

HOUSE PRICE PREDICTION USING REGRESSION TECHNIQUE



A MINI PROJECT REPORT

Submitted by:
Raganeer Verma
(Intern FlipRobo)

INTRODUCTION

House Price Is a very important topic for real state business. Thousands of houses are sold every day and There are some questions every buyer asks himself like: -

1. What is the actual price that this house deserves?
2. Am I paying a fair price?

Lots of confusion are coming in mind before buying a new house, So in this project, we proposed a machine learning model to predict a house price based on historical data of property markets.

Data Description: -

Our dataset contains two files:

1. Train data file for tanning a model (it contains 1460 rows and 81 columns)
2. Test data file for predict a model (it contains 1459 rows and 80 columns)
Test data does not have Target column Sale Price.

| Id | MSSubClass | MSZoning | LotFrontage | LotArea | Street | Alley | LotShape | LandContour | Utilities | LotConfig | LandSlope | Neighborhood | Condition1 | Condi |
|----|------------|----------|-------------|---------|--------|-------|----------|-------------|-----------|-----------|-----------|--------------|------------|-------|
| 1 | 60 | RL | 65.0 | 8450 | Pave | NaN | Reg | | Lvl | AllPub | Inside | Gtl | CollgCr | Norm |
| 2 | 20 | RL | 80.0 | 9600 | Pave | NaN | Reg | | Lvl | AllPub | FR2 | Gtl | Veenker | Feedr |
| 3 | 60 | RL | 68.0 | 11250 | Pave | NaN | IR1 | | Lvl | AllPub | Inside | Gtl | CollgCr | Norm |
| 4 | 70 | RL | 60.0 | 9550 | Pave | NaN | IR1 | | Lvl | AllPub | Corner | Gtl | Crawfor | Norm |
| 5 | 60 | RL | 84.0 | 14260 | Pave | NaN | IR1 | | Lvl | AllPub | FR2 | Gtl | NoRidge | Norm |

#details of test data
df_test.head()

| Id | MSSubClass | MSZoning | LotFrontage | LotArea | Street | Alley | LotShape | LandContour | Utilities | LotConfig | LandSlope | Neighborhood | Condition1 | Con |
|------|------------|----------|-------------|---------|--------|-------|----------|-------------|-----------|-----------|-----------|--------------|------------|-------|
| 1461 | 20 | RH | 80.0 | 11622 | Pave | NaN | Reg | | Lvl | AllPub | Inside | Gtl | Names | Feedr |
| 1462 | 20 | RL | 81.0 | 14267 | Pave | NaN | IR1 | | Lvl | AllPub | Corner | Gtl | Names | Norm |
| 1463 | 60 | RL | 74.0 | 13830 | Pave | NaN | IR1 | | Lvl | AllPub | Inside | Gtl | Gilbert | Norm |
| 1464 | 60 | RL | 78.0 | 9978 | Pave | NaN | IR1 | | Lvl | AllPub | Inside | Gtl | Gilbert | Norm |

Column name of dataset:

```

1 df.columns # All column name of dataset

Index(['Id', 'MSSubClass', 'MSZoning', 'LotFrontage', 'LotArea', 'Street',
      'Alley', 'LotShape', 'LandContour', 'Utilities', 'LotConfig',
      'LandSlope', 'Neighborhood', 'Condition1', 'Condition2', 'BldgType',
      'HouseStyle', 'OverallQual', 'OverallCond', 'YearBuilt', 'YearRemodAdd',
      'RoofStyle', 'RoofMatl', 'Exterior1st', 'Exterior2nd', 'MasVnrType',
      'MasVnrArea', 'ExterQual', 'ExterCond', 'Foundation', 'BsmtQual',
      'BsmtCond', 'BsmtExposure', 'BsmtFinType1', 'BsmtFinSF1',
      'BsmtFinType2', 'BsmtFinSF2', 'BsmtUnfSF', 'TotalBsmtSF', 'Heating',
      'HeatingQC', 'CentralAir', 'Electrical', '1stFlrSF', '2ndFlrSF',
      'LowQualFinSF', 'GrLivArea', 'BsmtFullBath', 'BsmtHalfBath', 'FullBath',
      'HalfBath', 'BedroomAbvGr', 'KitchenAbvGr', 'KitchenQual',
      'TotRmsAbvGrd', 'Functional', 'Fireplaces', 'FireplaceQu', 'GarageType',
      'GarageYrBlt', 'GarageFinish', 'GarageCars', 'GarageArea', 'GarageQual',
      'GarageCond', 'PavedDrive', 'WoodDeckSF', 'OpenPorchSF',
      'EnclosedPorch', '3SsnPorch', 'ScreenPorch', 'PoolArea', 'PoolQC',
      'Fence', 'MiscFeature', 'MiscVal', 'MoSold', 'YrSold', 'SaleType',
      'SaleCondition', 'SalePrice'],
      dtype='object')

```

Importing All Required Libraries:

Firstly, we import all the libraries which are used to predict our machine learning model.

```

1 #Importing required packages.
2 import pandas as pd
3 import numpy as np
4 import scipy.stats as st
5 #ploting libraries
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8 #feature engineering
9 from sklearn.preprocessing import StandardScaler, LabelEncoder
10 #train test split and cross validation
11 from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score
12 from scipy.stats import zscore
13 #metrics
14 from sklearn import metrics
15 from sklearn.metrics import confusion_matrix, classification_report, accuracy_score
16 from sklearn.metrics import r2_score, mean_squared_error
17 from sklearn.metrics import roc_curve, roc_auc_score
18 from sklearn.metrics import plot_roc_curve
19 #ML models for regression
20 from sklearn.linear_model import LinearRegression, Ridge, Lasso
21 from sklearn.svm import SVR
22 from sklearn.tree import DecisionTreeRegressor
23 from sklearn.ensemble import RandomForestRegressor
24 from sklearn.neighbors import KNeighborsRegressor
25 from sklearn.ensemble import AdaBoostRegressor
26 from sklearn.ensemble import ExtraTreesRegressor
27 from sklearn.ensemble import GradientBoostingRegressor

```

Data Analysis:

Before performing prediction task, we want to perform some Exploratory Data Analysis to gain a basic understanding of house price data.

