# Housing Price Project

Submitted by:

SNEHA SANTRA

## ABSTRACT

Houses are one of the necessary need of each and every person around the globe and therefore housing and real estate market is one of the markets which is one of the major contributors in the world's economy. It is a very large market and there are various companies working in the domain. Data science comes as a very important tool to solve problems in the domain to help the companies increase their overall revenue, profits, improving their marketing strategies and focusing on changing trends in house sales and purchases. Predictive modelling, Market mix modelling, recommendation systems are some of the machine learning techniques used for achieving the business goals for housing companies. Our problem is related to one such housing company.

# INTRODUCTION

- ## Business Problem Framing

  You are required to model the price of houses with the available independent variables. This model will then be used by the management to understand how exactly the prices vary with the variables. They can accordingly manipulate the strategy of the firm and concentrate on areas that will yield high returns. Further, the model will be a good way for the management to understand the pricing dynamics of a new market.

- ## Conceptual Background of the Domain Problem

  The project is sub-divided following section. These are:

  1. Loading necessary libraries

  2. Loading Dataset from a CSV file

  3. Summarization of Data to understand Dataset (Descriptive Statistics)

  4. Visualization of Data to understand Dataset (Plots, Graphs etc.)

  5. Processing the data for modeling

  6. skewness and outliers detection for better accuracy

  7. Build the model and select the right model and save it

- ## Review of Literature

  There are 81 columns including Mssubclass,Msxoining,lotarea,condition,lotshape,sale condition and main feature is sale price.

  There are 38 numerical features,17 discrete features,16 continious features and they are highly corelated with target variables.
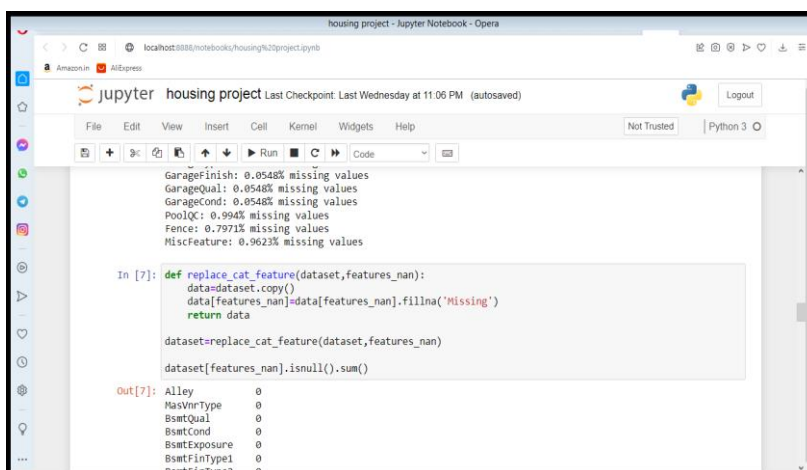
- ## Motivation for the Problem Undertaken

The company is looking at prospective properties to buy houses to enter the market. You are required to build a model using Machine Learning in order to predict the actual value of the prospective properties and decide whether to invest in them or not. For this company wants to know:

- Which variables are important to predict the price of variable?
- How do these variables describe the price of the house?

# <u>Analytical Problem Framing</u>

- ## Mathematical/ Analytical Modeling of the Problem
  Data Processing-There are some missing values in both test and train dataset so will fill it with nan values
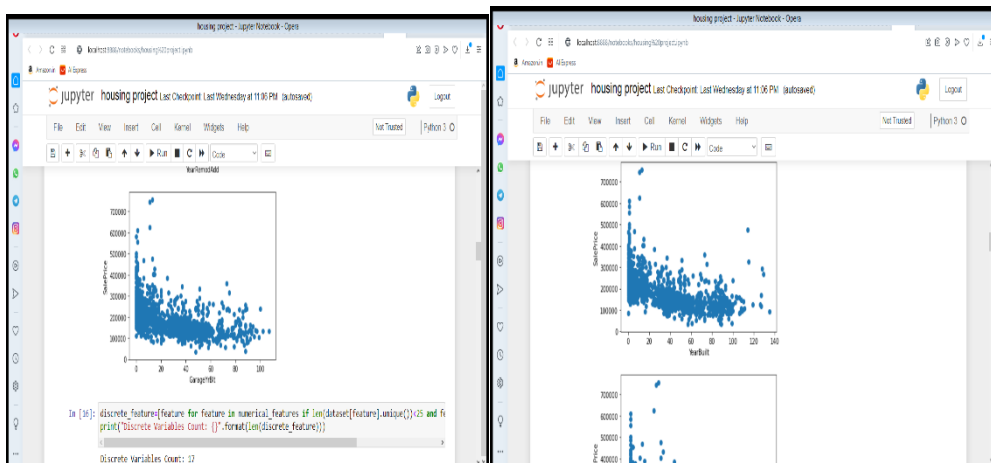
Now there is no missing values

