**FLIP ROBO**

# HOUSING PRICE PREDICTION MODEL

Submitted by:

BHAVNA PIPLANI

# **ACKNOWLEDGMENT**

This includes mentioning of all the references, research papers, data sources, professionals and other resources that helped and guided me in completion of the project.

# INTRODUCTION

## A. Problem Framing:

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.

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 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?

Description of the domain related concepts that will be useful for better understanding of the project.

Data Cleaning, Data Visualisation using different plotting methods like barplot, countplot, scatterplot, Data pre-processing using LabelEncoder, StandardScalar and Model Training

# B. Analytical Problem Framing

## a) Mathematical/ Analytical Modelling of the Problem

**Statistical modelling** is the process of applying statistical analysis to a dataset. A statistical model is a mathematical representation (or mathematical model) of observed data.

When data analysts apply various statistical models to the data they are investigating, they are able to understand and interpret the information more strategically. Rather than sifting through the raw data, this practice allows them to identify relationships between variables, make predictions about future sets of data, and visualize that data so that non-analysts and stakeholders can consume and leverage it.

most common techniques will fall into the following two groups:

- **Supervised learning**, including regression and classification models.
- **Unsupervised learning**, including clustering algorithms and association rules.

Some of the most common regression models include **LinearRegression,Lasso,Ridge,RandomForestRegressor, KNeighborsRegressor,GradientBoostingRegressor, AdaBoostRegressor**.

## b)  Data Sources and their formats

Data Source file was given in csv format with all the necessary variables for further Data Cleaning, Data pre-processing  and Model Training. Data Description can be seen in the following table:

| Variable | Definition |
|---|---|
| MSSubClass | Identifies the type of dwelling involved in the sale.<br>20     1-STORY 1946 & NEWER ALL STYLES<br>30     1-STORY 1945 & OLDER<br>40     1-STORY W/FINISHED ATTIC ALL AGES<br>45     1-1/2 STORY - UNFINISHED ALL AGES<br>50     1-1/2 STORY FINISHED ALL AGES<br>60     2-STORY 1946 & NEWER<br>70     2-STORY 1945 & OLDER<br>75     2-1/2 STORY ALL AGES<br>80     SPLIT OR MULTI-LEVEL<br>85     SPLIT FOYER<br>90     DUPLEX - ALL STYLES AND AGES<br>120     1-STORY PUD (Planned Unit Development) - 1946 & NEWE<br>150     1-1/2 STORY PUD - ALL AGES<br>160     2-STORY PUD - 1946 & NEWER<br>180     PUD - MULTILEVEL - INCL SPLIT LEV/FOYER<br>190     2 FAMILY CONVERSION - ALL STYLES AND AGES |
| MSZoning | Identifies the general zoning classification of the sale.<br>A Agriculture<br>C Commercial<br>FV     Floating Village Residential<br>I  Industrial<br>RH     Residential High Density<br>RL     Residential Low Density |

| | |
|---|---|
| | RP        Residential Low Density Park<br>RM        Residential Medium Density |
| LotFrontage | Linear feet of street connected to property |
| LotArea | Lot size in square feet |
| Street | Type of road access to property<br>   Grvl        Gravel<br>   Pave        Paved |
| Alley | Type of alley access to property<br>   Grvl        Gravel<br>   Pave        Paved<br>   NA           No alley access |
| LotShape | General shape of property<br><br>   Reg        Regular<br>   IR1        Slightly irregular<br>   IR2        Moderately Irregular<br>   IR3        Irregular |
| LandContour | Flatness of the property<br><br>   Lvl          Near Flat/Level<br>   Bnk         Banked - Quick and significant rise<br>                   from street grade    to building<br>   HLS         Hillside - Significant slope from side to side<br>   Low         Depression |
| Utilities | Type of utilities available<br><br>   AllPub         All public Utilities (E,G,W,& S)<br>   NoSewr        Electricity, Gas, and Water (Septic Tank)<br>   NoSeWa        Electricity and Gas Only<br>   ELO             Electricity only |
| LotConfig | Lot configuration<br><br>   Inside         Inside lot<br>   Corner         Corner lot<br>   CulDSac       Cul-de-sac<br>   FR2             Frontage on 2 sides of property<br>   FR3             Frontage on 3 sides of property |
| LandSlope | Slope of property<br><br>   Gtl             Gentle slope |

| | | |
|---|---|---|
| | Mod | Moderate Slope |
| | Sev | Severe Slope |
| Neighborhood | Physical locations within Ames city limits | |
| | Blmngtn | Bloomington Heights |
| | Blueste | Bluestem |
| | BrDale | Briardale |
| | BrkSide | Brookside |
| | ClearCr | Clear Creek |
| | CollgCr | College Creek |
| | Crawfor | Crawford |
| | Edwards | Edwards |
| | Gilbert | Gilbert |
| | IDOTRR | Iowa DOT and Rail Road |
| | MeadowV | Meadow Village |
| | Mitchel | Mitchell |
| | Names | North Ames |
| | NoRidge | Northridge |
| | NPkVill | Northpark Villa |
| | NridgHt | Northridge Heights |
| | NWAmes | Northwest Ames |
| | OldTown | Old Town |
| | SWISU | South & West of Iowa State University |
| | Sawyer | Sawyer |
| | SawyerW | Sawyer West |
| | Somerst | Somerset |
| | StoneBr | Stone Brook |
| | Timber | Timberland |
| | Veenker | Veenker |
| Condition1 | Proximity to various conditions | |
| | Artery | Adjacent to arterial street |
| | Feedr | Adjacent to feeder street |
| | Norm | Normal |
| | RRNn | Within 200' of North-South Railroad |
| | RRAn | Adjacent to North-South Railroad |
| | PosN | Near positive off-site feature--park, greenbelt, etc. |
| | PosA | Adjacent to postive off-site feature |
| | RRNe | Within 200' of East-West Railroad |
| | RRAe | Adjacent to East-West Railroad |