Statistical Analysis:

We use the describe function for statistical analysis and it's showing the mean, std, max, median, 25% and 75% value of numerical data and top frequent unique for categorical type data. We observed that data and analyse about missing value, unique value, skewed and outliers are present in our dataset.

| | count | unique | top | freq | mean | std | min | 25% | 50% | 75% | max |
|---------------|--------|--------|--------|------|--------------|--------------|---------|----------|----------|----------|----------|
| Id | 1460.0 | NaN | NaN | NaN | 730.5 | 421.610009 | 1.0 | 365.75 | 730.5 | 1095.25 | 1460.0 |
| MSSubClass | 1460.0 | NaN | NaN | NaN | 56.89726 | 42.300571 | 20.0 | 20.0 | 50.0 | 70.0 | 190.0 |
| MSZoning | 1460 | 5 | RL | 1151 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| LotFrontage | 1201.0 | NaN | NaN | NaN | 70.049958 | 24.284752 | 21.0 | 59.0 | 69.0 | 80.0 | 313.0 |
| LotArea | 1460.0 | NaN | NaN | NaN | 10516.828082 | 9981.264932 | 1300.0 | 7553.5 | 9478.5 | 11601.5 | 215245.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| MoSold | 1460.0 | NaN | NaN | NaN | 6.321918 | 2.703626 | 1.0 | 5.0 | 6.0 | 8.0 | 12.0 |
| YrSold | 1460.0 | NaN | NaN | NaN | 2007.815753 | 1.328095 | 2006.0 | 2007.0 | 2008.0 | 2009.0 | 2010.0 |
| SaleType | 1460 | 9 | WD | 1267 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| SaleCondition | 1460 | 6 | Normal | 1198 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| SalePrice | 1460.0 | NaN | NaN | NaN | 180921.19589 | 79442.502883 | 34900.0 | 129975.0 | 163000.0 | 214000.0 | 755000.0 |

64 rows x 12 columns

Check missing value Unique value and duplicate value in dataset:

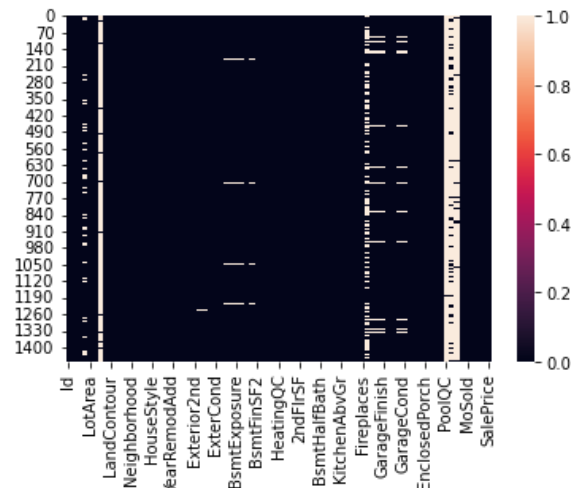
Lots of missing value is present but the column name is not shown there so we try to find out the column where missing value are present.

no duplicate present there because the row and column value is same.

```
In [18]: 1 df.duplicated().sum()
```

```
Out[18]: 0
```

| Unique Values | |
|---------------|------|
| Id | 1460 |
| MSSubClass | 15 |
| MSZoning | 5 |
| LotFrontage | 110 |
| LotArea | 1073 |
| ... | ... |
| MoSold | 12 |
| YrSold | 5 |
| SaleType | 9 |
| SaleCondition | 6 |
| SalePrice | 663 |



Here we check the percentage of nan values present in each feature:

I observe that sale price has 49.9% missing value because we join both train and test data but test data do not have this column that's why missing value present here so no need to impute this column because after cleaning treatment, I separate both the data again.

