

# Prediction of Purchase Amount for retail store customers

**Capstone 1- Springboard data science career track**

# Problem Statement:

---

**A retail Company wants to understand the customer purchase behaviour (specifically ,purchase amount) against various product of different categories during Black Friday sales.**

**Dataset contains the purchase summary of various customers for selected high volume products from last month.**

**The data set also contains customer demographics (age, gender, marital status, city\_type, stay\_in\_current\_city), product details (product\_id and product category) and Total purchase\_amount from last month.**

**Now, they want to build a model to predict the purchase amount of customer against various products which will help them to create personalized offer for customers against different products.**

# Dataset Information

---

The dataset contains about 537577 observations and 12 variable.

**FOLLOWING COLUMNS AND DATA TYPES ARE OBSERVED WHICH ARE CONVERTED TO SUITABLE DATATYPES FOR ML MODELS.**

COLUMN NAMES	DATATYPES
User_ID	537577 non-null int64
Product_ID	537577 non-null object
Gender	537577 non-null object
Age	537577 non-null object
Occupation	537577 non-null int64
City_Category	537577 non-null object
Stay_In_Current_City_Years	537577 non-null object
Marital_Status	537577 non-null int64
Product_Category_1	537577 non-null int64
Product_Category_2	370591 non-null float64
Product_Category_3	164278 non-null float64
Purchase	537577 non-null int64

# DATA WRANGLING

---

FOLLOWING COLUMNS AND DATA TYPES ARE OBSERVED WHICH ARE CONVERTED TO SUITABLE DATATYPES FOR ML MODELS.

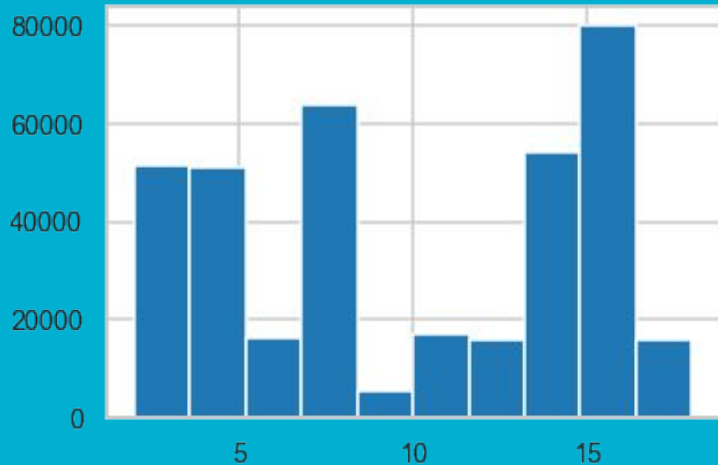
COLUMN NAMES	DATATYPES
User_ID	537577 non-null int64
Product_ID	537577 non-null object
Gender	537577 non-null object
Age	537577 non-null object
Occupation	537577 non-null int64
City_Category	537577 non-null object
Stay_In_Current_City_Years	537577 non-null object
Marital_Status	537577 non-null int64
Product_Category_1	537577 non-null int64

# MISSING VALUES

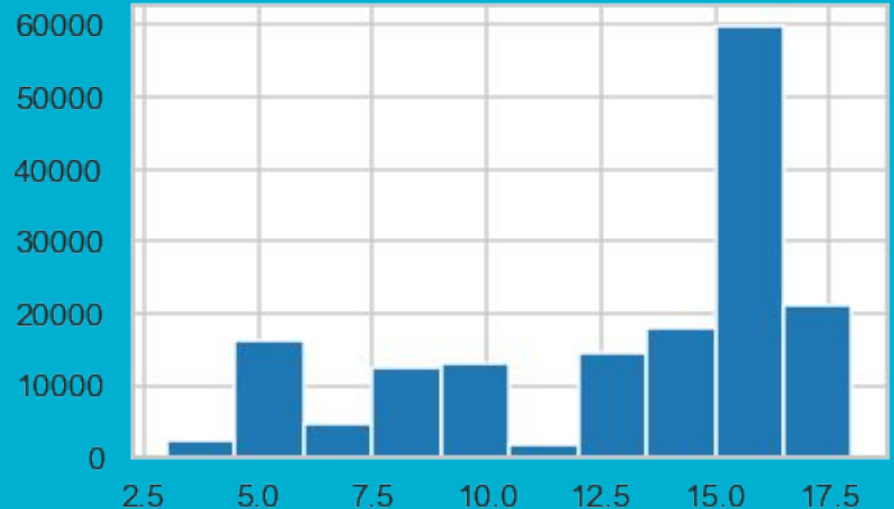
It is observed that there are missing values in both of the columns as mentioned below:

**histogram distribution for Product Category is non-normal distribution. However while imputing missing value ,we can replace the 'nan' value with 0 ,considering non-purchase of perticular item from the category.**

Product\_Category\_2      0.310627%

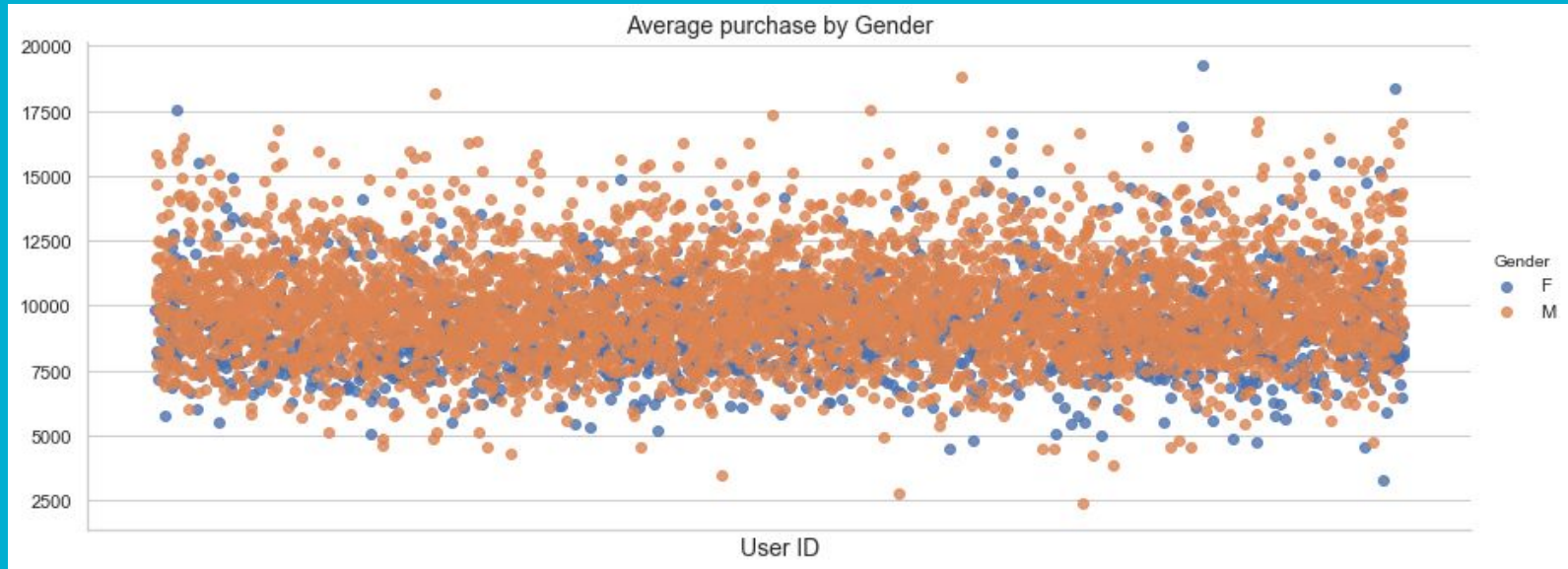


Product\_Category\_3      0.694410%



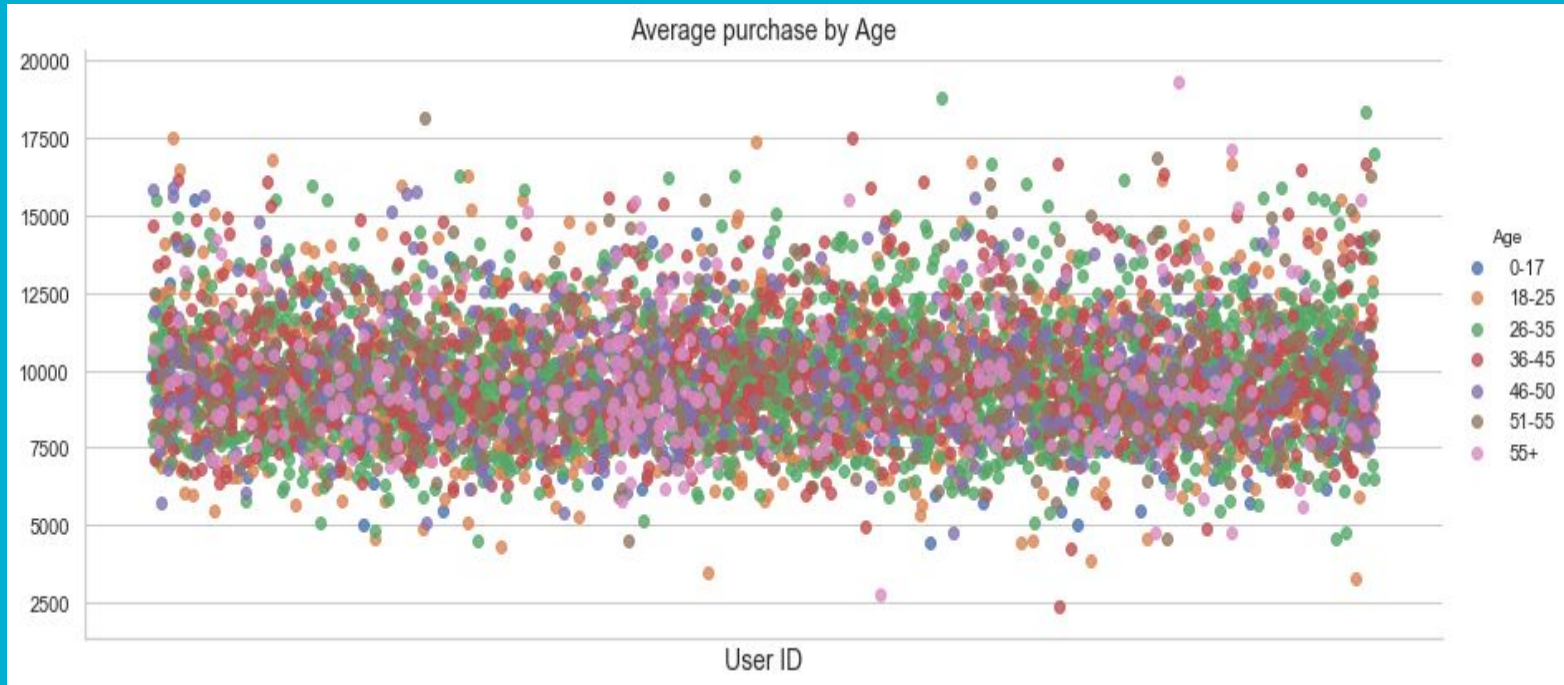
# DATA VISUALISATION

We can see the dominance of Male Buyers:



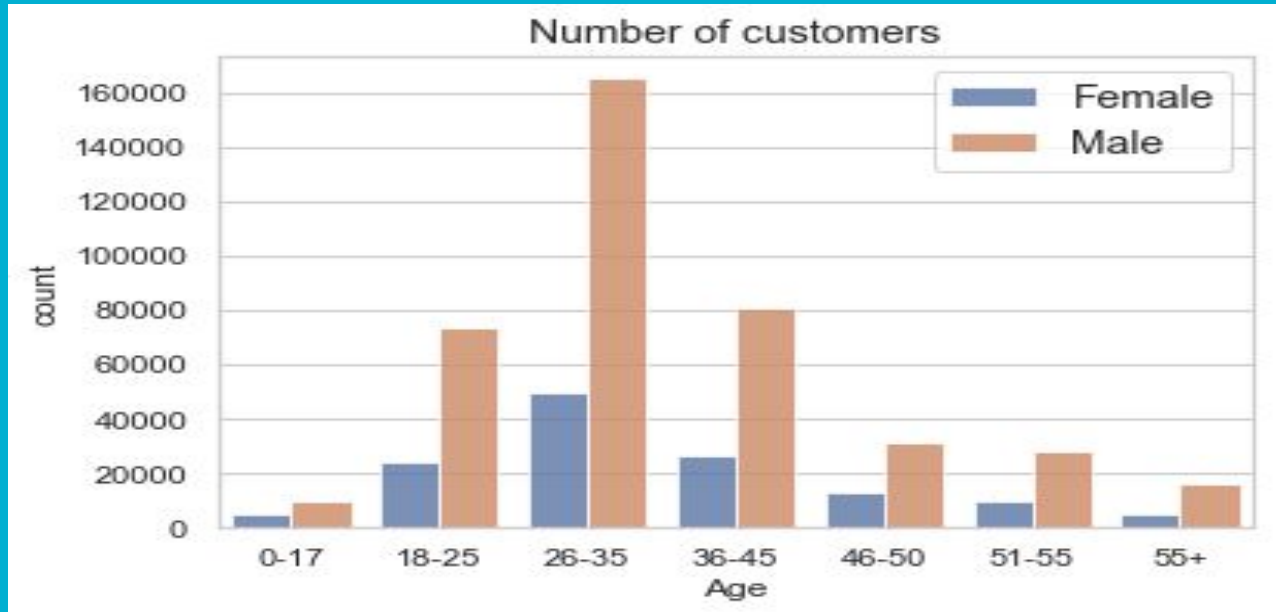
# Genderwise purchase distribution

---



# Gender Count between different Age Group

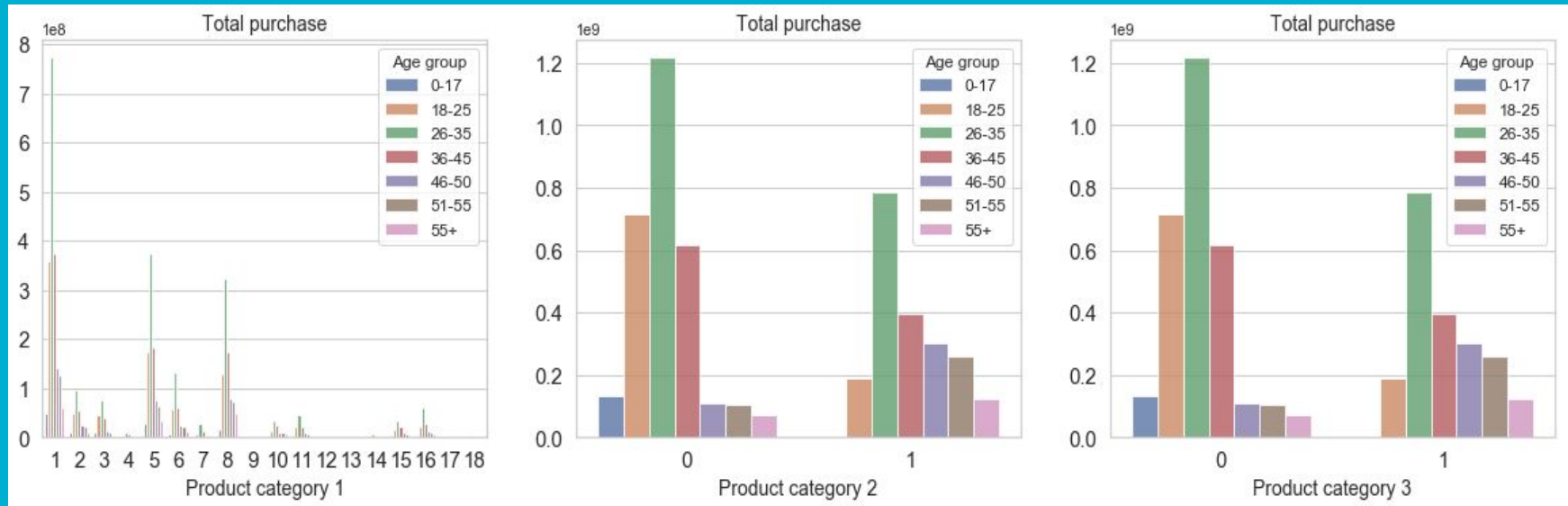
---







—We can observe that Items from Product Category1 are purchased more



# Observations based on Exploratory data analysis

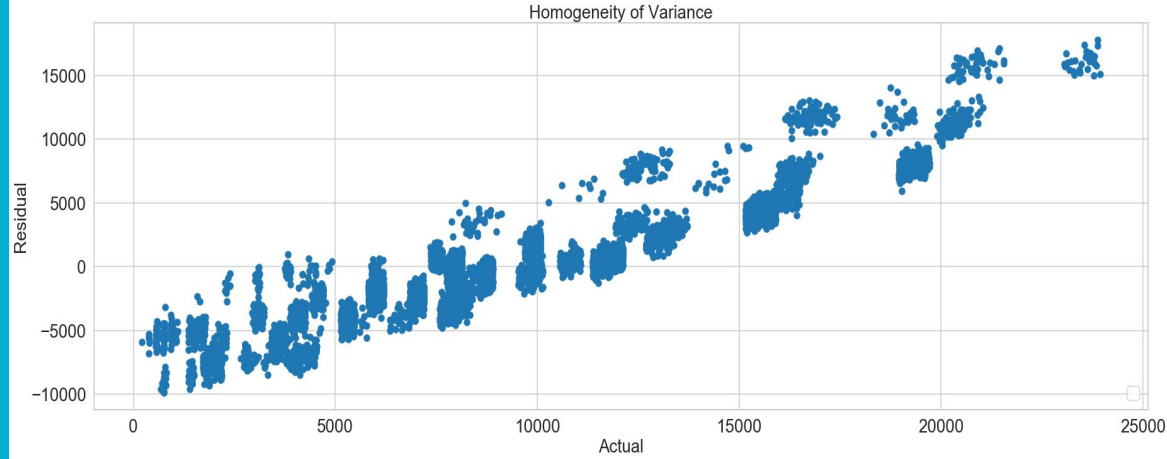
1. There are more numbers of male shoppers than female
2. The maximum selling item belongs to product category
3. More numbers of shoppers are single
4. The shoppers in the range 26-35
5. Customers from City B shopped the most