| | |
|---|---|
| Condition2 | Proximity to various conditions (if more than one is present)<br><br>Artery    Adjacent to arterial street<br>Feedr    Adjacent to feeder street<br>Norm    Normal<br>RRNn    Within 200' of North-South Railroad<br>RRAn    Adjacent to North-South Railroad<br>PosN    Near positive off-site feature--park, greenbelt, etc.<br>PosA    Adjacent to postive off-site feature<br>RRNe    Within 200' of East-West Railroad<br>RRAe    Adjacent to East-West Railroad |
| BldgType | Type of dwelling<br><br>1Fam    Single-family Detached<br>2FmCon.    Two-family Conversion; originally<br>        built as one-family dwelling<br>Duplx    Duplex<br>TwnhsE    Townhouse End Unit<br>TwnhsI    Townhouse Inside Unit |
| HouseStyle | Style of dwelling<br><br>1Story    One story<br>1.5Fin    One and one-half story: 2nd level finished<br>1.5Unf    One and one-half story: 2nd level unfinished<br>2Story    Two story<br>2.5Fin    Two and one-half story: 2nd level finished<br>2.5Unf    Two and one-half story: 2nd level unfinished<br>SFoyer    Split Foyer<br>SLvl    Split Level |
| OverallQual | Rates the overall material and finish of the house<br><br>10    Very Excellent<br>9    Excellent<br>8    Very Good<br>7    Good<br>6    Above Average<br>5    Average<br>4    Below Average<br>3    Fair<br>2    Poor<br>1    Very Poor |

| OverallCond | Rates the overall condition of the house |
|---|---|
| | 10     Very Excellent <br> 9     Excellent <br> 8     Very Good <br> 7     Good <br> 6     Above Average <br> 5     Average <br> 4     Below Average <br> 3     Fair <br> 2     Poor <br> 1     Very Poor |
| YearBuilt | Original construction date |
| YearRemodAdd | Remodel date (same as construction date if no remodeling or additions) |
| RoofStyle | Type of roof <br><br> Flat        Flat <br> Gable       Gable <br> Gambrel    Gabrel (Barn) <br> Hip          Hip <br> Mansard    Mansard <br> Shed        Shed |
| RoofMatl | Roof material <br><br> ClyTile      Clay or Tile <br> CompSh.    Standard (Composite) Shingle <br> Membran    Membrane <br> Metal        Metal <br> Roll         Roll <br> Tar&Grv     Gravel & Tar <br> WdShake    Wood Shakes <br> WdShngl    Wood Shingles |
| Exterior1st | Exterior covering on house <br><br> AsbShng      Asbestos Shingles <br> AsphShn      Asphalt Shingles <br> BrkComm     Brick Common <br> BrkFace       Brick Face <br> CBlock        Cinder Block |

| | |
|---|---|
| | CemntBd       Cement Board<br>HdBoard       Hard Board<br>ImStucc       Imitation Stucco<br>MetalSd       Metal Siding<br>Other       Other<br>Plywood       Plywood<br>PreCast       PreCast<br>Stone       Stone<br>Stucco       Stucco<br>VinylSd       Vinyl Siding<br>Wd Sdng       Wood Siding<br>WdShing       Wood Shingles |
| Exterior2nd | Exterior covering on house (if more than one material)<br><br>AsbShng       Asbestos Shingles<br>AsphShn       Asphalt Shingles<br>BrkComm       Brick Common<br>BrkFace       Brick Face<br>CBlock       Cinder Block<br>CemntBd       Cement Board<br>HdBoard       Hard Board<br>ImStucc       Imitation Stucco<br>MetalSd       Metal Siding<br>Other       Other<br>Plywood       Plywood<br>PreCast       PreCast<br>Stone       Stone<br>Stucco       Stucco<br>VinylSd       Vinyl Siding<br>Wd Sdng       Wood Siding<br>WdShing       Wood Shingles |
| MasVnrType | Masonry veneer type<br><br>BrkCmn       Brick Common<br>BrkFace       Brick Face<br>CBlock       Cinder Block<br>None       None<br>Stone       Stone |
| MasVnrArea | Masonry veneer area in square feet |
| ExterQual | Evaluates the quality of the material on the exterior |

| | |
|---|---|
| | Ex         Excellent<br>Gd         Good<br>TA        Average/Typical<br>Fa         Fair<br>Po        Poor |
| ExterCond | Evaluates the present condition of the material on the exterior<br><br>Ex         Excellent<br>Gd         Good<br>TA        Average/Typical<br>Fa         Fair<br>Po        Poor |
| Foundation | Type of foundation<br><br>BrkTil    Brick & Tile<br>CBlock   Cinder Block<br>PConc    Poured Contrete<br>Slab       Slab<br>Stone    Stone<br>Wood    Wood |
| BsmtQual | Evaluates the height of the basement<br><br>Ex         Excellent (100+ inches)<br>Gd         Good (90-99 inches)<br>TA        Typical (80-89 inches)<br>Fa         Fair (70-79 inches)<br>Po        Poor (<70 inches<br>NA       No Basement |
| BsmtCond | Evaluates the general condition of the basement<br><br>Ex         Excellent<br>Gd         Good<br>TA        Typical - slight dampness allowed<br>Fa         Fair - dampness or some cracking or settling<br>Po        Poor - Severe cracking, settling, or wetness<br>NA       No Basement |
| BsmtExposure | Refers to walkout or garden level walls<br><br>Gd         Good Exposure |

| | | |
|---|---|---|
| | Av | Average Exposure (split levels or foyers typically score average or above) |
| | Mn | Mimimum Exposure |
| | No | No Exposure |
| | NA | No Basement |
| BsmtFinType1 | Rating of basement finished area | |
| | GLQ | Good Living Quarters |
| | ALQ | Average Living Quarters |
| | BLQ | Below Average Living Quarters |
| | Rec | Average Rec Room |
| | LwQ | Low Quality |
| | Unf | Unfinshed |
| | NA | No Basement |
| BsmtFinSF1 | Type 1 finished square feet | |
| BsmtFinType2 | Rating of basement finished area (if multiple types) | |
| | GLQ | Good Living Quarters |
| | ALQ | Average Living Quarters |
| | BLQ | Below Average Living Quarters |
| | Rec | Average Rec Room |
| | LwQ | Low Quality |
| | Unf | Unfinished |
| | NA | No Basement |
| BsmtFinSF2 | Type 2 finished square feet | |
| BsmtUnfSF | Unfinished square feet of basement area | |
| TotalBsmtSF | Total square feet of basement area | |
| Heating | Type of heating | |
| | Floor | Floor Furnace |
| | GasA | Gas forced warm air furnace |
| | GasW | Gas hot water or steam heat |
| | Grav | Gravity furnace |
| | OthW | Hot water or steam heat other than gas |
| | Wall | Wall furnace |
| HeatingQC | Heating quality and condition | |
| | Ex | Excellent |
| | Gd | Good |
| | TA | Average/Typical |