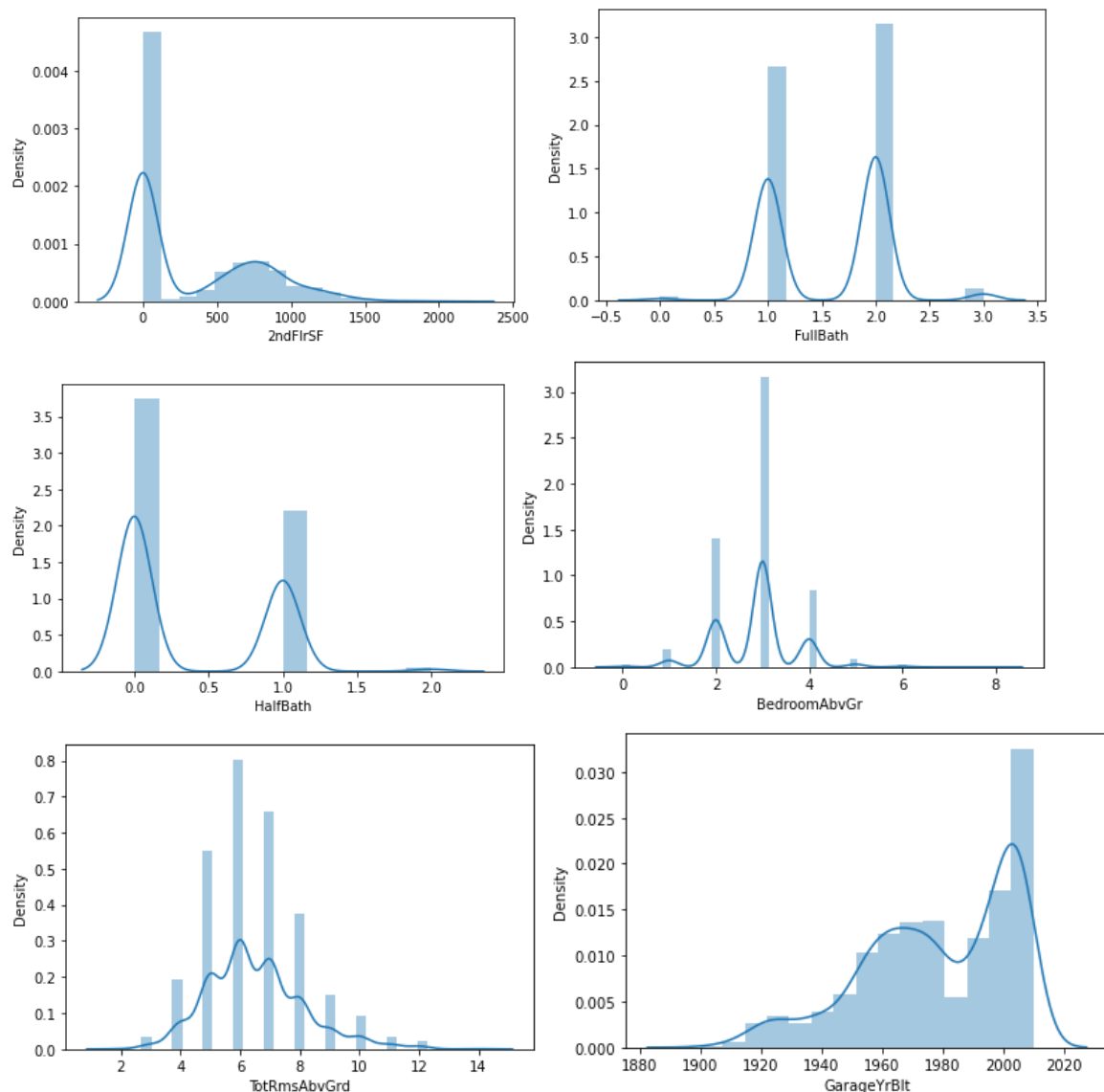
We might want to remove columns with more than 50% missing values

```
1 for features in feautres_na:
2     print(features,np.round(df[features].isnull().mean(),4), '%')
```

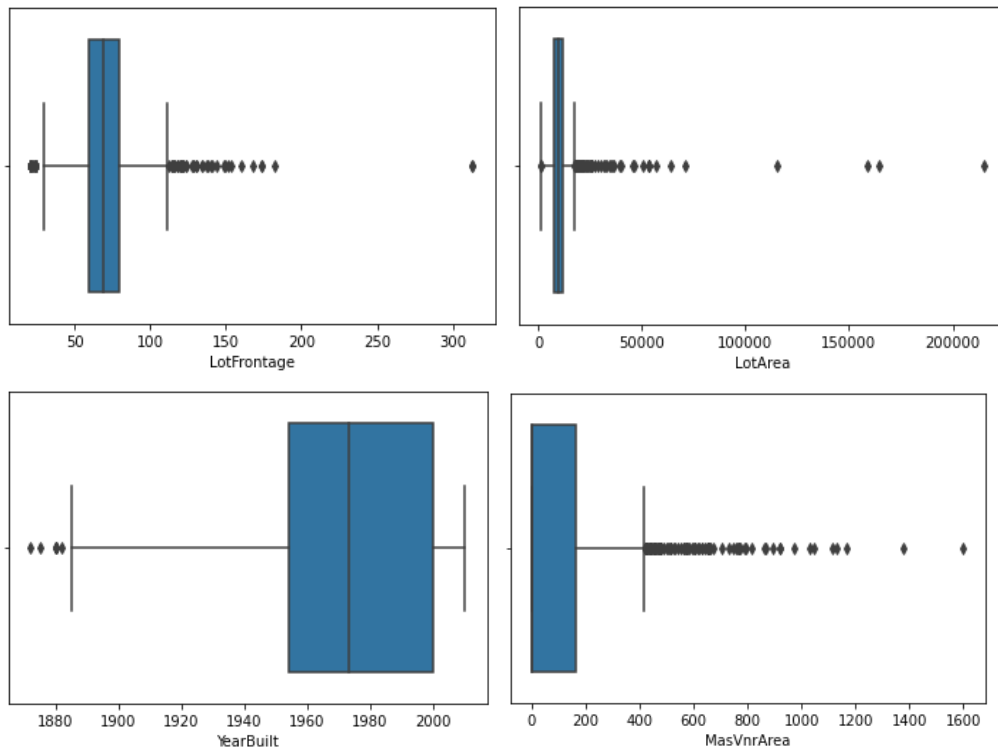
```
LotFrontage 0.1774 %
Alley 0.9377 %
MasVnrType 0.0055 %
MasVnrArea 0.0055 %
BsmtQual 0.0253 %
BsmtCond 0.0253 %
BsmtExposure 0.026 %
BsmtFinType1 0.0253 %
BsmtFinType2 0.026 %
Electrical 0.0007 %
FireplaceQu 0.4726 %
GarageType 0.0555 %
GarageYrBlt 0.0555 %
GarageFinish 0.0555 %
GarageQual 0.0555 %
GarageCond 0.0555 %
PoolQC 0.9952 %
Fence 0.8075 %
MiscFeature 0.963 %
```

Data Visualization: Data visualization is the graphical representation of information and data.