6. Customers who has resided in their city for 1 year shopped the most

# Machine Learning Models

---

**This is Supervised Model and we need to perform regression Analysis to predict Purchase Amount by Customer.**

**Based on below plot ,it is clear that data does not follow the linearity assumption such as homogeneity of variance which is essential for linear regression.**



# Tree Based Models

---

As our model fails to perform on Linear Regression with very low  $r^2$  of 0.107.

Its better to try tree based ensemble models.

I have tried following models:

1.Decision Tree

2.Random Forest

3.Gradientboosting Regressor

4.XGBoost

# ML Model building Process

---

As we have already done data preprocessing and data wrangling stage .

I have tried to do feature selection and removing constant feature if any.

I also stratified dataframe in order to remove any sampling biases while splitting

Dataframe using train\_test\_split.

Due to limiting computing capacity,the models are built on only 5% of data which

Accounts to be Train Data size -(18815,39)

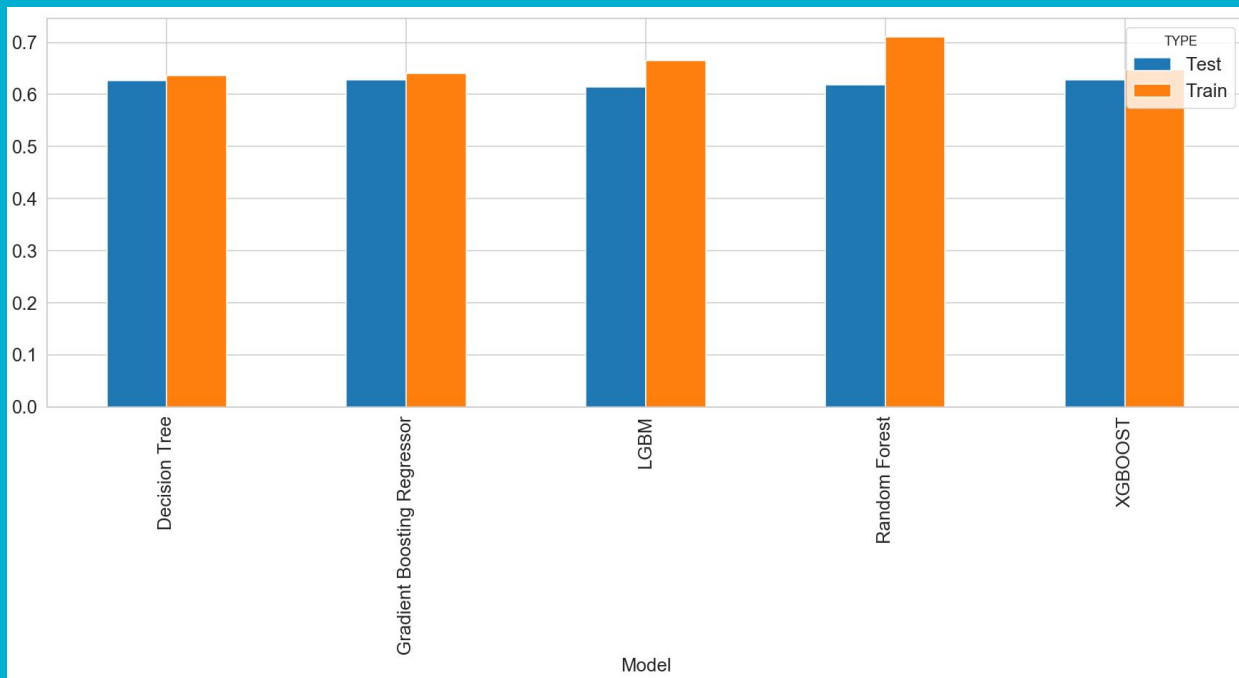
Test Data - (8064,39)