| | |
|---|---|
| | Fa          Fair<br>Po          Poor |
| CentralAir | Central air conditioning<br><br>N     No<br>Y     Yes |
| Electrical | Electrical system<br><br>SBrkr    Standard Circuit Breakers & Romex<br>FuseA    Fuse Box over 60 AMP and all Romex<br>              wiring (Average)<br>FuseF    60 AMP Fuse Box and mostly Romex<br>              wiring (Fair)<br>FuseP    60 AMP Fuse Box and mostly knob &<br>              tube wiring (poor)<br>Mix       Mixed |
| 1stFlrSF | First Floor square feet |
| 2ndFlrSF | Second floor square feet |
| LowQualFinSF | Low quality finished square feet (all floors) |
| GrLivArea | Above grade (ground) living area square feet |
| BsmtFullBath | Basement full bathrooms |
| BsmtHalfBath | Basement half bathrooms |
| FullBath | Full bathrooms above grade |
| HalfBath | Half baths above grade |
| Bedroom | Bedrooms above grade (does NOT include basement bedrooms) |
| Kitchen | Kitchens above grade |
| KitchenQual | Kitchen quality<br><br>Ex     Excellent<br>Gd    Good<br>TA    Typical/Average<br>Fa     Fair<br>Po     Poor |
| TotRmsAbvGrd | Total rooms above grade (does not include bathrooms) |
| Functional | Home functionality (Assume typical unless deductions are warranted)<br><br>Typ        Typical Functionality<br>Min1      Minor Deductions 1<br>Min2   Minor Deductions 2<br>Mod     Moderate Deductions |

| | | |
|---|---|---|
| | Maj1 | Major Deductions 1 |
| | Maj2 | Major Deductions 2 |
| | Sev | Severely Damaged |
| | Sal | Salvage only |
| Fireplaces | Number of fireplaces | |
| FireplaceQu | Fireplace quality | |
| | Ex | Excellent - Exceptional Masonry Fireplace |
| | Gd | Good - Masonry Fireplace in main level |
| | TA | Average - Prefabricated Fireplace in main living area or Masonry Fireplace in basement |
| | Fa | Fair - Prefabricated Fireplace in basement |
| | Po | Poor - Ben Franklin Stove |
| | NA | No Fireplace |
| GarageType | Garage location | |
| | 2Types | More than one type of garage |
| | Attchd | Attached to home |
| | Basment | Basement Garage |
| | BuiltIn | Built-In (Garage part of house – typically has room above garage) |
| | CarPort | Car Port |
| | Detchd | Detached from home |
| | NA | No Garage |
| GarageYrBlt | Year garage was built | |
| GarageFinish | Interior finish of the garage | |
| | Fin | Finished |
| | RFn | Rough Finished |
| | Unf | Unfinished |
| | NA | No Garage |
| GarageCars | Size of garage in car capacity | |
| GarageArea | Size of garage in square feet | |
| GarageQual | Garage quality | |
| | Ex | Excellent |
| | Gd | Good |
| | TA | Typical/Average |
| | Fa | Fair |
| | Po | Poor |

| | |
|---|---|
| | NA      No Garage |
| GarageCond | Garage condition<br><br>Ex      Excellent<br>Gd      Good<br>TA      Typical/Average<br>Fa      Fair<br>Po      Poor<br>NA      No Garage |
| PavedDrive | Paved driveway<br><br>Y Paved<br>P Partial Pavement<br>N Dirt/Gravel |
| WoodDeckSF | Wood deck area in square feet |
| OpenPorchSF | Open porch area in square feet |
| EnclosedPorch | Enclosed porch area in square feet |
| 3SsnPorch | Three season porch area in square feet |
| ScreenPorch | Screen porch area in square feet |
| PoolArea | Pool area in square feet |
| PoolQC | Pool quality<br><br>Ex      Excellent<br>Gd      Good<br>TA      Average/Typical<br>Fa      Fair<br>NA      No Pool |
| Fence | Fence quality<br><br>GdPrv  Good Privacy<br>MnPrv  Minimum Privacy<br>GdWo   Good Wood<br>MnWw Minimum Wood/Wire<br>NA      No Fence |
| MiscFeature | Miscellaneous feature not covered in other categories<br><br>Elev    Elevator<br>Gar2    2nd Garage (if not described in garage section)<br>Othr    Other<br>Shed    Shed (over 100 SF) |

|  |  |
|---|---|
|  | TenC  Tennis Court<br>NA    None |
| MiscVal | $Value of miscellaneous feature |
| MoSold | Month Sold (MM) |
| YrSold | Year Sold (YYYY) |
| SaleType | Type of sale<br><br>WD    Warranty Deed - Conventional<br>CWD   Warranty Deed - Cash<br>VWD   Warranty Deed - VA Loan<br>New    Home just constructed and sold<br>COD   Court Officer Deed/Estate<br>Con    Contract 15% Down payment regular terms<br>ConLw Contract Low Down payment and low interest<br>ConLI  Contract Low Interest<br>ConLD Contract Low Down<br>Oth    Other |
| SaleCondition | Condition of sale<br><br>Normal Normal Sale<br>Abnorml    Abnormal Sale - trade, foreclosure, short sale<br>AdjLand    Adjoining Land Purchase<br>Alloca     Allocation - two linked properties with<br>           separate deeds, typically condo with a garage unit<br>Family     Sale between family members<br>Partial    Home was not completed when<br>           last assessed (associated with New Homes) |

# c)   Explanatory Data Analysis

Hardware and Software Requirements and Tools Used

Anaconda Software – Jupiter Notebook

```
1  import pandas as pd
2  import numpy as np
3  import seaborn as sns
4  import matplotlib.pyplot as plt
5  import warnings
6  warnings.simplefilter("ignore")
```

