



2022

# HOUSE PRICE PREDICTION



Submitted by:  
**Bhushan Kumar Sharma,**  
**Intern Data Scientist**

## **ACKNOWLEDGMENT**

\*\*\*\*\*

I would like to express my special gratitude to “Flip Robo” team, who has given me this opportunity to deal with this dataset and it has helped me to improve my model building skills.

A huge thanks to my academic team “Datatrained” who is 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 to many other persons who has helped me directly or indirectly to complete the project.

**Following are the external references which I used:**

[www.geeksforgeeks.org](http://www.geeksforgeeks.org)

[www.stackoverflow.com](http://www.stackoverflow.com)

[www.w3school.com](http://www.w3school.com)

[www.google.com](http://www.google.com)

Datatrained Lectures

## . INTRODUCTION

### ➤ **Business 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 data is provided in the CSV file below.

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?

### ➤ **Conceptual Background of the Domain Problem**

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

for the management to understand the pricing dynamics of a new market.

Determining the listed price of a house is a challenging task, due to the many factors that price on the market. The focus of this project is developing machine learning models that can accurately predict the price of a house based on its features, in order to make informed purchases.

We implement and evaluate various learning methods on a dataset consisting of the sale prices of different makes and models across features. Our results show that GradientBoostingRegressor is giving best results.

### ➤ **Review of Literature**

- Two Datasets are provided, one is train dataset and second is test dataset.
- Prize will get decided based on various feature column
- After Pre-processing and EDA, data is used to build Machine learning model, as we have to predict prize of house, which clearly indicating this is a regression problem.
- Various pre-processing techniques are used in this project to give best data to our machine learning model, and at final it's giving good accuracy.

### ➤ **Motivation for the Problem Undertaken**

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

## ➤ Data Sources and their formats

This dataset is provided by fliprobo for building machine learning model to predict house price based on given feature column.

Dataset is in two part one is for building model and second is for which we have to predict price.

### **Train dataset:**

Dataset is having 1168 rows and 81 columns including target. In this particular datasets I have object, float and integer types of columns.

### **Test dataset:**

Dataset is having 292 rows and 80 columns and we have to predict price of house after creating ML model.

The information about features is as follows.

| train_df.head() |     |            |          |             |         |        |       |          |             |           |           |           |              |            |      |  |
|-----------------|-----|------------|----------|-------------|---------|--------|-------|----------|-------------|-----------|-----------|-----------|--------------|------------|------|--|
|                 | Id  | MSSubClass | MSZoning | LotFrontage | LotArea | Street | Alley | LotShape | LandContour | Utilities | LotConfig | LandSlope | Neighborhood | Condition1 | Cond |  |
| 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    | CuiDSac   | 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       |      |  |

| train_df.head() |      |            |          |            |             |             |           |              |           |          |             |             |            |          |  |  |
|-----------------|------|------------|----------|------------|-------------|-------------|-----------|--------------|-----------|----------|-------------|-------------|------------|----------|--|--|
|                 | 1    | Condition2 | BldgType | HouseStyle | OverallQual | OverallCond | YearBuilt | YearRemodAdd | RoofStyle | RoofMatl | Exterior1st | Exterior2nd | MasVnrType | MasVnrAr |  |  |
| m               | Norm | TwnhsE     | 1Story   | 6          | 5           | 1976        | 1976      | Gable        | CompShg   | Plywood  | Plywood     | Plywood     | None       | 0        |  |  |
| m               | Norm | 1Fam       | 1Story   | 8          | 6           | 1970        | 1970      | Flat         | Tar&Grv   | Wd Sdng  | Wd Sdng     | Wd Sdng     | None       | 0        |  |  |
| m               | Norm | 1Fam       | 2Story   | 7          | 5           | 1996        | 1997      | Gable        | CompShg   | MetalSd  | MetalSd     | MetalSd     | None       | 0        |  |  |
| m               | Norm | 1Fam       | 1Story   | 6          | 6           | 1977        | 1977      | Hip          | CompShg   | Plywood  | Plywood     | Plywood     | BrkFace    | 480      |  |  |
| m               | Norm | 1Fam       | 1Story   | 6          | 7           | 1977        | 2000      | Gable        | CompShg   | CemntBd  | CemntBd     | CemntBd     | Stone      | 126      |  |  |

| train_df.head() |            |           |           |            |          |          |              |              |            |              |            |           |     |  |  |  |
|-----------------|------------|-----------|-----------|------------|----------|----------|--------------|--------------|------------|--------------|------------|-----------|-----|--|--|--|
|                 | MasVnrArea | ExterQual | ExterCond | Foundation | BsmtQual | BsmtCond | BsmtExposure | BsmtFinType1 | BsmtFinSF1 | BsmtFinType2 | BsmtFinSF2 | BsmtUnfSF | Tot |  |  |  |
|                 | 0.0        | TA        | TA        | CBlock     | Gd       | TA       | No           | ALQ          | 120        | Unf          | 0          | 958       |     |  |  |  |
|                 | 0.0        | Gd        | Gd        | PConc      | TA       | Gd       | Gd           | ALQ          | 351        | Rec          | 823        | 1043      |     |  |  |  |
|                 | 0.0        | Gd        | TA        | PConc      | Gd       | TA       | Av           | GLQ          | 862        | Unf          | 0          | 255       |     |  |  |  |
|                 | 480.0      | TA        | TA        | CBlock     | Gd       | TA       | No           | BLQ          | 705        | Unf          | 0          | 1139      |     |  |  |  |
|                 | 126.0      | Gd        | TA        | CBlock     | Gd       | TA       | No           | ALQ          | 1246       | Unf          | 0          | 356       |     |  |  |  |

## HOUSE PRICE PREDICTION MODEL

| train_df.head() |         |           |            |            |          |          |              |           |              |              |          |          |          |  |
|-----------------|---------|-----------|------------|------------|----------|----------|--------------|-----------|--------------|--------------|----------|----------|----------|--|
| TotalBsmtSF     | Heating | HeatingQC | CentralAir | Electrical | 1stFlrSF | 2ndFlrSF | LowQualFinSF | GrLivArea | BsmtFullBath | BsmtHalfBath | FullBath | HalfBath | Bedrooms |  |
| 1078            | GasA    | TA        | Y          | SBrkr      | 958      | 0        | 0            | 958       | 0            | 0            | 2        | 0        |          |  |
| 2217            | GasA    | Ex        | Y          | SBrkr      | 2217     | 0        | 0            | 2217      | 1            | 0            | 2        | 0        |          |  |
| 1117            | GasA    | Ex        | Y          | SBrkr      | 1127     | 886      | 0            | 2013      | 1            | 0            | 2        | 1        |          |  |
| 1844            | GasA    | Ex        | Y          | SBrkr      | 1844     | 0        | 0            | 1844      | 0            | 0            | 2        | 0        |          |  |
| 1602            | GasA    | Gd        | Y          | SBrkr      | 1602     | 0        | 0            | 1602      | 0            | 1            | 2        | 0        |          |  |

| train_df.head() |              |             |              |            |            |             |            |             |              |            |            |            |             |
|-----------------|--------------|-------------|--------------|------------|------------|-------------|------------|-------------|--------------|------------|------------|------------|-------------|
| BedroomAbvGr    | KitchenAbvGr | KitchenQual | TotRmsAbvGrd | Functional | Fireplaces | FireplaceQu | GarageType | GarageYrBlt | GarageFinish | GarageCars | GarageArea | WoodDeckSF | OpenPorchSF |
| 2               | 1            | TA          | 5            | Typ        | 1          | TA          | Attchd     | 1977.0      | RFn          | 2          | 4          |            |             |
| 4               | 1            | Gd          | 8            | Typ        | 1          | TA          | Attchd     | 1970.0      | Unf          | 2          | 6          |            |             |
| 3               | 1            | TA          | 8            | Typ        | 1          | TA          | Attchd     | 1997.0      | Unf          | 2          | 4          |            |             |
| 3               | 1            | TA          | 7            | Typ        | 1          | TA          | Attchd     | 1977.0      | RFn          | 2          | 5          |            |             |
| 3               | 1            | Gd          | 8            | Typ        | 1          | TA          | Attchd     | 1977.0      | Fin          | 2          | 5          |            |             |

| train_df.head() |            |            |            |            |             |               |           |             |          |        |       |             |        |        |
|-----------------|------------|------------|------------|------------|-------------|---------------|-----------|-------------|----------|--------|-------|-------------|--------|--------|
| GarageArea      | GarageQual | GarageCond | PavedDrive | WoodDeckSF | OpenPorchSF | EnclosedPorch | 3SsnPorch | ScreenPorch | PoolArea | PoolQC | Fence | MiscFeature | MoSold | YrSold |
| 440             | TA         | TA         | Y          | 0          | 205         | 0             | 0         | 0           | 0        | Nan    | Nan   | Nan         |        |        |
| 621             | TA         | TA         | Y          | 81         | 207         | 0             | 0         | 224         | 0        | Nan    | Nan   | Nan         |        |        |
| 455             | TA         | TA         | Y          | 180        | 130         | 0             | 0         | 0           | 0        | Nan    | Nan   | Nan         |        |        |
| 546             | TA         | TA         | Y          | 0          | 122         | 0             | 0         | 0           | 0        | Nan    | MnPrv | Nan         |        |        |
| 529             | TA         | TA         | Y          | 240        | 0           | 0             | 0         | 0           | 0        | Nan    | Nan   | Nan         |        |        |

| train_df.head() |               |           |             |          |        |       |             |         |        |        |          |               |           |  |
|-----------------|---------------|-----------|-------------|----------|--------|-------|-------------|---------|--------|--------|----------|---------------|-----------|--|
| OpenPorchSF     | EnclosedPorch | 3SsnPorch | ScreenPorch | PoolArea | PoolQC | Fence | MiscFeature | MiscVal | MoSold | YrSold | SaleType | SaleCondition | SalePrice |  |
| 205             | 0             | 0         | 0           | 0        | Nan    | Nan   | Nan         | 0       | 2      | 2007   | WD       | Normal        | 128000    |  |
| 207             | 0             | 0         | 224         | 0        | Nan    | Nan   | Nan         | 0       | 10     | 2007   | WD       | Normal        | 268000    |  |
| 130             | 0             | 0         | 0           | 0        | Nan    | Nan   | Nan         | 0       | 6      | 2007   | WD       | Normal        | 269790    |  |
| 122             | 0             | 0         | 0           | 0        | Nan    | MnPrv | Nan         | 0       | 1      | 2010   | COD      | Normal        | 190000    |  |
| 0               | 0             | 0         | 0           | 0        | Nan    | Nan   | Nan         | 0       | 6      | 2009   | WD       | Normal        | 215000    |  |

## Data Information

```
train_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1168 entries, 0 to 1167
Data columns (total 81 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Id               1168 non-null    int64  
 1   MSSubClass        1168 non-null    int64  
 2   MSZoning          1168 non-null    object  
 3   LotFrontage       954 non-null    float64 
 4   LotArea           1168 non-null    int64  
 5   Street            1168 non-null    object  
 6   Alley              77 non-null    object  
 7   LotShape           1168 non-null    object  
 8   LandContour        1168 non-null    object  
 9   Utilities          1168 non-null    object  
 10  LotConfig          1168 non-null    object  
 11  LandSlope          1168 non-null    object  
 12  Neighborhood       1168 non-null    object  
 13  Condition1         1168 non-null    object  
 14  Condition2         1168 non-null    object  
 15  BldgType           1168 non-null    object  
 16  HouseStyle         1168 non-null    object  
 17  OverallQual        1168 non-null    int64  
 18  OverallCond        1168 non-null    int64  
 19  YearBuilt          1168 non-null    int64  
 20  YearRemodAdd       1168 non-null    int64  
 21  RoofStyle           1168 non-null    object  
 22  RoofMatl            1168 non-null    object  
 23  Exterior1st         1168 non-null    object  
 24  Exterior2nd         1168 non-null    object  
 25  MasVnrType          1161 non-null    object  
 26  MasVnrArea          1161 non-null    float64 
 27  ExterQual           1168 non-null    object  
 28  ExterCond           1168 non-null    object  
 29  Foundation          1168 non-null    object  
 30  BsmtQual            1138 non-null    object  
 31  BsmtCond            1138 non-null    object  
 32  BsmtExposure        1137 non-null    object  
 33  BsmtFinType1         1138 non-null    object  
 34  BsmtFinSF1           1168 non-null    int64  
 35  BsmtFinType2         1137 non-null    object  
 36  BsmtFinSF2           1168 non-null    int64  
 37  BsmtUnfSF            1168 non-null    int64  
 38  TotalBsmtSF          1168 non-null    int64  
 39  Heating              1168 non-null    object  
 40  HeatingQC             1168 non-null    object  
 41  CentralAir           1168 non-null    object
```

## HOUSE PRICE PREDICTION MODEL

---

```
42 Electrical      1168 non-null   object
43 1stFlrSF       1168 non-null   int64
44 2ndFlrSF       1168 non-null   int64
45 LowQualFinSF  1168 non-null   int64
46 GrLivArea      1168 non-null   int64
47 BsmtFullBath  1168 non-null   int64
48 BsmtHalfBath  1168 non-null   int64
49 FullBath       1168 non-null   int64
50 HalfBath        1168 non-null   int64
51 BedroomAbvGr   1168 non-null   int64
52 KitchenAbvGr   1168 non-null   int64
53 KitchenQual    1168 non-null   object
54 TotRmsAbvGrd  1168 non-null   int64
55 Functional     1168 non-null   object
56 Fireplaces      1168 non-null   int64
57 FireplaceQu    617 non-null    object
58 GarageType      1104 non-null   object
59 GarageYrBlt    1104 non-null   float64
60 GarageFinish   1104 non-null   object
61 GarageCars      1168 non-null   int64
62 GarageArea      1168 non-null   int64
63 GarageQual     1104 non-null   object
64 GarageCond      1104 non-null   object
65 PavedDrive     1168 non-null   object
66 WoodDeckSF     1168 non-null   int64
67 OpenPorchSF    1168 non-null   int64
68 EnclosedPorch  1168 non-null   int64
69 3SsnPorch      1168 non-null   int64
70 ScreenPorch    1168 non-null   int64
71 PoolArea       1168 non-null   int64
72 PoolQC         7 non-null     object
73 Fence          237 non-null   object
74 MiscFeature    44 non-null    object
75 MiscVal        1168 non-null   int64
76 MoSold         1168 non-null   int64
77 YrSold         1168 non-null   int64
78 SaleType       1168 non-null   object
79 SaleCondition  1168 non-null   object
80 SalePrice       1168 non-null   int64
dtypes: float64(3), int64(35), object(43)
memory usage: 739.2+ KB
```

Keys: As we can see above all column are not having same not null values, mean dataset is having null value and we need to over them to remove it.

Target variable is SalePrice

**Data Dtype:**

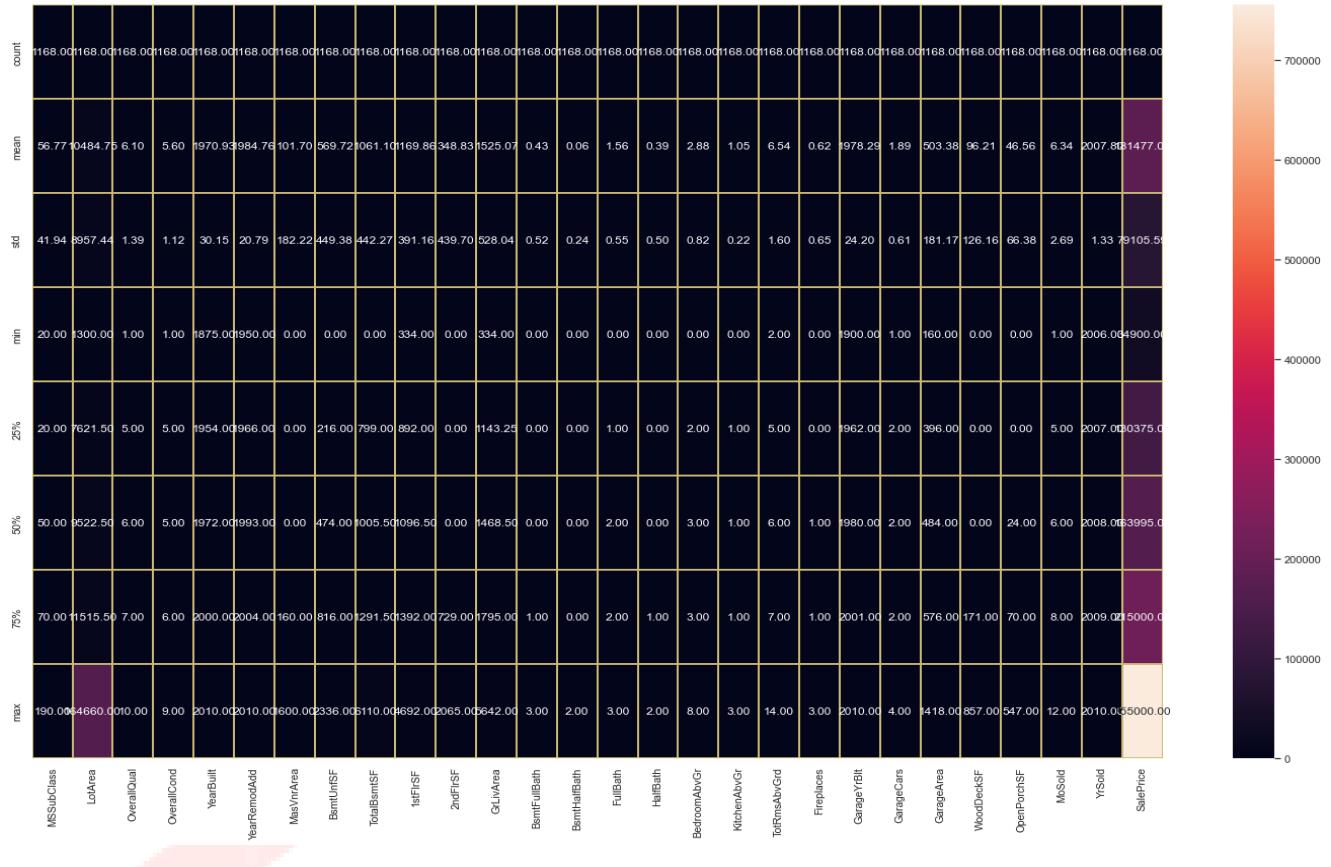
|              |         |              |        |               |         |
|--------------|---------|--------------|--------|---------------|---------|
| Id           | int64   | ExterQual    | object | Functional    | object  |
| MSSubClass   | int64   | ExterCond    | object | Fireplaces    | int64   |
| MSZoning     | object  | Foundation   | object | FireplaceQu   | object  |
| LotFrontage  | float64 | BsmtQual     | object | GarageType    | object  |
| LotArea      | int64   | BsmtCond     | object | GarageYrBlt   | float64 |
| Street       | object  | BsmtExposure | object | GarageFinish  | object  |
| Alley        | object  | BsmtFinType1 | object | GarageCars    | int64   |
| LotShape     | object  | BsmtFinSF1   | int64  | GarageArea    | int64   |
| LandContour  | object  | BsmtFinType2 | object | GarageQual    | object  |
| Utilities    | object  | BsmtFinSF2   | int64  | GarageCond    | object  |
| LotConfig    | object  | BsmtUnfSF    | int64  | PavedDrive    | object  |
| LandSlope    | object  | TotalBsmtSF  | int64  | WoodDeckSF    | int64   |
| Neighborhood | object  | Heating      | object | OpenPorchSF   | int64   |
| Condition1   | object  | HeatingQC    | object | EnclosedPorch | int64   |
| Condition2   | object  | CentralAir   | object | 3SsnPorch     | int64   |
| BldgType     | object  | Electrical   | object | ScreenPorch   | int64   |
| HouseStyle   | object  | 1stFlrSF     | int64  | PoolArea      | int64   |
| OverallQual  | int64   | 2ndFlrSF     | int64  | PoolQC        | object  |
| OverallCond  | int64   | LowQualFinSF | int64  | Fence         | object  |
| YearBuilt    | int64   | GrLivArea    | int64  | MiscFeature   | object  |
| YearRemodAdd | int64   | BsmtFullBath | int64  | MiscVal       | int64   |
| RoofStyle    | object  | BsmtHalfBath | int64  | MoSold        | int64   |
| RoofMatl     | object  | FullBath     | int64  | YrSold        | int64   |
| Exterior1st  | object  | HalfBath     | int64  | SaleType      | object  |
| Exterior2nd  | object  | BedroomAbvGr | int64  | SaleCondition | object  |
| MasVnrType   | object  | KitchenAbvGr | int64  | SalePrice     | int64   |
| MasVnrArea   | float64 | KitchenQual  | object | dtype: object |         |
|              |         | TotRmsAbvGrd | int64  |               |         |

Total 81 columns we have received by data, in which 80 are feature columns and last one is target column (saleprize).

As we to have find the house price and target variable is in numeric form so it is any regression problem.

## Analytical Problem Framing

- Describe dataset:



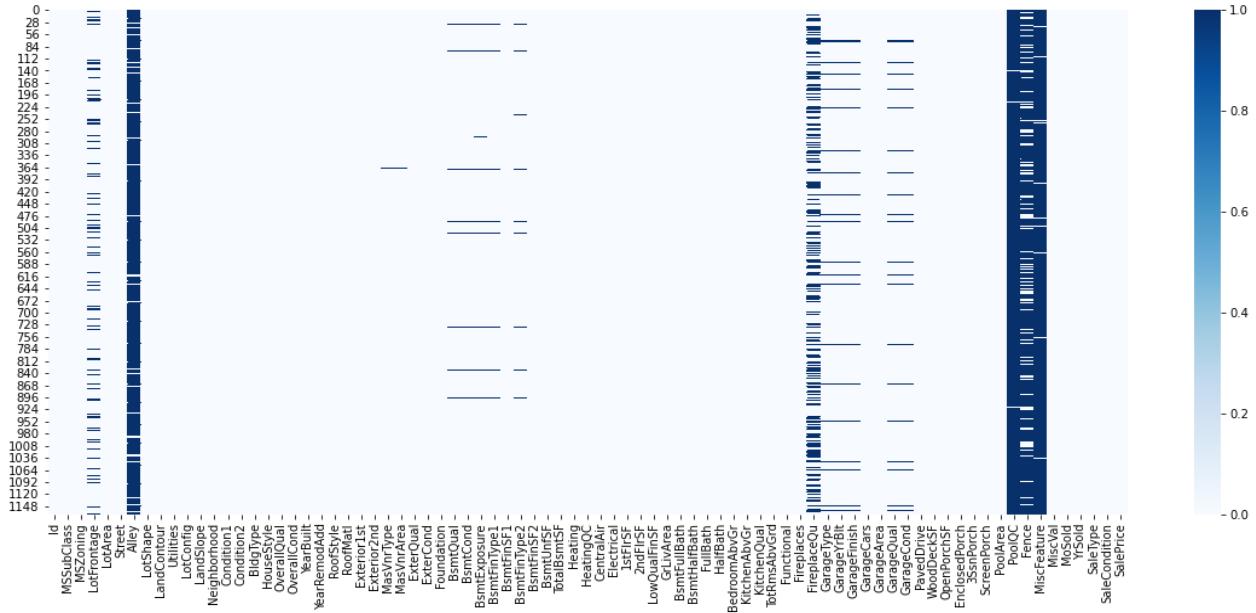
- As every column is showing high difference between mean and 50% percentile, which indicating the skewness of the column
- Very high difference found between min and max value of kilometre and prize columns.
- Mean and 50 percentile value is having some value difference, mean these columns are having skewness
- Count for each column is same mean no null value present.
- 25% and 75% is having difference, which also indicating , column spread.

The statistical figure (above mentioned) get to know by the df.describe() so many information, like min, max, standard deviation, 25 percentile, 50th percentile, 75 percentile.