Skewness: distribution plot for each column that contain continues variable. We seen that Skewness is present so we need to remove skewness using any Transform Technique.



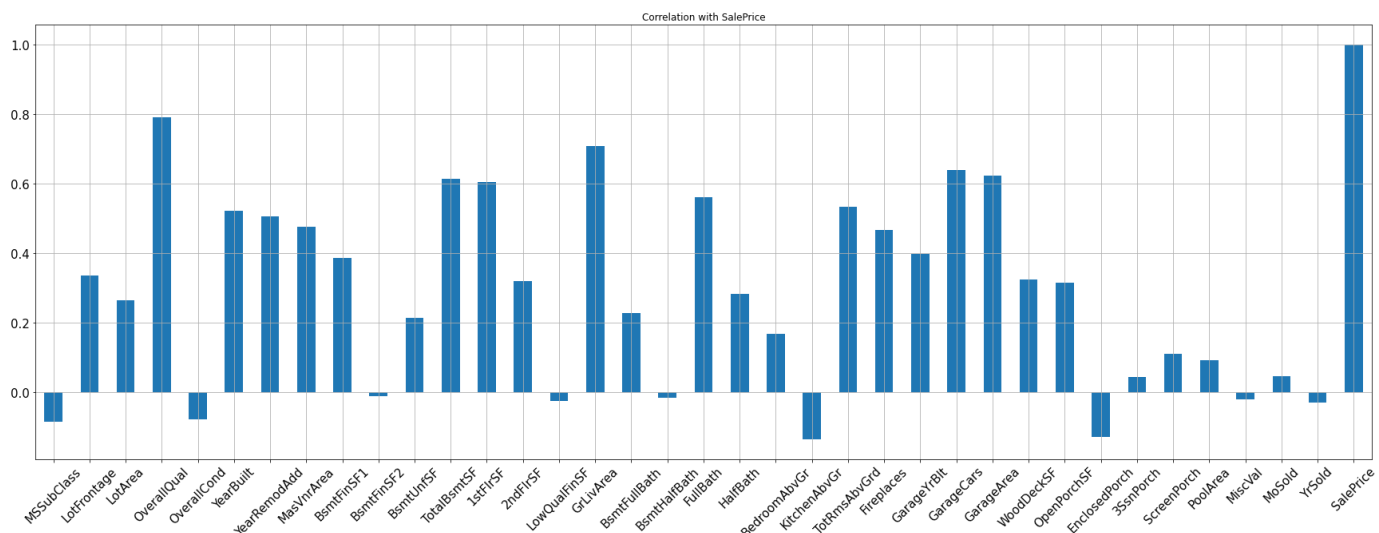
Outliers: Use box plot for visualizing outliers are present in our data. We observed that some outliers are there so we can remove it by using z-score or IQR technique.