Importing the dataset from the source file provided.

```
1  #importing train dataset
2  ds=pd.read_csv("housing_train.csv",sep='\t')
3  ds.head(5)
```

| | Id | MSSubClass | MSZoning | LotFrontage | LotArea | Street | Alley | LotShape | LandContour | Utilities | ... | PoolArea | PoolQC |
|---|-----|------------|----------|-------------|---------|--------|-------|----------|-------------|-----------|-----|----------|--------|
| 0 | 127 | 120 | RL | NaN | 4928 | Pave | NaN | IR1 | Lvl | AllPub | ... | 0 | NaN |
| 1 | 889 | 20 | RL | 95.0 | 15865 | Pave | NaN | IR1 | Lvl | AllPub | ... | 0 | NaN |
| 2 | 793 | 60 | RL | 92.0 | 9920 | Pave | NaN | IR1 | Lvl | AllPub | ... | 0 | NaN |
| 3 | 110 | 20 | RL | 105.0 | 11751 | Pave | NaN | IR1 | Lvl | AllPub | ... | 0 | NaN |
| 4 | 422 | 20 | RL | NaN | 16635 | Pave | NaN | IR1 | Lvl | AllPub | ... | 0 | NaN |

After importing the dataset , checking the datatypes of each column.

```
1  ds.dtypes #checking the datatypes of variables
```

```
Id                    int64
MSSubClass            int64
MSZoning             object
LotFrontage         float64
LotArea               int64
                     ...
MoSold                int64
YrSold                int64
SaleType             object
SaleCondition        object
SalePrice             int64
Length: 81, dtype: object
```

Train dataset has 1168 rows and 81 columns

```
1  ds.shape  # checking no of columns and rows
```

```
(1168, 81)
```

After checking the no. of rows and columns count, we will check the null values if any, column count, datatypes of columns.

```
1  dss=ds.columns[ds.isnull().any()]  # extracting all the variables having null values
2  dss
```

```
Index(['LotFrontage', 'Alley', 'MasVnrType', 'MasVnrArea', 'BsmtQual',
       'BsmtCond', 'BsmtExposure', 'BsmtFinType1', 'BsmtFinType2',
       'FireplaceQu', 'GarageType', 'GarageYrBlt', 'GarageFinish',
       'GarageQual', 'GarageCond', 'PoolQC', 'Fence', 'MiscFeature'],
      dtype='object')
```

The dataset is not clean. Data is missing in many of the features . There are missing or null values in the dataset. Nan values will be replaced by values mentioned in data description provided.

After data cleaning, we will see the correlation of features with the target variable.

```
plt.figure(figsize=(50,50))
sns.heatmap(ds.corr(),annot=True,cmap='YlGnBu')
```

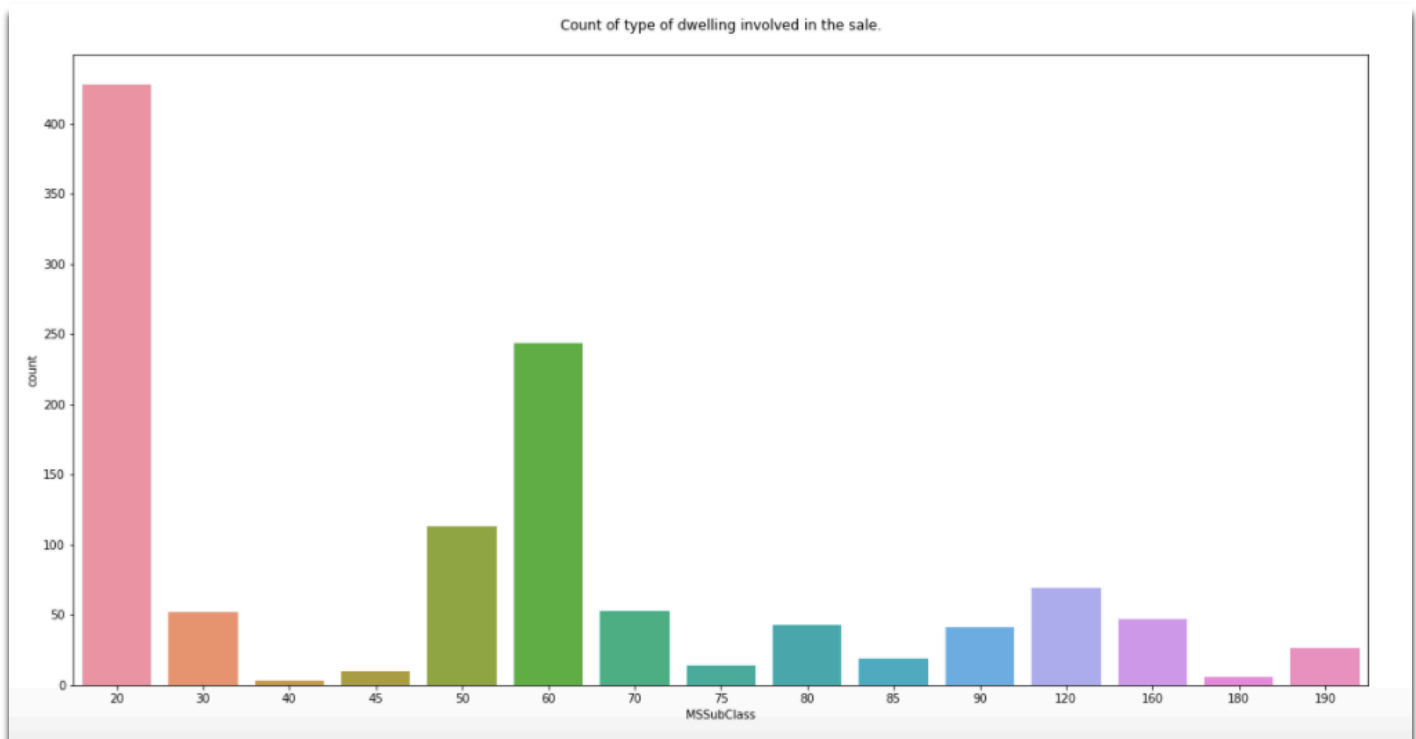<matplotlib.axes._subplots.AxesSubplot at 0x7fd8562efd60>

## key observations here

1. OverallQual is highly correlated with target variable ActualPrice.
2. Garagecars,GarageArea are highly correlated with each other.
3. GarageCars,garagearea,TotalBsmtSF, 1FirSF are highly correlated with target variable ActualPrice.

```python
1  plt.figure(figsize=(20,10))
2  sns.countplot(ds["MSSubClass"])
3  plt.title("Count of type of dwelling involved in the sale.\n ")
4  plt.show()
```



Count of type of dwelling involved in the sale.