---

Gridsearch and hyperparametrs tuning are tried to achieve optimum model performance.

# R<sup>2</sup> Values of Various Algorithms:

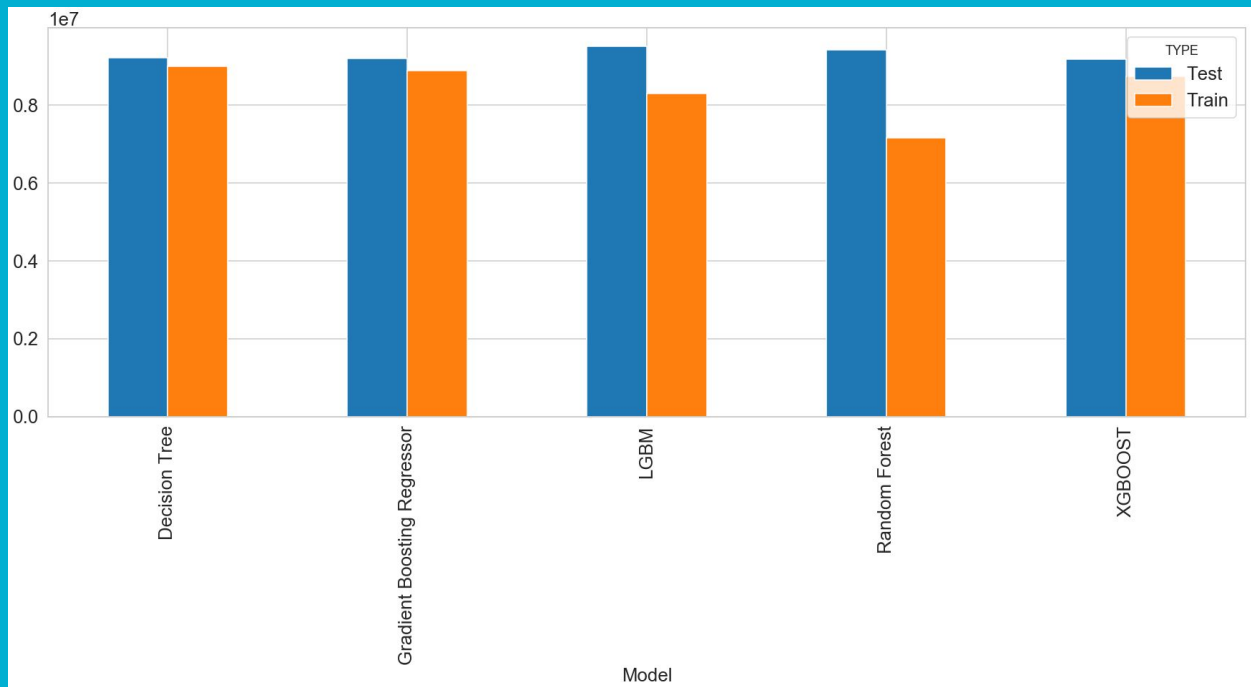
---





# MSE Score for various Algorithms:

---



# Conclusion

---

All the algorithms works equally though the  $R^2$  is not very good but we have learned

How the performance of ensembles is improved dramatically from 0.107  $r^2$  values to 0.64 approx. when compared with linear model.

Future work:

By working on whole dataset ,we may try to achieve more robust  $R^2$  and better performance by trying different combinations of hyperparameter values.