## ➤ Data Pre-processing:

### i. Checked Null values:



```
train_null_value_per = train_df.isnull().sum() / train_df.shape[0] * 100
train_null_value_per[train_null_value_per > 15]
```

|             |           |
|-------------|-----------|
| LotFrontage | 18.321918 |
| Alley       | 93.407534 |
| FireplaceQu | 47.174658 |
| PoolQC      | 99.400685 |
| Fence       | 79.708904 |
| MiscFeature | 96.232877 |

dtype: float64

```
# As these columns are having very high null values , therefore need to delete these columns from both
# datasets train and test
print(train_df.shape)
print(test_df.shape)
train_df.drop(columns = ['LotFrontage', 'Alley', 'FireplaceQu', 'PoolQC', 'Fence', 'MiscFeature'], inplace = True)
test_df.drop(columns = ['LotFrontage', 'Alley', 'FireplaceQu', 'PoolQC', 'Fence', 'MiscFeature'], inplace = True)
print(train_df.shape)
print(test_df.shape)

(1168, 81)
(292, 80)
(1168, 75)
(292, 74)
```

**Keys:** As these above mentioned columns are showing null value which is more than 15%, therefore I have decided to remove these column. As they are containing high null values

## HOUSE PRICE PREDICTION MODEL

- ii. Checked relation of columns then imputed values accordingly, maximum object type columns null values are replace with mode of the column.

```
train_df[train_df['MasVnrType'] == 'None']['MasVnrArea'].head(10)

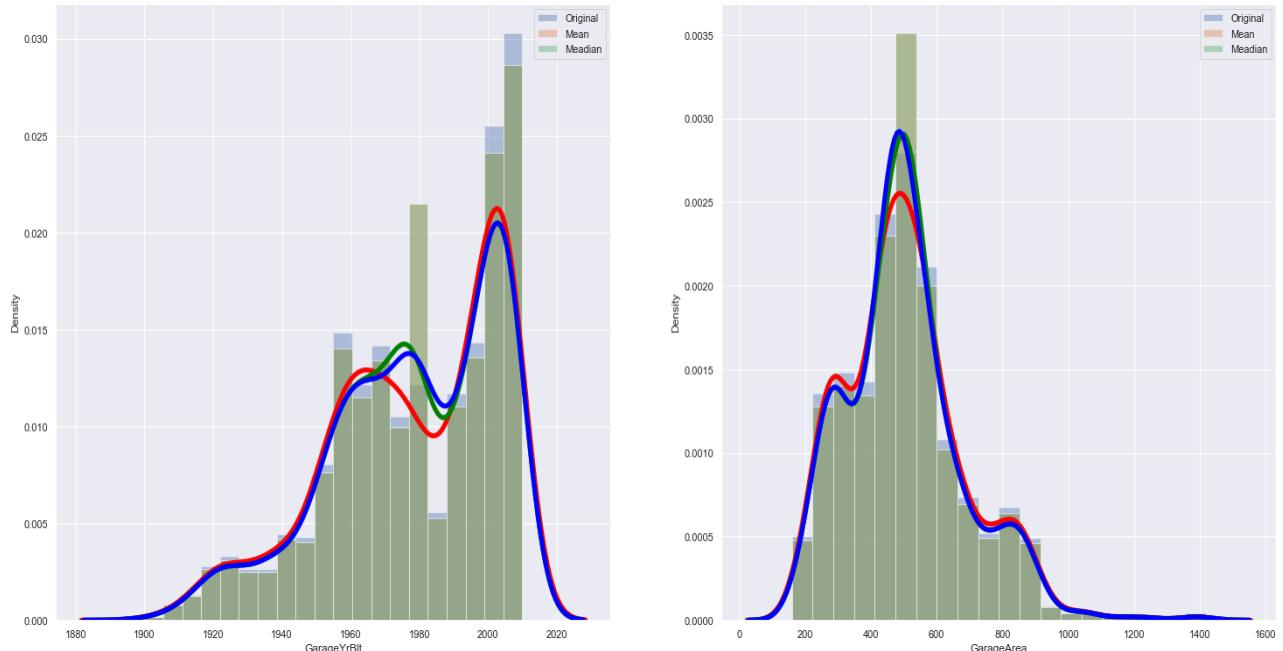
0    0.0
1    0.0
2    0.0
5    0.0
8    0.0
9    0.0
10   0.0
12   0.0
14   0.0
15   0.0
Name: MasVnrArea, dtype: float64

# As we can when MasVnrType column is None then MasVnrArea column has to be 0.0, therefore we will impute
# values accordingly
train_df['MasVnrType'].mode()
# mode of MasVnrType is 'None', and we can impute mode value here as it is a categorical column

0    None
dtype: object
```

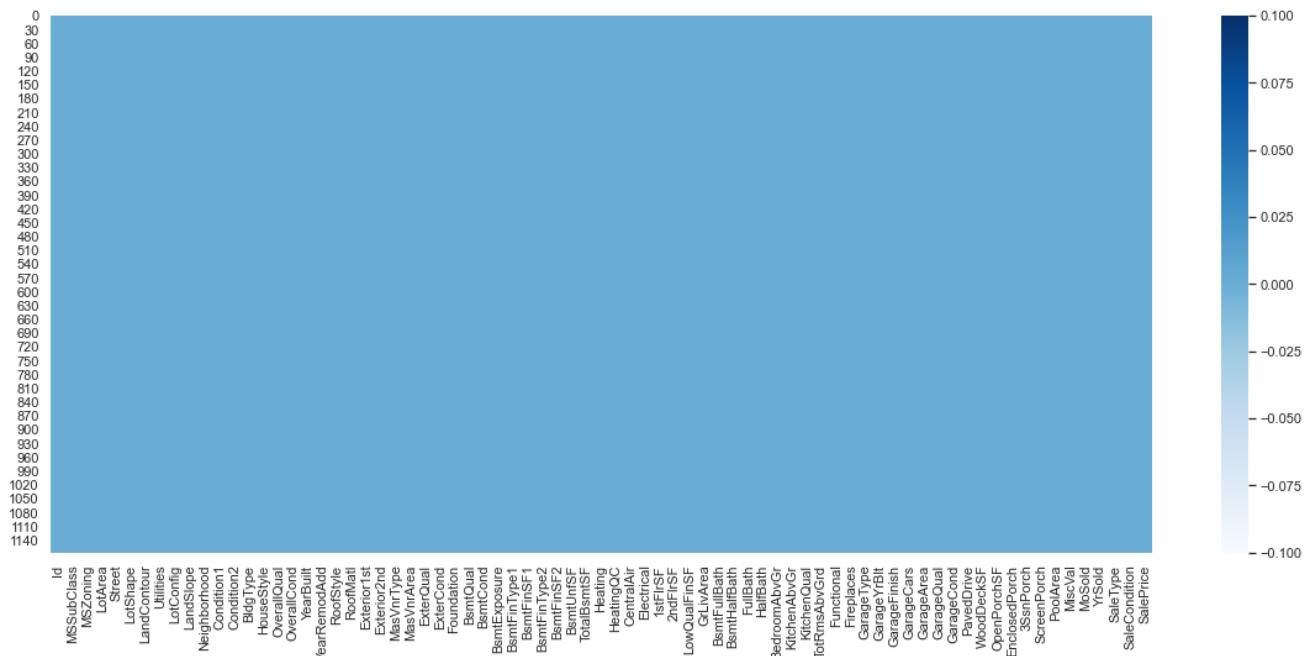
Note: This same process applied for each object type column, first checked relation, if found then imputed value accordingly else, mode value imputed

- iii. Numerical type columns, here, I have checked impact of mean and median imputation then impute value.



## HOUSE PRICE PREDICTION MODEL

- iv. Again checked Null values



- v. Few columns values has been divided into some range group and replace their value with range group.

Column names are given below:

```
# BsmtFinSF1  
# BsmtFinSF2  
# LowQualFinSF  
# EnclosedPorch  
# 3SsnPorch  
# ScreenPorch  
# PoolArea  
# MiscVal
```

**Applying range opeartion for BsmtFinSF1**

```
print('Min: ', train_df['BsmtFinSF1'].min())  
print('Max: ', train_df['BsmtFinSF1'].max())  
print('Length: ', len(train_df['BsmtFinSF1']))
```

```
Min:  0  
Max:  5644  
Length:  1168
```

```
BsmtFinSF1_class = []

for i in train_df['BsmtFinSF1']:
    if i in range (0, 101):
        BsmtFinSF1_class.append('0-100')

    elif i in range (101, 901):
        BsmtFinSF1_class.append('101-900')

    elif i in range(901,1801):
        BsmtFinSF1_class.append('901-1800')

    elif i in range(1801, 2701):
        BsmtFinSF1_class.append('1801-2700')

    elif i in range(2701, 3601):
        BsmtFinSF1_class.append('2701-3600')

    elif i in range(3601, 4501):
        BsmtFinSF1_class.append('3601-4500')

    elif i in range(4501, 5401):
        BsmtFinSF1_class.append('4501-5400')

    elif i in range(5401, 6000):
        BsmtFinSF1_class.append('5400+ ')

train_df['BsmtFinSF1_Class'] = BsmtFinSF1_class
```

### *Applying range opeartion for BsmtFinSF2*

```
print( 'Min: ', train_df['BsmtFinSF2'].min())
print( 'Max: ', train_df['BsmtFinSF2'].max())
print('Lenght: ', len(train_df['BsmtFinSF2']))
```

```
Min:  0
Max:  1474
Lenght:  1168
```

```

BsmtFinSF2_class = []

for i in train_df['BsmtFinSF2']:
    if i in range(0, 101):
        BsmtFinSF2_class.append('0-100')
    elif i in range(101, 401):
        BsmtFinSF2_class.append('101-400')
    elif i in range(401, 801):
        BsmtFinSF2_class.append('401-800')
    elif i in range(801, 1201):
        BsmtFinSF2_class.append('801-1200')
    elif i in range(1201, 2000):
        BsmtFinSF2_class.append('1200+')
len(BsmtFinSF2_class)
train_df['BsmtFinSF2_class'] = BsmtFinSF2_class

```

: 1168

#### *Applying range opeartion for LowQualFinSF*

```

: print('Min: ', train_df['LowQualFinSF'].min())
print('Max: ', train_df['LowQualFinSF'].max())
print('Length: ', len(train_df['LowQualFinSF']))

```

Min: 0  
 Max: 572  
 Length: 1168

```

: # for training dataset
LowQualFinSF_class = []

for i in train_df['LowQualFinSF']:
    if i in range(0, 1):
        LowQualFinSF_class.append('Zero')

    elif i in range(1, 151):
        LowQualFinSF_class.append('1-150')

    elif i in range(151, 301):
        LowQualFinSF_class.append('151-300')

    elif i in range(301, 451):
        LowQualFinSF_class.append('301-450')

    elif i in range(451, 601):
        LowQualFinSF_class.append('451-600')

```

```

train_df['LowQualFinSF_class'] = LowQualFinSF_class

```

### *Applying range operation for EnclosedPorch*

```
: print( 'Min: ', train_df['EnclosedPorch'].min())
  print( 'Max: ', train_df['EnclosedPorch'].max())
  print('Length: ', len(train_df['EnclosedPorch']))
```

```
Min:  0
Max:  552
Length:  1168
```

```
# for training dataset
EnclosedPorch_class = []

for i in train_df['EnclosedPorch']:
    if i in range (0, 1):
        EnclosedPorch_class.append('Zero')

    elif i in range (1, 151):
        EnclosedPorch_class.append('1-150')

    elif i in range(151,301):
        EnclosedPorch_class.append('151-300')

    elif i in range(301, 451):
        EnclosedPorch_class.append('301-450')

    elif i in range(451, 601):
        EnclosedPorch_class.append('451-600')

train_df['EnclosedPorch_class'] = EnclosedPorch_class
```

### *Applying range operation for 3SsnPorch*

```
: print( 'Min: ', train_df['3SsnPorch'].min())
  print( 'Max: ', train_df['3SsnPorch'].max())
  print('Length: ', len(train_df['3SsnPorch']))
```

```
Min:  0
Max:  508
Length:  1168
```

```
# for training dataset
SsnPorch_class = []

for i in train_df['3SsnPorch']:
    if i in range (0, 1):
        SsnPorch_class.append('Zero')

    elif i in range (1, 151):
        SsnPorch_class.append('1-150')

    elif i in range(151,301):
        SsnPorch_class.append('151-300')

    elif i in range(301, 451):
        SsnPorch_class.append('301-450')

    elif i in range(451, 601):
        SsnPorch_class.append('451-600')

train_df['3SsnPorch_class'] = SsnPorch_class
```

### *Applying range opeartion for ScreenPorch*

```
print( 'Min: ', train_df['ScreenPorch'].min())
print( 'Max: ', train_df['ScreenPorch'].max())
print('Lenght: ', len(train_df['ScreenPorch']))
```

Min: 0  
Max: 480  
Lenght: 1168

```
# for training dataset
ScreenPorch_class = []

for i in train_df['ScreenPorch']:
    if i in range (0, 1):
        ScreenPorch_class.append('Zero')

    elif i in range (1, 151):
        ScreenPorch_class.append('1-150')

    elif i in range(151,301):
        ScreenPorch_class.append('151-300')

    elif i in range(301, 451):
        ScreenPorch_class.append('301-450')

    elif i in range(451, 601):
        ScreenPorch_class.append('450+')
```

```
train_df['ScreenPorch_class'] = ScreenPorch_class
```

### *Applying range operation for PoolArea*

```
print( 'Min: ', train_df['PoolArea'].min())
print( 'Max: ', train_df['PoolArea'].max())
print('Length: ', len(train_df['PoolArea']))
```

```
Min:  0
Max:  738
Length:  1168
```

```
# for training dataset
PoolArea_class = []

for i in train_df['PoolArea']:
    if i in range (0, 1):
        PoolArea_class.append('Zero')

    elif i in range (1, 151):
        PoolArea_class.append('1-150')

    elif i in range(151,301):
        PoolArea_class.append('151-300')

    elif i in range(301, 451):
        PoolArea_class.append('301-450')

    elif i in range(451, 601):
        PoolArea_class.append('450-600')

    elif i in range(601, 751):
        PoolArea_class.append('601-750')
```

```
: train_df['PoolArea_class'] = PoolArea_class
```

### *Applying range operation for MiscVal*

```
: print( 'Min: ', train_df['MiscVal'].min())
print( 'Max: ', train_df['MiscVal'].max())
print('Length: ', len(train_df['MiscVal']))
```

```
Min:  0
Max:  15500
Length:  1168
```

```

: # for training dataset
MiscVal_class = []

for i in train_df['MiscVal']:
    if i in range (0, 1):
        MiscVal_class.append('Zero')

    elif i in range (1, 1001):
        MiscVal_class.append('1-150')

    elif i in range(1001,6001):
        MiscVal_class.append('151-300')

    elif i in range(6001, 11001):
        MiscVal_class.append('301-450')

    elif i in range(11001, 16000):
        MiscVal_class.append('11001 -16000')

len(MiscVal_class)
:
```

```
train_df['MiscVal_class'] = MiscVal_class
```

## FLIP FLOP

- vi. One record of Exterior1st and Exterior2nd column was wrong updated, which was creating problem at the time of encoding

That wrong value replace with mode of that column

```

test_df['Exterior1st'] = test_df['Exterior1st'].replace ('CBlock', train_df['Exterior1st'].mode()[0])
test_df['Exterior2nd'] = test_df['Exterior2nd'].replace ('CBlock', train_df['Exterior2nd'].mode()[0])

print(test_df['Exterior1st'].unique())
print('_____')
print(test_df['Exterior2nd'].unique())

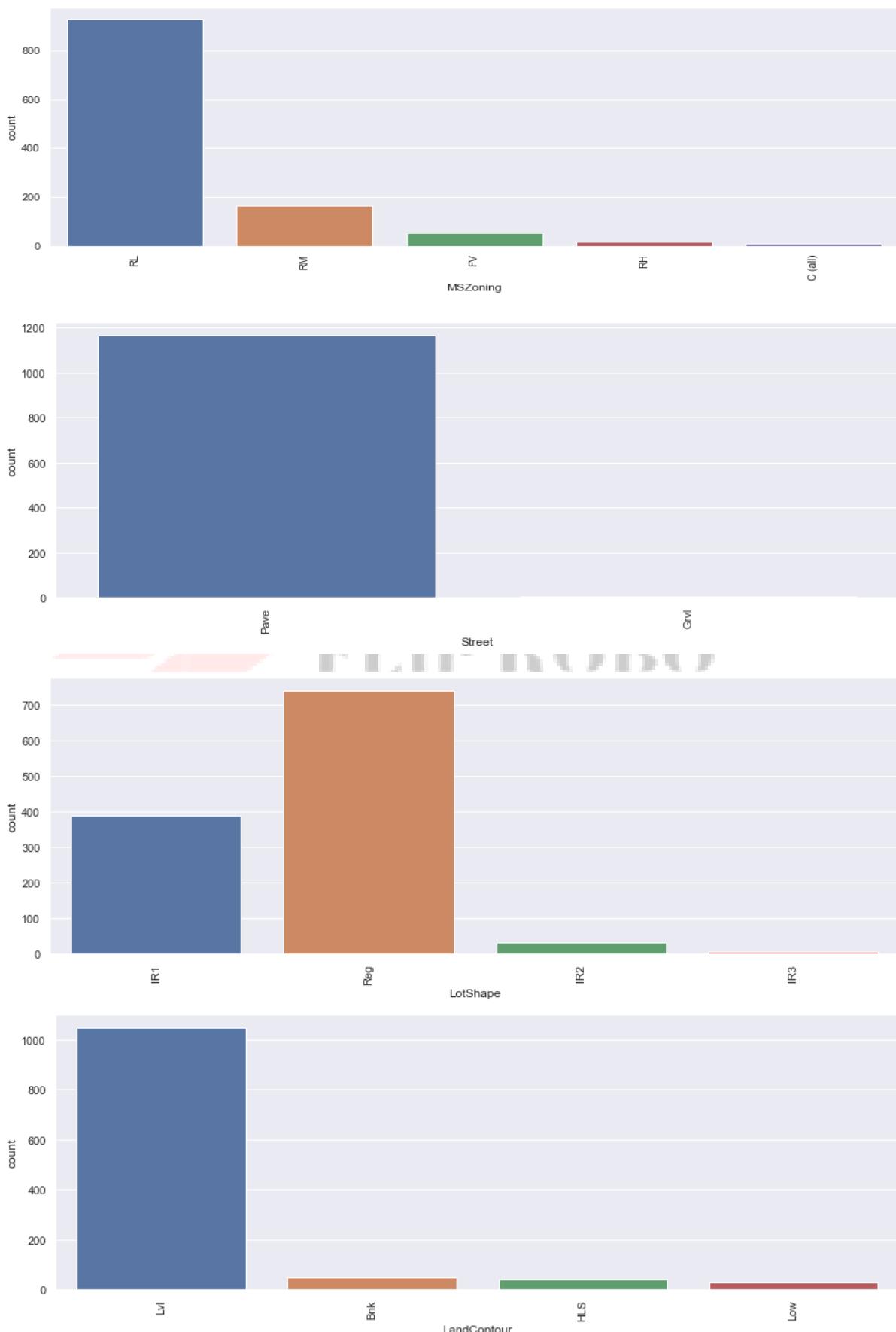
['VinylSd' 'HdBoard' 'MetalSd' 'CemntBd' 'Wd Sdng' 'Plywood' 'BrkFace'
 'AsbShng' 'WdShing' 'BrkComm' 'Stucco']

['VinylSd' 'HdBoard' 'MetalSd' 'CmentBd' 'Wd Sdng' 'Plywood' 'Wd Shng'
 'Brk Cmn' 'ImStucc' 'AsbShng' 'BrkFace' 'Stone' 'Stucco']
```

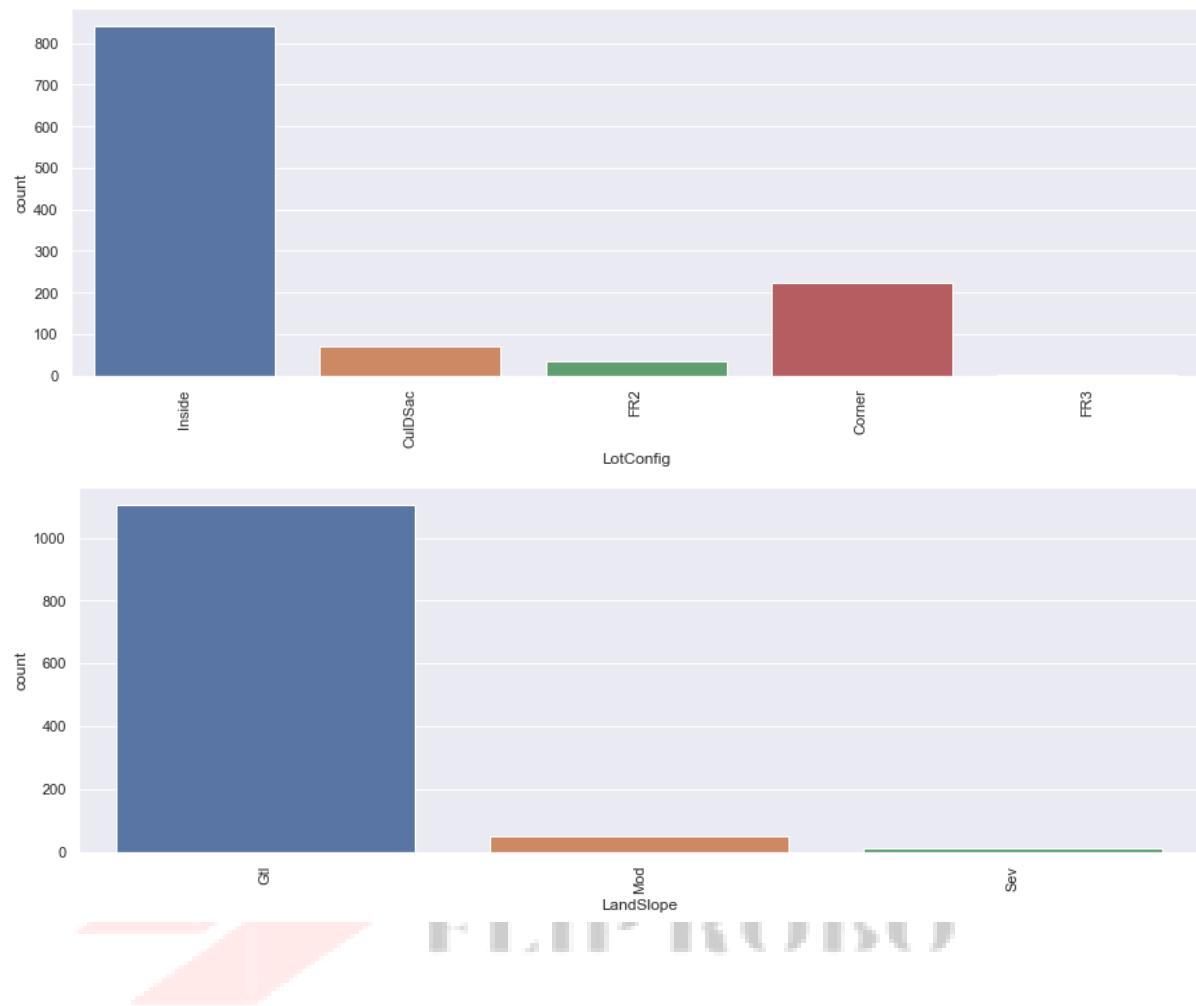
## Visualization

# HOUSE PRICE PREDICTION MODEL

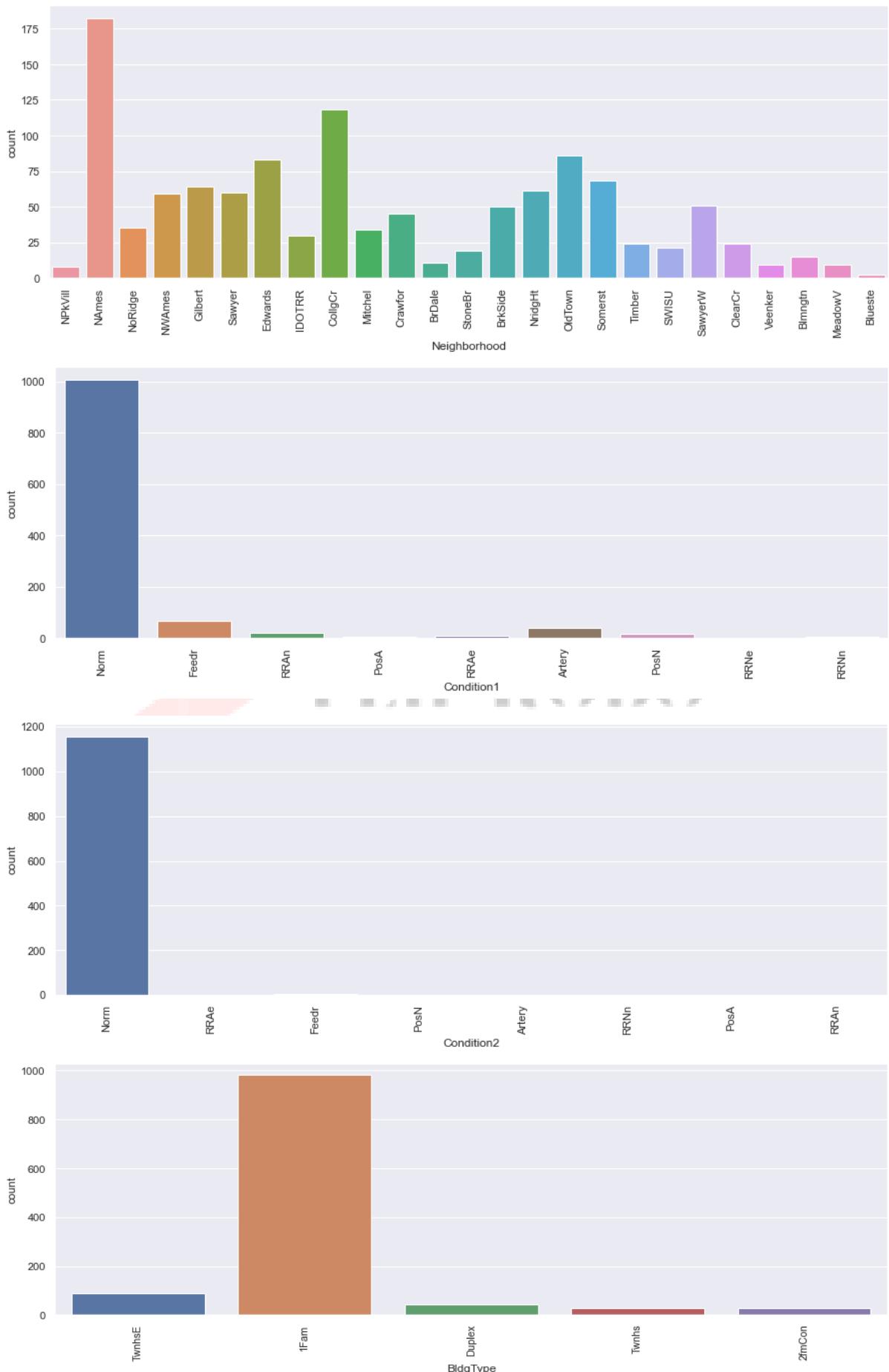
## ➤ Bar Chart:



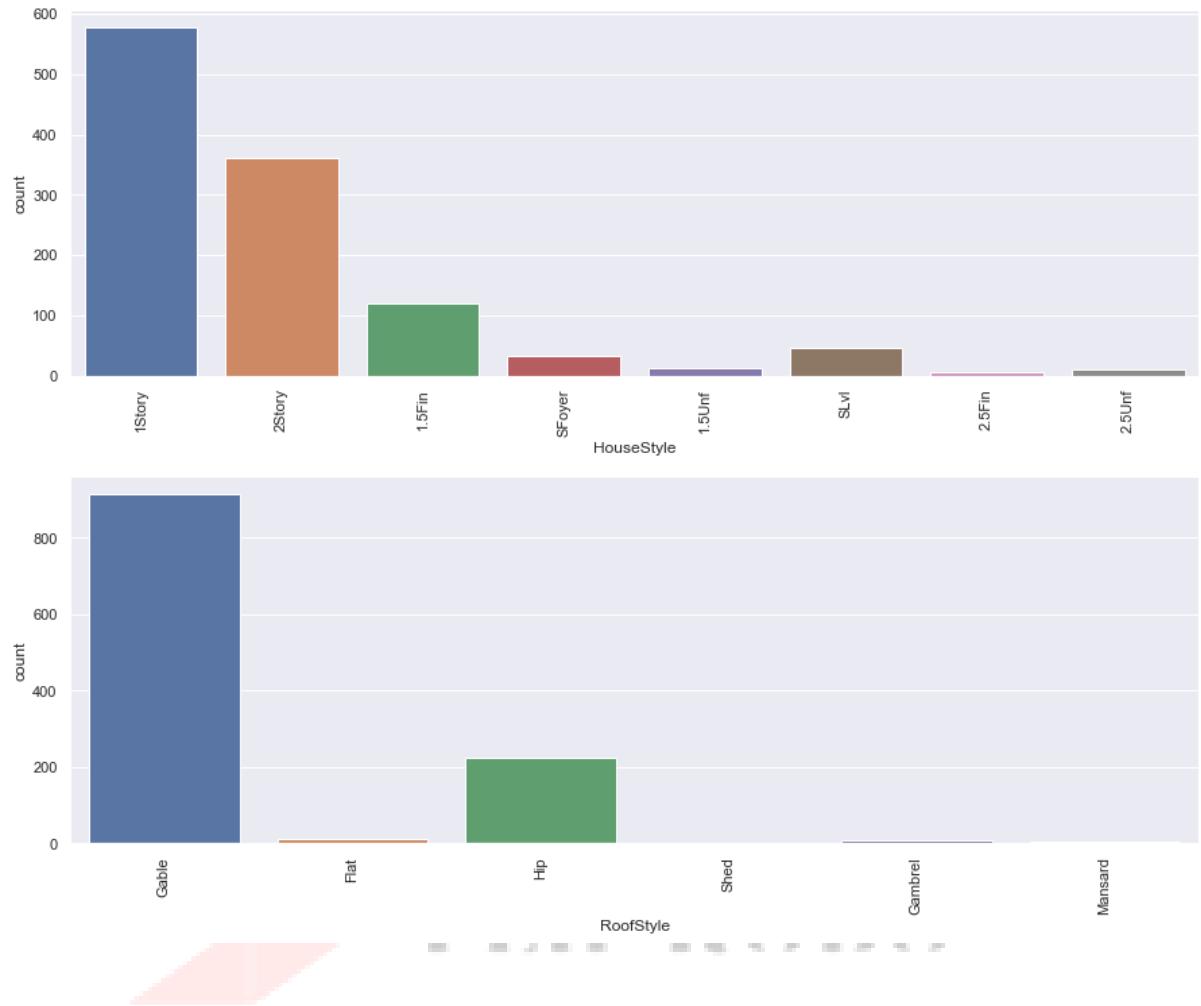
## HOUSE PRICE PREDICTION MODEL



# HOUSE PRICE PREDICTION MODEL



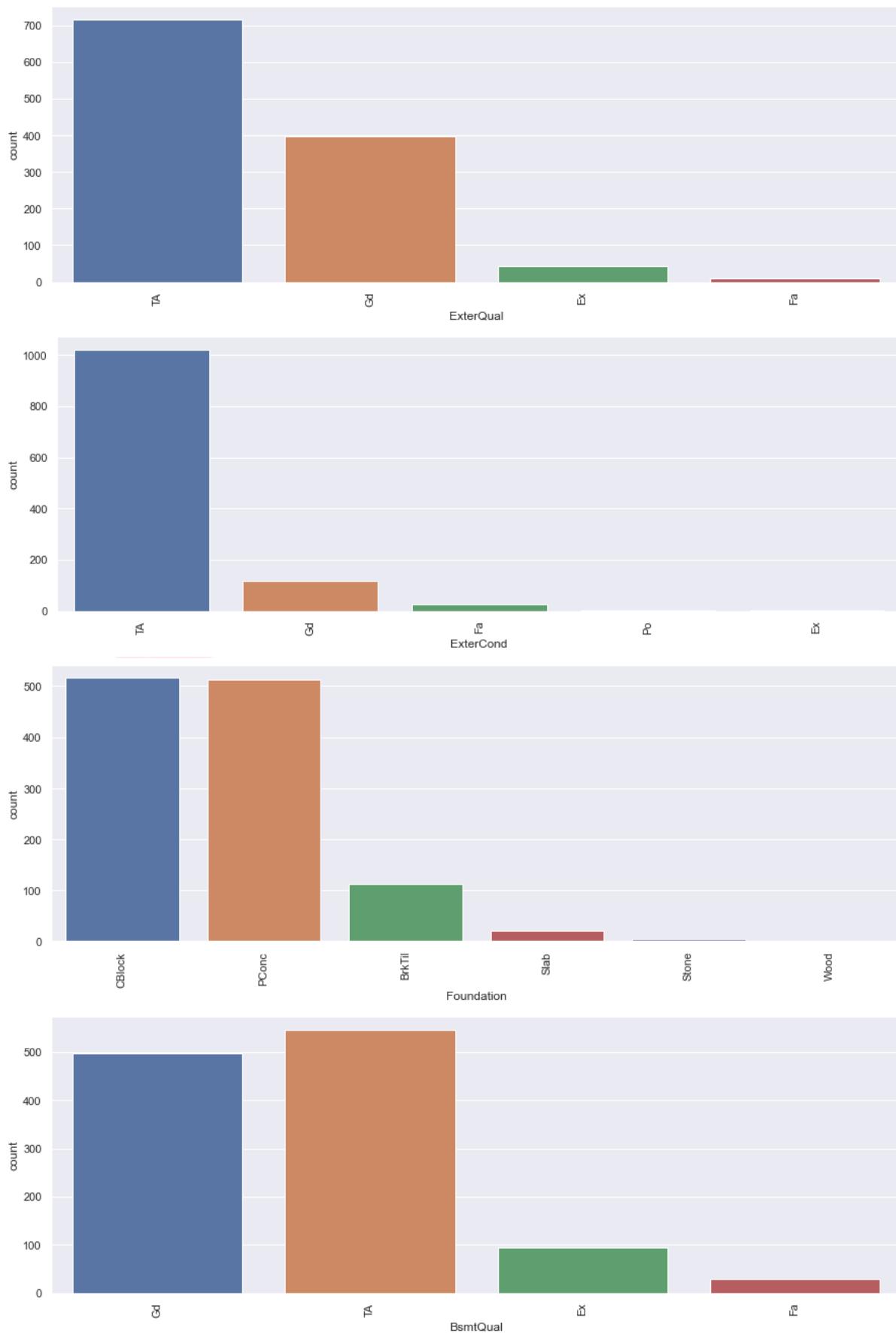
## HOUSE PRICE PREDICTION MODEL



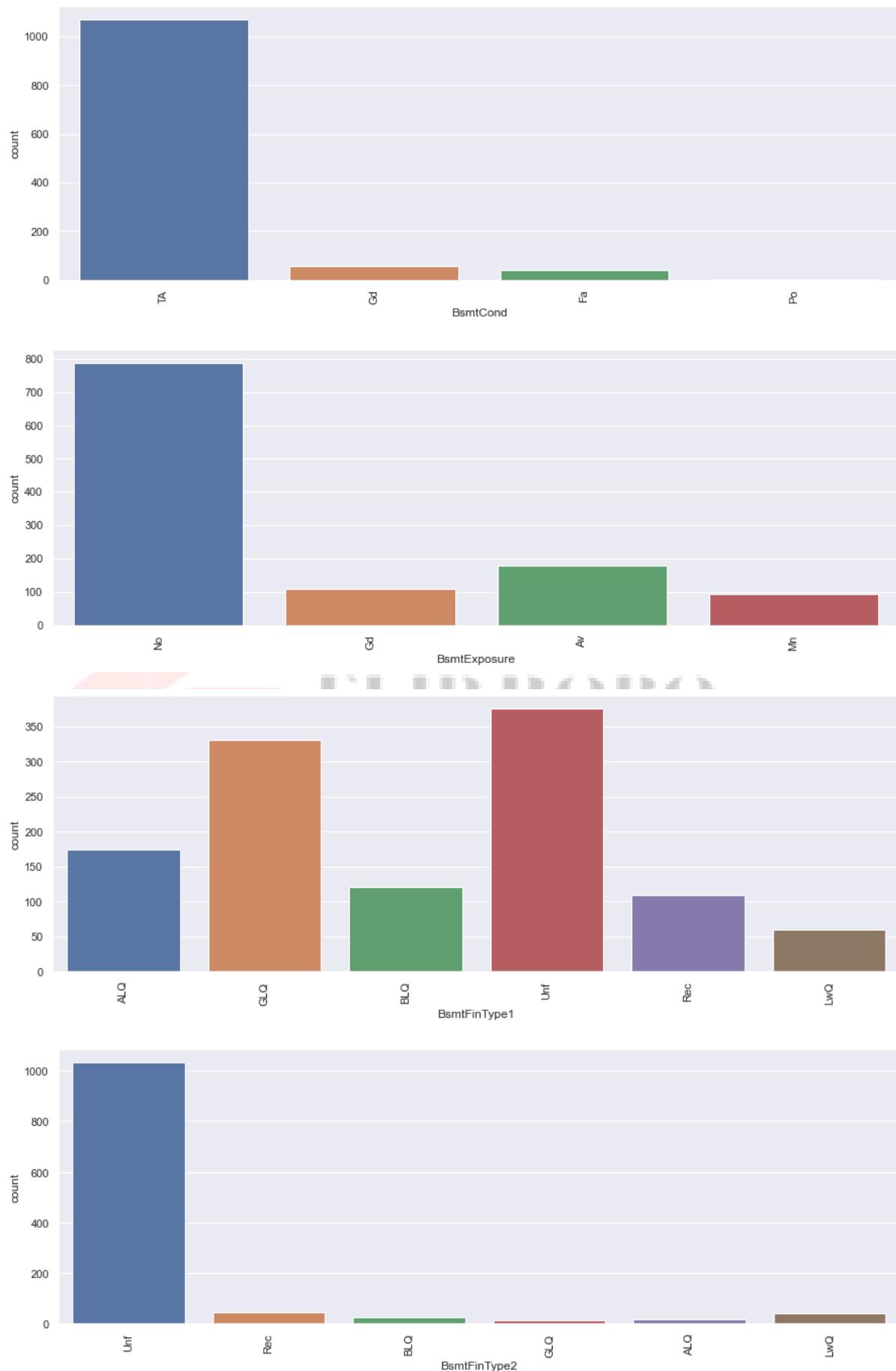
## HOUSE PRICE PREDICTION MODEL



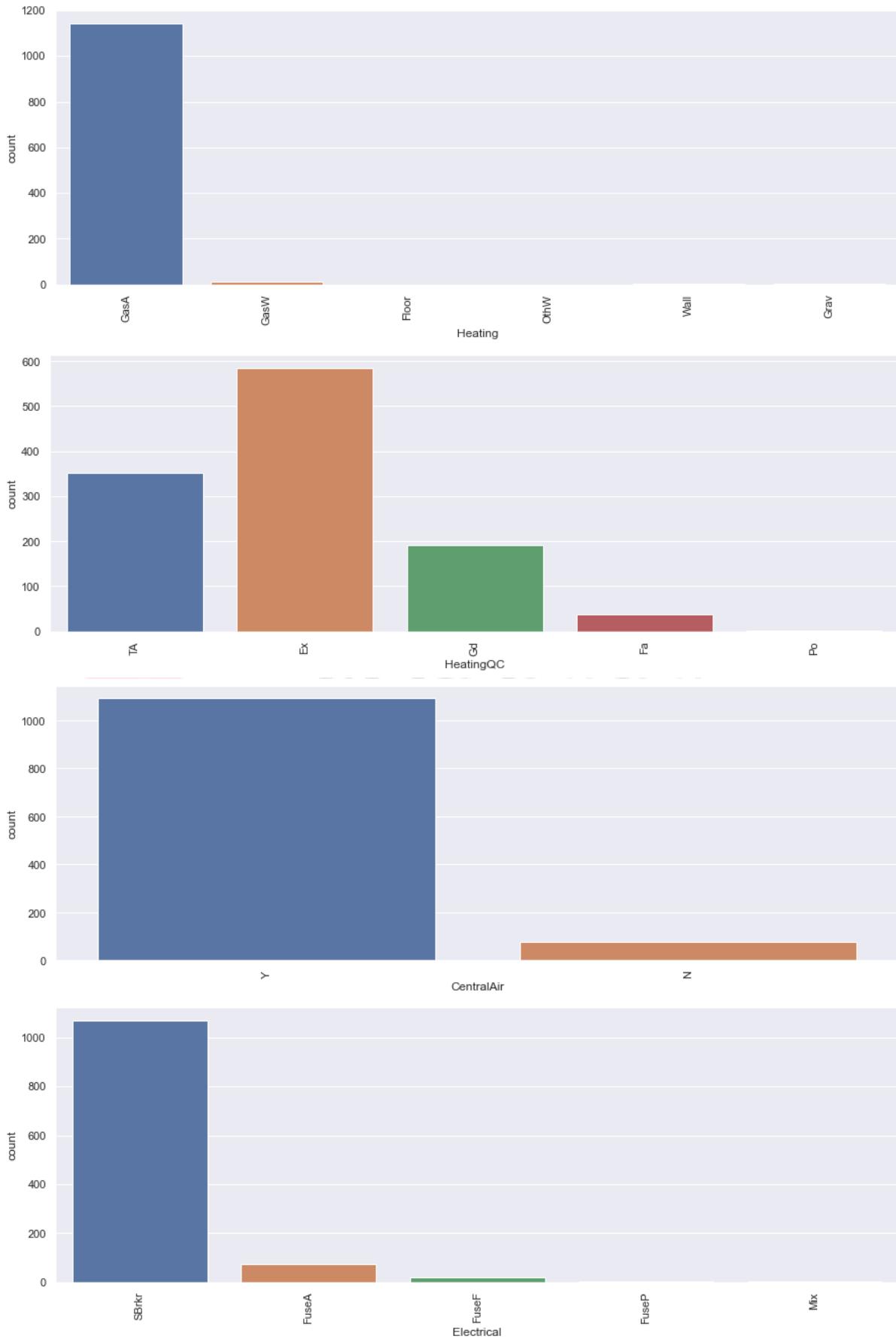
## HOUSE PRICE PREDICTION MODEL



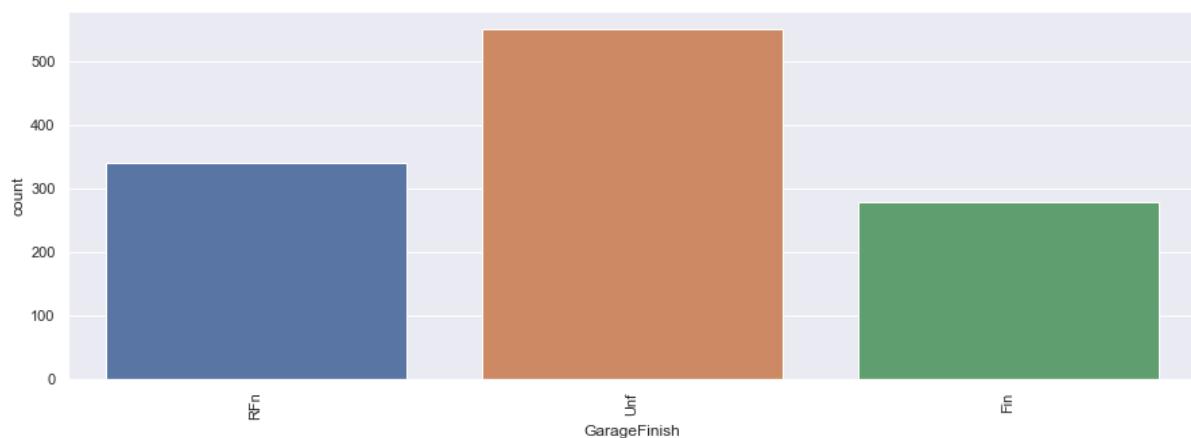
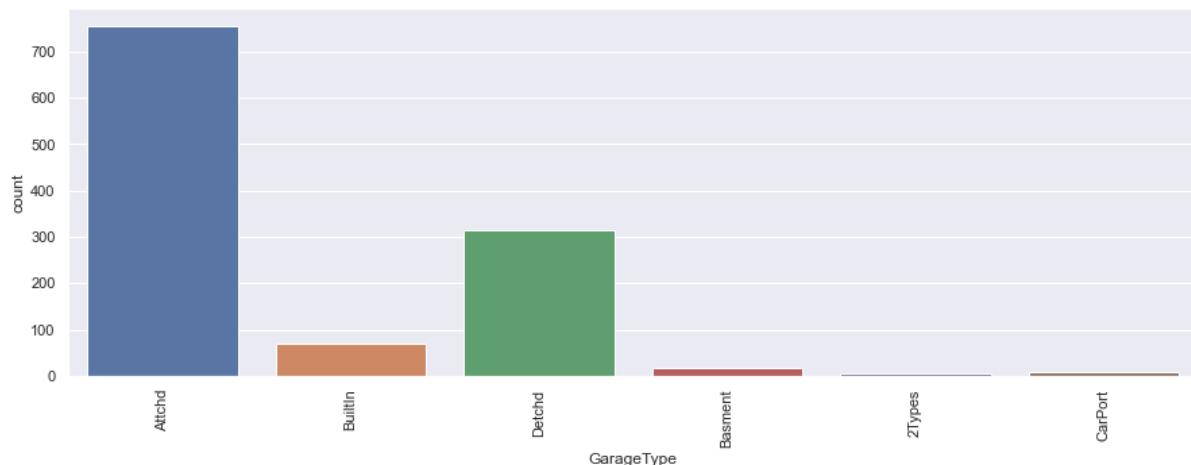
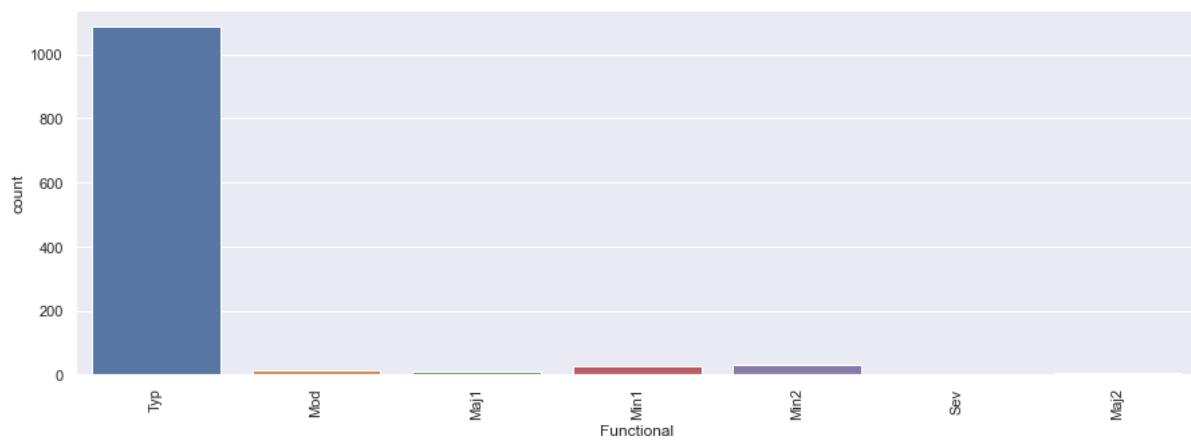
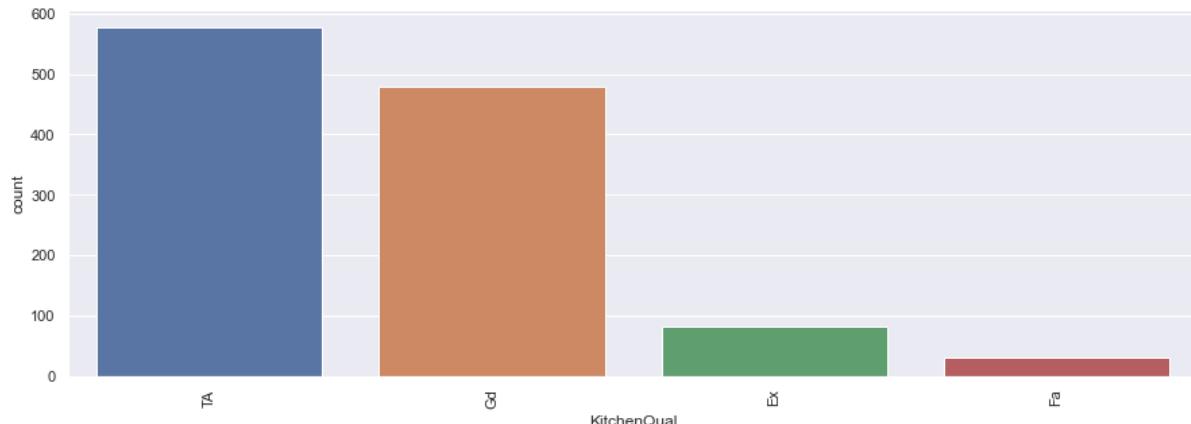
## HOUSE PRICE PREDICTION MODEL