Above countplot shows that
1. 1-STORY 1946 & NEWER ALL STYLES are maximum which are in sale.
2. 1-STORY W/FINISHED ATTIC ALL AGES is least one which is in sale.

```
1  plt.figure(figsize=(10,5))
2  sns.countplot(ds["MSZoning"])
3  plt.title("Count of type of dwelling involved in the sale.\n ")
4  plt.show()
```



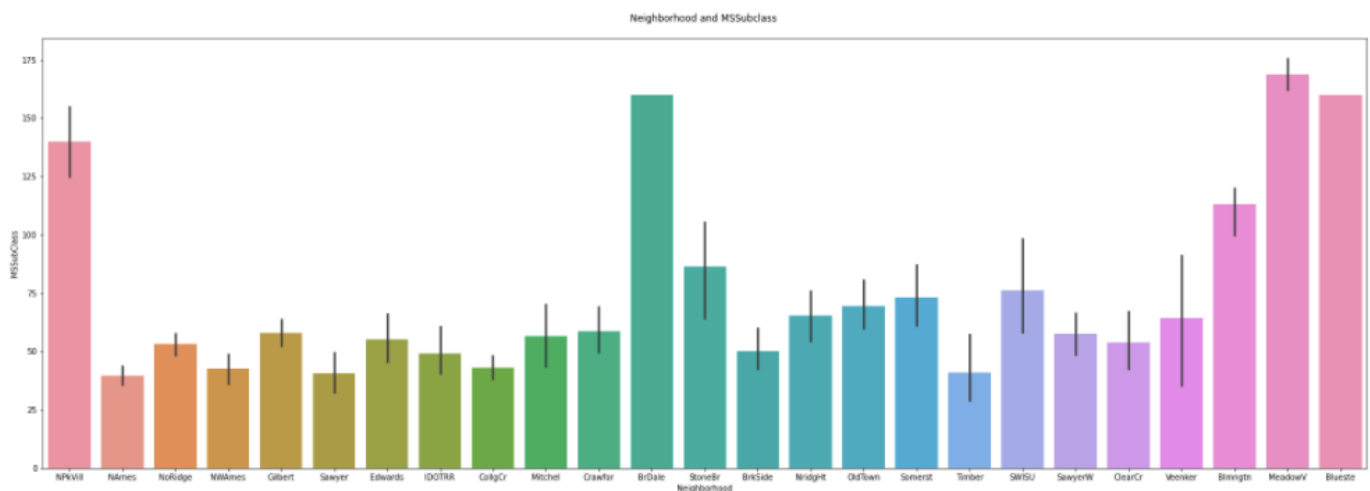Count of type of dwelling involved in the sale.

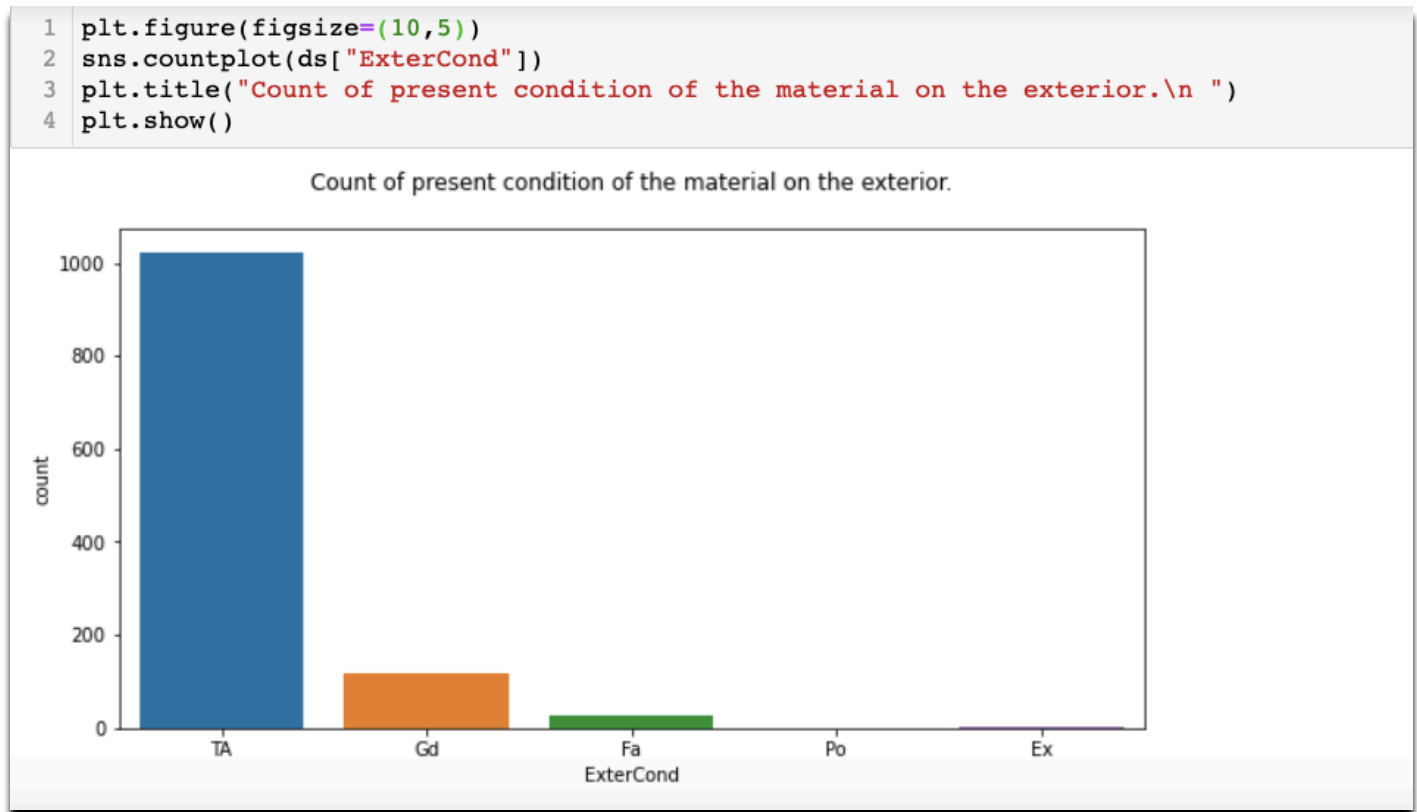Above countplot shows that Residential Low Density property is maximum on sale. and Commercial is least on sale.

```
1  plt.figure(figsize=(30, 10))
2  sns.barplot(x='Neighborhood',y='MSSubClass',data=ds)
3  plt.title("Neighborhood and MSSubclass \n")
4  plt.show()
```