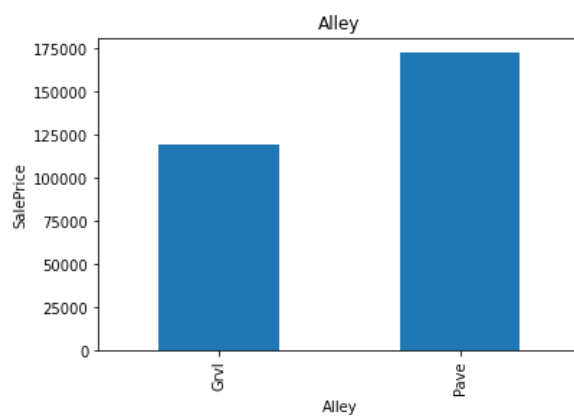
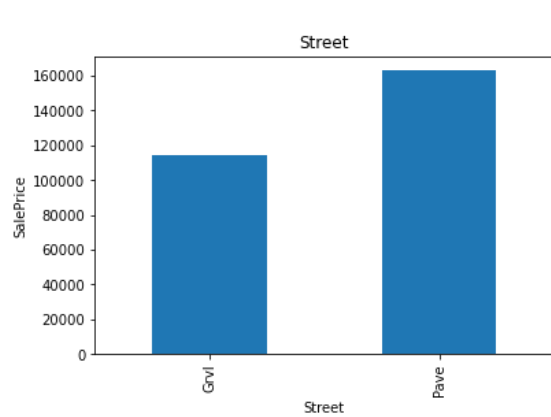
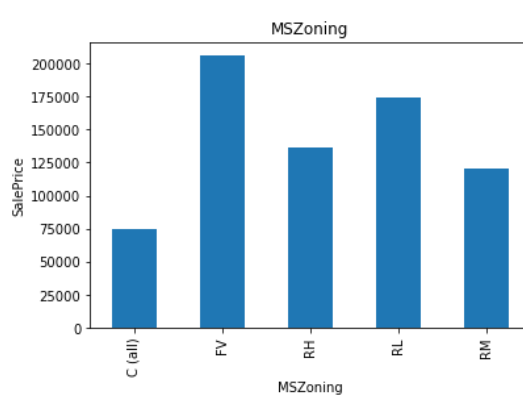
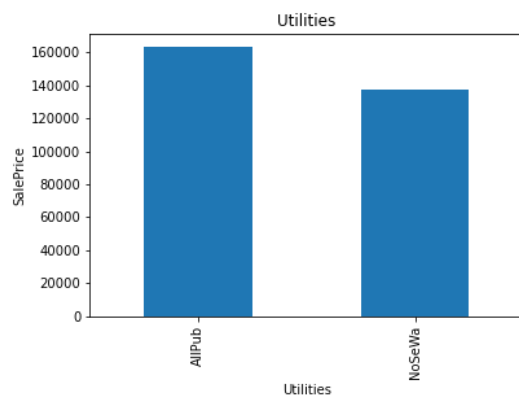
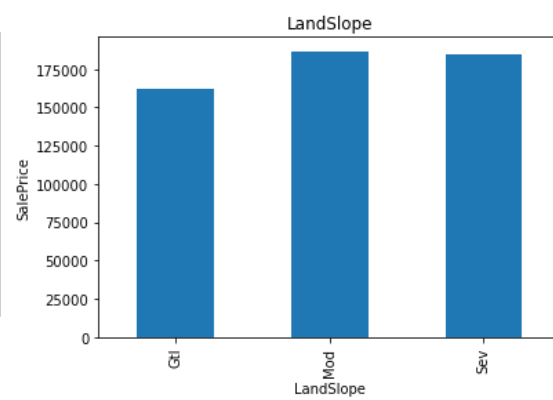
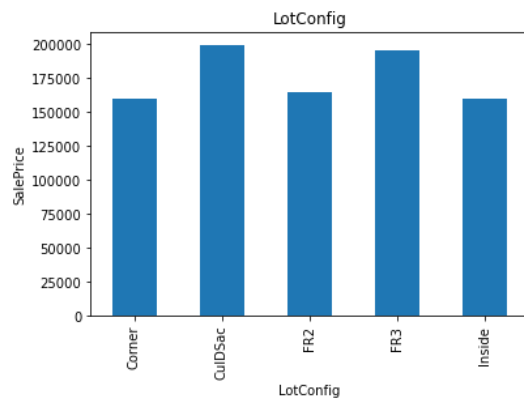
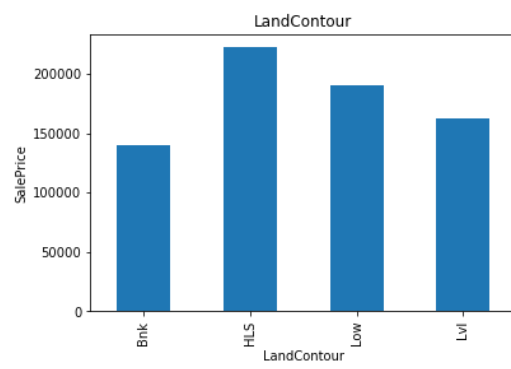
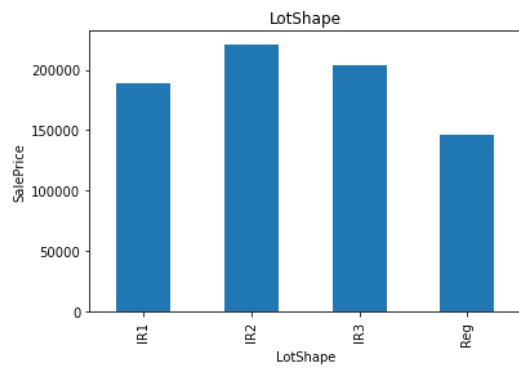
Correlation:

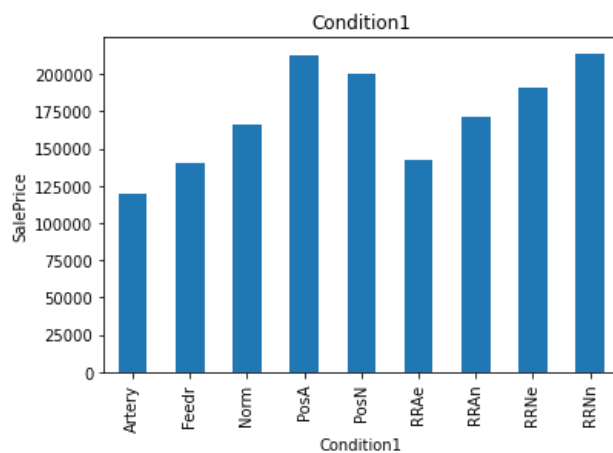
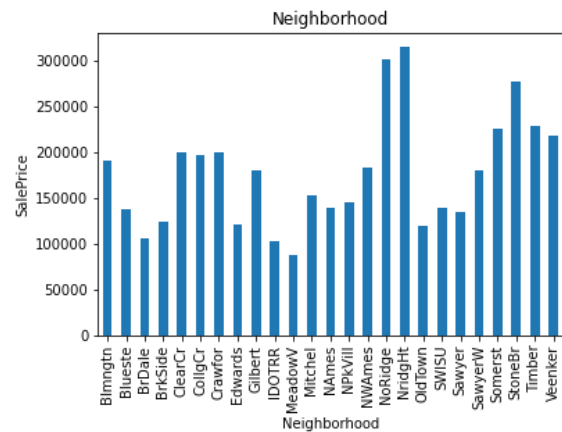
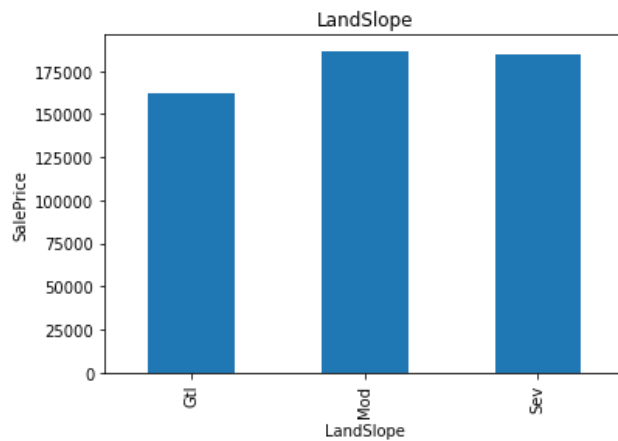
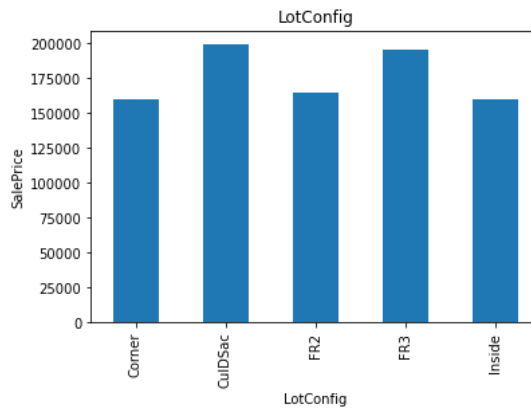
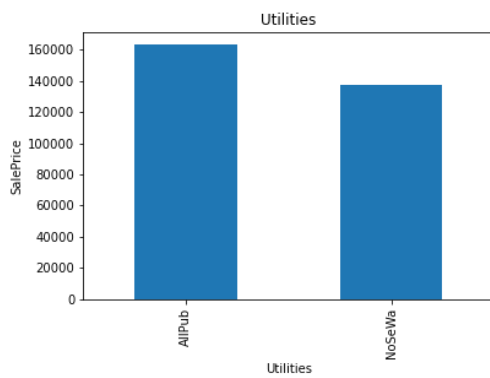
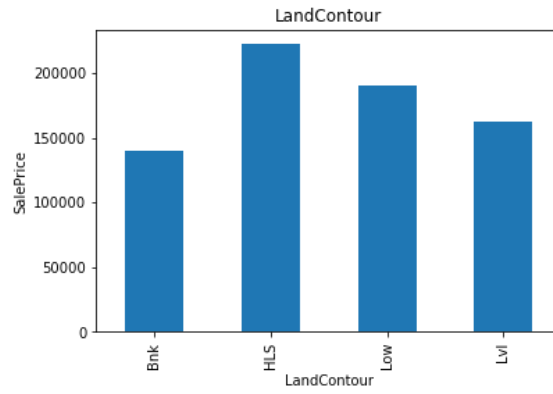
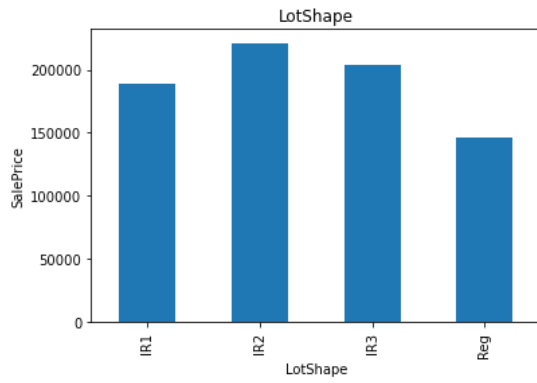
Checking the Correlation between the target value and rest columns.

We observed that most of the column is positively correlated with target variable.