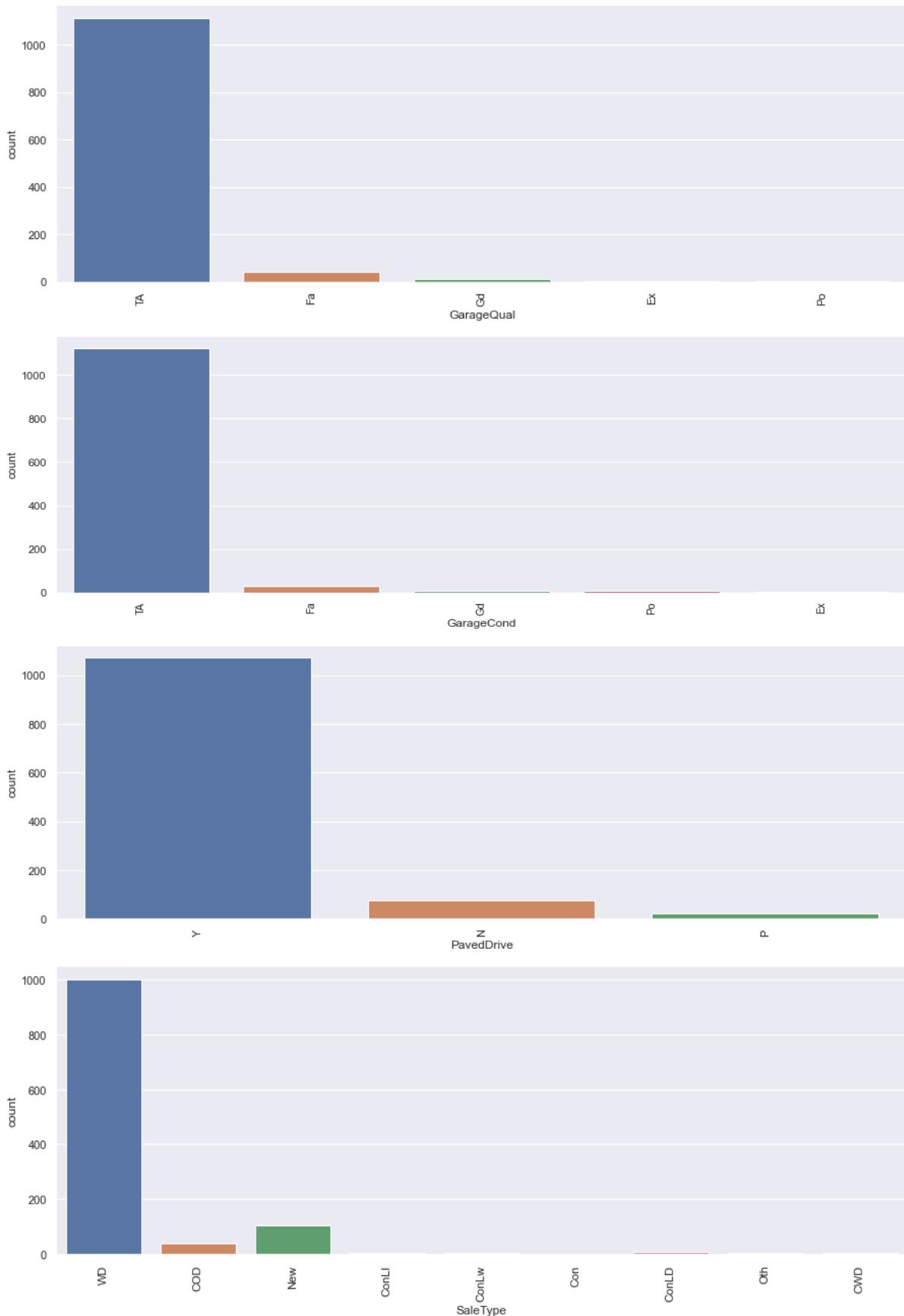
## HOUSE PRICE PREDICTION MODEL



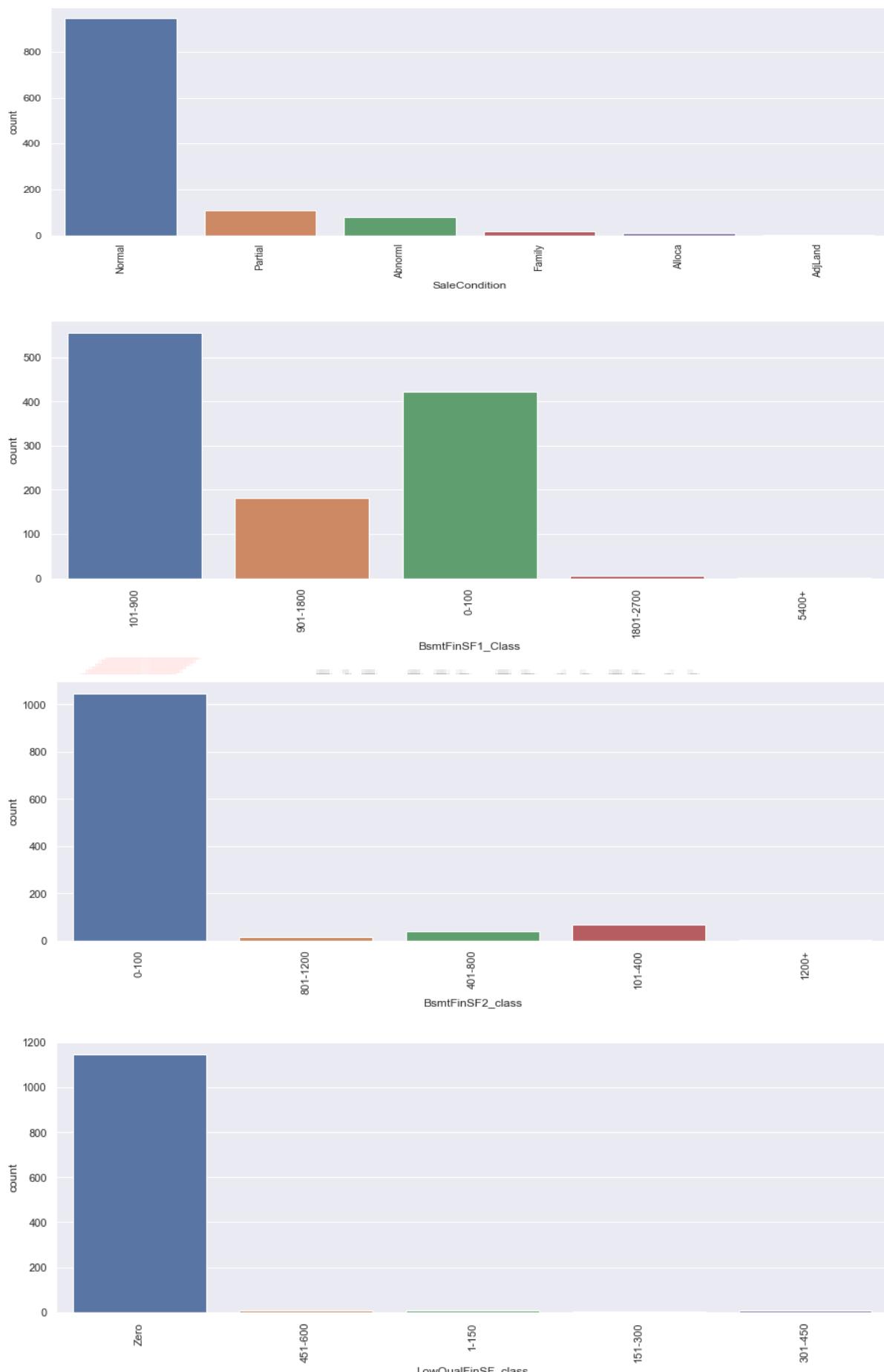
## HOUSE PRICE PREDICTION MODEL



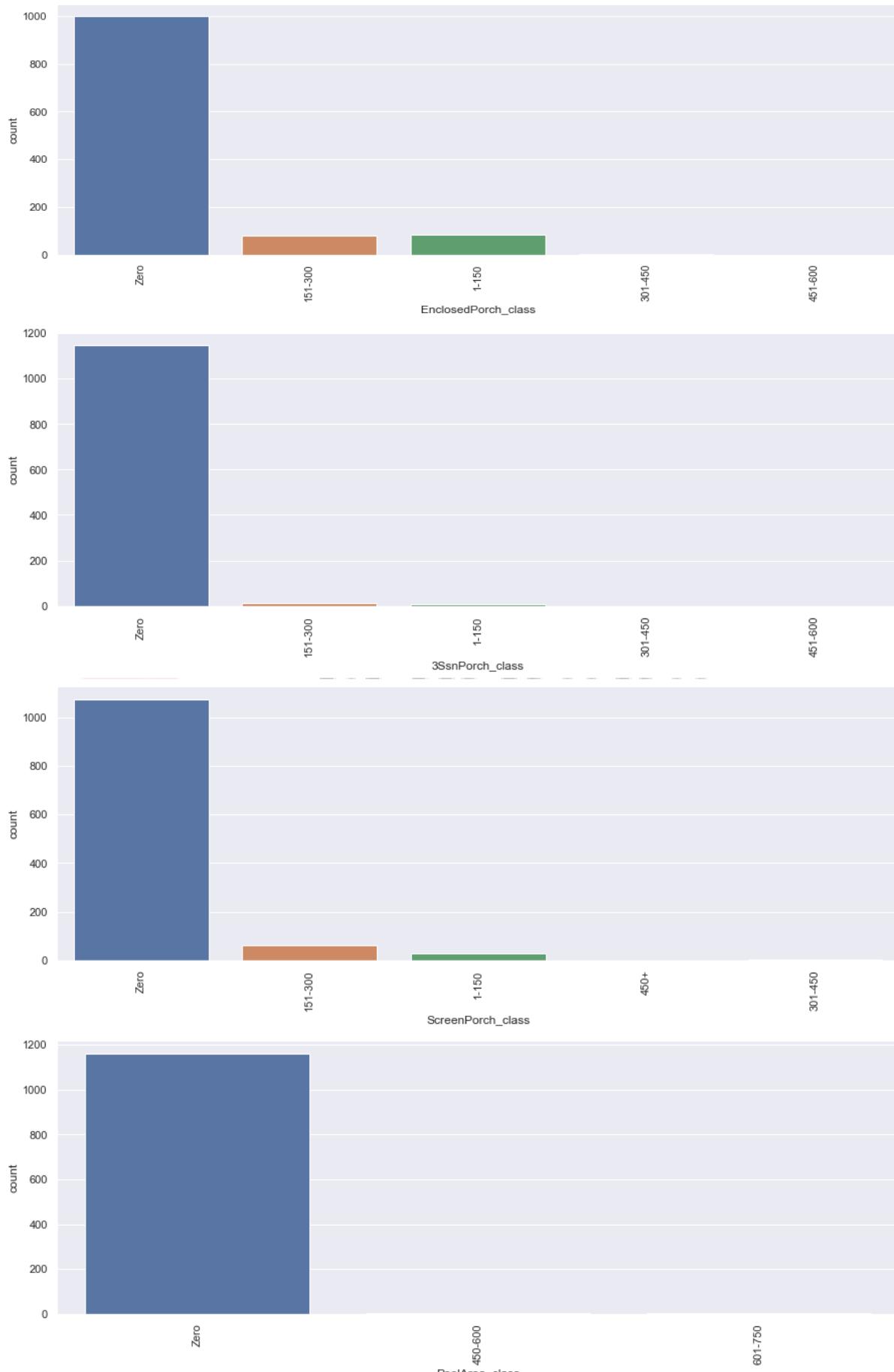
## HOUSE PRICE PREDICTION MODEL



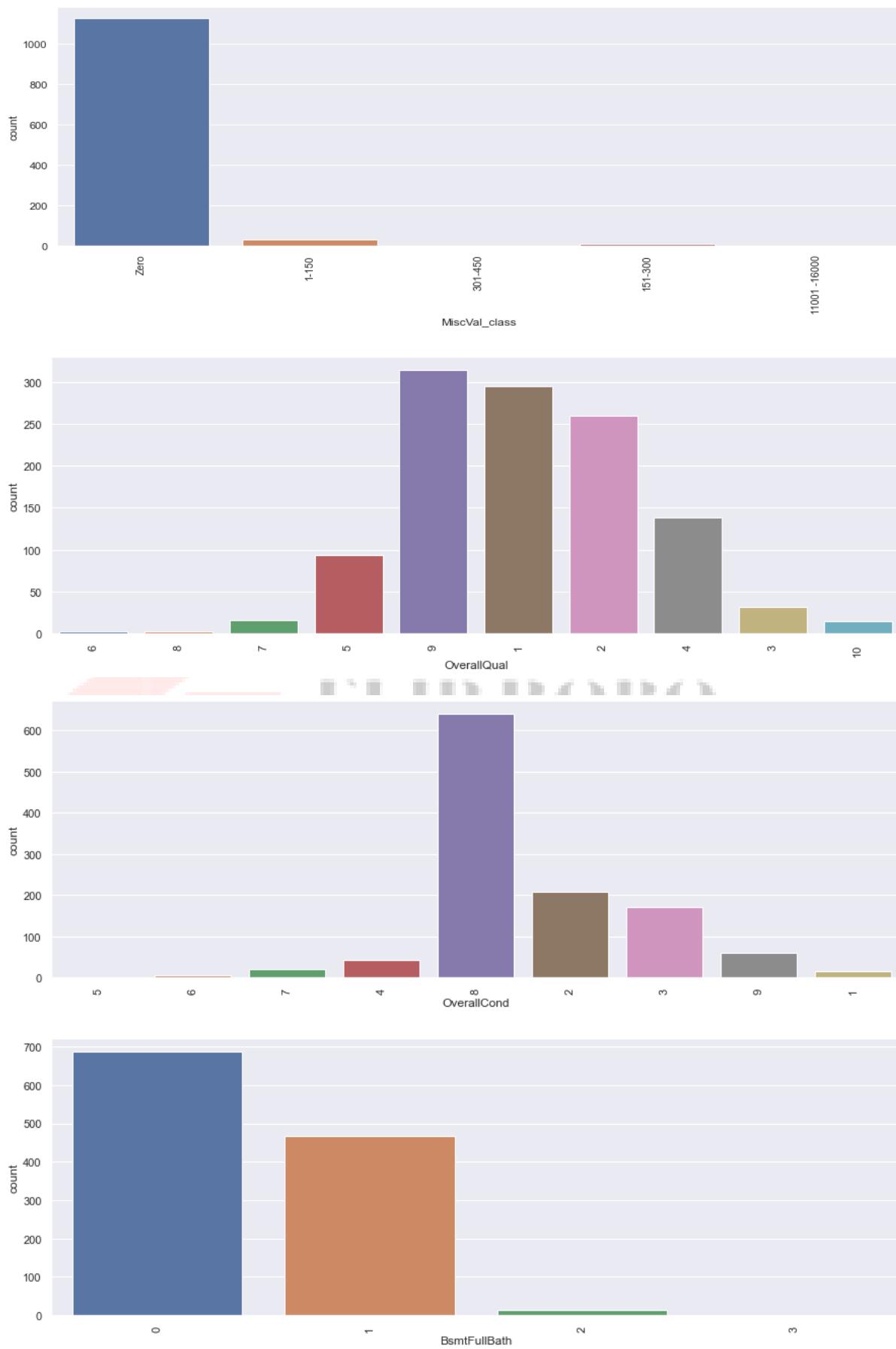
## HOUSE PRICE PREDICTION MODEL



## HOUSE PRICE PREDICTION MODEL



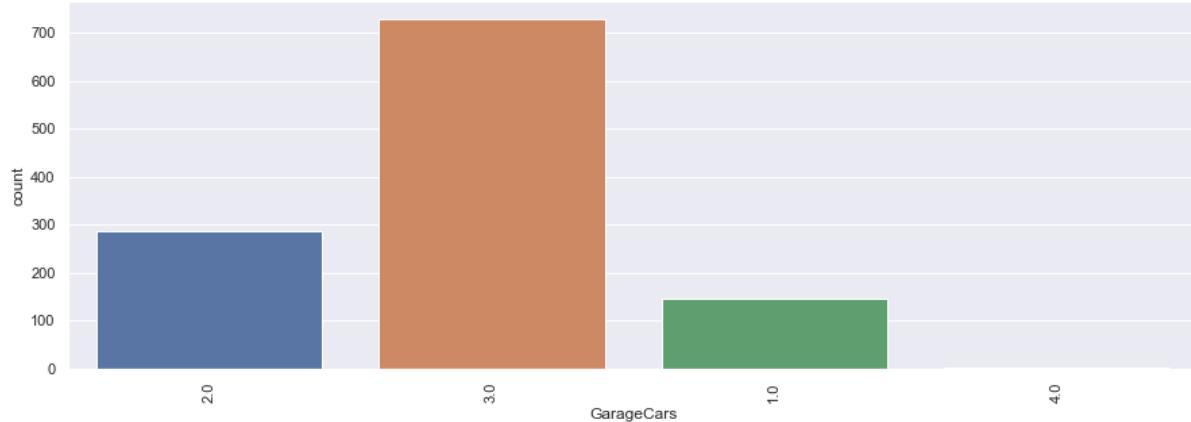
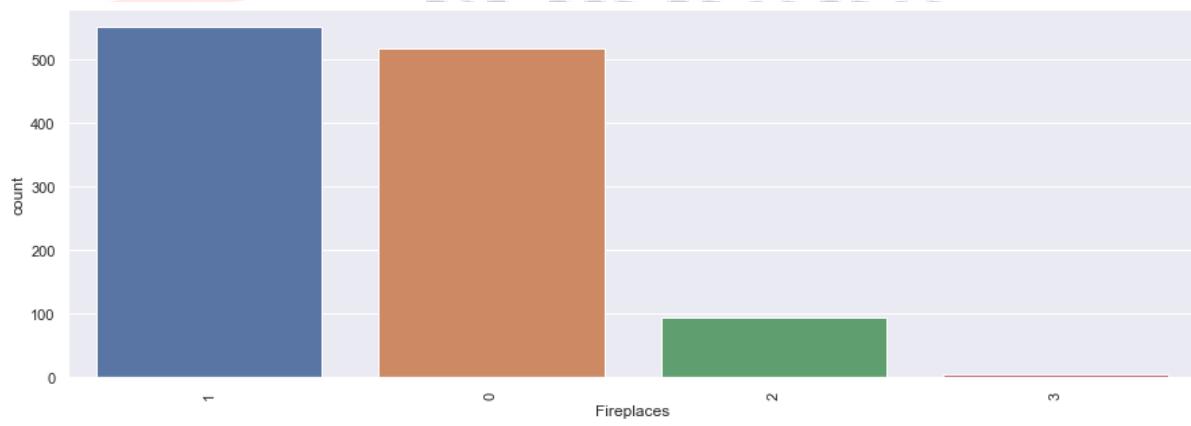
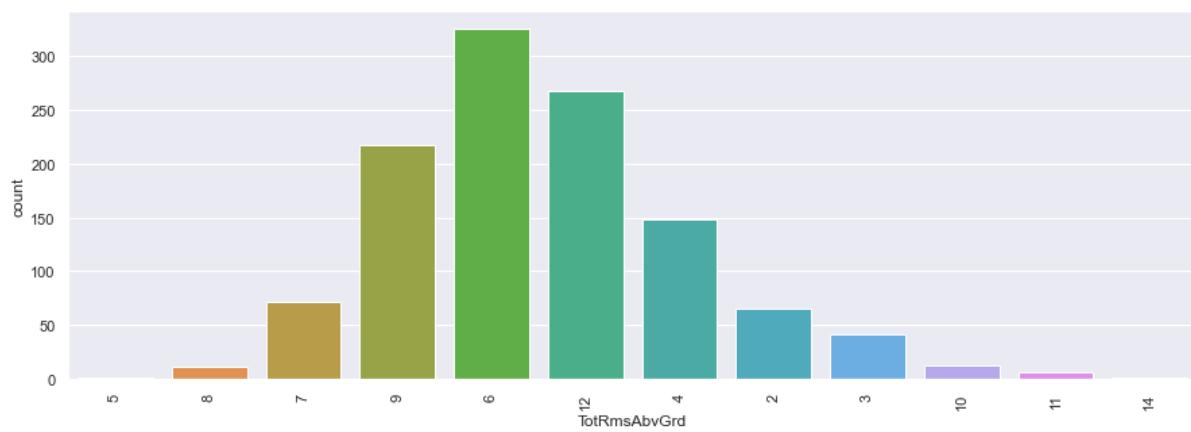
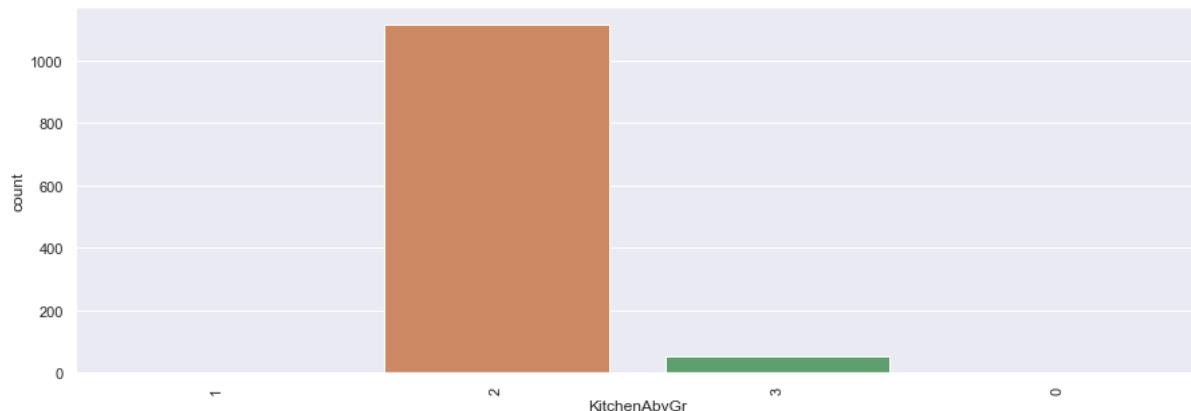
# HOUSE PRICE PREDICTION MODEL



## HOUSE PRICE PREDICTION MODEL

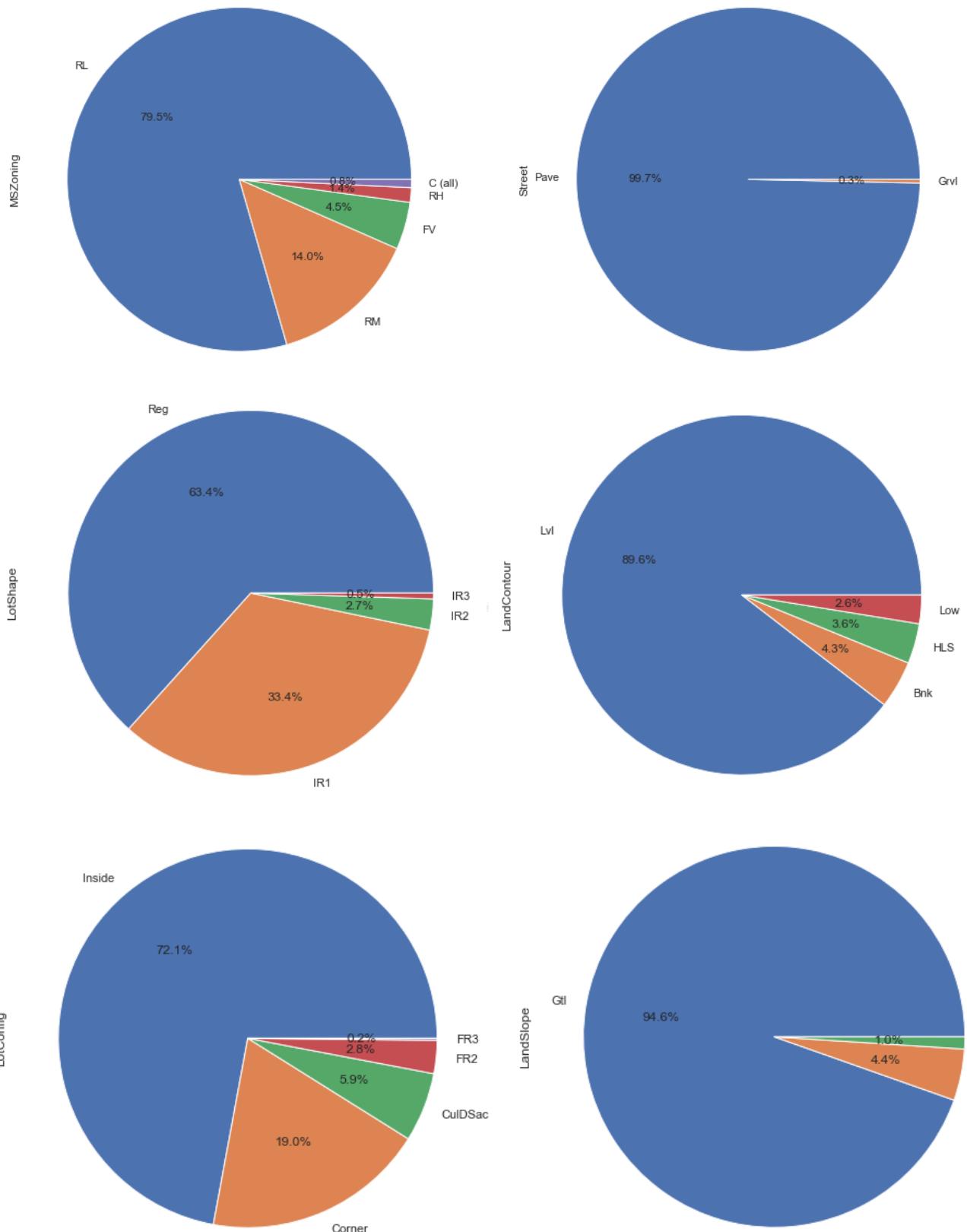


## HOUSE PRICE PREDICTION MODEL

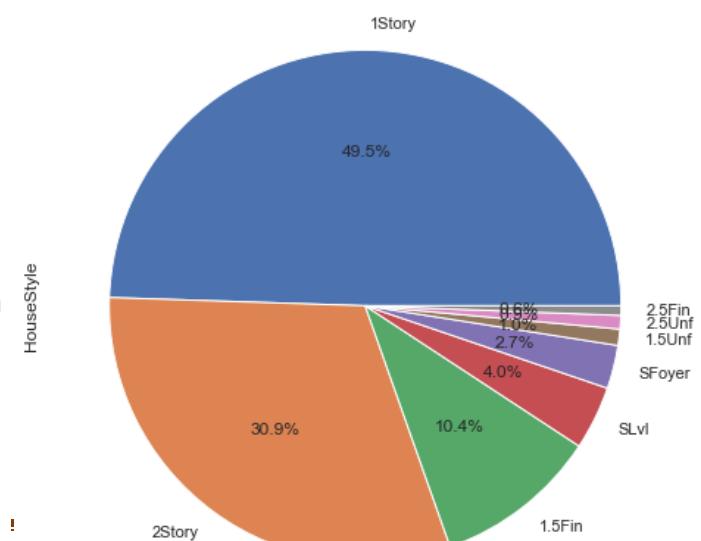
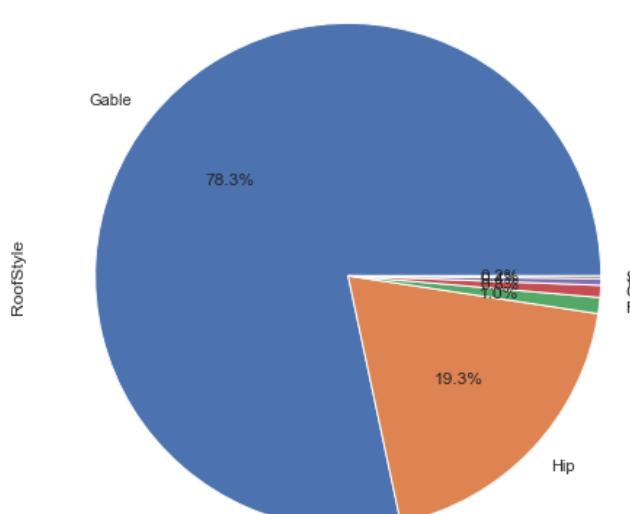
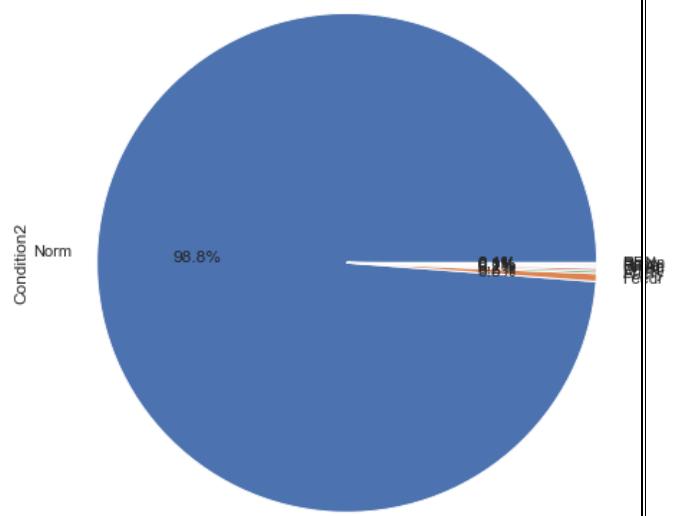
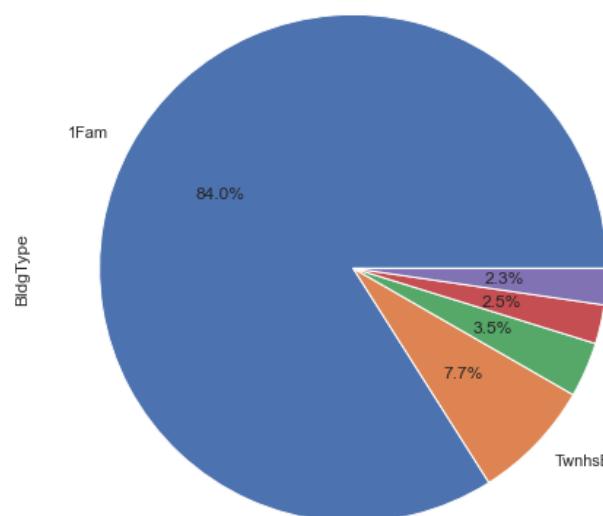
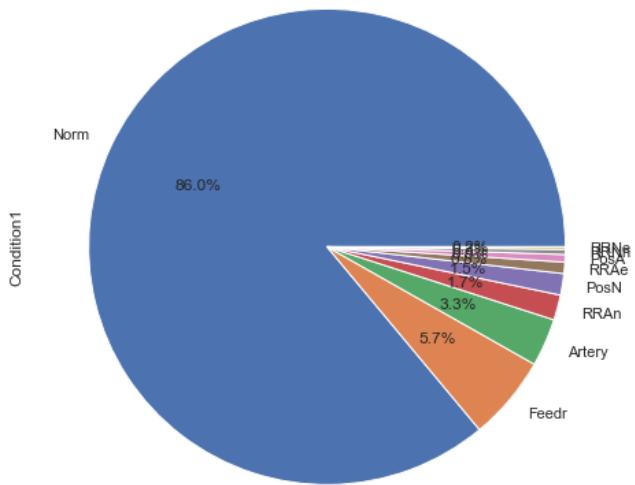
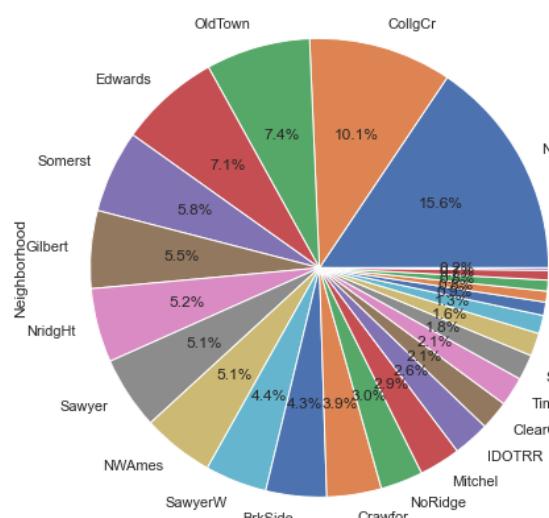