Neighborhood and MSSubclass

above barplot shows the neighbourhood of all MSSubclass dwelling involved in sale. like "Meadow Village" is in neighbourhood of "2-STORY PUD - 1946 & NEWER"
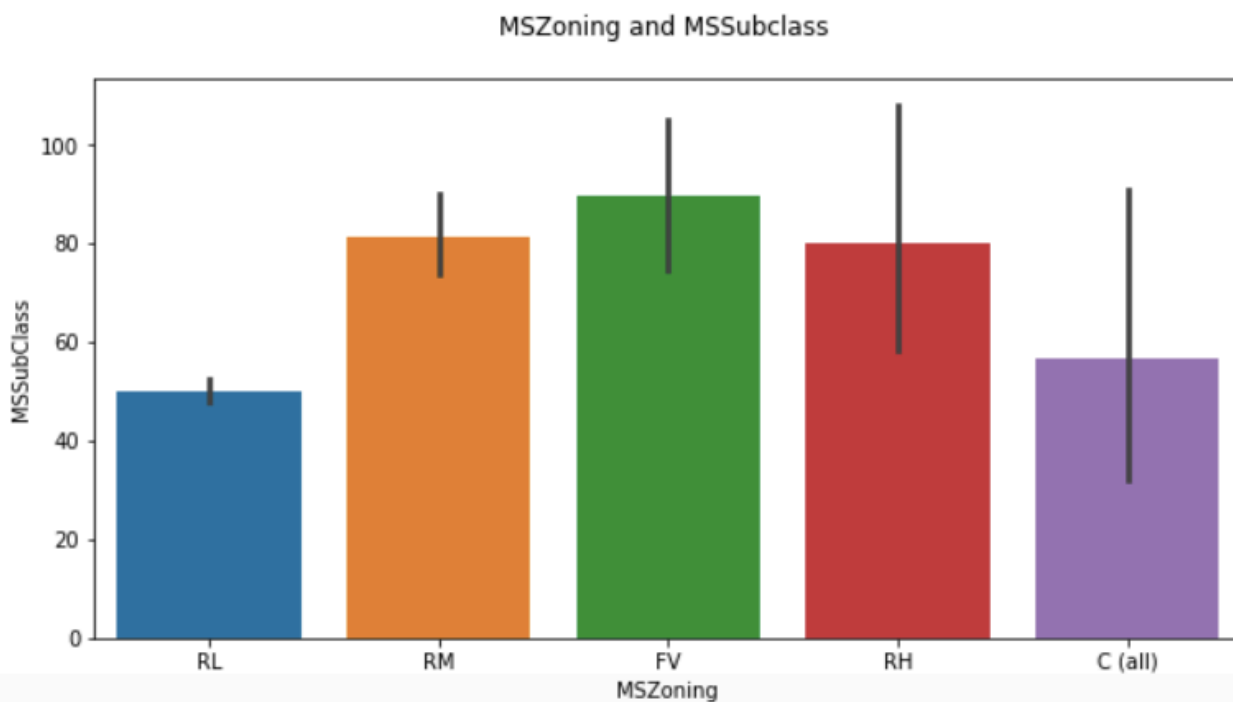
```
1  plt.figure(figsize=(10,5))
2  sns.countplot(ds["ExterCond"])
3  plt.title("Count of present condition of the material on the exterior.\n ")
4  plt.show()
```



Count of present condition of the material on the exterior.

Above countplot Evaluates the present condition of the material on the exterior is Average/Typical

```
1  plt.figure(figsize=(10,5))
2  sns.barplot(x='MSZoning',y='MSSubClass',data=ds)
3  plt.title("MSZoning and MSSubclass \n")
4  plt.show()
```



MSZoning and MSSubclass

Above countplot shows Maximum dwelling involved in sale are "DUPLEX - ALL STYLES AND AGES" having general zoning classification as "Floating Village Residential"
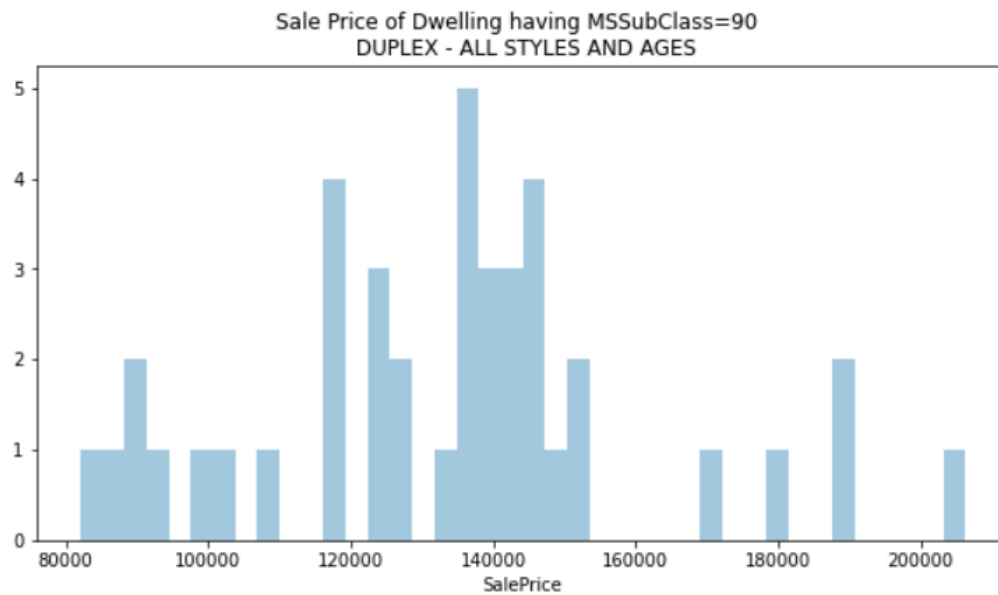
```
1  plt.figure(figsize=(10,5))
2  sns.distplot(ds[ds['MSSubClass']==90]['SalePrice'],kde=False,bins=40)
3  plt.title('Sale Price of Dwelling having MSSubClass=90  \n DUPLEX - ALL STYLES AND AGES')
```

