



Project Report
On
House Price Prediction Analysis

Submitted by:

AJIT MADAME

ACKNOWLEDGMENT

I would like to express my special gratitude to “Flip Robo” team, who has given me this opportunity to deal with a beautiful dataset and it has helped me to improve my analyzation skills. And I want to express my huge gratitude to Ms. Gulshana Chaudhary (SME Flip Robo), she is the person who has helped me to get out of all the difficulties I faced while doing the project.

A huge thanks to my academic team “Data trained” who are the reason behind what I am today. Last but not least my parents who have been my backbone in every step of my life. And also thank you for many other persons who has helped me directly or indirectly to complete the project.

Contents:

1. Introduction

- Business Problem Framing
- Conceptual Background of the Domain Problem
- Review of literature
- Motivation for the Problem undertaken

2. Analytical Problem Framing

- Mathematical/ Analytical Modelling of the Problem
- Data Sources and their formats
- Data Pre-processing Done
- Data Input – Logic – Output Relationships
- Underlying Assumptions
- Hardware, Software and Tools Used

3. Data Analysis and Visualization

- Univariate Visualization
- Bivariate Visualizations
- Multivariate Visualization

4. Model Developments and Evaluation

- Features selections
- The model algorithms used
- Interpretation of the result
- Hyperparameter tuning

5. Conclusions

- Key Finding and conclusions
- Limitation of this works and scope for future works

1.INTRODUCTION

1.1Business Problem Framing:

Houses are one of the necessary needs 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?**

1.2 Conceptual Background of the Domain Problem

Housing Attributes: Studying the structural, locational, and economic attributes of housing properties is crucial in understanding their mutually inclusive relationships with their pricing.

House is one of human life's most essential needs, along with other fundamental needs such as food, water, and much more. Demand for houses grew rapidly over the years as people's living standards improved. While there are people who make their house as an investment and property, yet most people around the world are buying a house as their shelter or as their livelihood.

Predictive modelling, Market mix modelling, recommendation systems are some of the machine learning techniques used for achieving the business goals for housing companies.

1.3 Review of Literature

Trends in housing prices indicate the current economic situation and also are a concern to the buyers and sellers. Many factors have an impact on house prices, such as the number of bedrooms and bathrooms. House price depends upon its location as well. A house with great accessibility to highways, schools, and malls would have a greater price as compared to a house with no such accessibility. Predicting house prices manually is a difficult task and generally not

very accurate, hence there are many systems developed for house price prediction.

From studying the papers and analyzing the research work it is learnt that locational attributes and structural attributes are prominent factors in predicting house prices. Studies suggest that there exists a close relationship between House pricing and locational attributes such as distance from the closest shopping center, train station, position offering views of hills or shore, the neighborhood in which the property is situated etc. Structural attributes of the house like lot size, lot shape, quality and condition of the house, garage capacity, rooms, Lot frontage, number of bedrooms, bathrooms, overall finishing of the house etc. play a big role in influencing the house price. Neighborhood qualities can be included in deciding house price. Factors like efficiency of public education, community social status, the socio-cultural demographics improve the worth of a property.

1.4 Motivation for the Problem Undertaken

Everyone wishes to buy and live in a house which suits their lifestyle and which provides amenities according to their needs. There are many factors that are to be taken into consideration like area, location, view etc. for prediction of house price. It is very difficult to predict house price as it is constantly changing and quite often the prices are exaggerated for which people who want to buy houses, and various real estate agencies who want to invest in properties, find it difficult to buy or sell houses. For this reason, I build an advanced automated Machine Learning model.

2.Analytical Problem Framing

2.1 Mathematical/ Analytical Modelling of the Problem

Various Regression analysis techniques were used to build predictive models to understand the relationships that exist between Housing sales prices and various Housing property attributes. The Regression analysis models were used to predict the Sale price value for changes in Housing property attributes.

In order to forecast house price, predictive models such as Linear regression Model, Random Forest Regression model, Gradient Boosting Regression Model, Extreme Gradient Boost Regression, K- Nearest Neighbors Regression model were used to predict how the values of Sale Price depended on the independent variables of various Housing property attributes.

2.2 Data Sources and their formats

The dataset is provided by Flip Robo which is in the format csv.

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.

```
1 data.head()
```

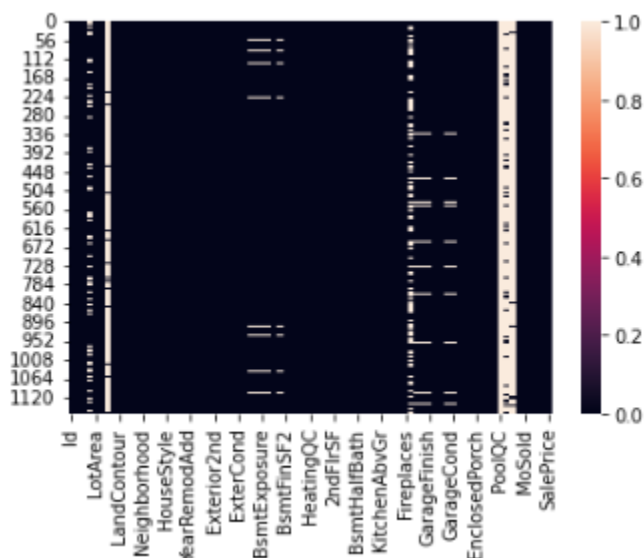
	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	Utilities	LotConfig	LandSlope	Neighborhood	Condition1
0	127	120	RL	NaN	4928	Pave	NaN	IR1	Lvl	AllPub	Inside	Gtl	NPKVill	Norm
1	889	20	RL	95.0	15865	Pave	NaN	IR1	Lvl	AllPub	Inside	Mod	NAmes	Norm
2	793	60	RL	92.0	9920	Pave	NaN	IR1	Lvl	AllPub	CulDSac	Gtl	NoRidge	Norm
3	110	20	RL	105.0	11751	Pave	NaN	IR1	Lvl	AllPub	Inside	Gtl	NWAmes	Norm
4	422	20	RL	NaN	16635	Pave	NaN	IR1	Lvl	AllPub	FR2	Gtl	NWAmes	Norm

Condition1	Condition2	BldgType	HouseStyle	OverallQual	OverallCond	YearBuilt	YearRemodAdd	RoofStyle	RoofMatl	Exterior1st	Exterior2nd	MasVnrType
Norm	Norm	TwnhsE	1Story	6	5	1976	1976	Gable	CompShg	Plywood	Plywood	None
Norm	Norm	1Fam	1Story	8	6	1970	1970	Flat	Tar&Grv	Wd Sdng	Wd Sdng	None
Norm	Norm	1Fam	2Story	7	5	1996	1997	Gable	CompShg	MetalSd	MetalSd	None
Norm	Norm	1Fam	1Story	6	6	1977	1977	Hip	CompShg	Plywood	Plywood	BrkFace
Norm	Norm	1Fam	1Story	6	7	1977	2000	Gable	CompShg	CemntBd	CmentBd	Stone

Training Dataset contains 1168 entries and 81 variables, while Test Dataset contains 292 entries and 80 variables.

2.3 Data Pre-processing Done

Checking Null Values



- We can see, Alley, Pool Quality, Fence and Miscellaneous Features are having more 90% null values so I will drop this column.
- In Fire place quality feature are having more than 55% of missing data so I will drop this feature.
- The ID columns from test and train datasets were also dropped since they don't contribute to building a good model for predicting the target variable values.

2.4 Data Inputs - Logic - Output Relationships

The Datasets consist mainly of object data type variables and a few float and int data type variables. The relationships between the independent variables and dependent variable were analyzed.

Features like Lot area, Lot Frontage, Overall Quality, Overall Condition, Basement Finishing, Total Basement Surface Area, first and 2nd Floor square feet, Garage capacity, Total rooms have a positive linear relationship, therefore increase in their values leads to increase in Sale Price.

2.5 Underlying Assumptions

```
1 data.describe()
```

	MSSubClass	LotFrontage	LotArea	OverallQual	OverallCond	YearBuilt	YearRemodAdd	MasVnrArea	BsmtFinSF1	BsmtFinSF2	BsmtUnfSF
count	1168.000000	1168.000000	1168.000000	1168.000000	1168.000000	1168.000000	1168.000000	1168.000000	1168.000000	1168.000000	1168.000000
mean	56.767979	70.988470	10484.749144	6.104452	5.595890	1970.930651	1984.758562	102.310078	444.726027	46.647260	569.721747
std	41.940650	22.437056	8957.442311	1.390153	1.124343	30.145255	20.785185	182.047152	462.664785	163.520016	449.375525
min	20.000000	21.000000	1300.000000	1.000000	1.000000	1875.000000	1950.000000	0.000000	0.000000	0.000000	0.000000
25%	20.000000	60.000000	7621.500000	5.000000	5.000000	1954.000000	1966.000000	0.000000	0.000000	0.000000	216.000000
50%	50.000000	70.988470	9522.500000	6.000000	5.000000	1972.000000	1993.000000	0.000000	385.500000	0.000000	474.000000
75%	70.000000	79.250000	11515.500000	7.000000	6.000000	2000.000000	2004.000000	160.000000	714.500000	0.000000	816.000000
max	190.000000	313.000000	164660.000000	10.000000	9.000000	2010.000000	2010.000000	1600.000000	5644.000000	1474.000000	2336.000000

- Big difference between max value and 75% in Sale Price, MSSubClass, LotFrontage, LotArea, BsmtFinSF1, BsmtFinSF2, etc. indicates presence of outliers.
- Higher std than mean in columns: MasVnrArea, BsmtFinSF1, BsmtFinSF2, WoodDeckSF, OpenPorchSF, EnclosedPorch, 3SsnPorch etc. indicates presence of skewness.