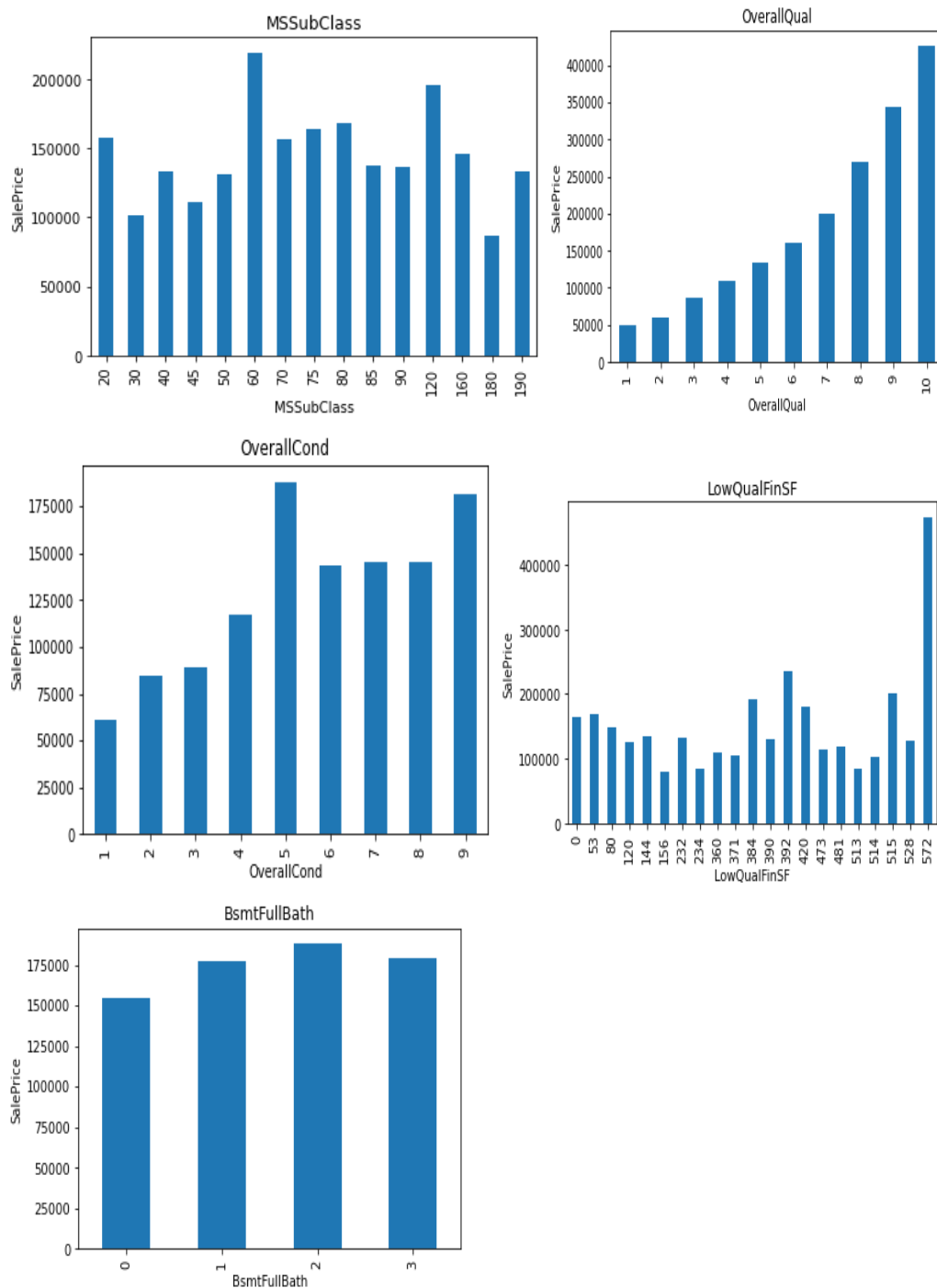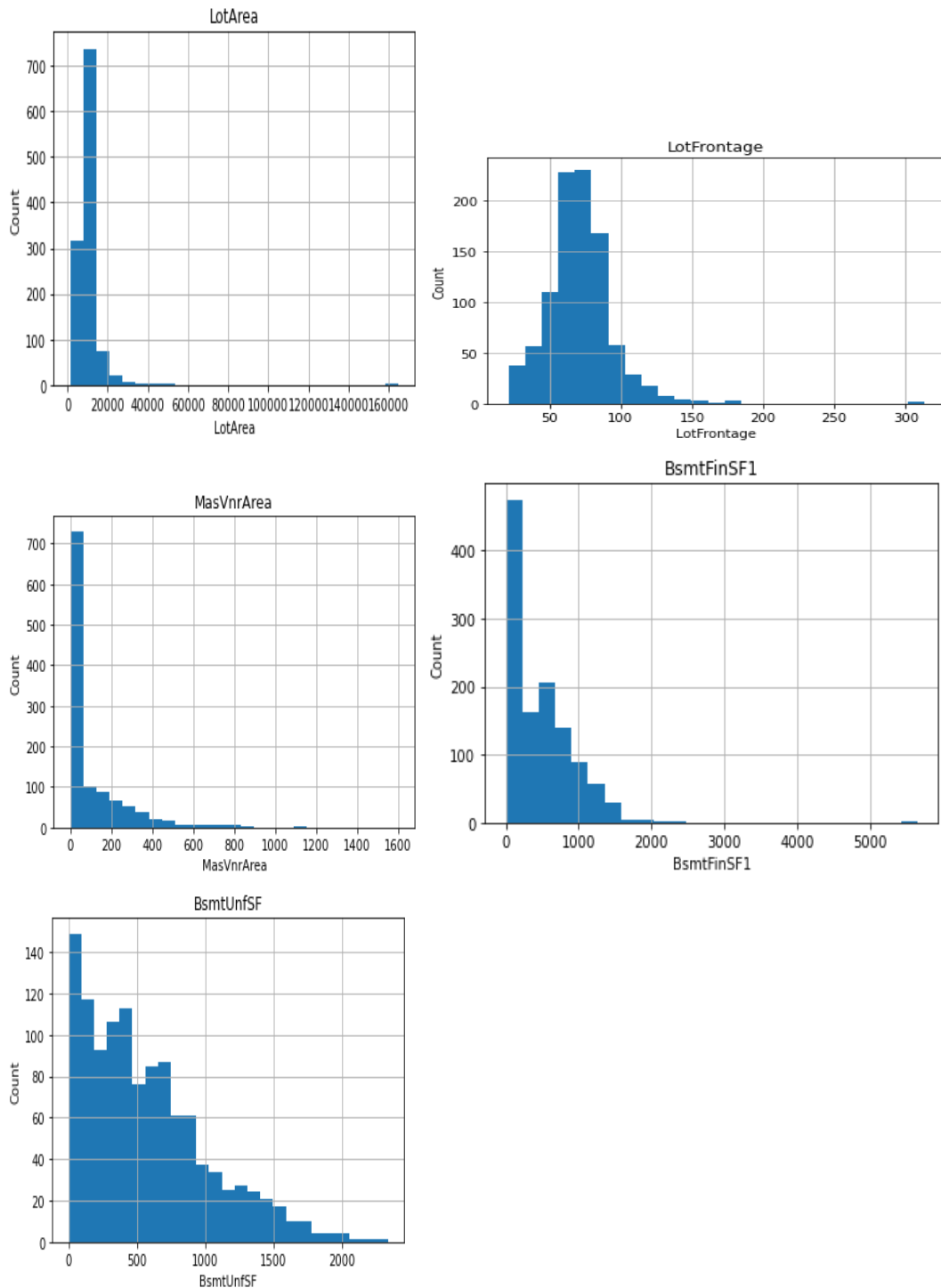And we drop some columns to scalling the dataset .