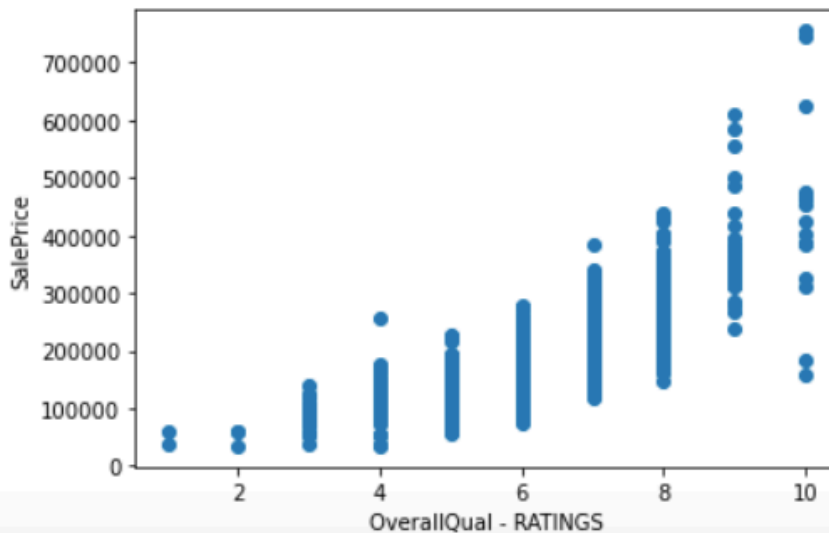Text(0.5, 1.0, 'Sale Price of Dwelling having MSSubClass=90  \n DUPLEX - ALL STYLES AND AGES')



Above distplot shows the saleprice of dwelling where MSSubClass=90 i.e DUPLEX - ALL STYLES AND AGES

```
1  plt.scatter(ds["OverallQual"],ds["SalePrice"])
2  plt.xlabel("OverallQual - RATINGS")
3  plt.ylabel("SalePrice")
4  plt.show()
```



Above scatterplot shows that Highest rating have highest prices.

## d)  Pre-Processing Pipeline

Maximum columns have the numerical values other than 18 features. For moving further for model training, we need to transform these nominal values into numerical values by encoding the data.

**Label Encoding**

```
1  # changing the nominal value to integer for training model
2  from sklearn.preprocessing import LabelEncoder
3  le=LabelEncoder()
4  list1=['MSZoning','Street', 'Alley', 'LotShape', 'LandContour', 'Utilities', 'LotConfig', 'LandSlope',
5        'Neighborhood','Condition1', 'Condition2', 'BldgType', 'HouseStyle', 'RoofStyle','RoofMatl',
6        'Exterior1st','Exterior2nd','MasVnrType','ExterQual','ExterCond','Foundation','BsmtQual',
7        'BsmtCond', 'BsmtExposure','BsmtFinType1','BsmtFinType2','Heating','HeatingQC','CentralAir',
8        'Electrical','KitchenQual','Functional','FireplaceQu','GarageType','GarageFinish','GarageQual',
9        'GarageCond','PavedDrive','PoolQC','Fence','MiscFeature','SaleType','SaleCondition']
10 for val in ds_cat:
11     ds[val]=le.fit_transform(ds[val].astype(str))
```
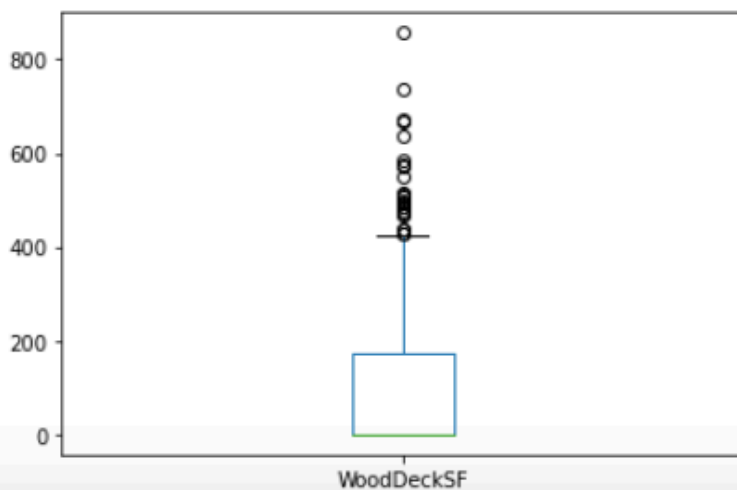
After transforming the nominal variables, the dataset is all numerical.

## Handling the outliers¶

```
1  ds["WoodDeckSF"].plot.box()
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f8bd3f8ed60>



As many type of dwelling may have lotarea, basementfin1 , alley area etc and many may not have . that may be making a huge difference in mean and median and 75% and max values . So, not removing the outliers

We will check the distribution of skewness in the data. If it is there, we will remove it.

```
1  ds.skew()  # checking if the data is skewed
```

```
Id                  0.026526
MSSubClass          1.422019
MSZoning           -1.796785
LotFrontage         2.710383
LotArea            10.659285
                      ...
MoSold              0.220979
YrSold              0.115765
SaleType           -3.660513
SaleCondition      -2.671829
SalePrice           1.953878
Length: 81, dtype: float64
```

**little skewness is there so we will remove it**

```
1  # seperatng the target variable
2  ds_x=ds.drop(columns=['SalePrice'])
3  y_t=pd.DataFrame(ds['SalePrice'])
4  print(ds_x.shape, y_t.shape)
```

```
(1168, 80) (1168, 1)
```

```
1  from sklearn.preprocessing import power_transform
```

```
1  ds_x=power_transform(ds_x,method='yeo-johnson')
```

Once skewness removal has been done, We will do scaling of the dataset. Scaling is a technique to standardize the independent features present in the data in a fixed range. If feature scaling is not done, then a machine learning algorithm tends to weigh greater values, higher and consider smaller values as the lower values, regardless of the unit of the values.

**scaling the dataset**

```python
from sklearn.preprocessing import StandardScaler
```

```python
#scaling the dataset
from sklearn.preprocessing import StandardScaler
sc=StandardScaler()
scaledX=sc.fit_transform(ds_x)
scaledX.shape
```