2.6 Hardware, Software and Tool Used

Hardware Used:

Processor – Intel core i3

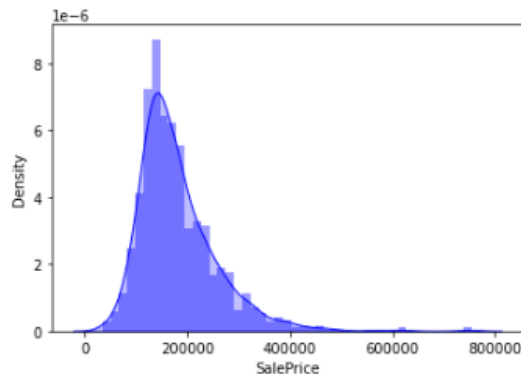
Physical Memory – 8 GB

Software Used:

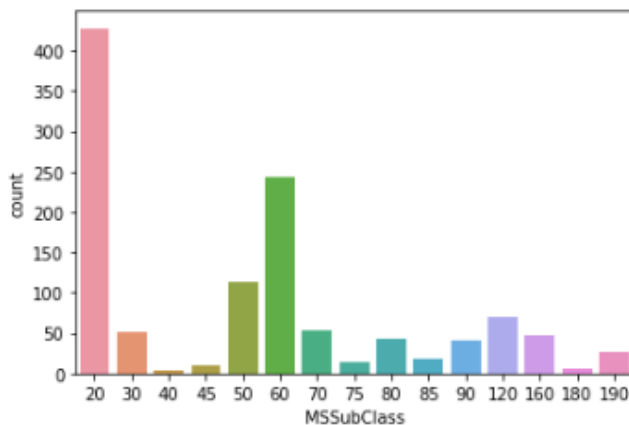
- Windows 10 Operating System
- Anaconda Package and Environment Manager
- Jupyter Notebook
- Python Libraries used: In Which Pandas, Seaborn, Matplotlib, Numpy and Scipy
- sklearn for Modelling Machine learning algorithms, Data Encoding, Evaluation metrics, Data Transformation, Data Scaling, Component analysis, Feature selection etc.

3.Data Analysis and Visualization

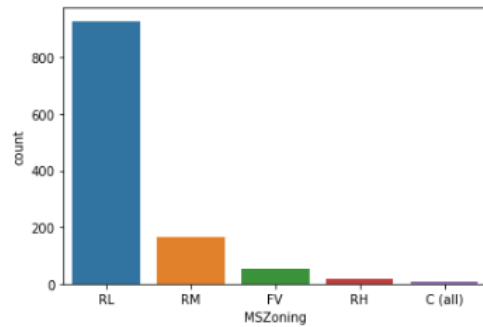
3.1 Univariate Visualization



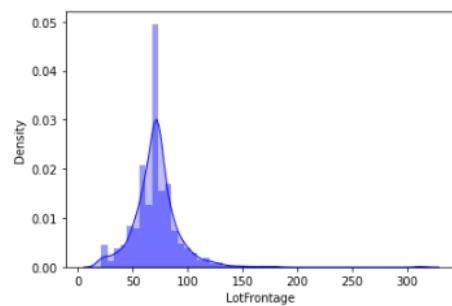
- SalePrice is not normally distributed. It contains some outliers.
- SalePrice are positively or right skewed.



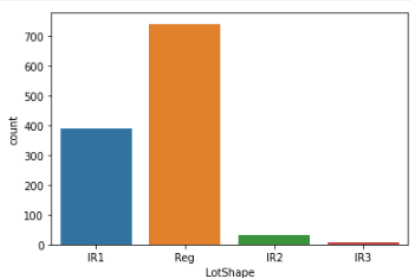
- It is dwelling involved while Selling the Houses.
- We can see, 1-STORY 1946 & NEWER ALL STYLES i.e. 20 has maximum count followed by 2-STORY 1946 & NEWER i.e. 60.
- 40 and 180 has very few counts than others.



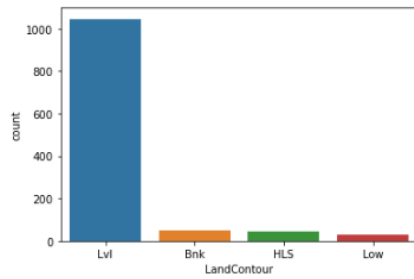
- We can see, residual with low density has maximum zone where houses are more number for selling followed by residual medium density.
- Very few houses are for sale in commercial zone.



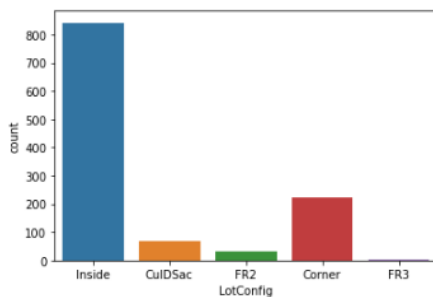
- LotFrontage means Linear feet of street connected to property. It is a right skewed column.
- It may have some outliers.
- We can see, Maximum LotFrontage are in 30 to 110.



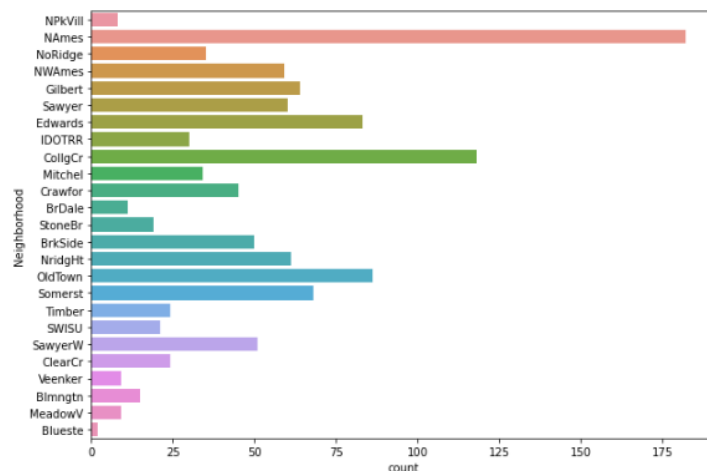
- Regular Shape of property contains maximum count followed by Slightly irregular.
- Irregular shape of property contains lowest count followed by Moderately Irregular.



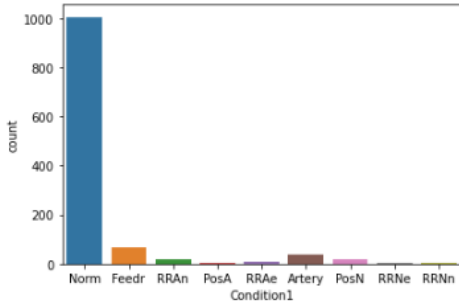
- Most Properties have Near Flat/Level LandContour.
- Remaining contain almost same count of landcontour.



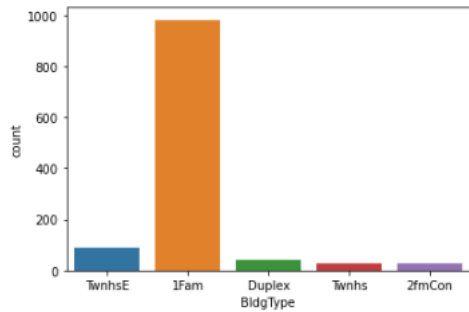
- Inside lot is the most common Lot configuration and it has maximum count than others.
- Frontage on 3 sides of property has lowest count. It means it rarely used configuration.



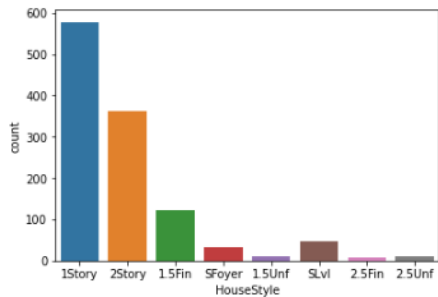
- Most Housing properties are situated in Neighborhoods of North Ames, followed by College Creek, Edwards and Old Town.
- Very Few houses in neighborhood of Bluestem.



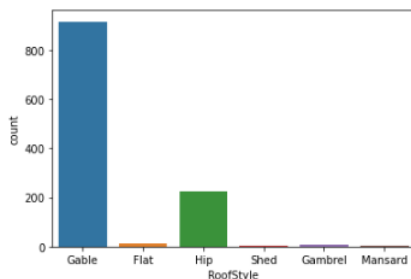
- Most Housing properties are in proximity to Normal conditions
- Very few properties are in proximity to Within 200' of North-South Railroad



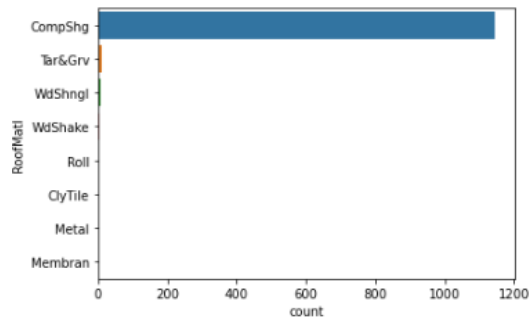
- Single-family Detached type of family dwelling are most common properties than others.
- Townhouse End Unit type of family dwelling are very few properties than others.



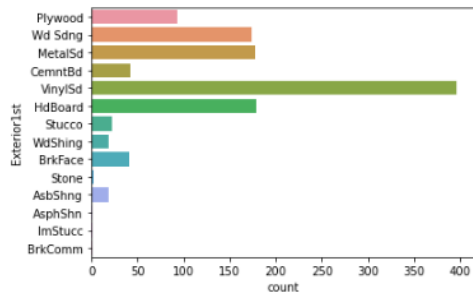
- Most Housing properties 1 storied and 2 storied.
- There are less housing properties wo and one-half story: 2nd level finished.



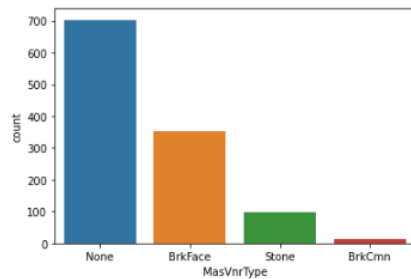
- Most Houses have Gable roof style.
- Shed, Gambrel and Mansard has very low count.



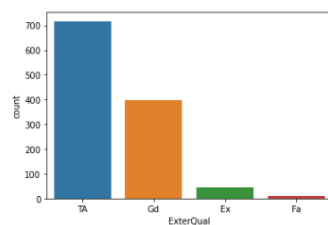
- Most Houses have roofs made of Standard (Composite) Shingle.
- Almost remaining all are having lowest count than CompShg.