# EDA-



This is the relation between houseprice and yearsold.



This is the relation between saleprice vs yearbuild and garageYrBlt

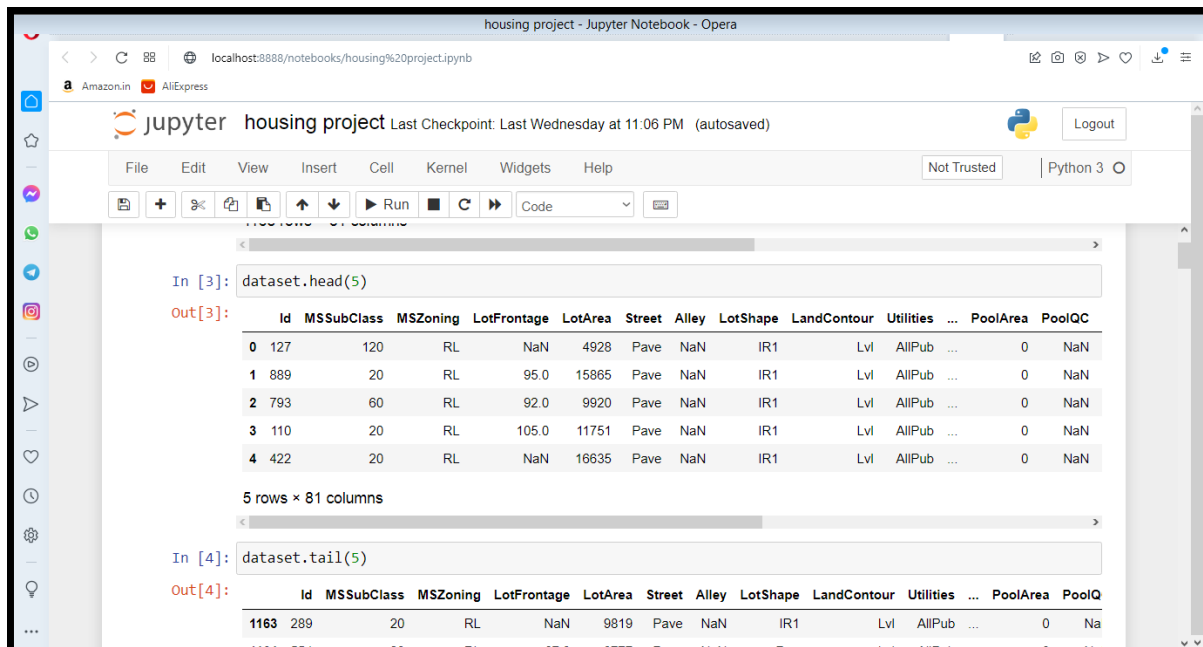This is the relation between salePrice and discrete features.

This is the plot/graph between count and continious features.

- ## Data Sources and their formats

A US-based housing company named **Surprise Housing** has decided to enter the Australian market. The company uses data analytics to purchase houses at a price below their actual values and flip them at a higher price. For the same purpose, the company has collected a data set from the sale of houses in Australia. The data is provided in the CSV file below.

This data consist of text and train data which is scrapped and There are 81 columns and 1460 entries. For the classification problem under

consideration we have used some object and some float and some int columns for input variable and our output column is **sale price and sale condition** that is also object type.
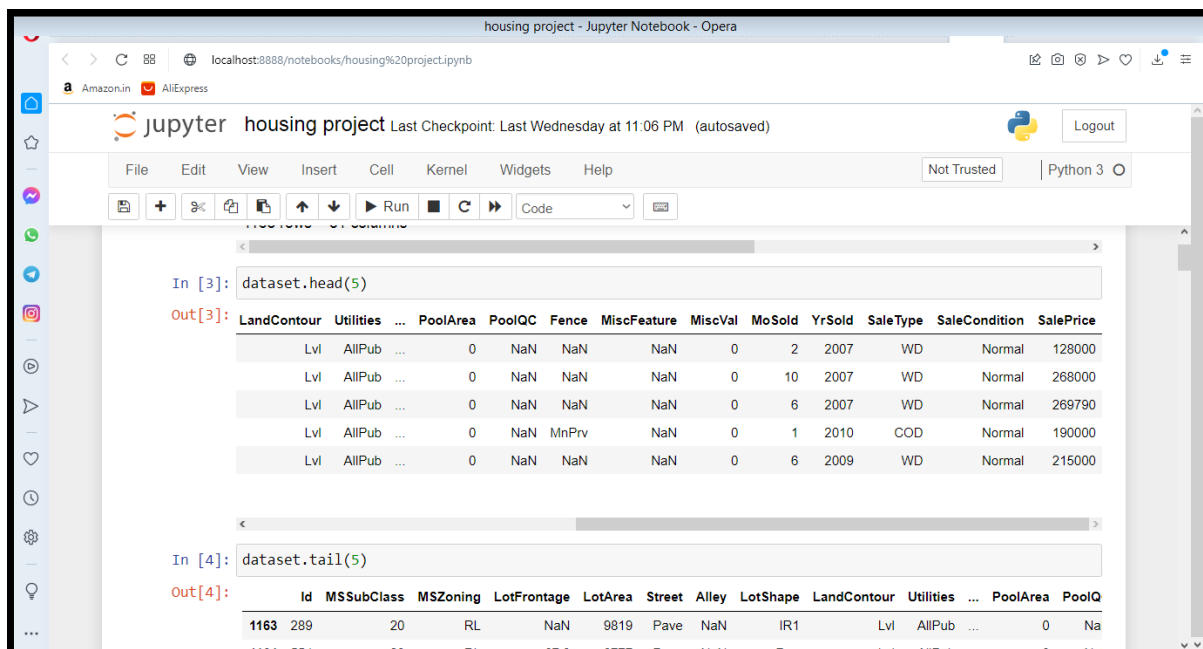




