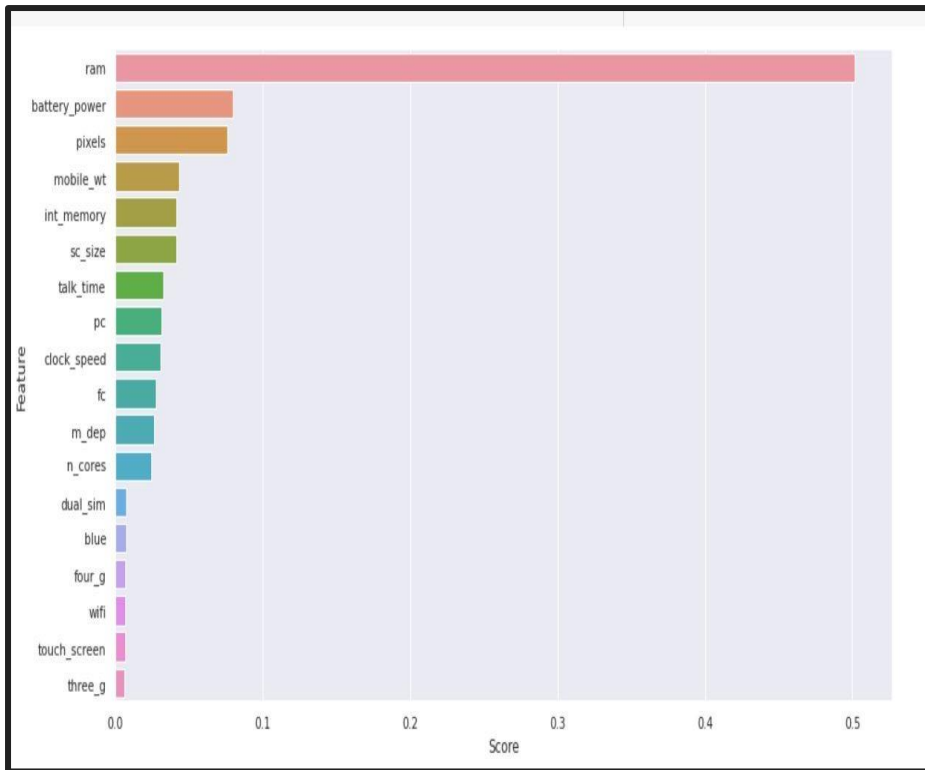
Evaluation of models:



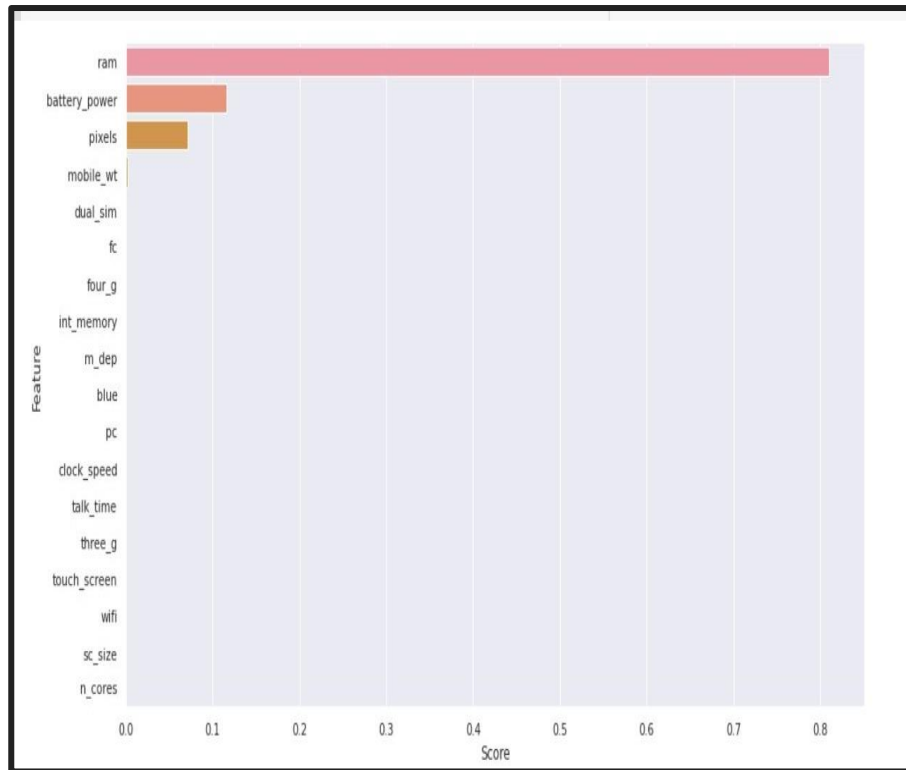


Feature Importance

Random forest

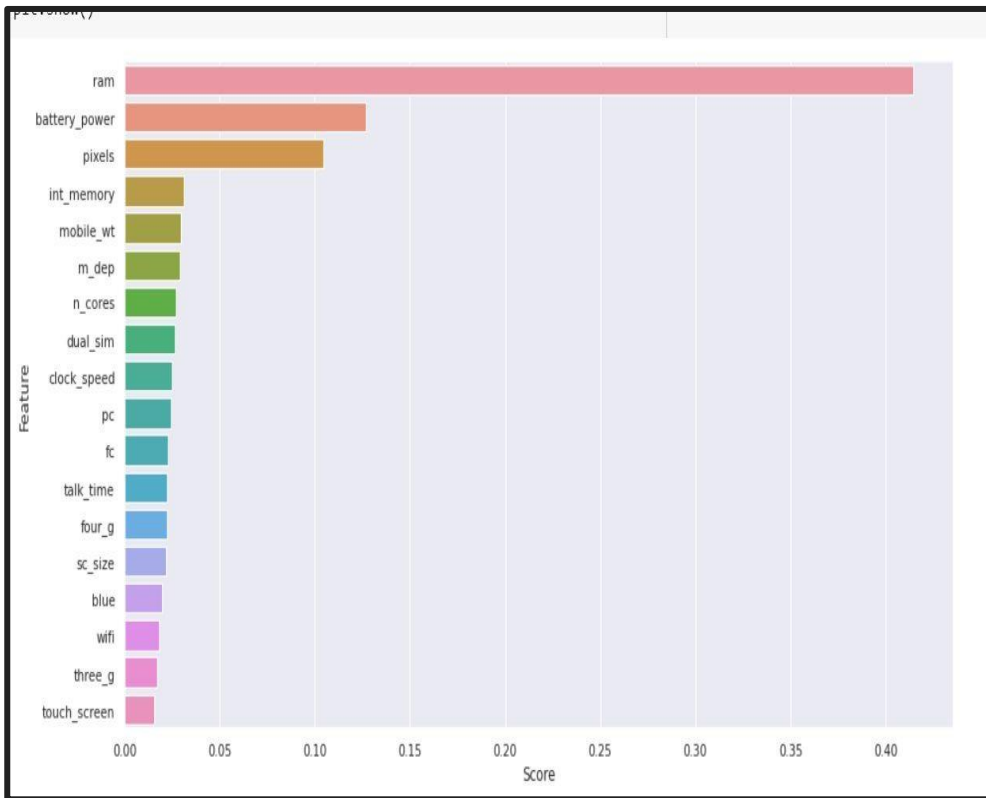


Decision Tree

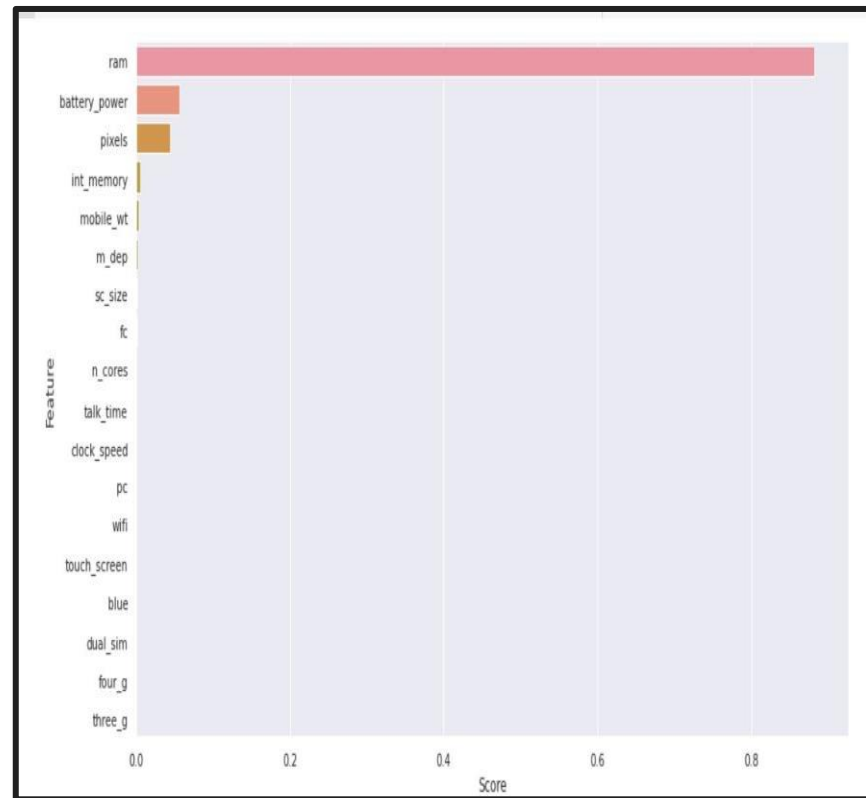


❖ Feature Importance

XGBoost



Gradient Boosting



❖ Challenges

- We performed “Hypothesis driven EDA” based on domain, but unlikely most of our hypothesis got rejected by our data.
- Most of the models are not able to get good accuracy for each class of target variable.
- We hit a ceiling at 94% accuracy using a single model.

Conclusions:

- We Started with Data understanding, data wrangling, basic EDA where we found the relationships, trends between price range and other independent variables.
- We selected the best features for predictive modeling by using K best feature selection method using Chi square statistic.
- Implemented various classification algorithms, out of which the SVM(Support vector machine) algorithm gave the best performance with 94% train accuracy and 90 % test accuracy.
- XG boost is the second best good model which gave good performance 98% train accuracy and 89% test accuracy score.
- KNN gave very worst model performance.
- We checked for the feature importance's of each model. RAM, Battery Power, Px_height and px_width,Screen size,mobile weight contributed the most while predicting the price range.

THANK YOU