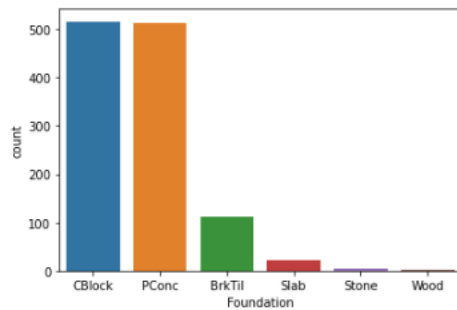
- Vinyl Siding is the most common exterior covering used followed by HD Board.
- Asphalt Shingles has lowest count of exterior are used.



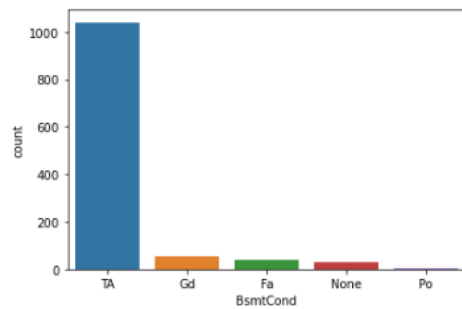
- Most Houses don't have a Masonry veneer type.
- while some have Brick Face Masonry veneer type.



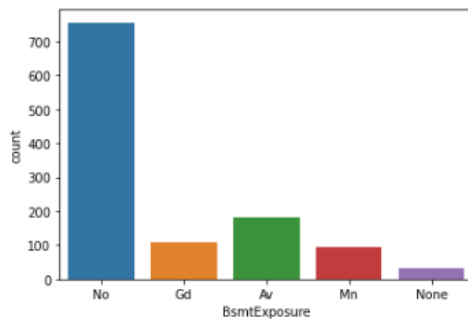
- The quality of the material on the exterior is most commonly average/typical.
- The quality of material on exterior is very few are fair.



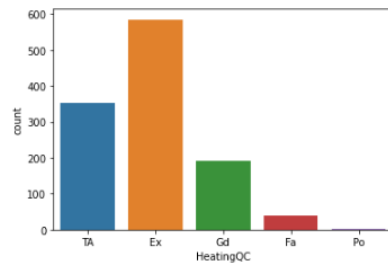
- Two of the most common foundation types are Cinder Block and Poured Contrete.
- Very few are used Wood type for foundation.



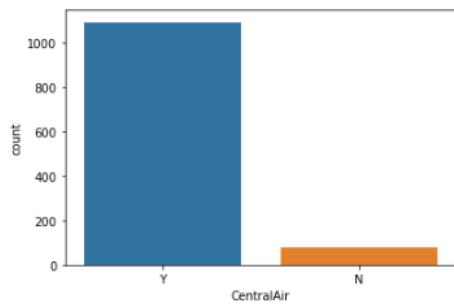
- The general condition of the basement is commonly Typical with slight dampness.
- Some of Basement have very Poor condition.



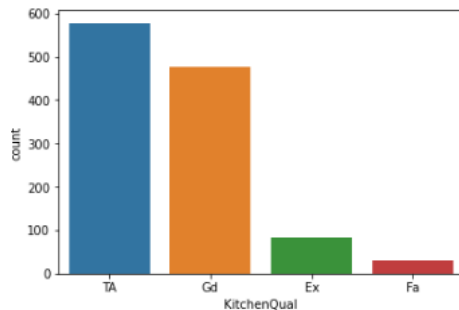
- Basements most commonly have no exposure.
- But Some Basement has Average Exposure.



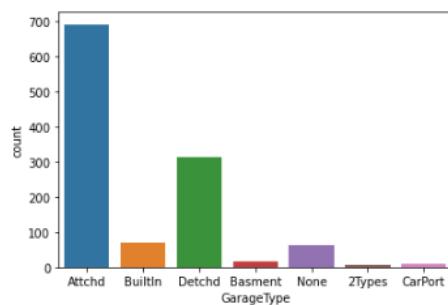
- Most houses have Excellent Heating quality and condition followed by Average/Typical.
- Very low almost don't have poor quality of Heating condition of houses.



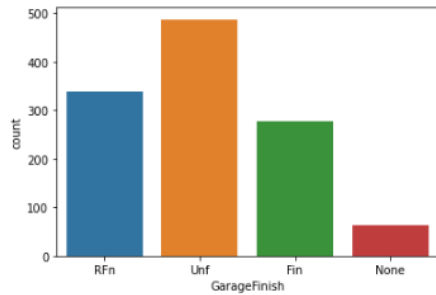
- Most houses have Central air conditioning.
- Very few houses don't have central air conditioning.



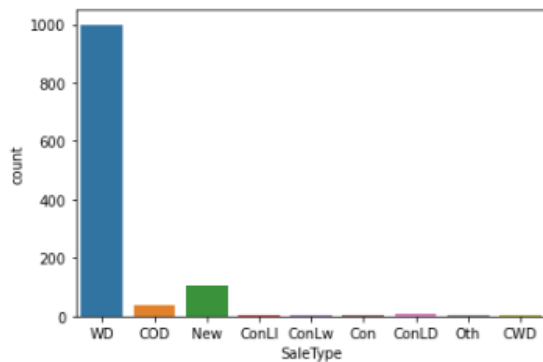
- Most houses have Most houses have Typical/Average and Good Kitchen quality.
- Very few have fair quality of kitchen.



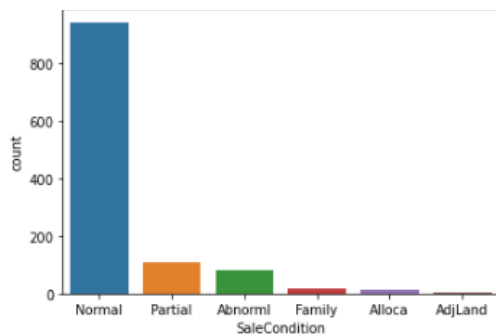
- Most houses have a Garage Attached to home.
- Many houses have a Detached from home.



- Most houses have an Unfinished garage.
- But so many houses have a Rough Finished garage.

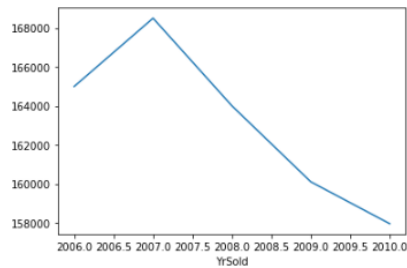


- Almost all are having Warranty Deed - Conventional is the most common Type of sale.

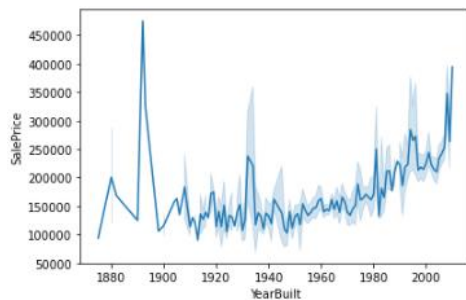


- Almost all are having Condition of sale is Normal Sale.

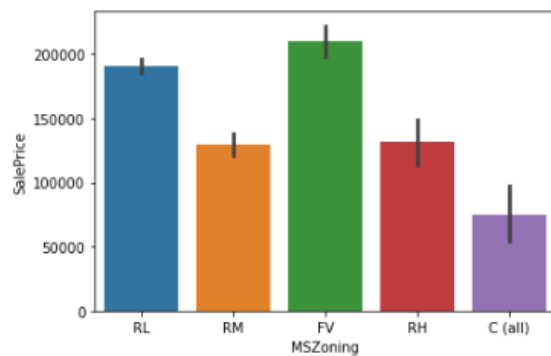
3.2 Bivariate Visualization



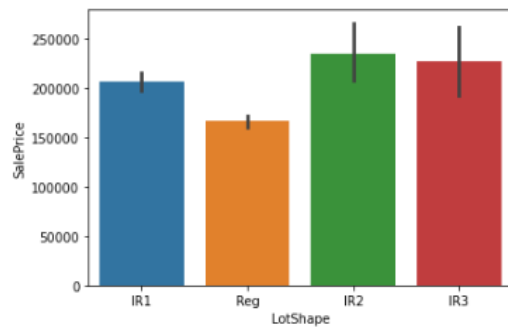
- We can see, in 2007 housing sales price is in peak. But after 2007 Hoses price is drastically dropped.



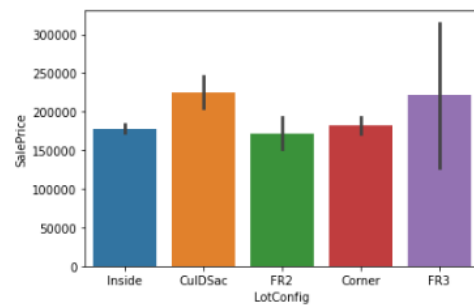
- We can see, houses price are increases after the year of built 1990.
- But in year of 1880 to 1990, sale price of house is high.



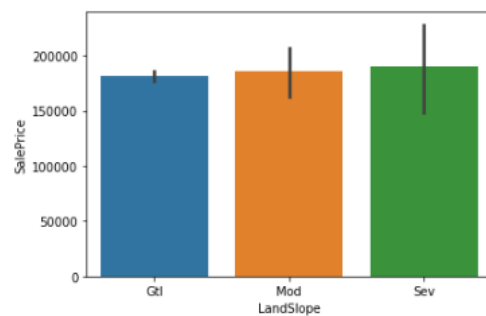
- Floating Village Residential has maximum sale Price of houses followed by Residential Low Density.
- Commercial houses are having low Sale Price comparatively.



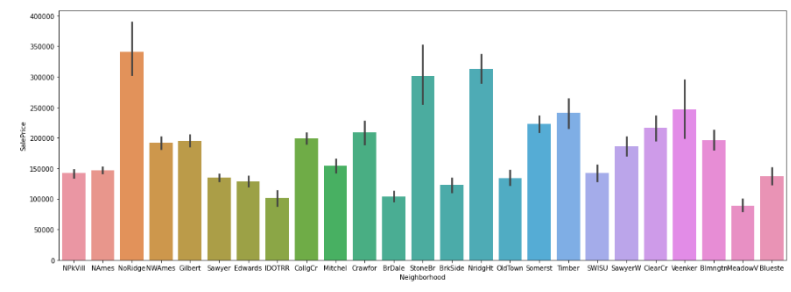
- Moderately Irregular type of shape houses has maximum sale price followed by Irregular.
- But Regular type of Lotshape has low sale price than others.