- Data Preprocessing Done

What were the steps followed for the cleaning of the data? What were the assumptions done and what were the next actions steps over that?
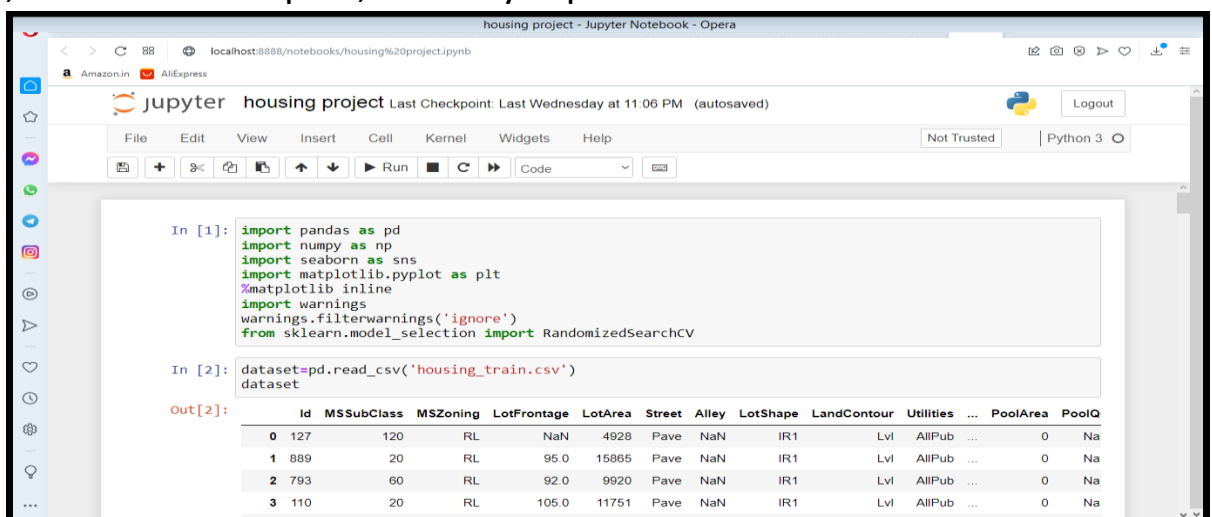
- Data Inputs- Logic- Output Relationships

Describe the relationship behind the data input, its format, the logic in between and the output. Describe how the input affects the output.

- State the set of assumptions (if any) related to the problem under consideration

Here, you can describe any presumptions taken by you.

- Hardware and Software Requirements and Tools Used

Here we use lots of liaberies like pandas,numpy,matplot,seaborn, and we use python language for the coding purpose and import some other metrics liaberies also for model building like sklearn metrics ,classification report,accuracy report etc.
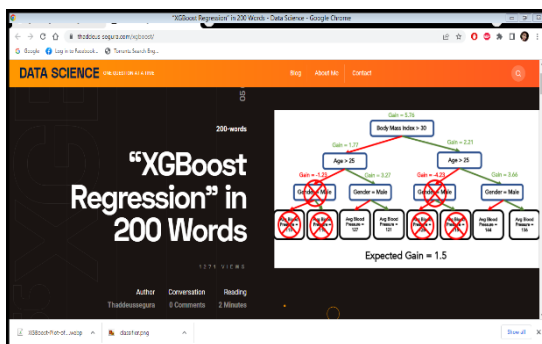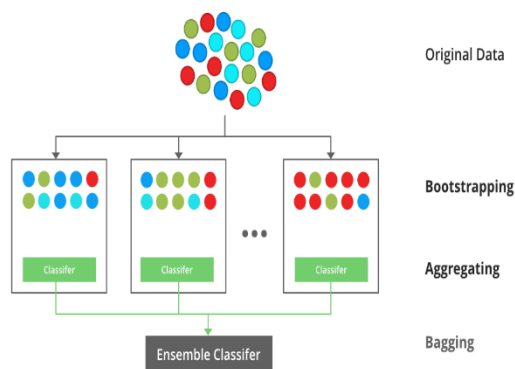