```
(1168, 80)
```

Libraries and packages used for model training are listed below

**Data Modelling**

```python
# importing our libraries
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression,Lasso,Ridge,ElasticNet
from sklearn.ensemble import RandomForestRegressor,AdaBoostRegressor,GradientBoostingRegressor
from sklearn.neighbors import KNeighborsRegressor
```

```python
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import cross_val_score
from sklearn.metrics import r2_score,mean_absolute_error,mean_squared_error
```

## e) Model Development and Evaluation

Identification of possible problem-solving approaches (methods)

most common techniques will fall into the following two groups:

a. Supervised learning, including regression and classification models.
b. Unsupervised learning, including clustering algorithms and association rules.

For this dataset, I will be using regression model because the output variable is continuous.

28

## Testing of Identified Approaches (Algorithms)

Here, In this project I will be using LinearRegression, Lasso, Ridge, RandomForestRegressor,KNeighborsRegressor,GradientBoostingRegressor, AdaBoostRegressor algorithms

## Finding the best random state :

finding the best random state

```
 1  best_rstate=0
 2  accuracy=0
 3  for i in range(30,200):
 4      x_train,x_test,y_train,y_test=train_test_split(scaledX,y_t,test_size=.22,random_state=i)
 5      mod=LinearRegression()
 6      mod.fit(x_train,y_train)
 7      predlr=mod.predict(x_test)
 8      tempaccu=r2_score(y_test,predlr)
 9      if(tempaccu>accuracy):
10          accuracy=tempaccu
11          best_rstate=i
12
13  print("Best Accuracy",accuracy*100, "Random state",best_rstate)
```

```
Best Accuracy 86.00361454154653 Random state 148
```

## Using the best random state:

using the best random state

```
 1  x_train,x_test,y_train,y_test=train_test_split(scaledX,y_t,test_size=.22,random_state=148)
```

```
 1  x_train.shape , x_test.shape
```

```
((911, 80), (257, 80))
```

```
 1  y_train.shape , y_test.shape
```

```
((911, 1), (257, 1))
```

Using different algorithms, we will try to find the best model.

## Finding the best model

```python
#using algorithms in for loops
model=[LinearRegression(),Lasso(),Ridge(),RandomForestRegressor(),KNeighborsRegressor(),GradientBoostingRegressor()
for m in model:
    m.fit(x_train,y_train)
    y_pred=m.predict(x_test)
    r2score=r2_score(y_test,y_pred)
    cvscore=cross_val_score(LinearRegression(),x_train,y_train,cv=5).mean()
    print(m , "\nAccuracy Score of " ,r2score*100, "Cross Val Score", {cvscore*100})
    print("*********************************************************************\n")
```

```
LinearRegression()
Accuracy Score of  86.00361454154653 Cross Val Score {67.181131669369}
*********************************************************************

Lasso()
Accuracy Score of  86.00664904863633 Cross Val Score {67.181131669369}
*********************************************************************

Ridge()
Accuracy Score of  86.00405781614823 Cross Val Score {67.181131669369}
*********************************************************************

RandomForestRegressor()
Accuracy Score of  90.69848499135716 Cross Val Score {67.181131669369}
*********************************************************************

KNeighborsRegressor()
Accuracy Score of  78.47326833612966 Cross Val Score {67.181131669369}
*********************************************************************

GradientBoostingRegressor()
Accuracy Score of  91.76786443443163 Cross Val Score {67.181131669369}
*********************************************************************
```

Here, **GradientBoostingRegressor** has performed the best with accuracy score round(2) is 0.91

Doing a **GridSearchCV** is a great way to do hyper parameters tuning.

**Hyperparameter Tuning**

```
1  # GradientBoostingRegressor is best performing model so finding its best parameter
2  from sklearn.model_selection import GridSearchCV
```

```
1  GBR=GradientBoostingRegressor()
2  GBR.fit(x_train,y_train)
3  y_pred=GBR.predict(x_test)
4  r2score=r2_score(y_test,y_pred)
5  cvscore=cross_val_score(GBR,x_train,y_train,cv=5).mean()
6  print( "\nAccuracy Score of ",GBR ,"is",r2score*100,"and", "Cross Val Score is", {cvscore*100})
7  print("********************************************************************************\n")
8  search_grid={'n_estimators':[5, 6, 7, 8, 9, 10, 11, 12, 13, 15],'learning_rate':[.001,0.01,.1],'max_depth':[1,2,4])
9  search=GridSearchCV(estimator=GBR,param_grid=search_grid,scoring='neg_mean_squared_error',n_jobs=1)
```

```
Accuracy Score of  GradientBoostingRegressor() is 91.88182151052608 and Cross Val Score is {70.03890946820167}
********************************************************************************
```

Now, as the model is performing good with the score of 91% , we will save the predicted_model .

## Saving the model- Serialization

```
1  # saving the prediction model
2
3  import pickle
4  filename="Housingprice.pkl"
5  pickle.dump(GBR,open(filename,'wb'))
```

```
1  # load the model
2  fitted_model=pickle.load(open("Housingprice.pkl",'rb'))
```

After doing serialization, we will fit the model on test data of 292 rows given.

```
1  # predictions over test data (houseprice_test.csv)
2  predictions=fitted_model.predict(scaled_df)
```

```
1  predictions=predictions.astype(int)
```

```
1  ds_pred=pd.DataFrame(data=predictions,columns=['SalePrice'])
2  ds_pred
```

|     | SalePrice |
| --- | --- |
| 0 | 371689 |
| 1 | 212453 |
| 2 | 231748 |
| 3 | 188676 |
| 4 | 218146 |
| ... | ... |
| 287 | 246127 |
| 288 | 141279 |
| 289 | 157507 |
| 290 | 159595 |
| 291 | 95635 |

292 rows × 1 columns

# **CONCLUSION**

Key Findings and Conclusions of the Study

1.    OverallQual is highly correlated with target variable SalePrice.

2.    Garagecars,GarageArea are highly correlated with each other.

3.    GarageCars, garagearea, TotalBsmtSF, 1FirSF are highly correlated with target variable SalePrice.

4.    It was found that removing outliers will be loss of more of the data. So, I decided not to remove them.

5.    GradientBoostingRegressor was the best fit model.