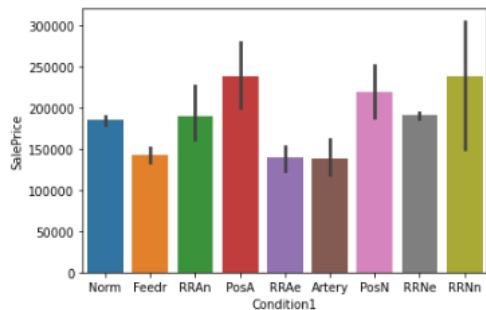
- Cul-de-sac of Lot configuration has highest sale price than others followed by Frontage on 3 sides of property.
- Inside of Lot configuration has lowest sale price than others.



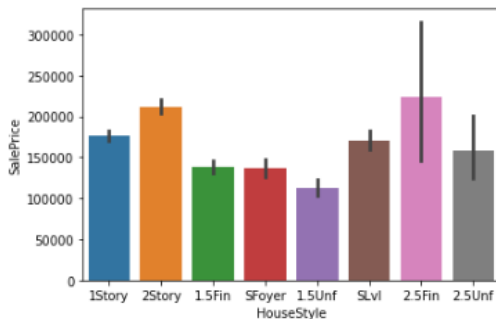
- Severe Slope has highest sale price than others.
- But all are having almost same sale price so slope is not affecting more to sale price.



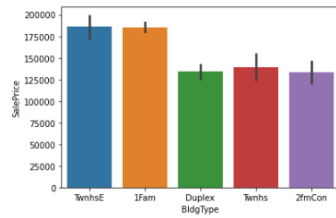
- Northridge has maximum sale price; it means that Northridge neighborhood has maximum sale price of houses followed by Northridge Heights.
- The houses are Meadow Village neighborhood has lowest sale price of houses than others.



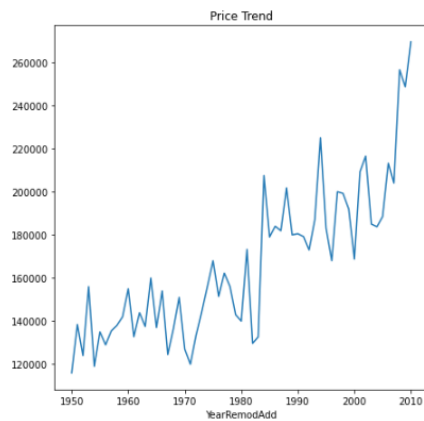
- The houses are proximity to condition Adjacent to positive off-site feature this has maximum sale price followed by Adjacent to North-South Railroad.
- Hoses have adjacent to arterial street has lowest sale price.



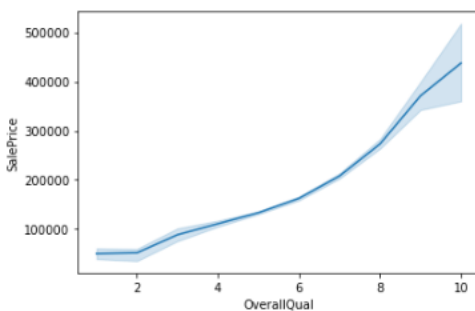
- Two and one-half story: 2nd level finished this type of houses dwelling has maximum sale price than others.
- One and one-half story: 2nd level finished it has lowest sale price.



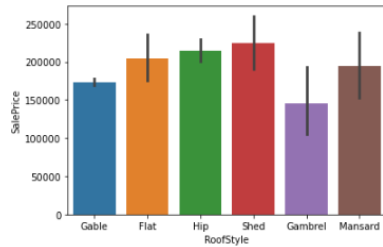
- Townhouse End Unit, tghis type of dwelling has maximum sale price than others followed by ownhouse Inside Unit.
- Duplex has lowest sale than others.



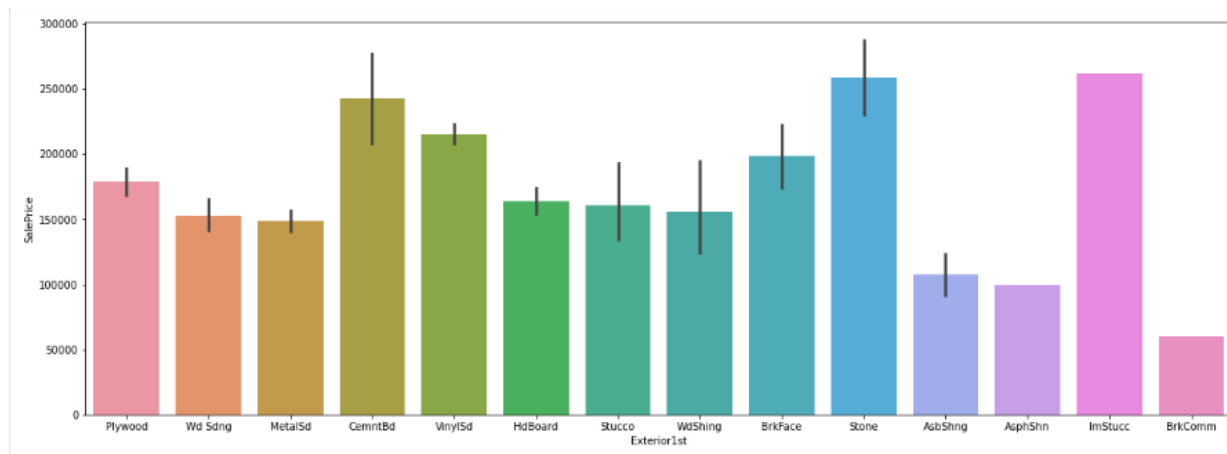
- It is observed that Sales value is higher for houses which were remodeled more recently.



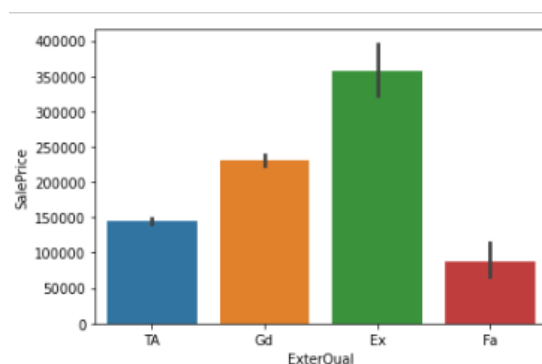
- We can see, Overall quality has linear relationship with Sale Price.
- If quality of houses is increasing then price of houses is also increasing.



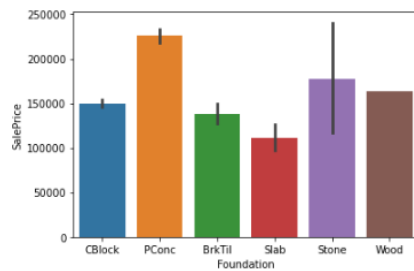
- Shed type of roof style of houses has maximum sale price than others followed by Hip type of roof style.
- Gambrel type of roof style of houses has lowest sale price.



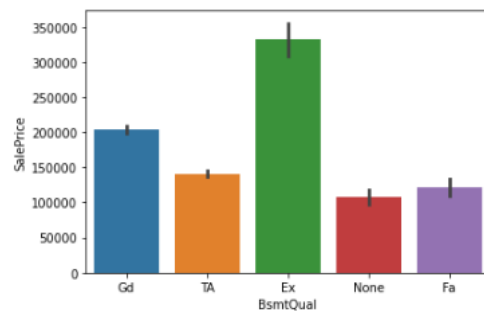
- Imitation Stucco type Exterior covering on house has maximum sale price followed by stone.
- Brick Common type Exterior covering on house has lowest sale price followed by Asphalt Shingles.



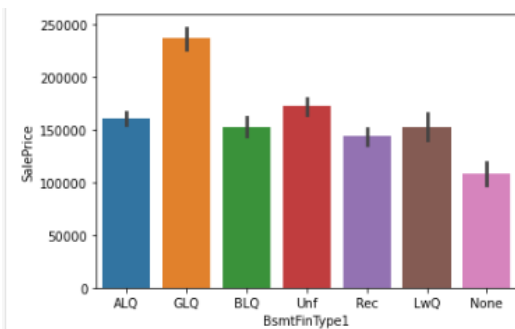
- Excellent quality of exterior has maximum sale price followed by good quality.
- Fair quality houses have lowest sale price than others.



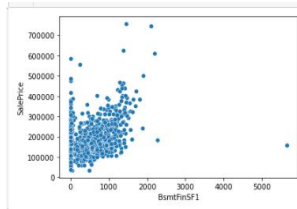
- Poured Concrete type of foundation houses has maximum sale price than others followed by Stone.
- Slab type of foundation has lowest sale price than others.



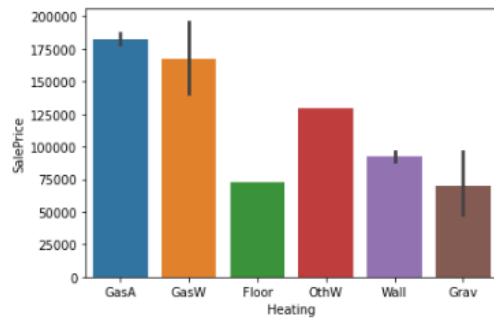
- Excellent quality of Basement has more sale price than others followed by good quality of Basement.
- None of type or no Basement houses has lowest sale price than others.