# Model/s Development and Evaluation

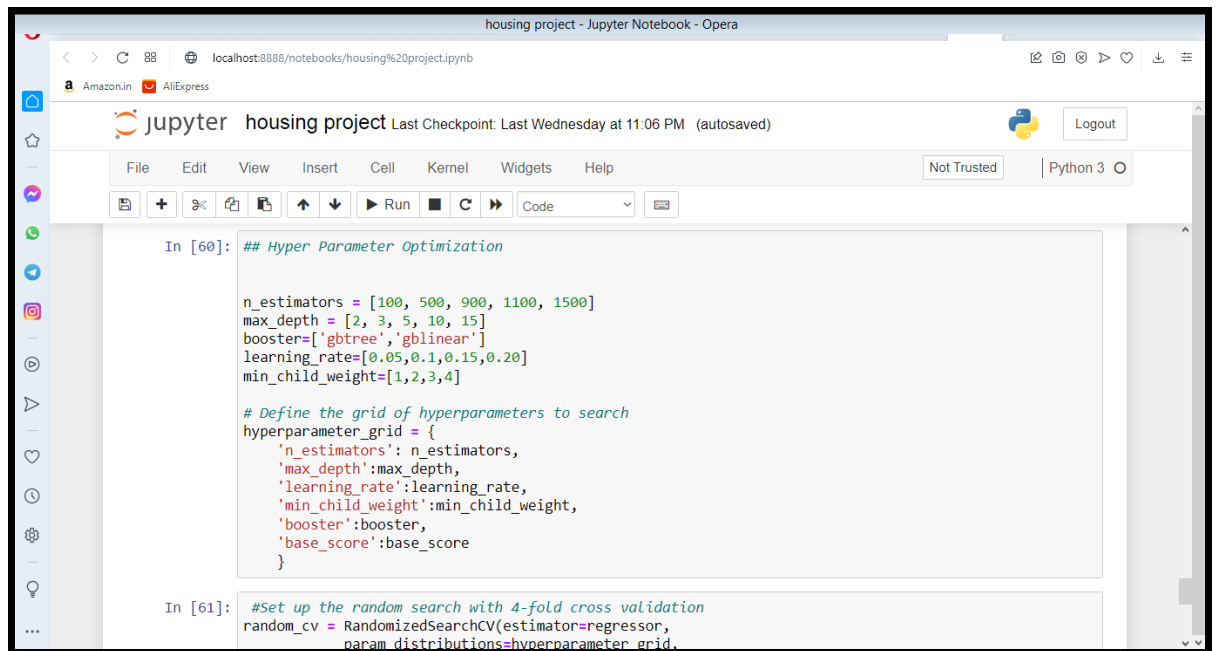- Identification of possible problem-solving approaches (methods)

  We use xgboost classifier and regressor for the model building as it is lassoregression problem and finally we save the classifier model for gd accuracy.

- Testing of Identified Approaches (Algorithms)

  We use xgboost classifier and regressor





- Run and Evaluate selected models

This is hyper parameter optimization and 4-fold cross validation for set up the random search and we choose the best parameter

- Key Metrics for success in solving problem under consideration

    We choose xgboost classifier as the accuracy of the hyper parameter accuracy of classifier is best .



# CONCLUSION

- Key Findings and Conclusions of the Study

After the Final Submission of test data, my accuracy score was  90%

Feature engineering helped me increase my accuracy.

Amazingly xgboost classifier worked better than all other Ensemble models.

- ### Learning Outcomes of the Study in respect of Data Science

Data science comes as a very important tool to solve problems in the domain to help the companies increase their overall revenue, profits, improving their marketing strategies and focusing on changing trends in house sales and purchases. Predictive modelling, Market mix modelling, recommendation systems are some of the machine learning techniques used for achieving the business goals for housing companies. Our problem is related to one such housing company.