## HOUSE PRICE PREDICTION MODEL

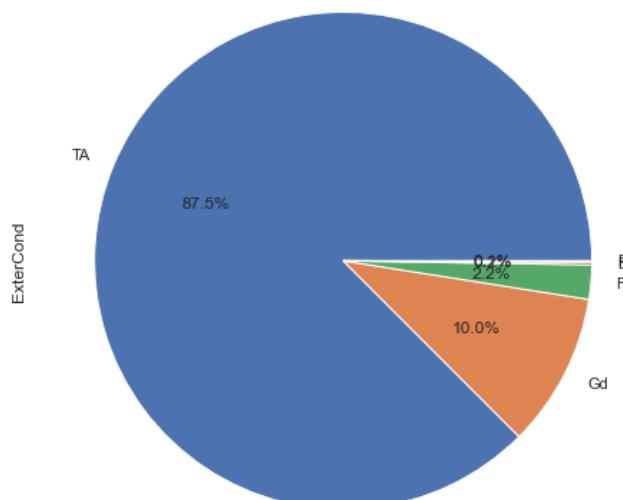
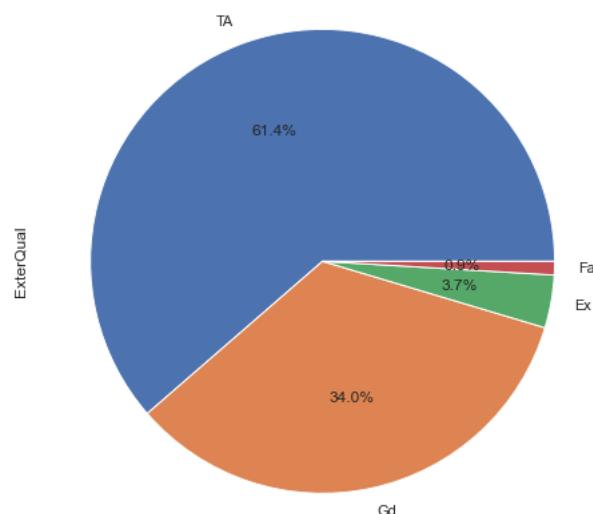
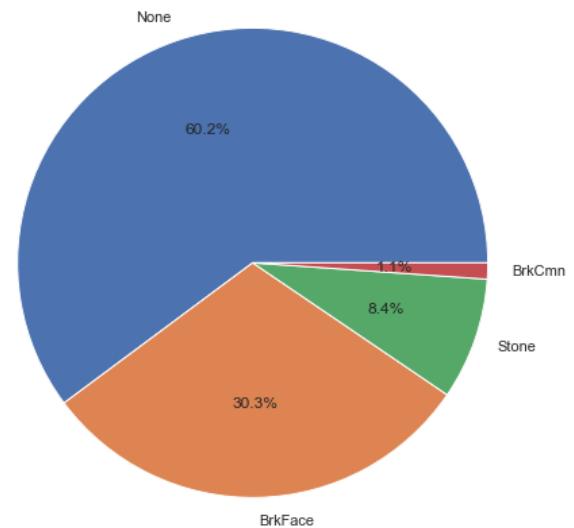
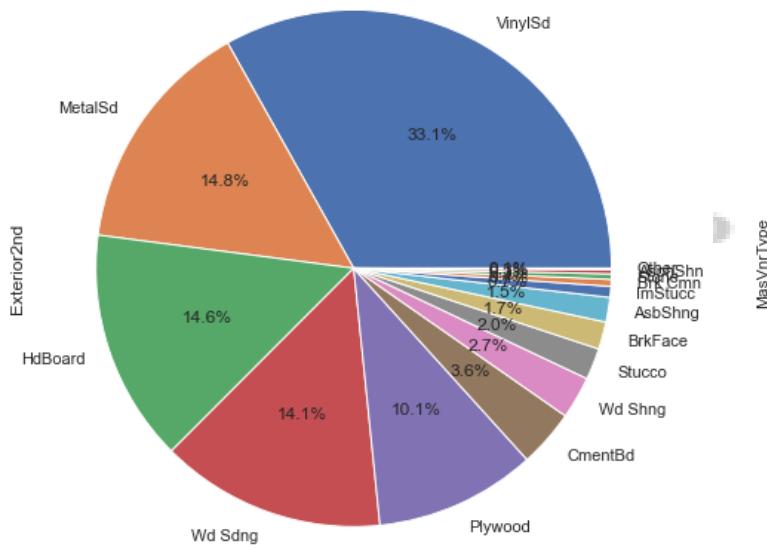
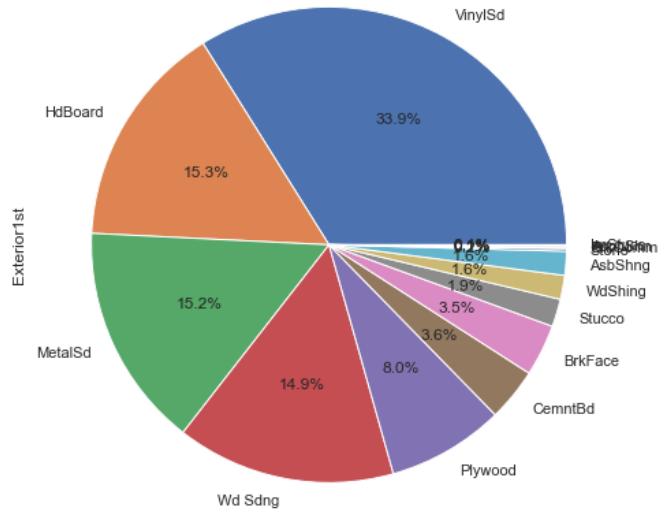
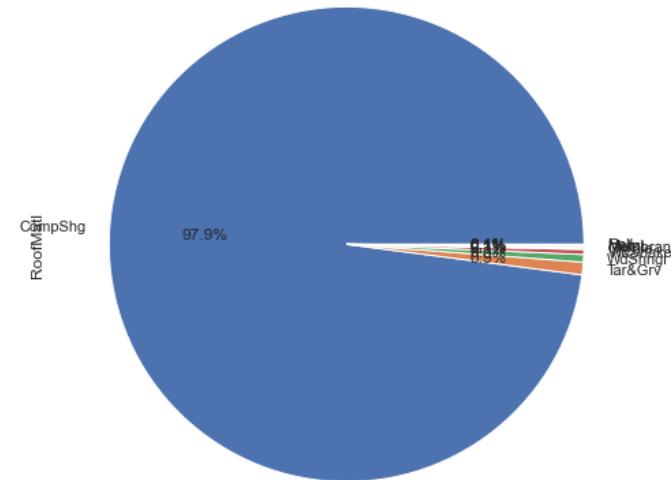
### ➤ Pie Chart:



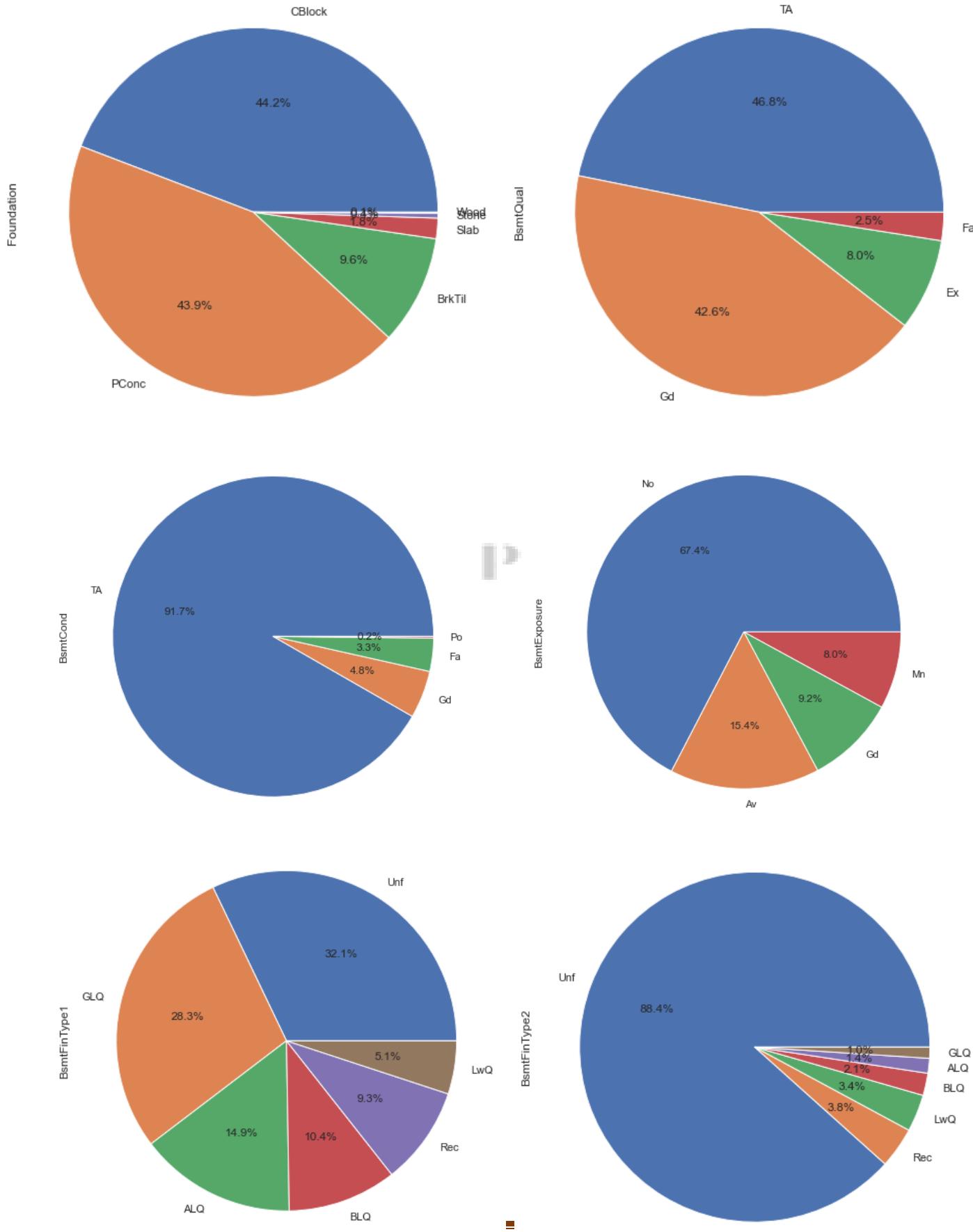
# HOUSE PRICE PREDICTION MODEL



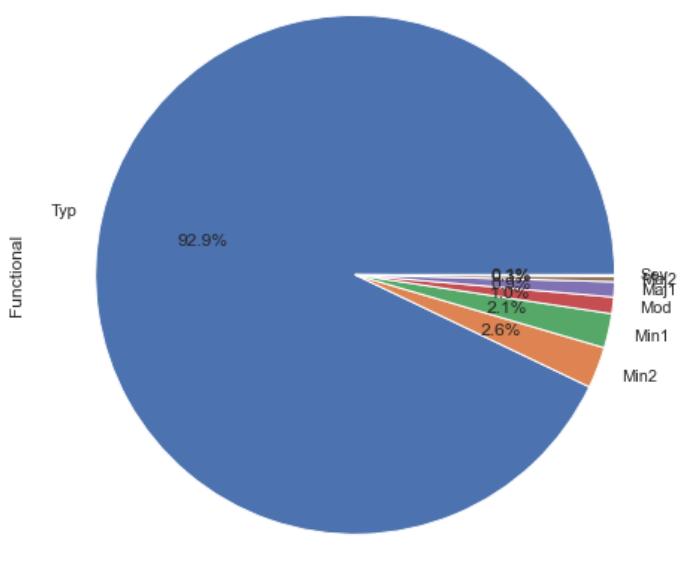
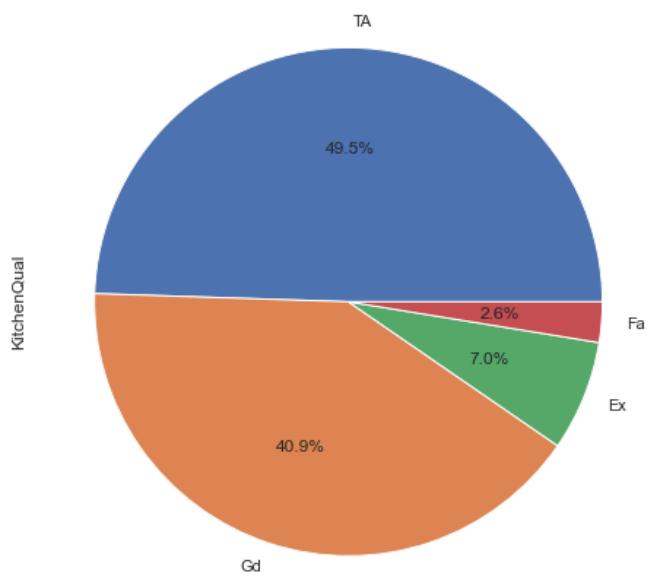
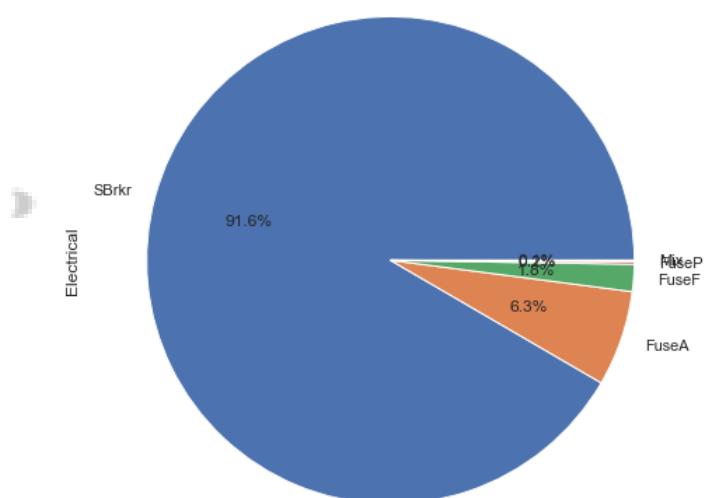
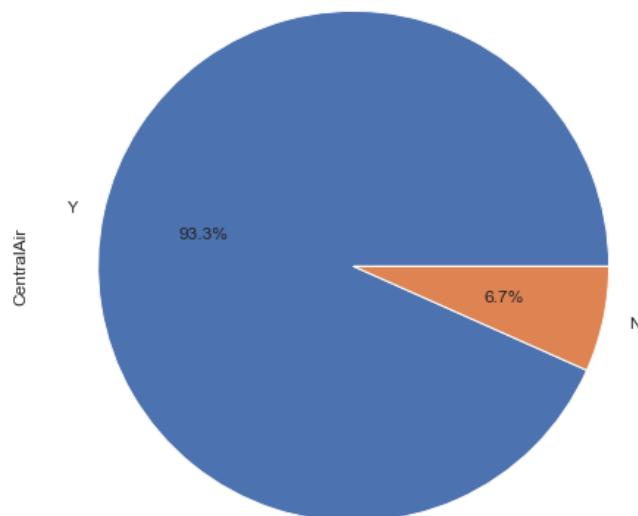
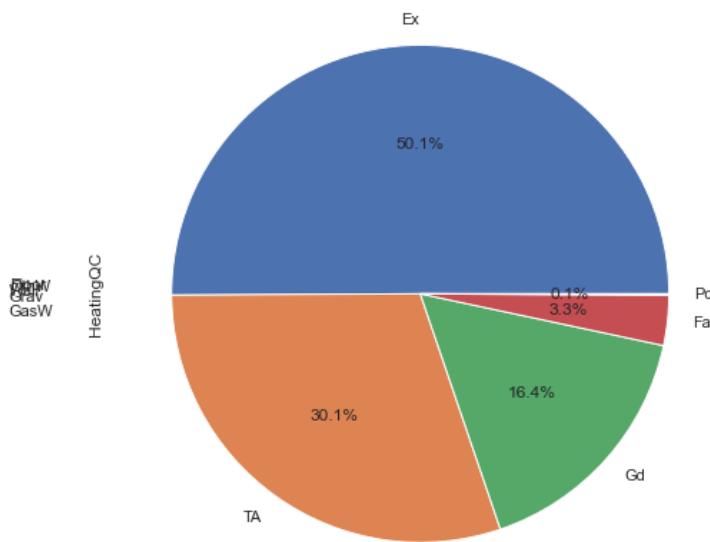
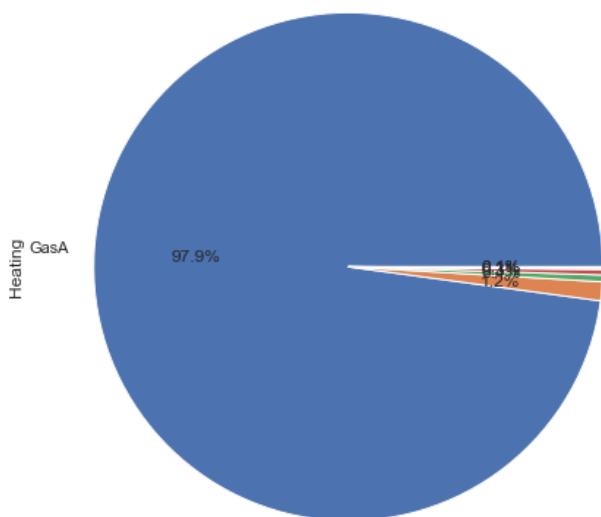
## HOUSE PRICE PREDICTION MODEL



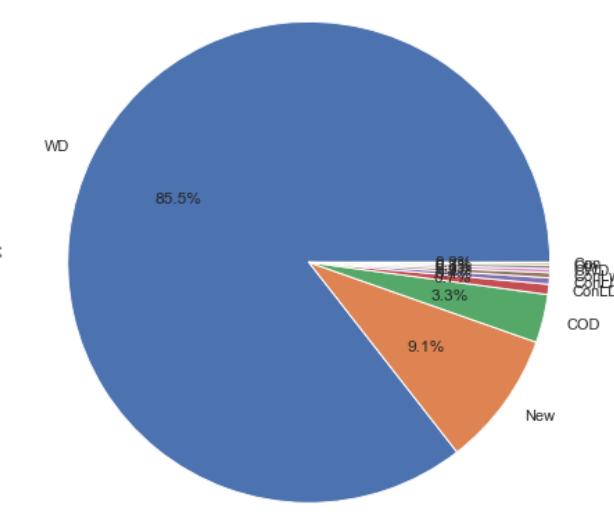
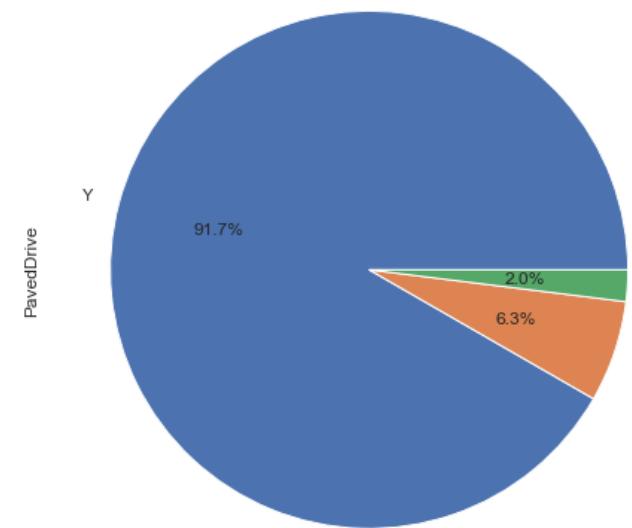
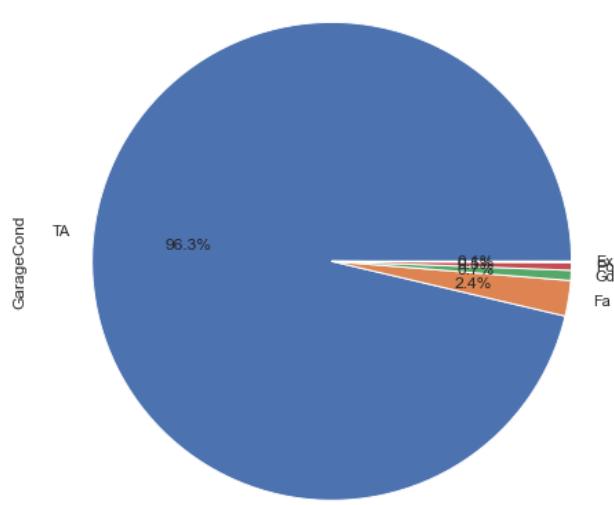
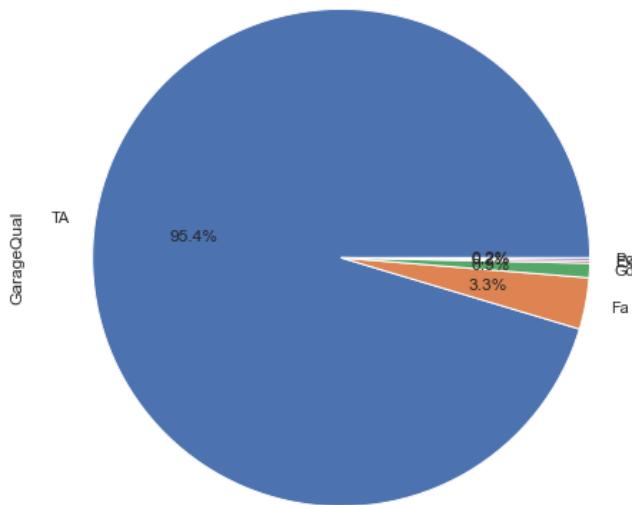
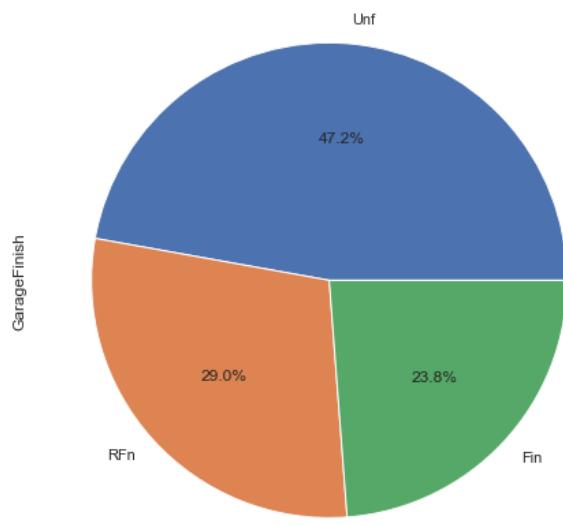
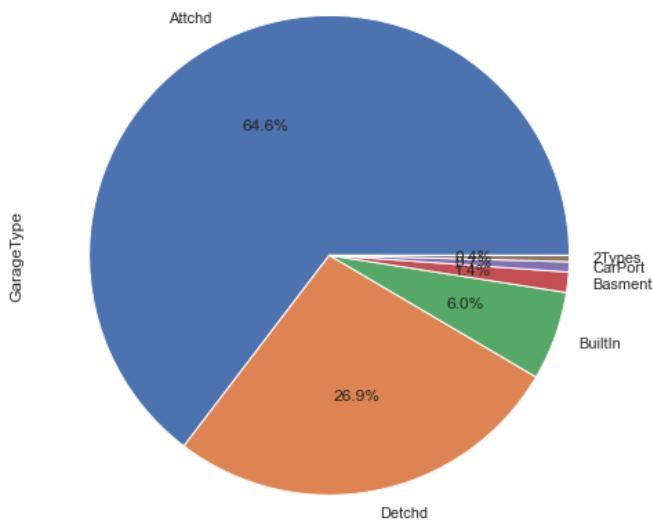
# HOUSE PRICE PREDICTION MODEL