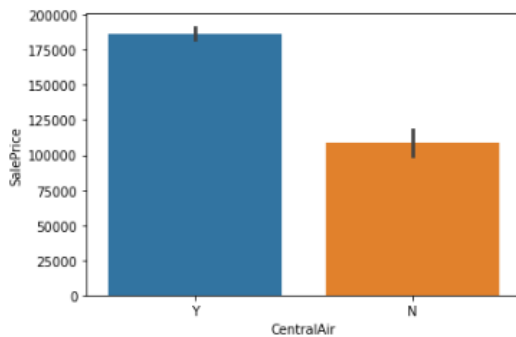
- Good Living Quarters type of Rating of basement finished area has highest sale price than others followed Unfinished.
- No Basement has lowest sale price.



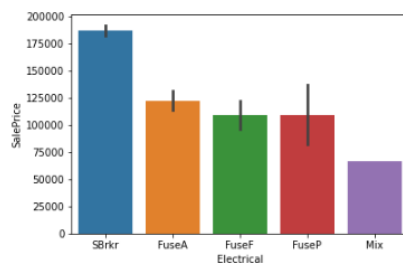
- It has linear relationship with each other's.
- Maximum sale price of Type 1 finished square feet is lies in 0 to 2000.



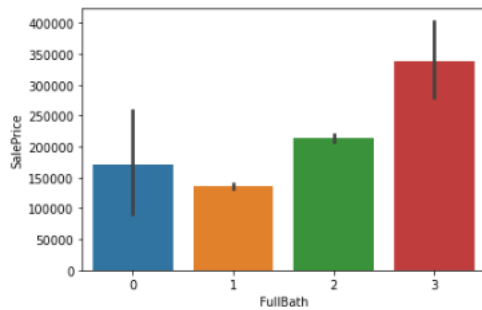
- Gas forced warm air furnace of heating houses has more sale price than others followed by Gas hot water or steam heat.
- Floor Furnace used heating has less sale price than others.



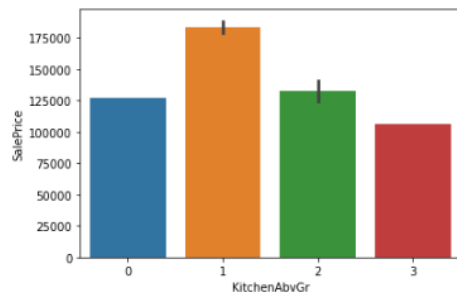
- Those houses have cetralAir conditioning that house has maximum sale price than others.
- Those houses are not used central air conditioning has minimum sale price than others.



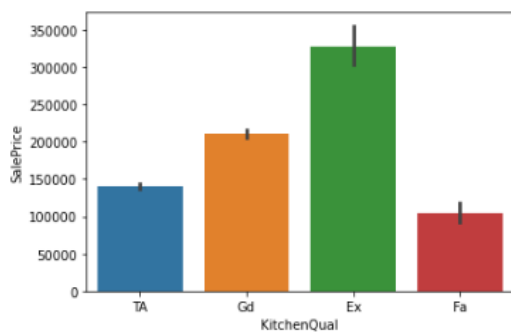
- Standard Circuit Breakers & Romex type electrical system has higher sale price than others.
- Mix type of electrical system has lower sale price than others.



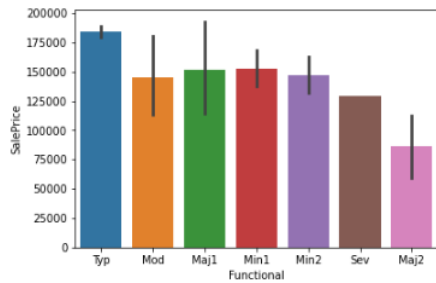
- Type 3 of FullBath has maximum sale price than others.
- type 1 FullBath has lowest sale price than others.



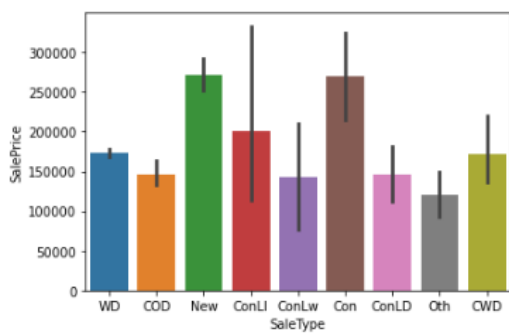
- 1 type of kitchen above grade has maximum sale price followed by 2.
- All Kitchen has equally distributed up to 100,000.



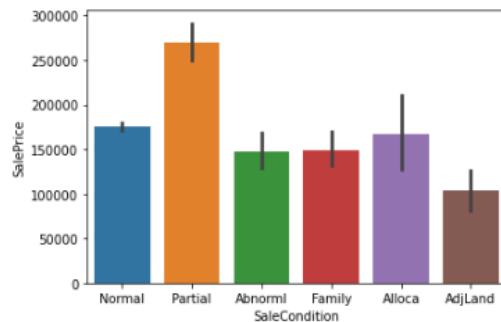
- Excellent type quality of Kitchen has maximum sale price followed by good quality of Kitchen.
- Fair quality of Kitchen has lowest sale price than others.



- Typical Functionality of garage near to house has maximum sale price than others.
- Major Deductions 2 has lowest sale price.

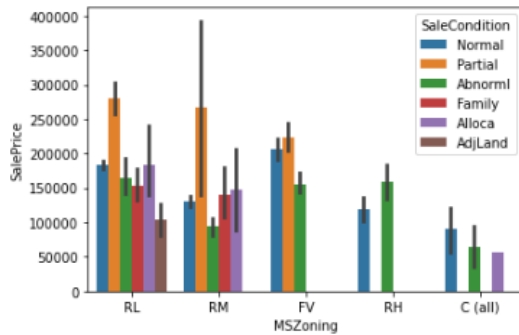


- Contract 15% Down payment regular terms sale type has more sale price than others.
- Others type of sale type has lower sale price.

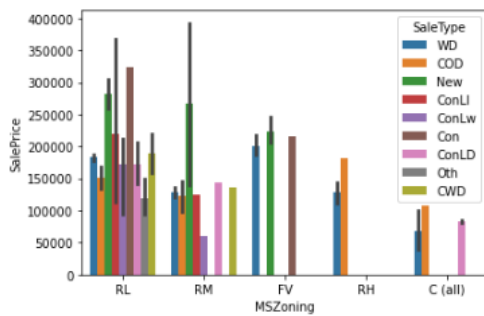


- Partial type of sale condition has maximum house sale price than others.
- Adjoining Land Purchase of sale condition house of sale price.

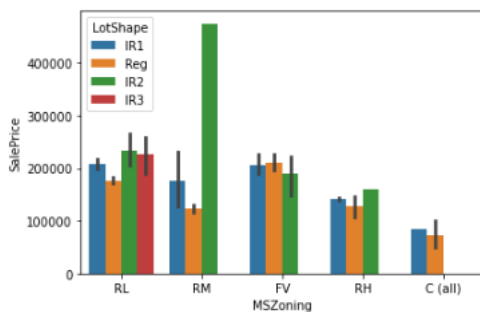
3.3 Multivariate Visualization



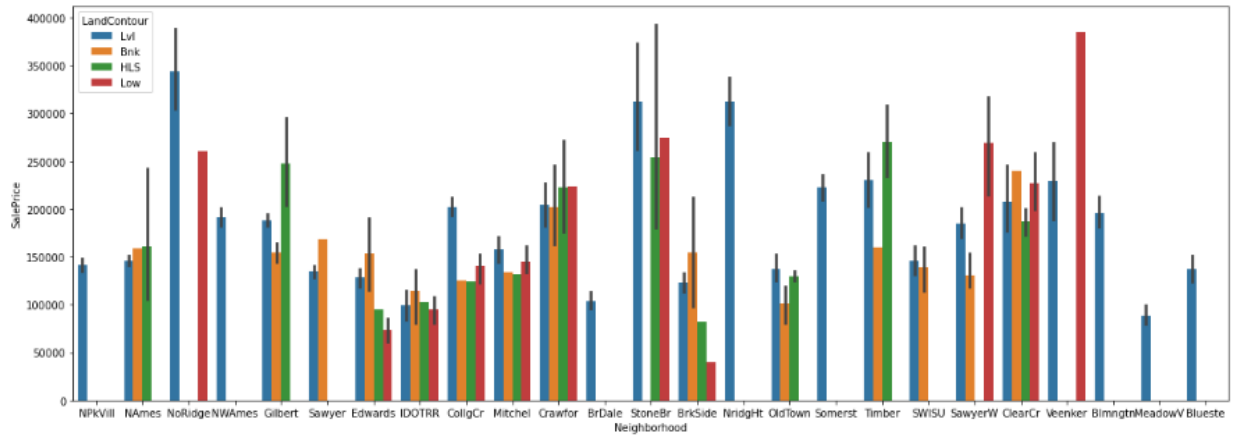
- New Homes are the most popular in all types of zoning.
- Partial sale type of condition has highest sale price.



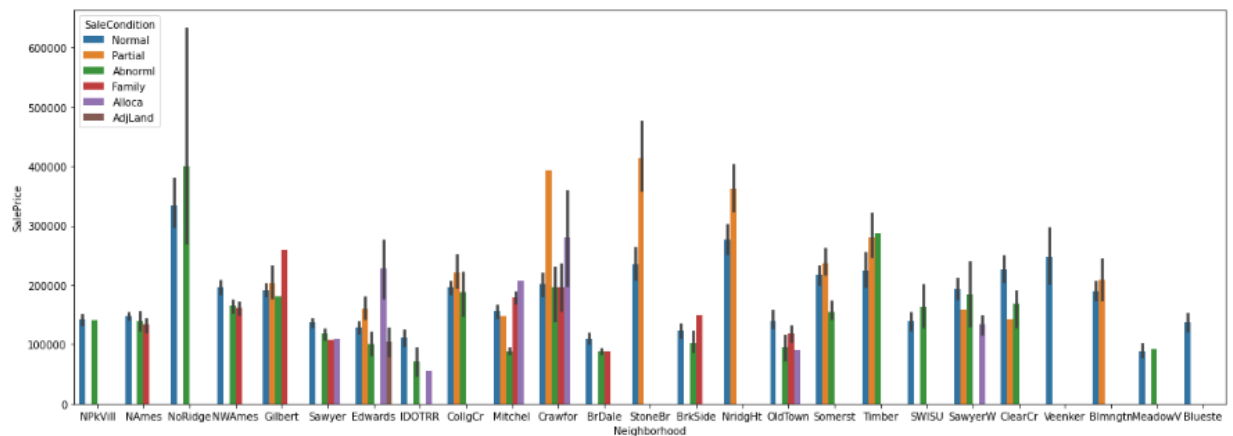
- Low interest contract are the most popular sale types in low density.
- New Houses are top among other sale type and it has good sale price.