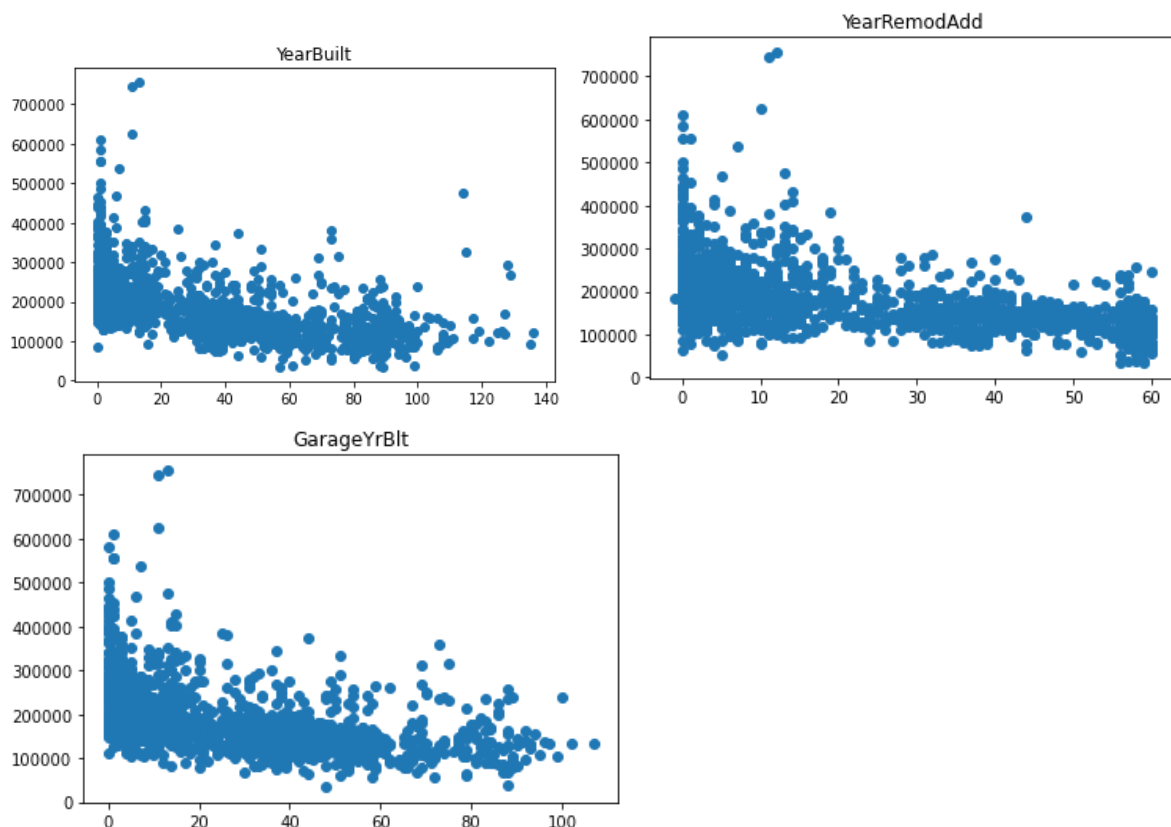


Bivariate Analysis: Data visualization between Discrete and Continues features and Target variable





Data visualization between Year Sold and Target variable:



Feature Engineering:

Handling missing value

```
1 for feature in categorical_features:
2     df[feature] = df[feature].fillna(df[feature].mode()[0]) #train
3     df_test[feature] = df_test[feature].fillna(df_test[feature].mode()[0]) #test
```

```
1 for feature in discrete_features:
2     df[feature] = df[feature].fillna(df[feature].mode()[0]) #train
3     df_test[feature] = df_test[feature].fillna(df_test[feature].mode()[0]) #test
```

```
1 for feature in year_features:
2     df[feature] = df[feature].fillna(df[feature].mode()[0]) #train
3     df_test[feature] = df_test[feature].fillna(df_test[feature].mode()[0]) #test
```

```
1 # as sales price is absent in test data so we want to remove it from the feature list
2 continous_numerical_features = [feature for feature in continous_features if feature not in 'SalePrice']
```

```
1 for feature in continous_numerical_features:
2     df[feature] = df[feature].fillna(df[feature].mean()) #train
3     df_test[feature] = df_test[feature].fillna(df_test[feature].mean()) #test
```

Label Encoder: Convert numerical variable into categorical variable.

```
1 from sklearn.preprocessing import LabelEncoder
2 le=LabelEncoder()
```

```
1 #for changing all catigorical value into numeric value
2
3 final_df=final_df.apply(le.fit_transform)
```

```
1 final_df.head()
```

| | MSSubClass | MSZoning | LotFrontage | LotArea | Street | LotShape | LandContour | Utilities | LotConfig | LandSlope | Neighborhood | Condition1 | Condition2 | Bldg1 |
|---|------------|----------|-------------|---------|--------|----------|-------------|-----------|-----------|-----------|--------------|------------|------------|-------|
| 0 | 5 | 3 | 41 | 619 | 1 | 3 | 3 | 0 | 4 | 0 | 5 | 2 | 2 | |
| 1 | 0 | 3 | 58 | 895 | 1 | 3 | 3 | 0 | 2 | 0 | 24 | 1 | 2 | |
| 2 | 5 | 3 | 44 | 1266 | 1 | 0 | 3 | 0 | 4 | 0 | 5 | 2 | 2 | |
| 3 | 6 | 3 | 36 | 883 | 1 | 0 | 3 | 0 | 0 | 0 | 6 | 2 | 2 | |
| 4 | 5 | 3 | 62 | 1670 | 1 | 0 | 3 | 0 | 2 | 0 | 15 | 2 | 2 | |

Drop Unnecessary columns:

Drop the column which is not useful for prediction purpose

```
1 for feature in more_than_50_percent_misssing_value_features:
2     df.drop([feature],axis = 1, inplace = True)
3     df_test.drop([feature],axis = 1, inplace = True)
```

We will also drop the id column as it is not useful for prediction

```
1 df.drop(['Id'],axis = 1, inplace = True)
2 df_test.drop(['Id'],axis = 1, inplace = True)
```

Scaling using Standard Scaler:

By using standard scaler, we arrange the value between 0 and 1.

```
1 from sklearn.preprocessing import StandardScaler
2 scaler=StandardScaler()
3 x=pd.DataFrame(scaler.fit_transform(x), columns=x.columns)
```

```
1 x.head()
```

| | MSSubClass | MSZoning | LotFrontage | LotArea | Street | LotShape | LandContour | Utilities | LotConfig | LandSlope | Neighborhood | Condition1 | Condition2 | Bldg1 |
|---|------------|-----------|-------------|-----------|----------|-----------|-------------|-----------|-----------|-----------|--------------|------------|------------|-------|
| 0 | 0.177030 | -0.045532 | -0.290681 | -0.508928 | 0.064238 | 0.750731 | 0.314667 | -0.02618 | 0.604670 | -0.225716 | -1.206215 | -0.036289 | -0.03174 | |
| 1 | -0.983043 | -0.045532 | 0.546875 | 0.000388 | 0.064238 | 0.750731 | 0.314667 | -0.02618 | -0.628316 | -0.225716 | 1.954302 | -1.188074 | -0.03174 | |
| 2 | 0.177030 | -0.045532 | -0.142877 | 0.685012 | 0.064238 | -1.378933 | 0.314667 | -0.02618 | 0.604670 | -0.225716 | -1.206215 | -0.036289 | -0.03174 | |
| 3 | 0.409045 | -0.045532 | -0.537022 | -0.021756 | 0.064238 | -1.378933 | 0.314667 | -0.02618 | -1.861302 | -0.225716 | -1.039872 | -0.036289 | -0.03174 | |
| 4 | 0.177030 | -0.045532 | 0.743947 | 1.430533 | 0.064238 | -1.378933 | 0.314667 | -0.02618 | -0.628316 | -0.225716 | 0.457215 | -0.036289 | -0.03174 | |

Model Building: Our target variable is Continues type of variable so we can use many Regressions model.