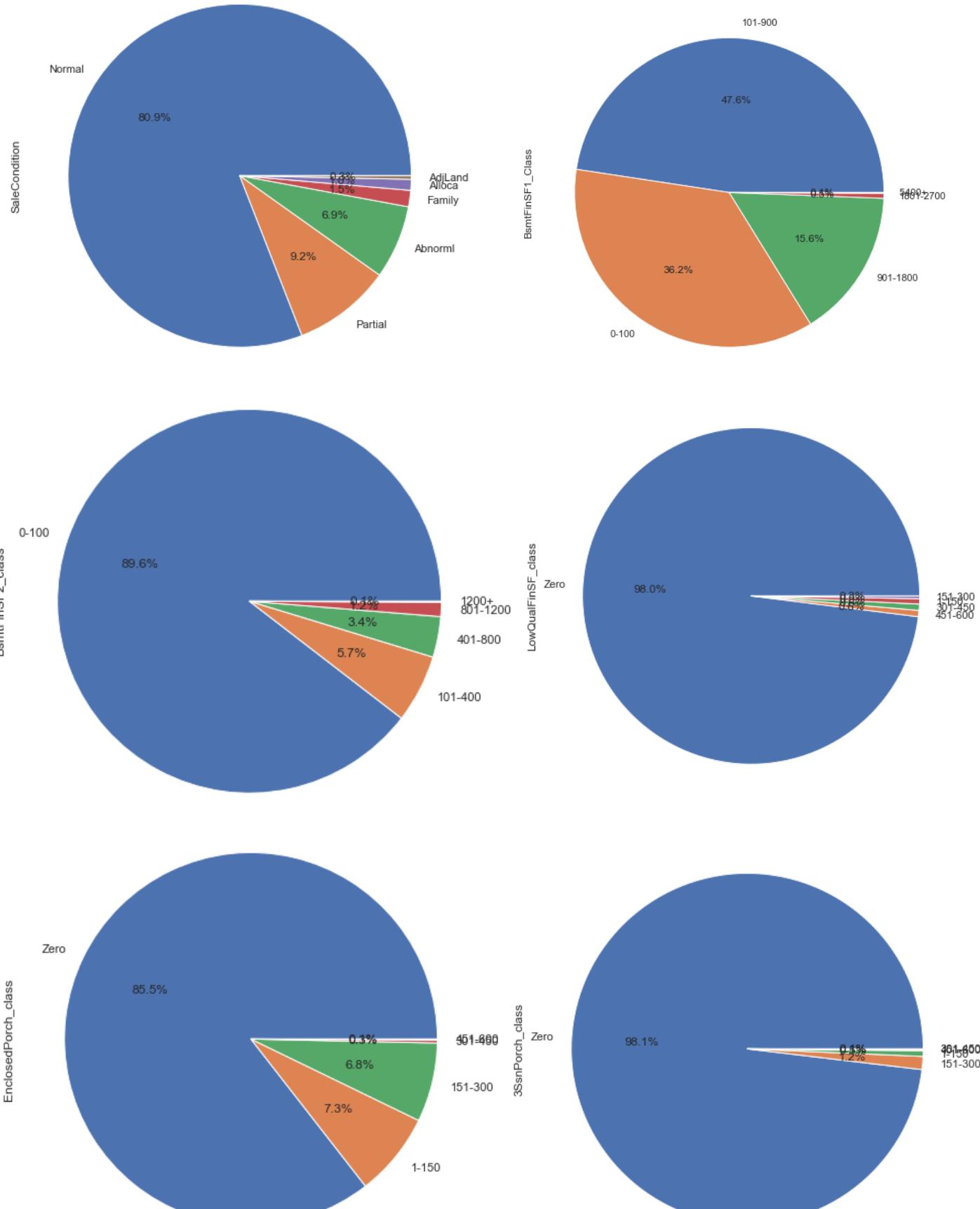
# HOUSE PRICE PREDICTION MODEL



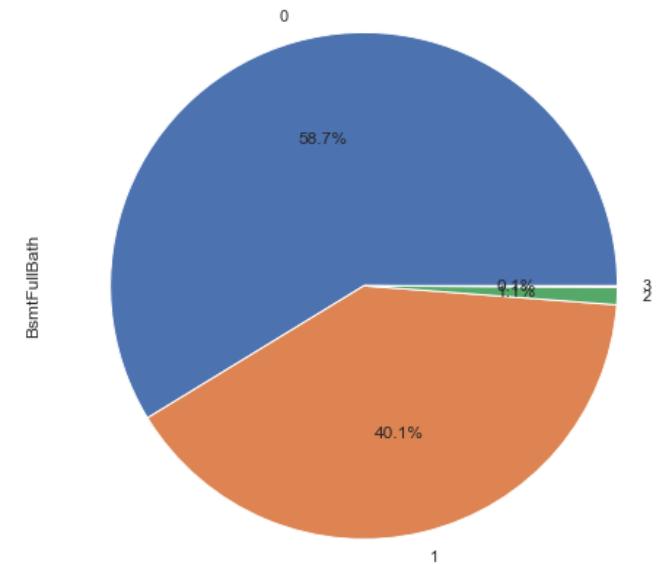
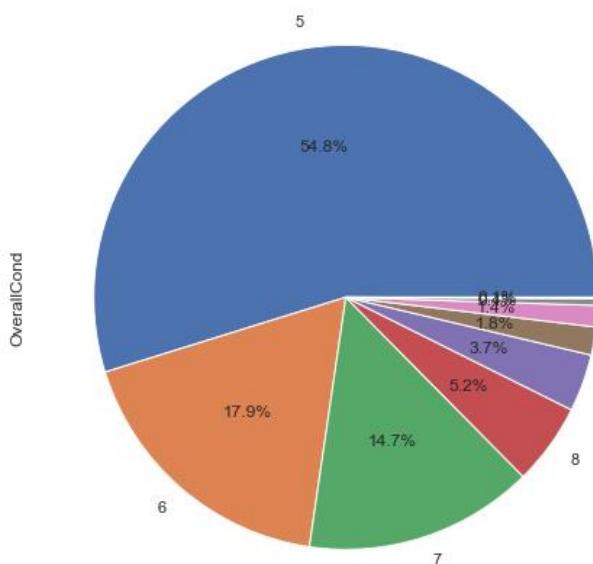
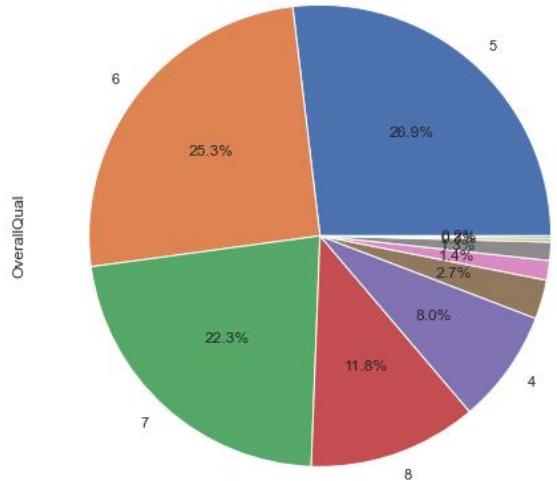
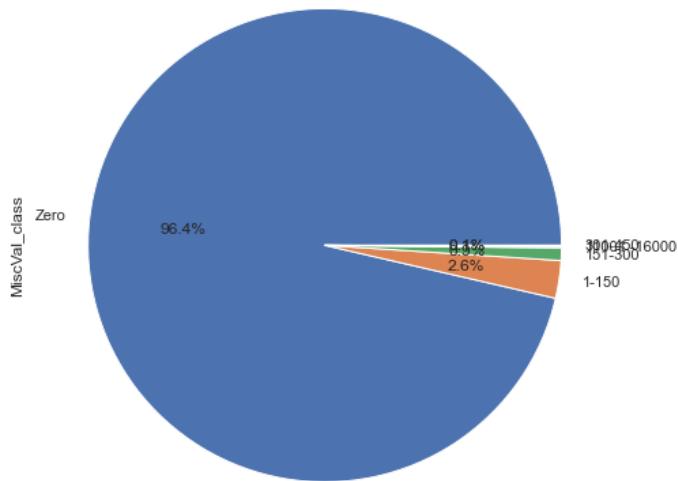
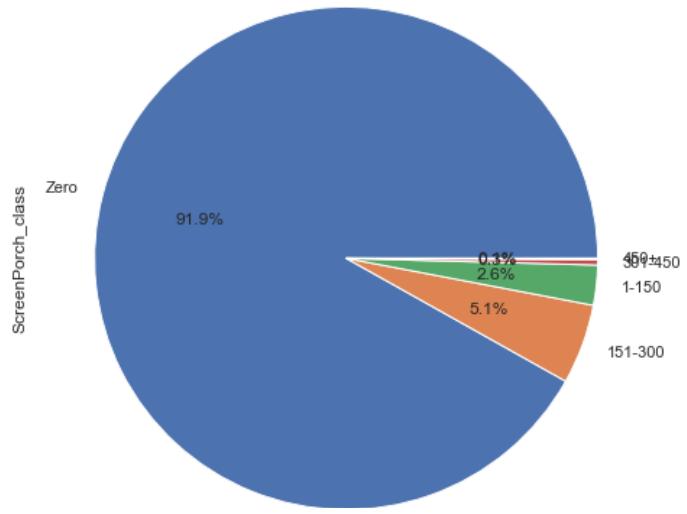
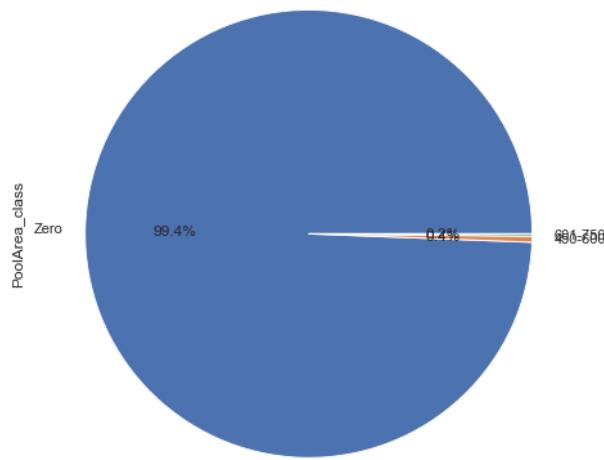
## HOUSE PRICE PREDICTION MODEL



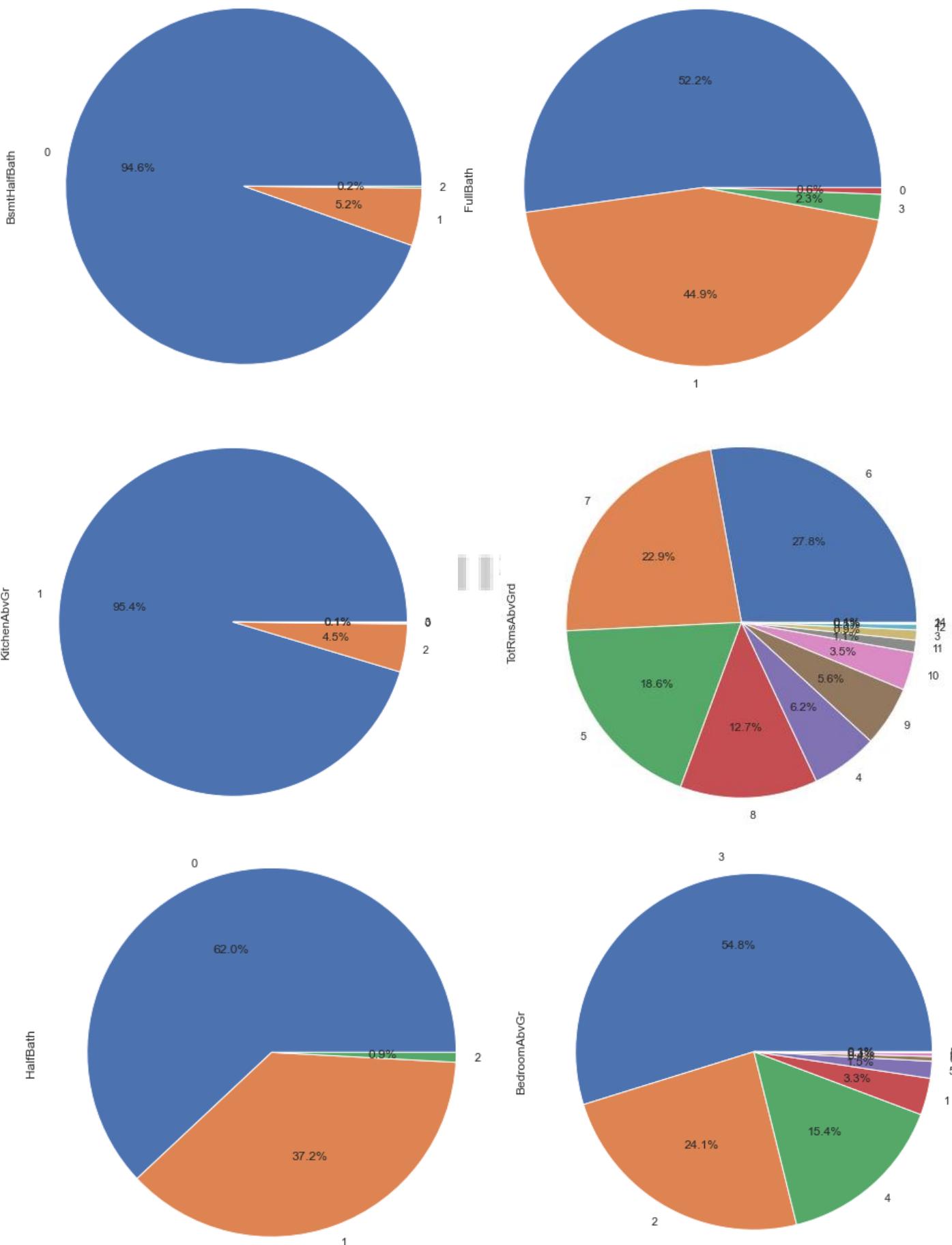
## HOUSE PRICE PREDICTION MODEL



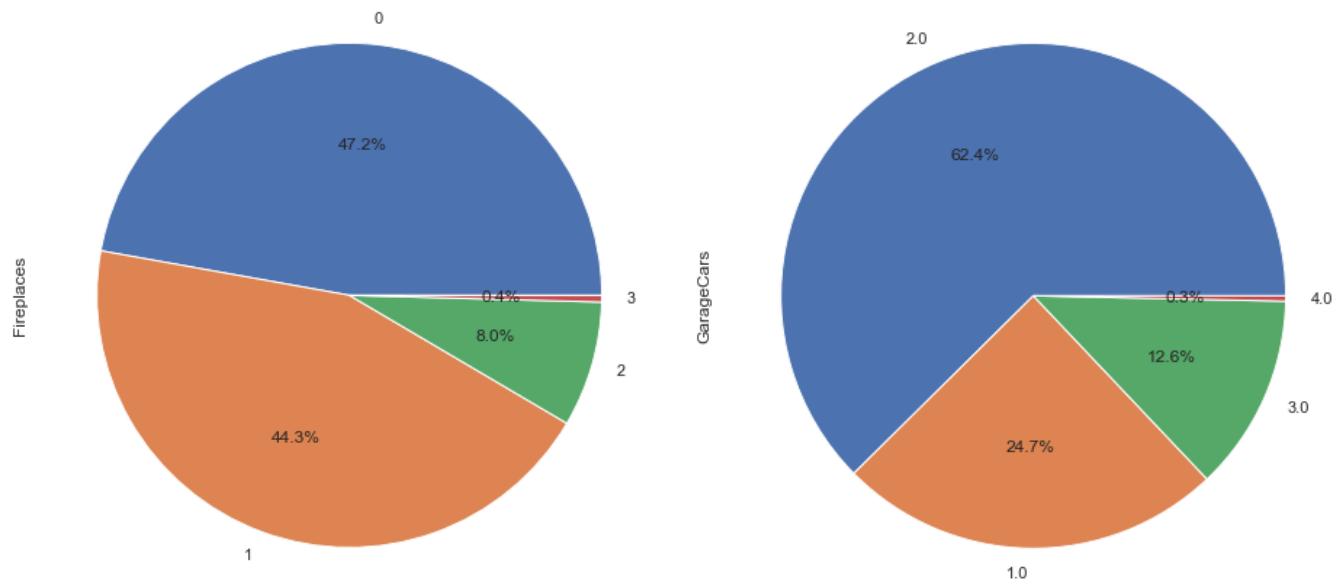
## HOUSE PRICE PREDICTION MODEL



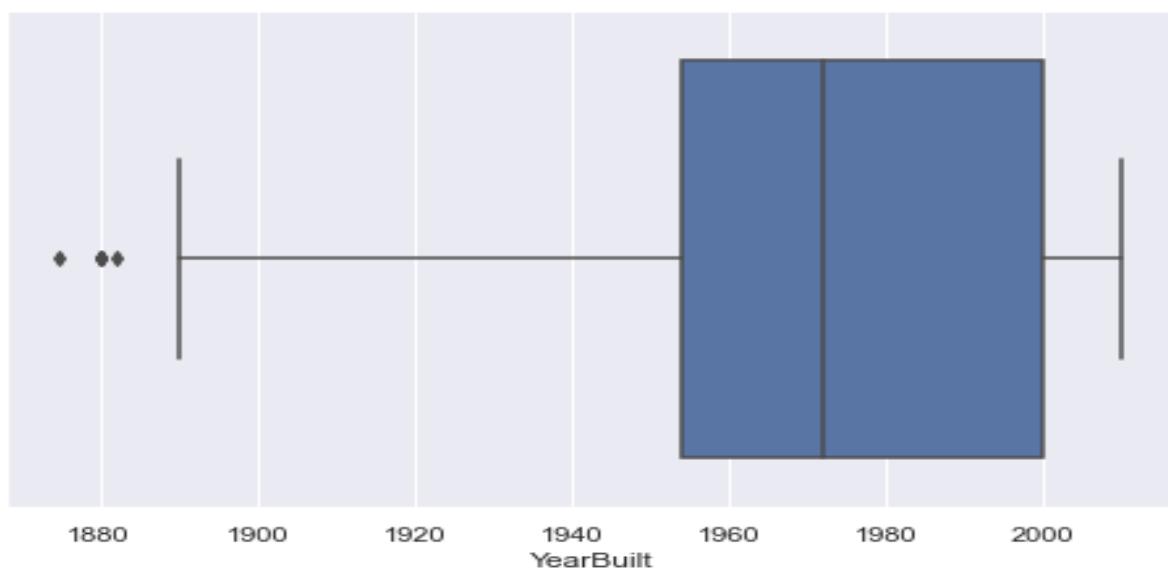
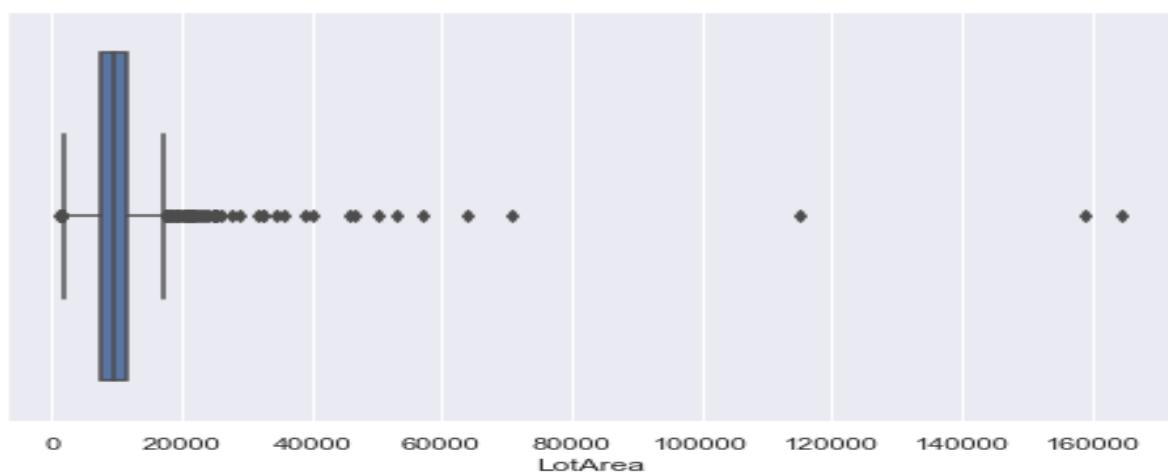
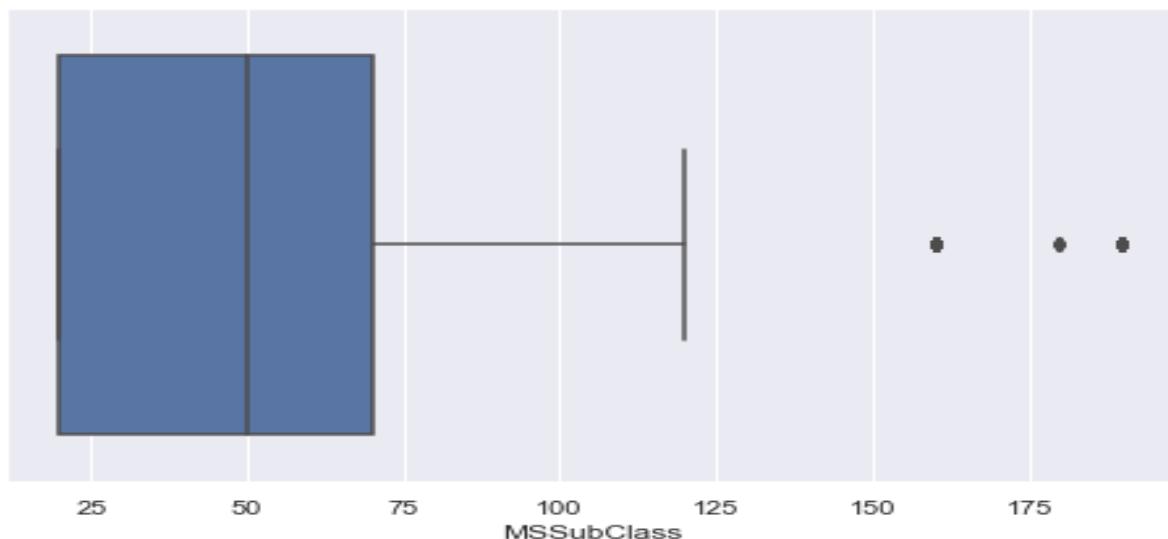
## HOUSE PRICE PREDICTION MODEL



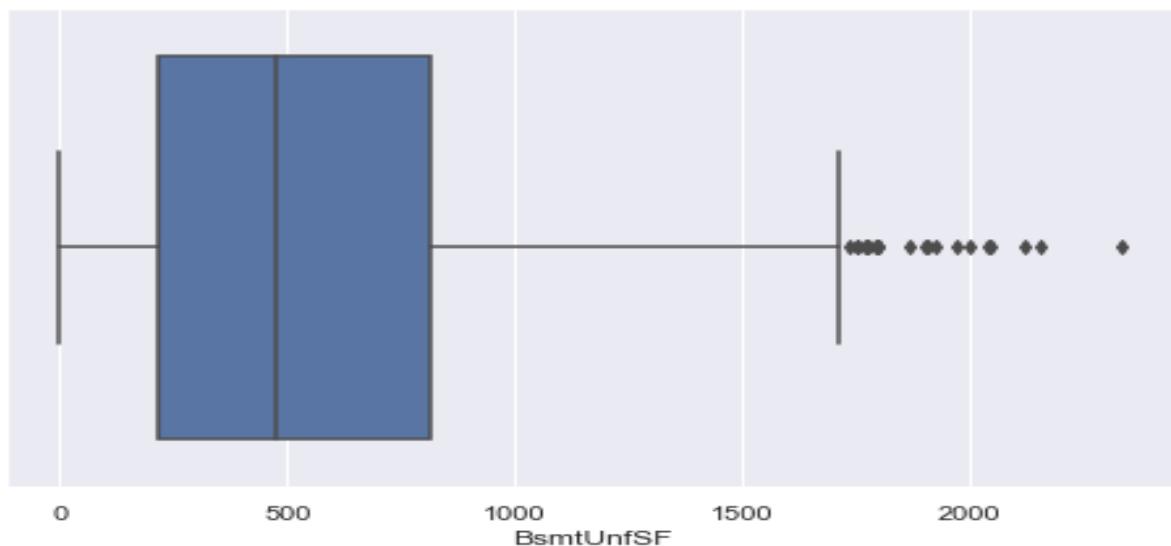
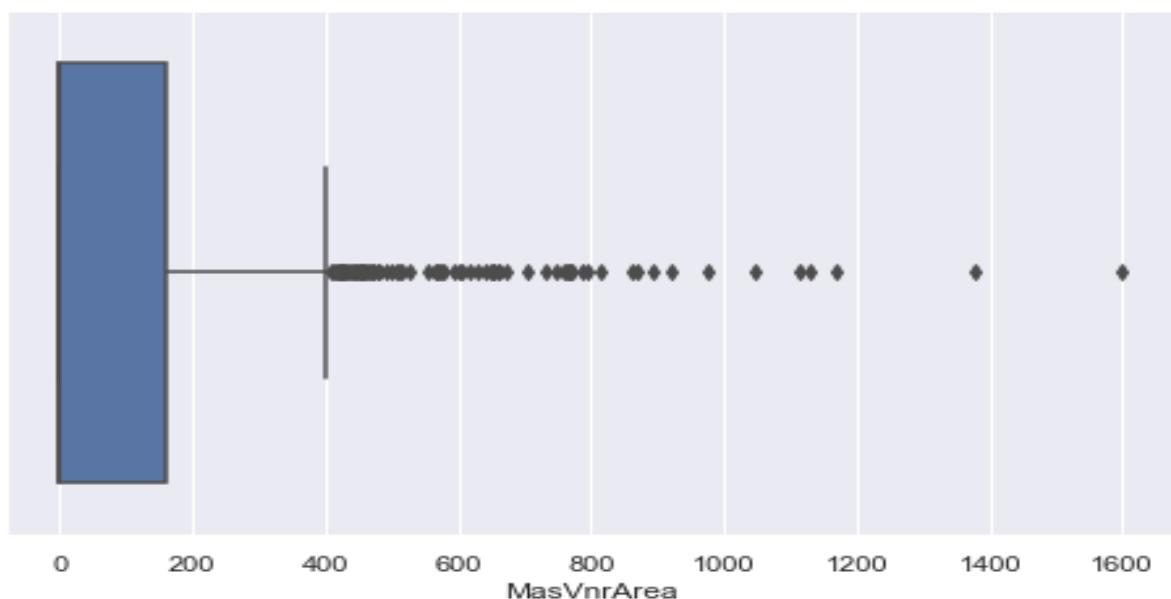
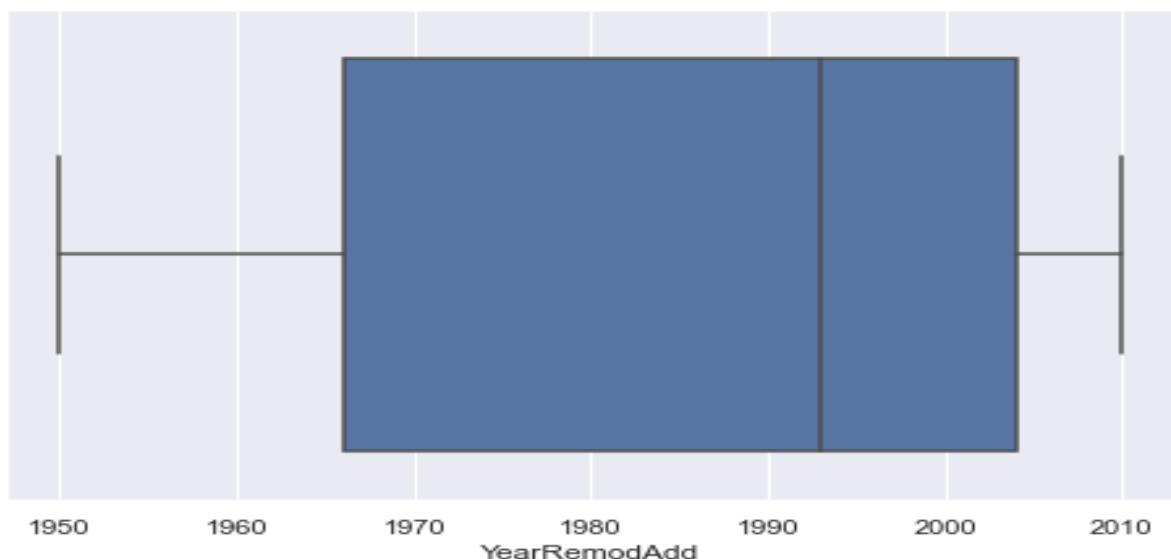
## HOUSE PRICE PREDICTION MODEL



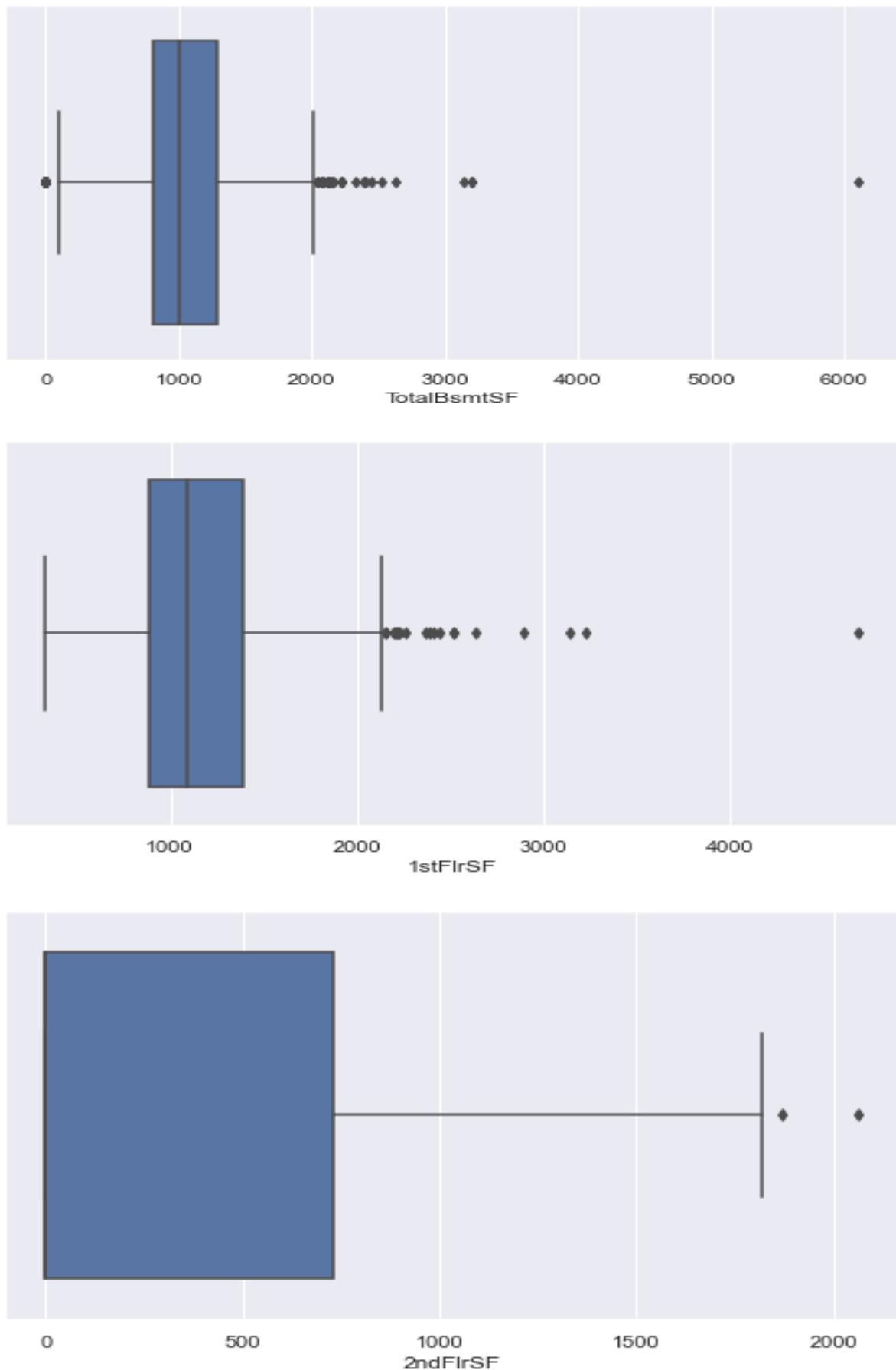
## Outliers:



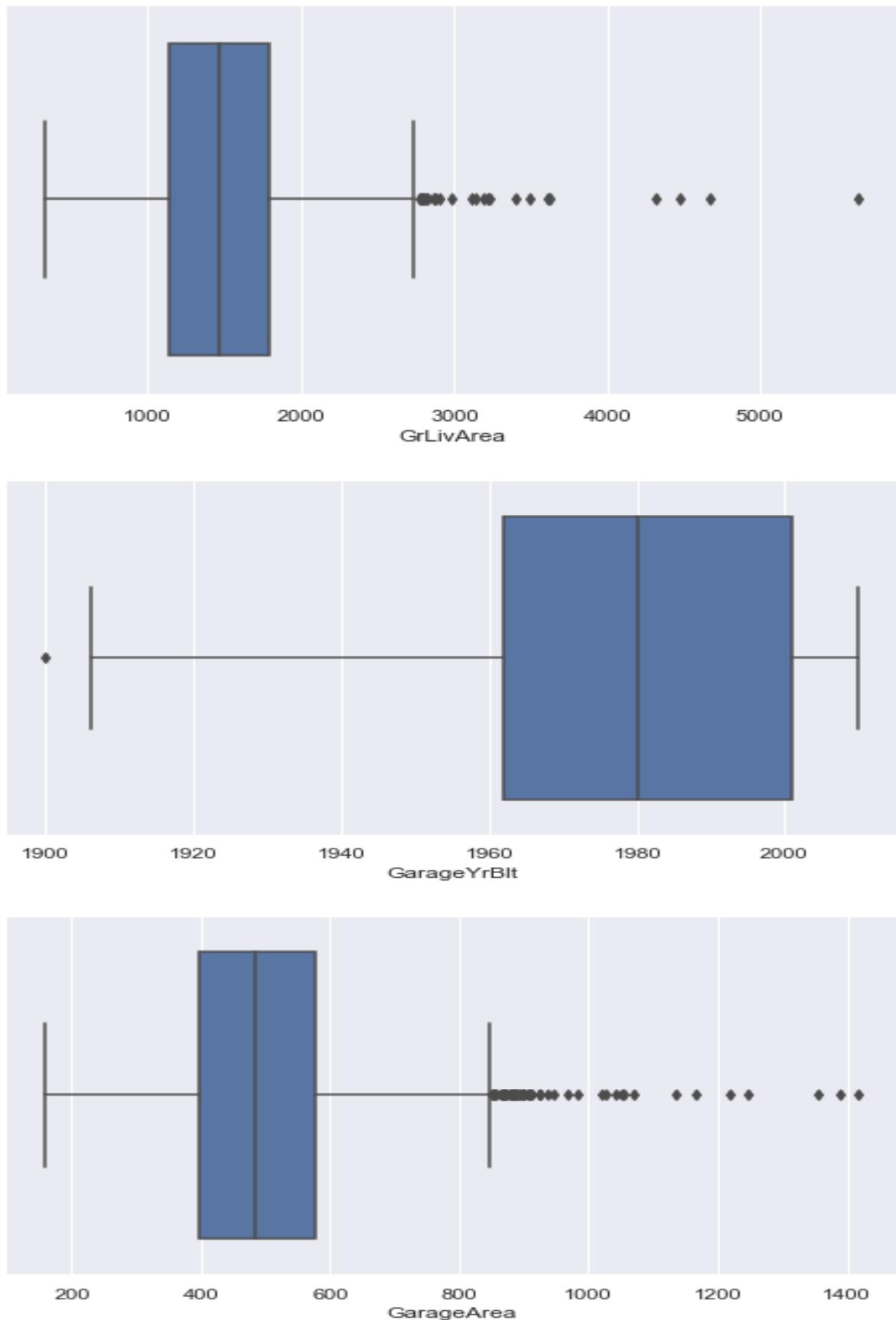
## HOUSE PRICE PREDICTION MODEL



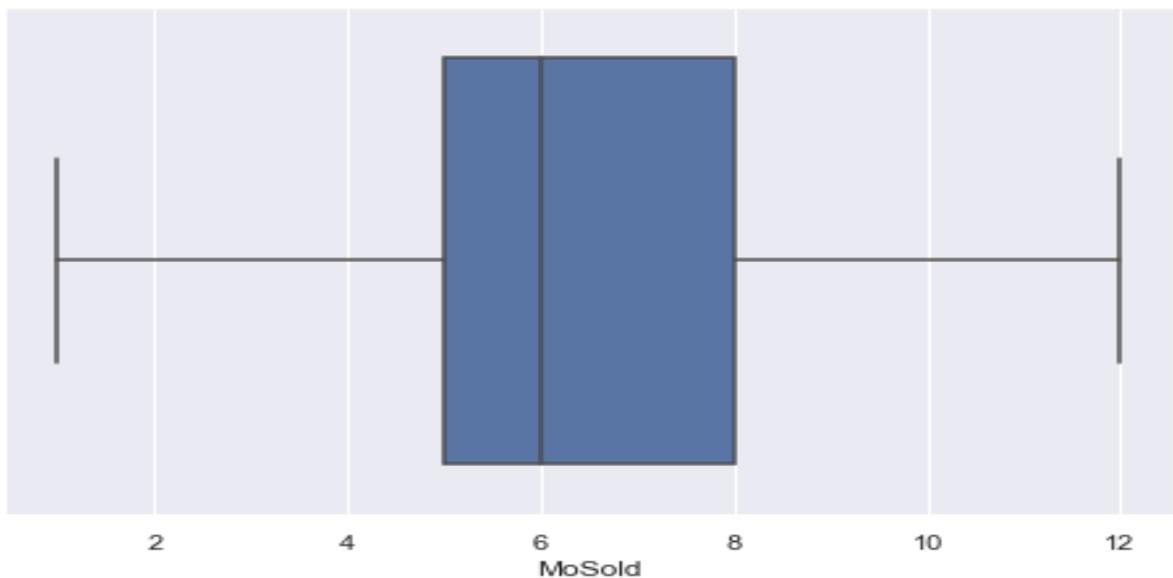
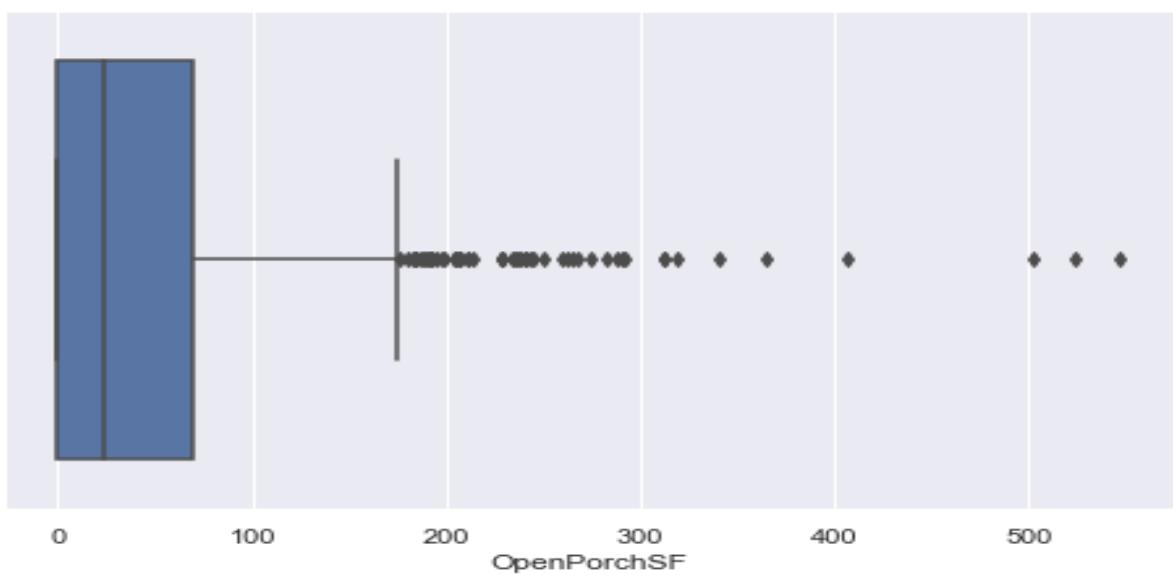
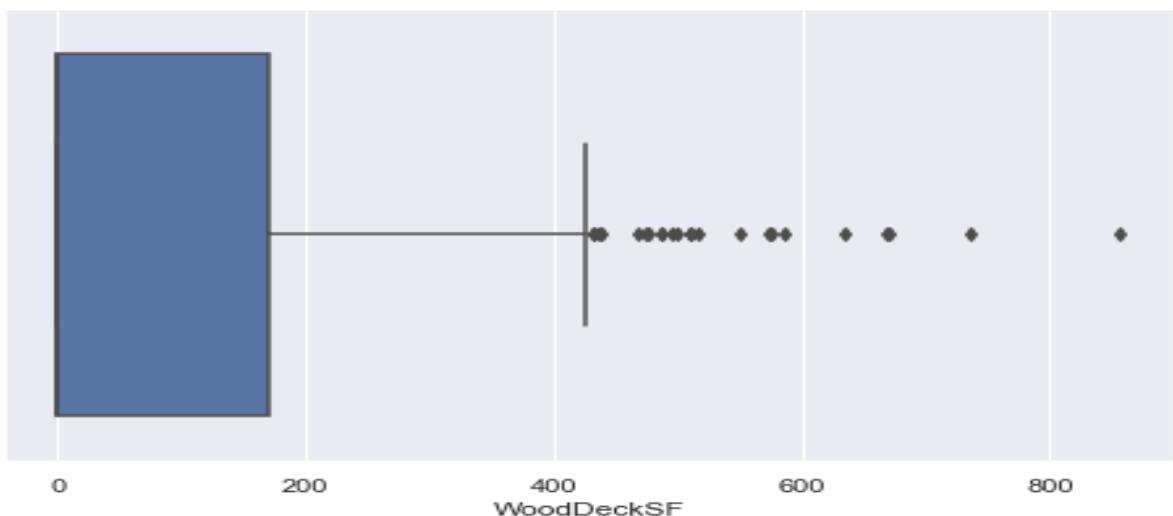
## HOUSE PRICE PREDICTION MODEL



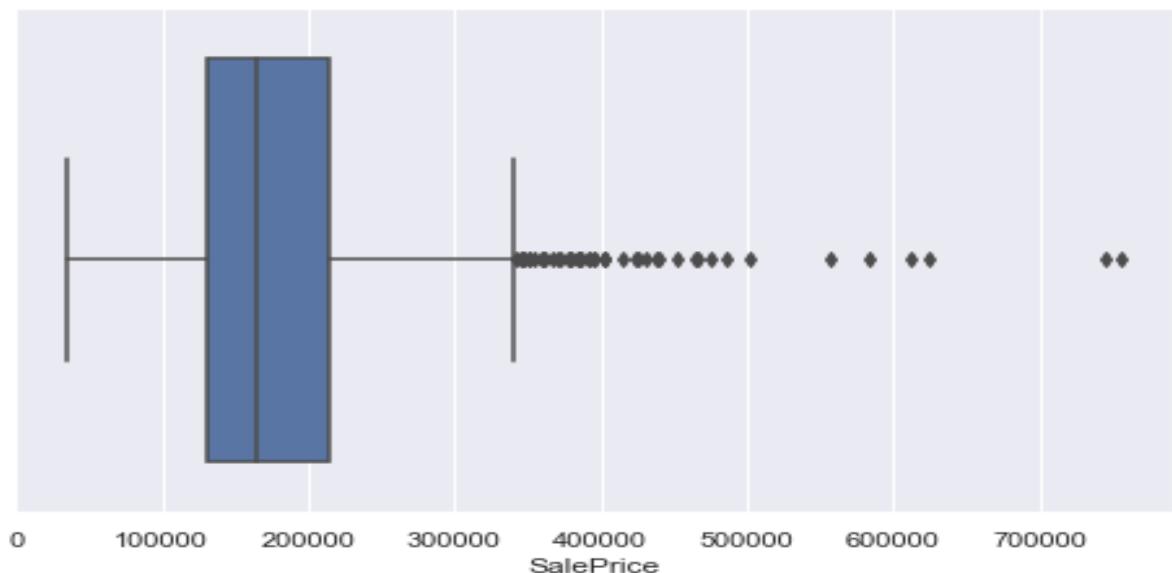
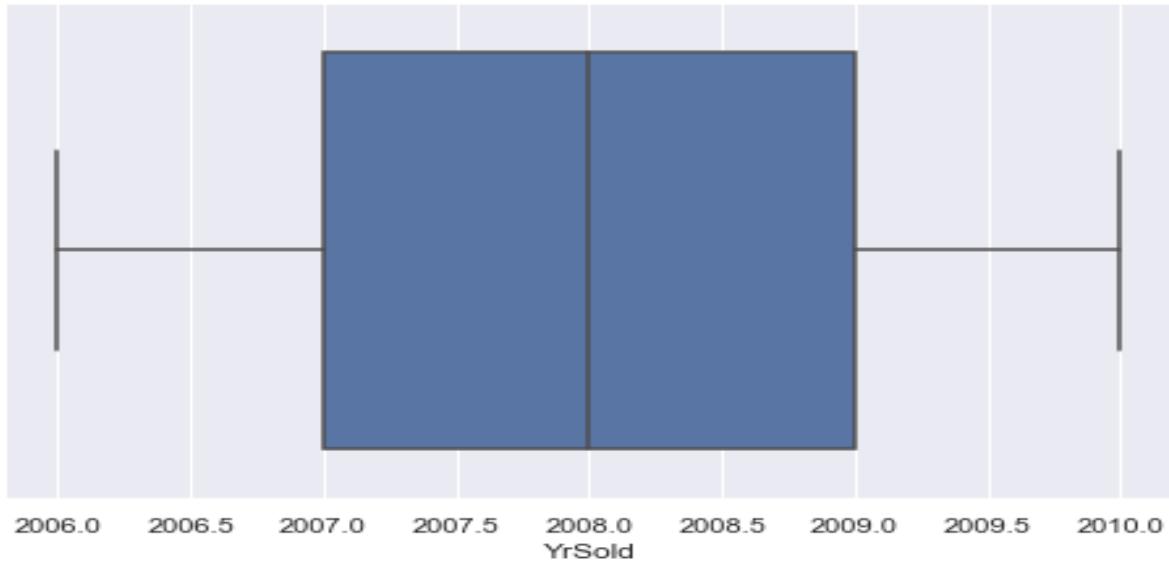
## HOUSE PRICE PREDICTION MODEL



## HOUSE PRICE PREDICTION MODEL



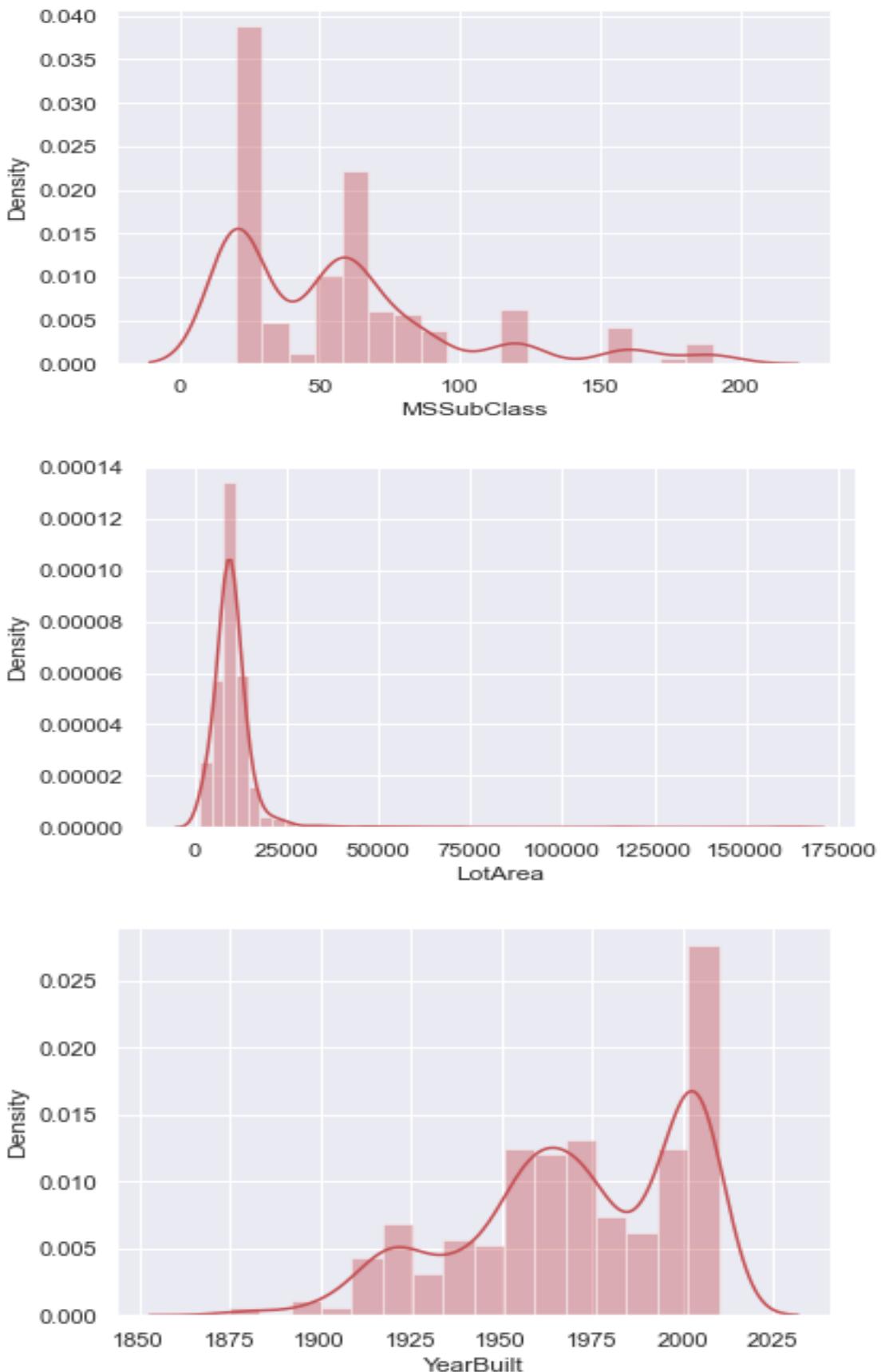
## HOUSE PRICE PREDICTION MODEL



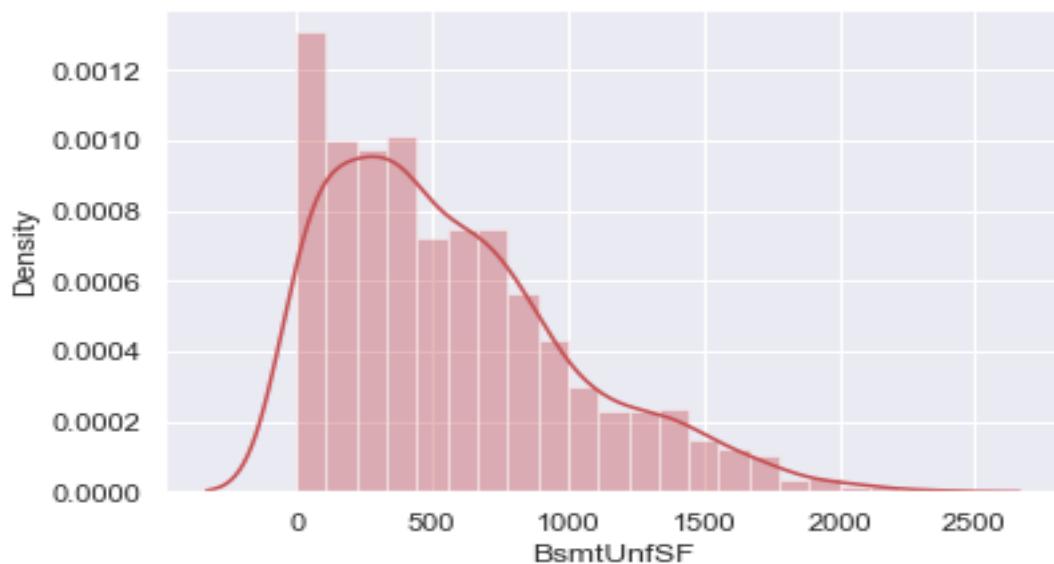
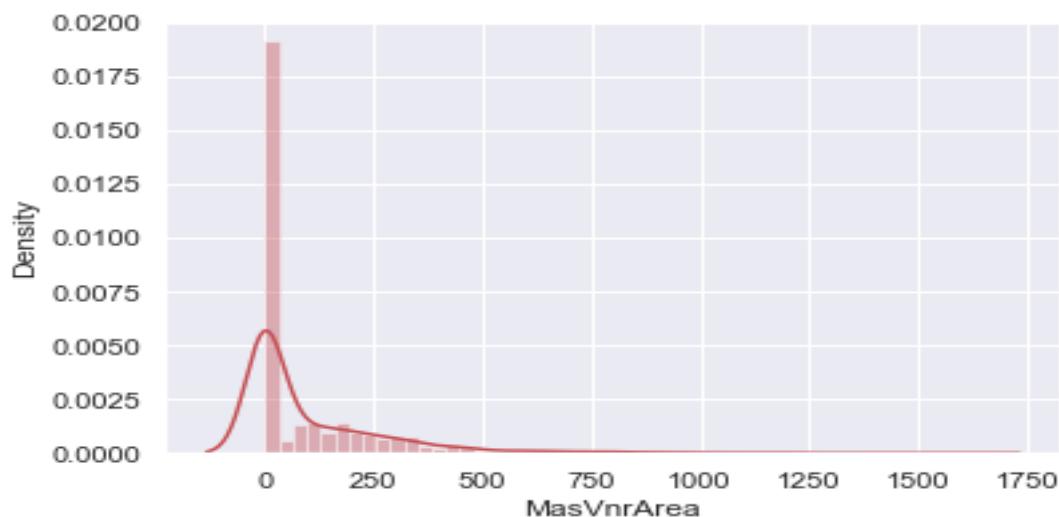
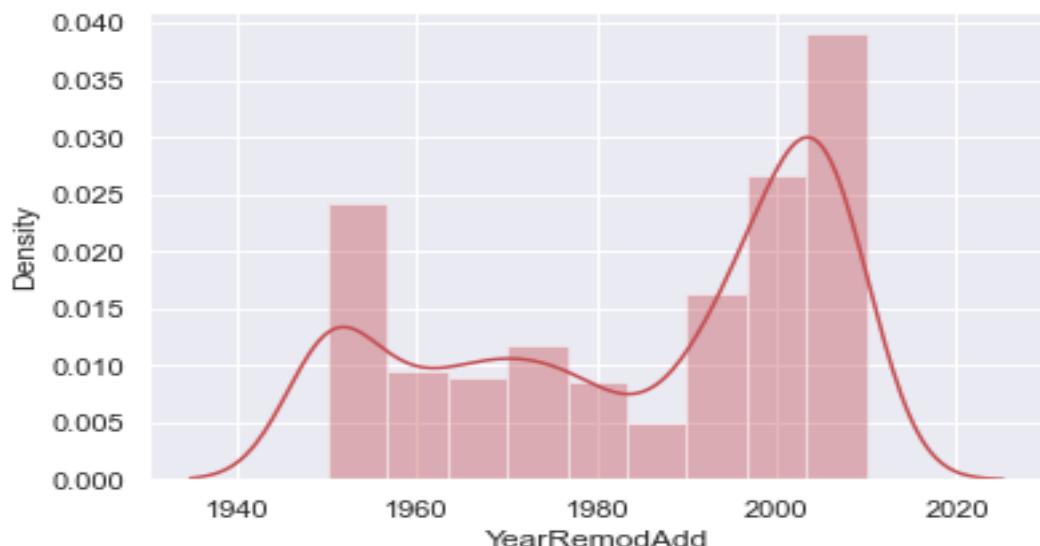
**Keys:** Columns are having outliers, but when I performed IQR and zscore method to remove outliers these both method was giving high loss of data.