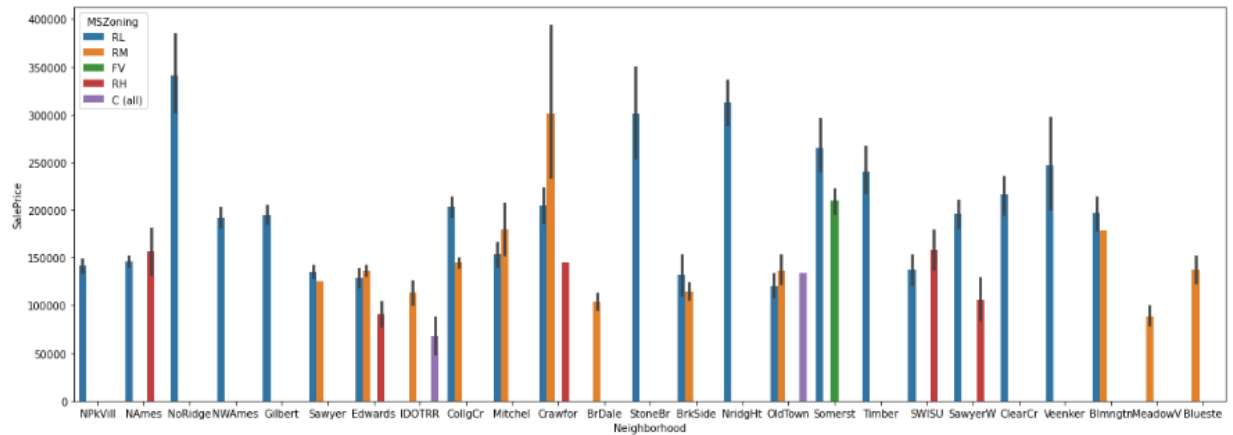
- Partially irregular and irregular plot shapes are most popular in low and medium residential zones.



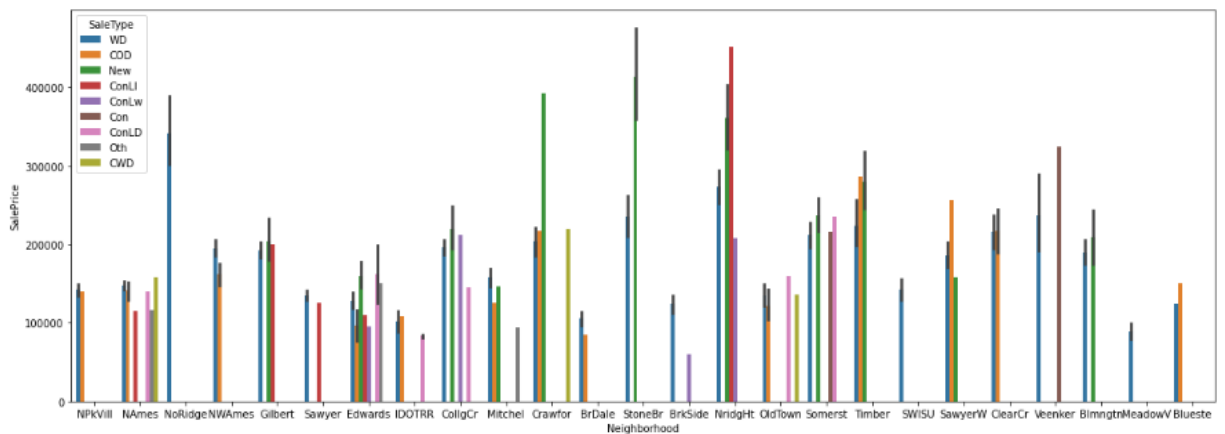
- Most housing properties established in levelled regions in North Ridge sell for the highest.
- Most Housing properties in levelled regions of Stone Brook sell for highest followed by banked region and hillsides.
- Houses in levelled region of NorthRidge heights sell for the most while housing properties in depressed regions of Veenker sell for the highest prices.



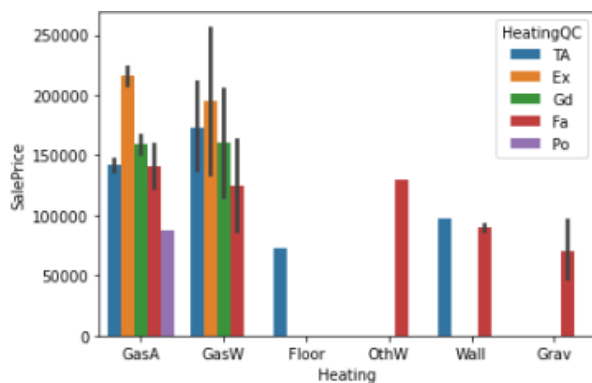
- Most housing properties that are newly established in Crawford,Stone Brook, Timberlane,North Ridge Heights,Bloomington Heights sell for the highest.
- Most Housing properties in North Ridge sell for trade, foreclosure, short sale and normal sale in North Ridge.



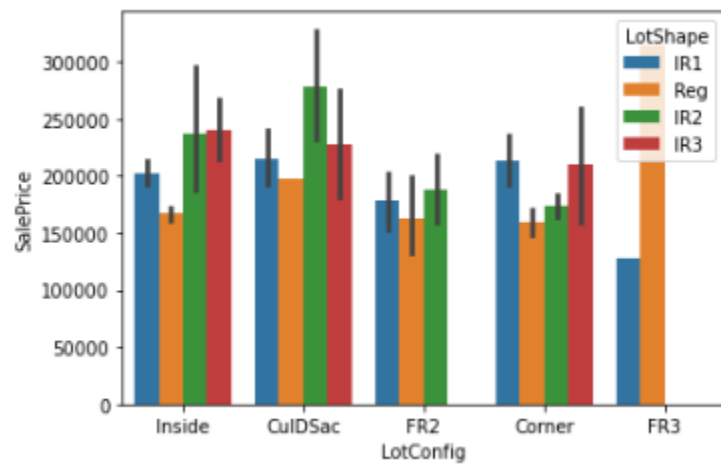
- Most houses sold in North Ridge, North Ridge Heights, Somerset, Timber Lane, Veenker, Bloomington Heights are in low density residential zones.
- North Ames has more houses sold in High density residential zones, while Crawford has more houses sold in medium density residential zones.



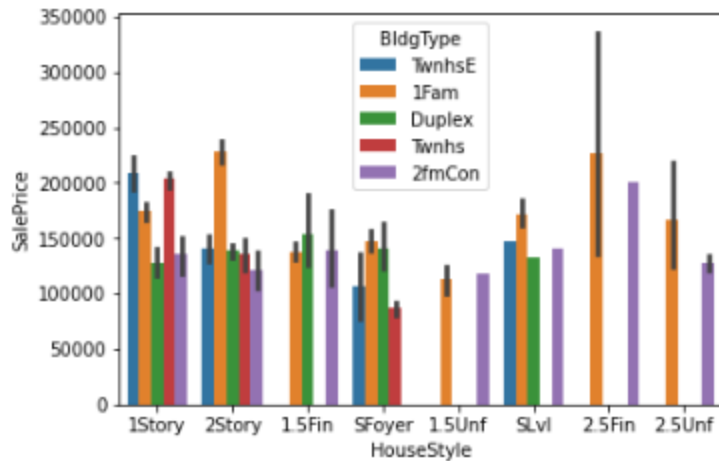
- Warranty Deed - Conventional, Home just constructed and sold, Contract Low Interest Court Officer Deed/Estate are the most common sale types.



- Excellent quality of Gas forced warm air furnace and Gas hot water heating systems fetches the highest amount of money.



- 3-sided Frontage properties with Regular plot shape sell for the highest.



- Two and one-half story: 2nd level finished housing properties sell for the highest.

Encoding the Categorical Columns

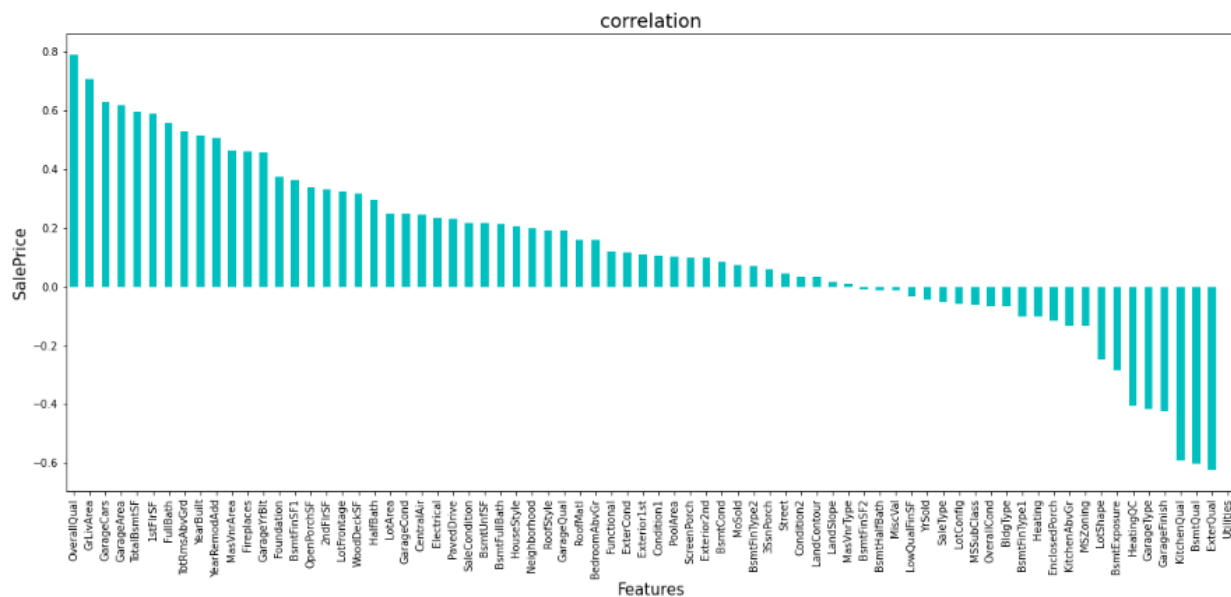
Before Proceeding with finding the correlations of the columns, the data of the categorical columns needs to be encoded using Label Encoder.

```
1 from sklearn.preprocessing import LabelEncoder
```

```
1 le=LabelEncoder()
```

```
1 for col in data[data.columns[data.dtypes == 'object']]:  
2     data[col] = le.fit_transform(data[col])
```

Visualizing correlation of feature columns with label column.