LinearRegression, Ridge, Lasso, DecisionTreeRegressor, RandomForestRegressor, KNeighborsRegressor, AdaBoostRegressor, ExtraTreesRegressor, GradientBoostingRegressor.

Now we use many regressions algorithm but the performance of ExtraTreesRegressor is best so this is our final model and perform hyperparameter tuning for this model.

The r2 value and cross validation score of ExtraTreesRegressor model is 90% .

```
1 # Choosing Extra Trees Regressor
2
3 param = {'criterion' : ['mse', 'mae'],
4          'min_samples_split' : [2, 3],
5          'n_jobs' : [-1, 1]
6          }
```

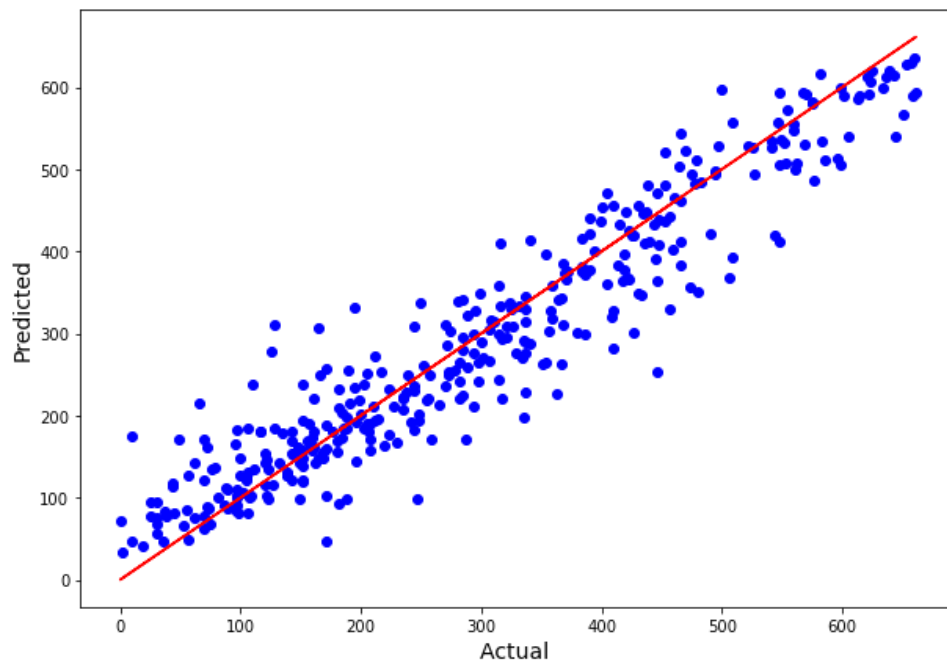
```
1 GSCV = GridSearchCV(ExtraTreesRegressor(),param, cv=5)
2 print(GSCV.fit(x_train,y_train))
3 print(GSCV.best_params_)
```

```
GridSearchCV(cv=5, estimator=ExtraTreesRegressor(),
              param_grid={'criterion': ['mse', 'mae'],
                           'min_samples_split': [2, 3], 'n_jobs': [-1, 1]})
{'criterion': 'mse', 'min_samples_split': 2, 'n_jobs': -1}
```

```
1 final_model = RandomForestRegressor(criterion = 'mse',min_samples_split = 2, n_jobs=-1)
2 final_model.fit(x_train,y_train)
3 predfm = final_model .predict(x_test)
4 scr = cross_val_score(final_model ,x,y,cv=5)
5 print("r2 score= ",r2_score(y_test,predfm)*100)
6 print('Cross validation score for RandomForestRegressor Model is ',scr.mean()*100)
7 print("MSE= ",mean_squared_error(y_test,predfm))
8 print("RMSE= ",np.sqrt(mean_squared_error(y_test,predfm)))
```

```
r2 score= 90.01634760346646
Cross validation score for RandomForestRegressor Model is 89.22800059636457
MSE= 2948.367931780822
RMSE= 54.29887597161493
```

Best fit Line for Final Model:



Saving and Loading Final Model:

```
1 import pickle
2 filename_1="house_price.pkl"
3 pickle.dump(final_model,open(filename_1,"wb"))
```

```
1 loaded_model=pickle.load(open("house_price.pkl","rb"))
2 result=loaded_model.score(x_test,y_test)
3 print(result*100)
```

90.01634760346646

Prediction and Conclusion:

Now we predict our saved model by using test data.

| | | |
|---|---|--|
| 1 | conclusion=pd.DataFrame([loaded_model.predict(x_test)[:],predfm[:]],index=["predicted","original"]) | |
| 1 | conclusion.T | |

| | predicted | original |
|-----|-----------|----------|
| 0 | 180.92 | 180.92 |
| 1 | 310.60 | 310.60 |
| 2 | 236.21 | 236.21 |
| 3 | 321.58 | 321.58 |
| 4 | 181.14 | 181.14 |
| ... | ... | ... |
| 360 | 421.84 | 421.84 |
| 361 | 556.73 | 556.73 |
| 362 | 110.51 | 110.51 |
| 363 | 249.68 | 249.68 |
| 364 | 411.04 | 411.04 |

365 rows × 2 columns

Conclusion/ Results

As a result, we can say, ExtraTreesRegressor is the best fit to this dataset, gives 90% accuracy. Gradient, lasso, ridge Models working very well on this Dataset. Ensembled all the models together for better results.

In [113]:



import pickle

.