- OverallQual,GrLiveArea,GarageCars,GarageArea,TotalBsmtSF,1stFlrSF,FullBath,TotRmsAbvGrd,MasVnrArea,FirePlaces have the strongest positive correlation with SalePrice.
- WhileBsmtQual,ExterQual,KitchenQual,GarageFinish,House_age,Remod_age,HeatingQC,Garage_age have the strongest negative correlation with SalePrice.

4. Models Development and Evaluation

4.1 Features Selection

Features were first checked for presence of multicollinearity and based on the respective ANOVA f-score values, the feature columns were selected that would best predict the Target variable, to train and test machine learning models.

```
1 from statsmodels.stats.outliers_influence import variance_inflation_factor
```

```
1 vif = pd.DataFrame()
```

```
1 vif["Features"] = x.columns
2 vif['vif'] = [variance_inflation_factor(scaled_X,i) for i in range(scaled_X.shape[1])]
3 vif
```

0]:

	Features	vif
0	MSSubClass	5.007543
1	MSZoning	1.357159
2	LotFrontage	1.755825
3	LotArea	1.755646
4	Street	1.195391
5	LotShape	1.256674
6	LandContour	1.329278
7	Utilities	NaN
8	LotConfig	1.160932
9	LandSlope	1.593436
10	Neighborhood	1.263801
11	Condition1	1.45170

BsmtFinSF1, 1stFlrSF, 2ndFlrSF, GrLivArea, GarageCars, GarageArea, HouseAge exhibit high multicollinearity.

```
1 from sklearn.feature_selection import SelectKBest, f_classif
```

```
1 bestfeat = SelectKBest(score_func = f_classif, k = 'all')
2 fit = bestfeat.fit(x,y)
3 dfscores = pd.DataFrame(fit.scores_)
4 dfcolumns = pd.DataFrame(x.columns)
```

```
1 fit = bestfeat.fit(x,y)
2 dfscores = pd.DataFrame(fit.scores_)
3 dfcolumns = pd.DataFrame(x.columns)
4 dfcolumns.head()
5 featureScores = pd.concat([dfcolumns,dfscores],axis = 1)
6 featureScores.columns = ['Feature', 'Score']
7 print(featureScores.nlargest(75,'Score'))
```

	Feature	Score
14	OverallQual	5.303071
60	MiscVal	3.564855
24	ExterQual	3.514221
42	KitchenQual	2.617125
49	GarageCars	2.578547
38	FullBath	2.435854
27	BsmtQual	2.334113
50	GarageArea	2.316328
16	YearBuilt	2.133300
23	MasVnrArea	1.852976
4	Street	1.835751
3	LotArea	1.826320
17	YearRemodAdd	1.813783
48	GarageFinish	1.811582
47	GarageYrBlt	1.725406
32	Heating	1.707885
43	TotRmsAbvGrd	1.656866
1	MSZoning	1.640044
45	Fireplaces	1.591973
34	CentralAir	1.557680
46	GarageType	1.544438
26	Foundation	1.528516
55	OpenPorchSF	1.460290
5	LotShape	1.407526
33	HeatingQC	1.358939
39	HalfBath	1.337597
9	Neighborhood	1.281079

Using SelectKBest and f_classif for measuring the respective ANOVA f-score values of the columns, the best 70 features were selected.

Using StandardScaler, the features were scaled by resizing the distribution values so that mean of the observed values in each feature column is 0 and standard deviation is 1.

From sklearn.model_selection import train_test_split, It is used to split the data into train and test data. Training data comprised 70% of total data whereas test data comprised 30% based on the best random state that would result in best model accuracy.

Finding Best Random State For Each Model

```
: 1 # finding Best Random state
2 maxAccu=0
3 maxRS=0
4
5 for i in range(1, 1000):
6     X_train, X_test, y_train, y_test = train_test_split(x, y, test_size=0.3, random_state=i)
7     lr=LinearRegression()
8     lr.fit(X_train, y_train)
9     pred = lr.predict(X_test)
10    r2 = r2_score(y_test, pred)
11
12    if r2>maxAccu:
13        maxAccu=r2
14        maxRS=i
15
16 print("Best r2 score is", maxAccu,"on Random State", maxRS)
```

By Using Above codes, I had found best random state for each model.

4.2 The model algorithms used

Linear Regression:

```
: 1 # Splitting the data into train and test
2 X_train, X_test, y_train, y_test = train_test_split(x, y, train_size=0.7, test_size=0.3, random_state=135)
3
4
5
6 results=pd.DataFrame(columns=['MAE', 'MSE', 'RMSE', 'R2-score'])
7
8 for method,func in regressors.items():
9     model = func.fit(X_train,y_train)
10    pred = model.predict(X_test)
11    results.loc[method]= [np.round(mean_absolute_error(y_test,pred),3),
12                          np.round(mean_squared_error(y_test,pred),3),
13                          np.sqrt(mean_squared_error(y_test,pred)),
14                          np.round(r2_score(y_test,pred),3)
15
16 ]
```

1	results				
		MAE	MSE	RMSE	R2-score
Linear Regression		21049.16	9.239717e+08	30396.90339	0.87

Cross Validation of Model

```

1 y_pred = lr.predict(X_test)
2 from sklearn.model_selection import cross_val_score
3 lss = r2_score(y_test,y_pred)

1 for j in range(4,10):
2     isscore = cross_val_score(lr,x,y,cv=j)
3     lsc = isscore.mean()
4     print("At cv:-",j)
5     print('Cross validation score is:-',lsc*100)
6     print('accuracy_score is:-',lss*100)
7     print('\n')

```

At cv:- 4
Cross validation score is:- 76.95323033369756
accuracy_score is:- 87.2139667319613

At cv:- 5
Cross validation score is:- 76.3315189979477
accuracy_score is:- 87.2139667319613

At cv:- 6
Cross validation score is:- 77.21221871690848
accuracy_score is:- 87.2139667319613

At cv:- 7
Cross validation score is:- 75.726237523794
accuracy_score is:- 87.2139667319613

```

1 lsscore_selected = cross_val_score(lr,x,y,cv=8).mean()
2 print("The cv score is: ",lsscore_selected,"\nThe accuracy score is: ",lss)

```

The cv score is: 0.7732514097659455
The accuracy score is: 0.8721396673196131

Random Forest Regression:

```
: 1 # Splitting the data into train and test
2 X_train, X_test, y_train, y_test = train_test_split(x, y, train_size=0.7, test_size=0.3, random_state=135)
```

```
: 1 regressors = {
2
3     'Random Forest' : RandomForestRegressor(),
4
5 }
6
7 results=pd.DataFrame(columns=['MAE', 'MSE', 'RMSE', 'R2-score'])
8
9 for method,func in regressors.items():
10     model = func.fit(X_train,y_train)
11     pred = model.predict(X_test)
12     results.loc[method]= [np.round(mean_absolute_error(y_test,pred),3),
13                           np.round(mean_squared_error(y_test,pred),3),
14                           np.sqrt(mean_squared_error(y_test,pred)),
15                           np.round(r2_score(y_test,pred),3)
16
17 ]
```

```
: 1 results
```

```
:

```

	MAE	MSE	RMSE	R2-score
Random Forest	19674.147	7.932564e+08	28164.808191	0.888

Cross Validation of model

```
1 rf = RandomForestRegressor()
2 rf.fit(X_train,y_train)
3 y_pred = rf.predict(X_test)
4 lss = r2_score(y_test,y_pred)
```

```
1 for j in range(4,10):
2     isscore = cross_val_score(rf,x,y,cv=j)
3     lsc = isscore.mean()
4     print("At cv:-",j)
5     print('Cross validation score is:-',lsc*100)
6     print('accuracy_score is:-',lss*100)
7     print('\n')
```

```
At cv:- 4
Cross validation score is:- 80.51160233099495
accuracy_score is:- 88.87169946361256
```

```
At cv:- 5
Cross validation score is:- 81.15381458000097
accuracy_score is:- 88.87169946361256
```

```
At cv:- 6
Cross validation score is:- 81.78879590456685
accuracy_score is:- 88.87169946361256
```

```

1 lsscore_selected = cross_val_score(rf,x,y,cv=6).mean()
2 print("The cv score is: ",lsscore_selected,"\nThe accuracy score is: ",lss)

```

The cv score is: 0.8159460782039397

The accuracy score is: 0.8887169946361255

Gradient Boosting Regressor:

```

]: 1 # Splitting the data into train and test
2 X_train, X_test, y_train, y_test = train_test_split(x, y, train_size=0.7, test_size=0.3, random_state=135)

```

```

]: 1 regressors = {
2
3     'Gradient Boost Regressor' : GradientBoostingRegressor(),
4
5 }
6
7 results=pd.DataFrame(columns=['MAE', 'MSE', 'RMSE', 'R2-score'])
8
9 for method,func in regressors.items():
10     model = func.fit(X_train,y_train)
11     pred = model.predict(X_test)
12     results.loc[method]= [np.round(mean_absolute_error(y_test,pred),3),
13                           np.round(mean_squared_error(y_test,pred),3),
14                           np.sqrt(mean_squared_error(y_test,pred)),
15                           np.round(r2_score(y_test,pred),3)
16
17 ]

```

```
]: 1 results
```

```
]:
```

	MAE	MSE	RMSE	R2-score
Gradient Boost Regressor	19343.363	7.566516e+08	27507.30013	0.893

Cross Validation of Model

```

: 1 gbr = GradientBoostingRegressor()
2 gbr.fit(X_train,y_train)
3 y_pred = gbr.predict(X_test)
4 from sklearn.model_selection import cross_val_score
5 lss = r2_score(y_test,y_pred)

```

```

: 1 for j in range(4,10):
2     isscore = cross_val_score(gbr,x,y,cv=j)
3     lsc = isscore.mean()
4     print("At cv:-",j)
5     print('Cross validation score is:-',lsc*100)
6     print('accuracy_score is:-',lss*100)
7     print('\n')

```

At cv:- 4

Cross validation score is:- 80.79615011670003

accuracy_score is:- 89.26980978643795

At cv:- 5

Cross validation score is:- 81.67587269317079

accuracy_score is:- 89.26980978643795

At cv:- 6

Cross validation score is:- 81.7989606484789

accuracy_score is:- 89.26980978643795

At cv:- 7

Cross validation score is:- 80.22241734713116

```

1 lsscore_selected = cross_val_score(gbr,x,y,cv=8).mean()
2 print("The cv score is: ",lsscore_selected,"\nThe accuracy score is: ",lss)

```

The cv score is: 0.8275259580055516
The accuracy score is: 0.8926980978643795

KNeighbors Regressor:

```

: 1 # Splitting the data into train and test
2 X_train, X_test, y_train, y_test = train_test_split(x, y, train_size=0.7, test_size=0.3, random_state=88)

```

```

: 1 regressors = {
2     'KNN Regressor': KNeighborsRegressor()
3 }
4
5 results=pd.DataFrame(columns=['MAE', 'MSE', 'RMSE', 'R2-score'])
6
7
8 for method,func in regressors.items():
9     model = func.fit(X_train,y_train)
10    pred = model.predict(X_test)
11    results.loc[method]= [np.round(mean_absolute_error(y_test,pred),3),
12                          np.round(mean_squared_error(y_test,pred),3),
13                          np.sqrt(mean_squared_error(y_test,pred)),
14                          np.round(r2_score(y_test,pred),3)
15
16 ]

```

```

: 1 results

```

```

:

```

	MAE	MSE	RMSE	R2-score
KNN Regressor	24728.548	1.254821e+09	35423.452651	0.798

Cross Validation of Model

```

1 knn = KNeighborsRegressor()
2 knn.fit(X_train,y_train)
3 y_pred = knn.predict(X_test)
4 lss = r2_score(y_test,y_pred)

```

```

1 for j in range(4,10):
2     isscore = cross_val_score(knn,x,y,cv=j)
3     lsc = isscore.mean()
4     print("At cv:-",j)
5     print('Cross validation score is:-',lsc*100)
6     print('accuracy_score is:-',lss*100)
7     print('\n')

```

At cv:- 4
Cross validation score is:- 68.74970165184517
accuracy_score is:- 79.79925665605202

At cv:- 5
Cross validation score is:- 68.44350007773714
accuracy_score is:- 79.79925665605202

At cv:- 6
Cross validation score is:- 68.81050093745115
accuracy_score is:- 79.79925665605202

At cv:- 7
Cross validation score is:- 67.54686991175045
accuracy_score is:- 79.79925665605202

```

1 lsscore_selected = cross_val_score(knn,x,y,cv=9).mean()
2 print("The cv score is: ",lsscore_selected,"\nThe accuracy score is: ",lss)

```

The cv score is: 0.6896135094821572
The accuracy score is: 0.7979925665605202

XGBOOST Regressor:

```

: 1 # Splitting the data into train and test
2 X_train, X_test, y_train, y_test = train_test_split(x, y, train_size=0.7, test_size=0.3, random_state=248)

```

```

: 1 regressors = {
2
3     'XG Boost Regressor' : XGBRegressor()
4
5 }
6
7 results=pd.DataFrame(columns=['MAE', 'MSE', 'RMSE', 'R2-score'])
8
9 for method,func in regressors.items():
10     model = func.fit(X_train,y_train)
11     pred = model.predict(X_test)
12     results.loc[method]= [np.round(mean_absolute_error(y_test,pred),3),
13                           np.round(mean_squared_error(y_test,pred),3),
14                           np.sqrt(mean_squared_error(y_test,pred)),
15                           np.round(r2_score(y_test,pred),3)
16
17 ]

```

```

: 1 results

```

```

:

```

	MAE	MSE	RMSE	R2-score
XG Boost Regressor	19071.368	7.928610e+08	28157.787913	0.878

Cross Validation of model

```

1 xgb = XGBRegressor()
2 xgb.fit(X_train,y_train)
3
4 y_pred = xgb.predict(X_test)
5 lss = r2_score(y_test,y_pred)

```

```

1 for j in range(4,10):
2     isscore = cross_val_score(xgb,x,y,cv=j)
3     lsc = isscore.mean()
4     print("At cv:-",j)
5     print('Cross validation score is:-',lsc*100)
6     print('accuracy_score is:-',lss*100)
7     print('\n')

```

At cv:- 4
Cross validation score is:- 79.58022840101758
accuracy_score is:- 87.81513506055637

At cv:- 5
Cross validation score is:- 80.0349543348284
accuracy_score is:- 87.81513506055637

At cv:- 6
Cross validation score is:- 80.21661569357038
accuracy_score is:- 87.81513506055637

At cv:- 7
Cross validation score is:- 79.59236402424216
accuracy_score is:- 87.81513506055637


```

1 lsscore_selected = cross_val_score(xgb,x,y,cv=9).mean()
2 print("The cv score is: ",lsscore_selected,"\nThe accuracy score is: ",lss)

```

The cv score is: 0.8153907083689376
The accuracy score is: 0.8781513506055637

Regularization

```

: 1 from sklearn.model_selection import GridSearchCV
2 from sklearn.linear_model import Lasso

: 1 parameters = {'alpha':[0.0001,0.001,0.01,0.1,1,10],
2               'random_state':list(range(0,10))}
3
4 ls = Lasso()
5 clf = GridSearchCV(ls,parameters)
6 clf.fit(X_train,y_train)
7 clf.best_params_

: {'alpha': 10, 'random_state': 0}

```

```

: 1 ls = Lasso(alpha=10,random_state=0)
2 ls.fit(X_train,y_train)
3 ls_score_training = ls.score(X_train,y_train)
4 pred_ls = ls.predict(X_test)
5 ls_score_training*100

: 81.57014236088625

```

```

: 1 pred = r2_score(y_test,pred_ls)
2 pred*100

: 79.96281560503131

```

```

: 1 cv_score = cross_val_score(ls,x,y,cv = 4)
2 cv_mean = cv_score.mean()
3 cv_mean*100

```

: 76.99018112764352

Mean Squared Error and Root Mean Squared Error metrics were used to evaluate the Model performance. The advantage of MSE and RMSE being that it is easier to compute the gradient. As, we take square of the error, the effect of larger errors become more pronounced than smaller error, hence the model can now focus more on the larger errors.

Cross validation is a technique for assessing how the statistical analysis generalizes to an independent data set. It is a technique for evaluating machine learning models by training several models on subsets of the available input data and evaluating them on the complementary subset of the data. Using cross-validation, there are high chances that we can detect over-fitting with ease.

4.3 Interpretation of the results

Based on comparing Accuracy Score results with Cross Validation results, it is determined that Random Forest Regressor is the best model. It also has the lowest Root Mean Squared Error score.

4.4 Hyperparameter Tuning

```
1 from sklearn.model_selection import GridSearchCV
2 from sklearn.ensemble import RandomForestRegressor

1 # Splitting the data into train and test
2 X_train, X_test, y_train, y_test = train_test_split(x, y, train_size=0.7, test_size=0.3, random_state=135)

1 rf = RandomForestRegressor()

1 parameters = {'criterion':['mse','mae'],
2               'max_features':['auto','sqrt','log2'],
3               'min_samples_split':[2,4,6,8,10],
4               'min_samples_leaf':[1,3,4,5,6,7],
5               'max_depth':[5,10,15],
6               }

1 grd = GridSearchCV(rf,param_grid=parameters)

1 grd.fit(X_train,y_train)
2
3 grd.best_params_

{'criterion': 'mae',
 'max_depth': 15,
 'max_features': 'sqrt',
 'min_samples_leaf': 1,
 'min_samples_split': 2}
```

```

1 rf = RandomForestRegressor(min_samples_split=2, min_samples_leaf=1, max_features='auto', max_depth=25, criterion='mse')
2
3 rf.fit(X_train,y_train)
4 y_pred = rf.predict(X_test)
5 print(r2_score(y_test,y_pred))

```

0.8923695356137124

Based on the input parameter values and after fitting the train datasets The Random Forest Regressor model was further tuned based on the parameter values yielded from GridsearchCV. The Random Forest Regressor model displayed an accuracy of 89.26%

The Model Test on Testing Dataset

```

1 Prediction_accuracy = pd.DataFrame({'Predictions': mod.predict(x), 'Actual Values': y[0:292]})
2 Prediction_accuracy.head(30)

```

	Predictions	Actual Values
0	345389.76	128000
1	248879.40	268000
2	244240.28	269790
3	192829.35	190000
4	204121.40	215000
5	94568.16	219210
6	142821.08	121500
7	323740.84	155000
8	234132.16	140000
9	180091.54	118500
10	104211.90	119500
11	149221.20	237000
12	134270.00	201000
13	166900.15	126500
14	311552.05	135500
15	117117.61	165000
16	128407.50	120500
17	132008.38	194500

5. Conclusions

5.1 Key Finding and Conclusions

- Structural attributes of the house Structural attributes of the house like lot size, lot shape, quality and condition of the house, garage capacity, rooms, Lot frontage, number of bedrooms, bathrooms, overall finishing of the house etc. play a big role in influencing the house price.
- Neighborhood qualities can be included in deciding house price.
- Training dataset is small because of this Data loss is huge so I am not removing outliers.

5.2 Limitation of this works and scope for future works

The housing market is affected by economic status, interest rate, real income and population density change. In contract to these market – side considerations, the available inventory can decide house price. For cycle of rising demand and limited supply, house prices will go up and threat of insecurity will increase.

In short economic factors may affect the houses price but, in this dataset, no economic features are used.