Therefore, outliers not removed and worked on existing dataset

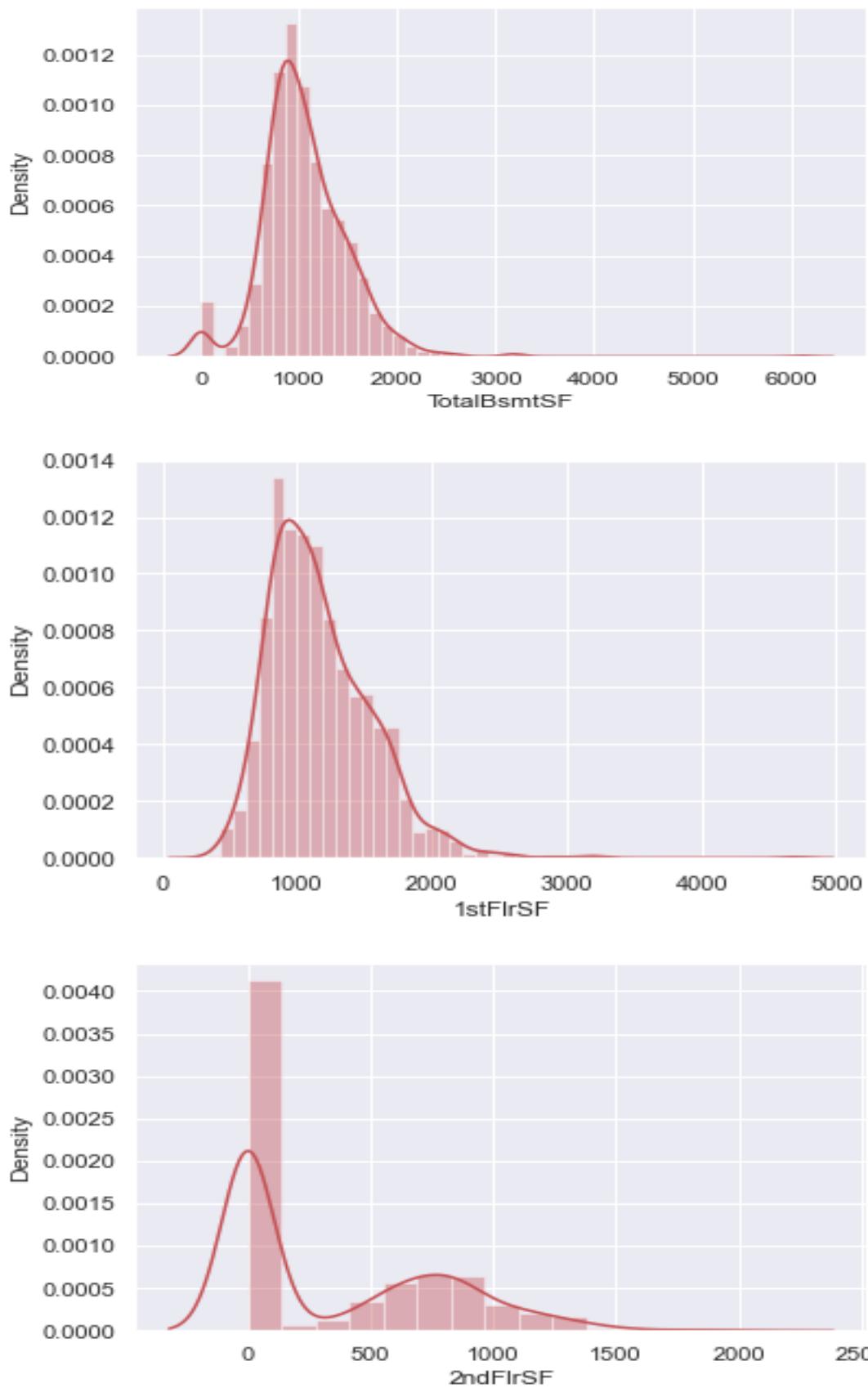
## Distribution and Skewness:



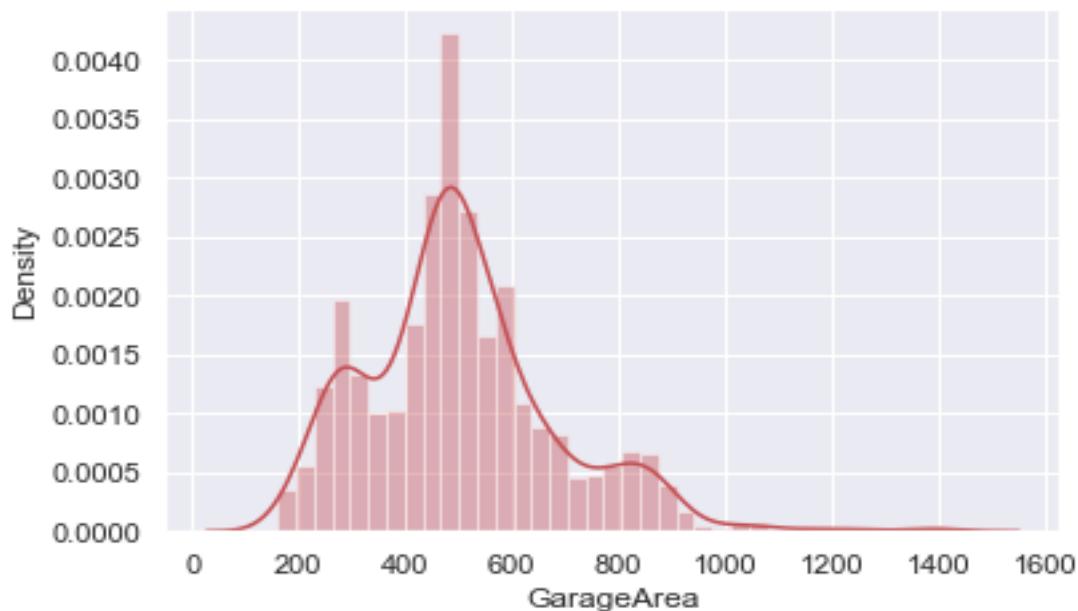
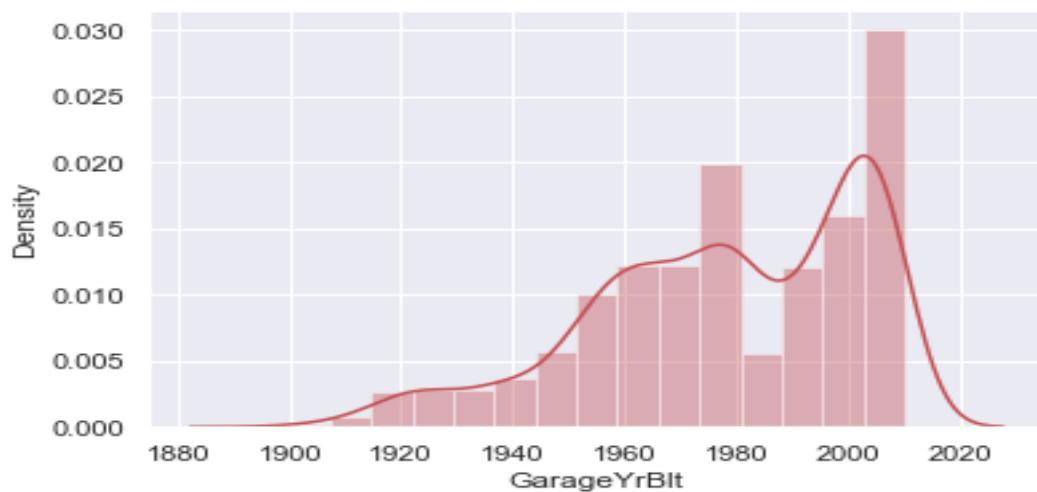
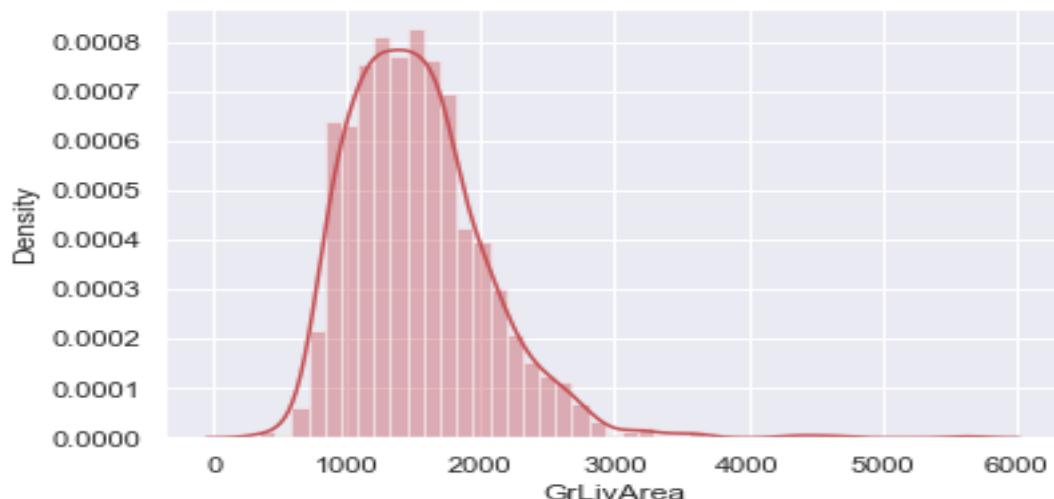
## HOUSE PRICE PREDICTION MODEL



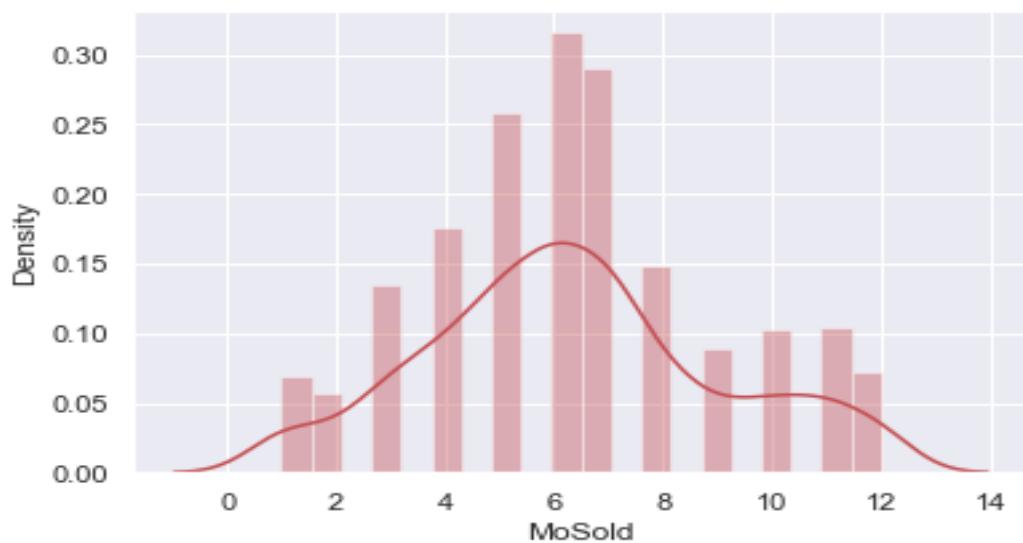
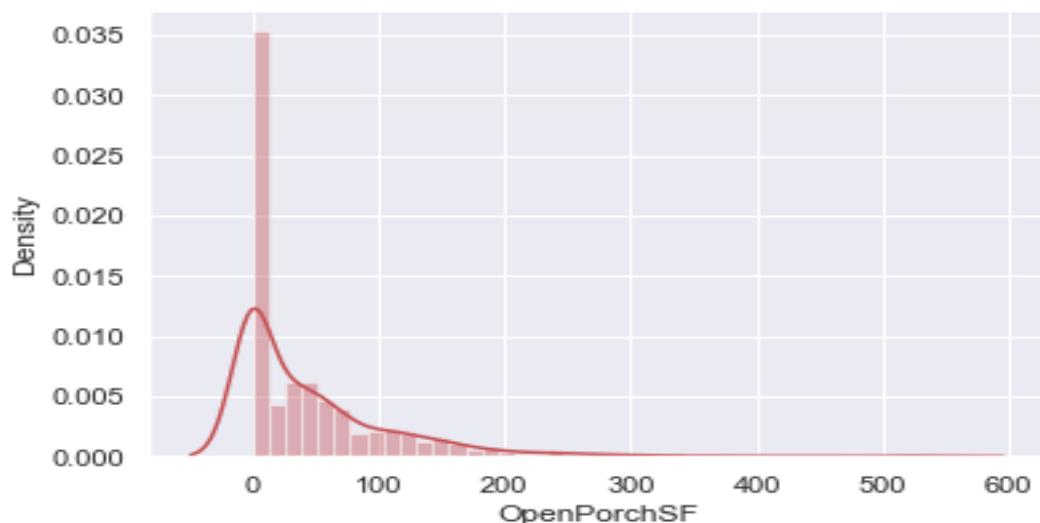
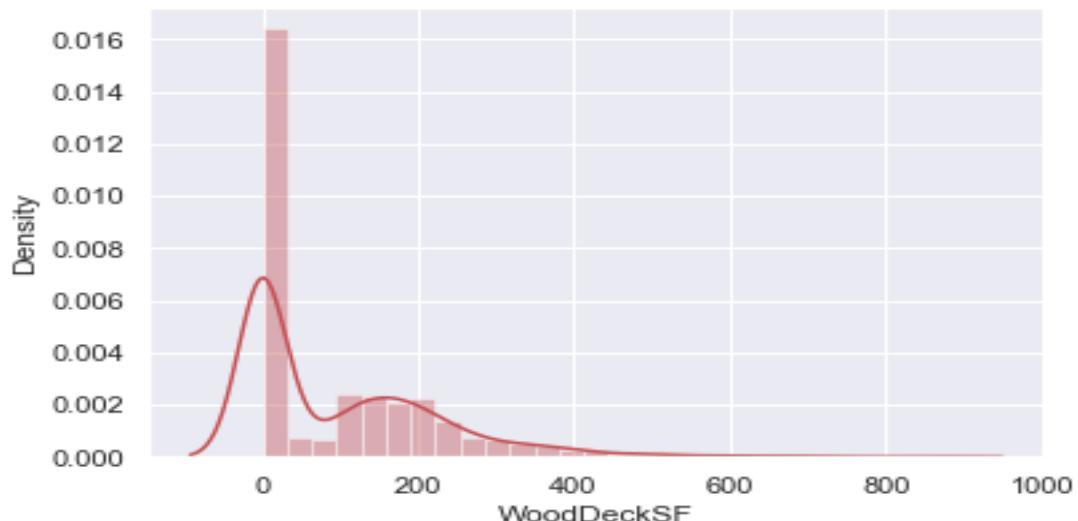
## HOUSE PRICE PREDICTION MODEL



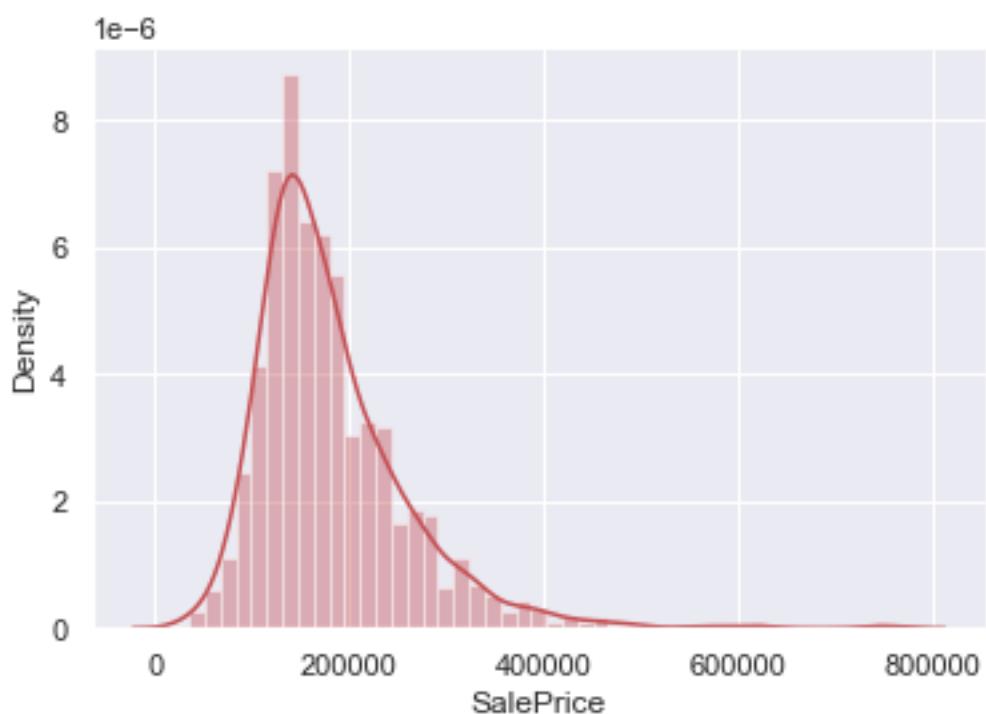
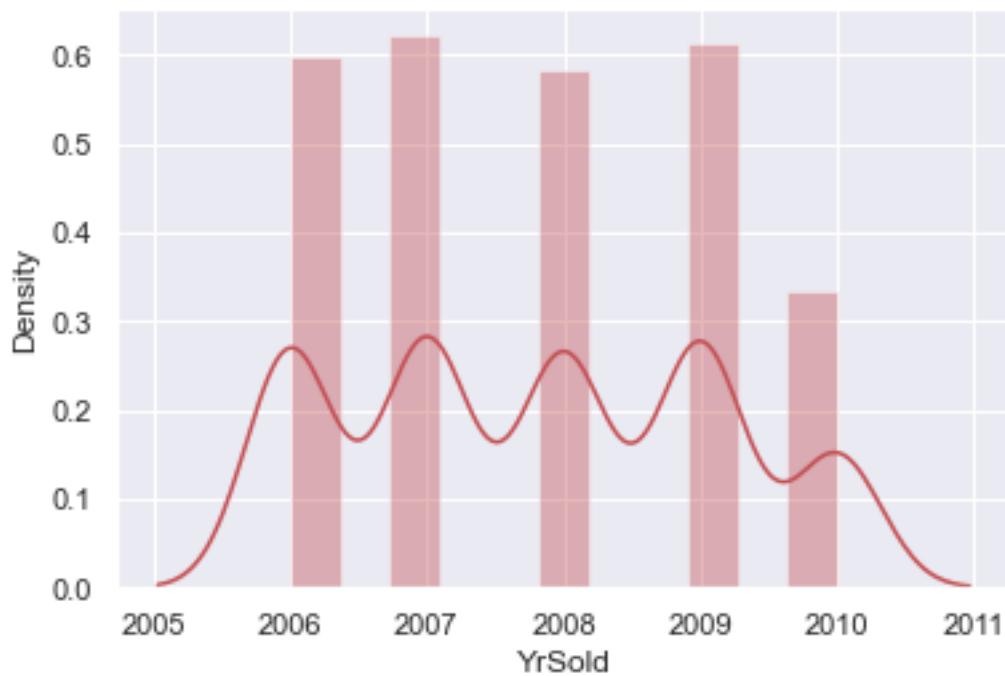
## HOUSE PRICE PREDICTION MODEL



## HOUSE PRICE PREDICTION MODEL



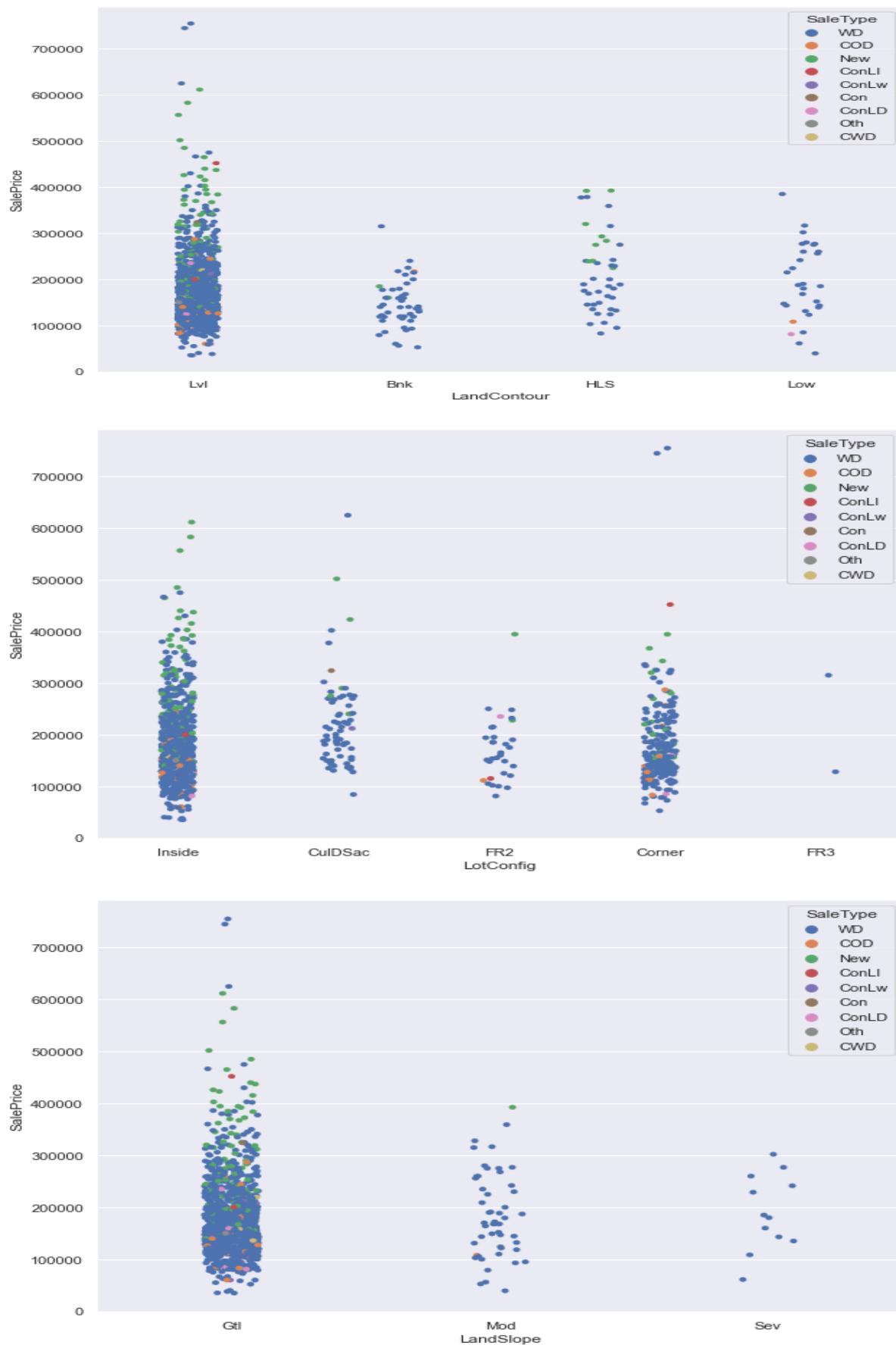
## HOUSE PRICE PREDICTION MODEL



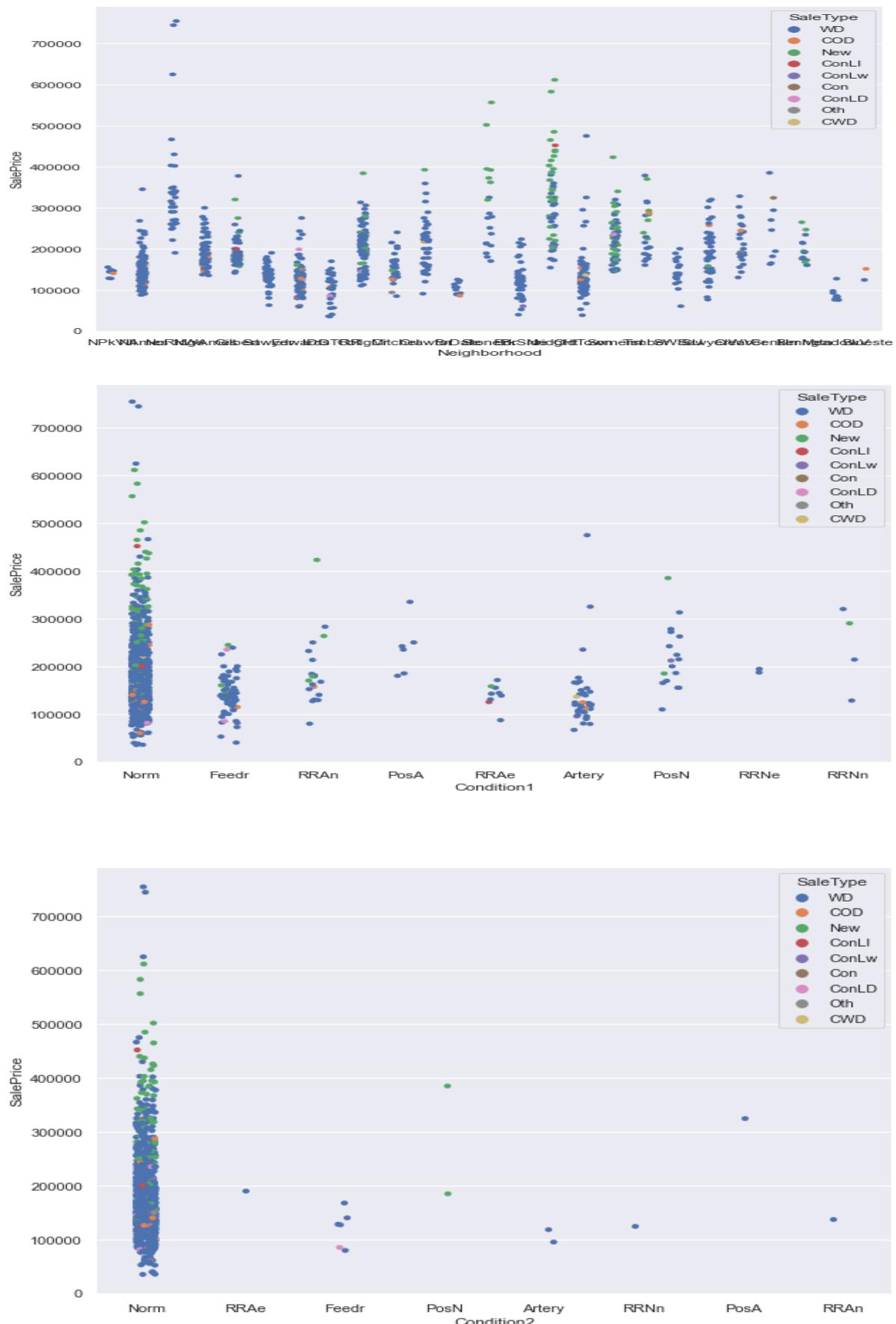
## Bi-variate Analysis:



## HOUSE PRICE PREDICTION MODEL



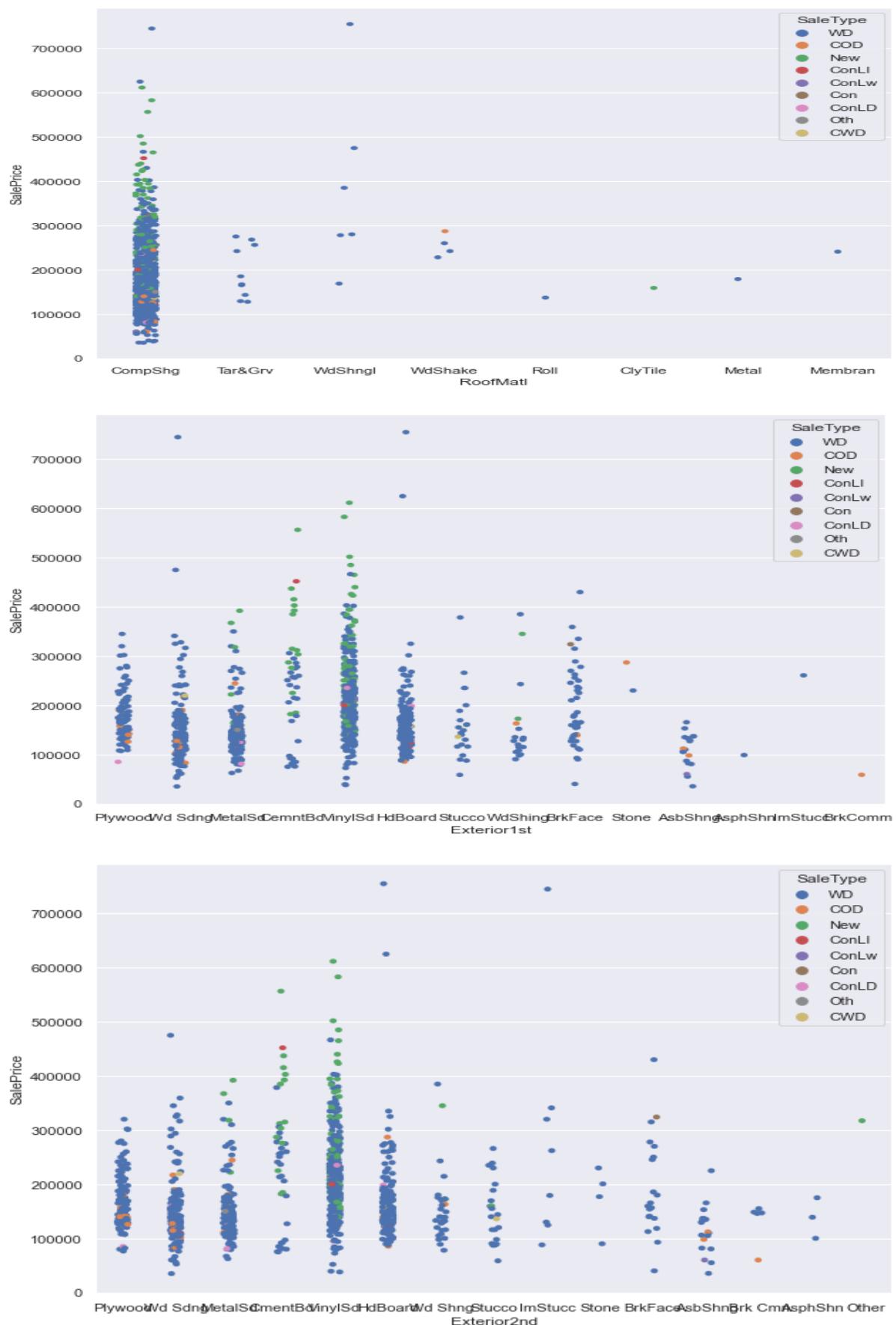
## HOUSE PRICE PREDICTION MODEL



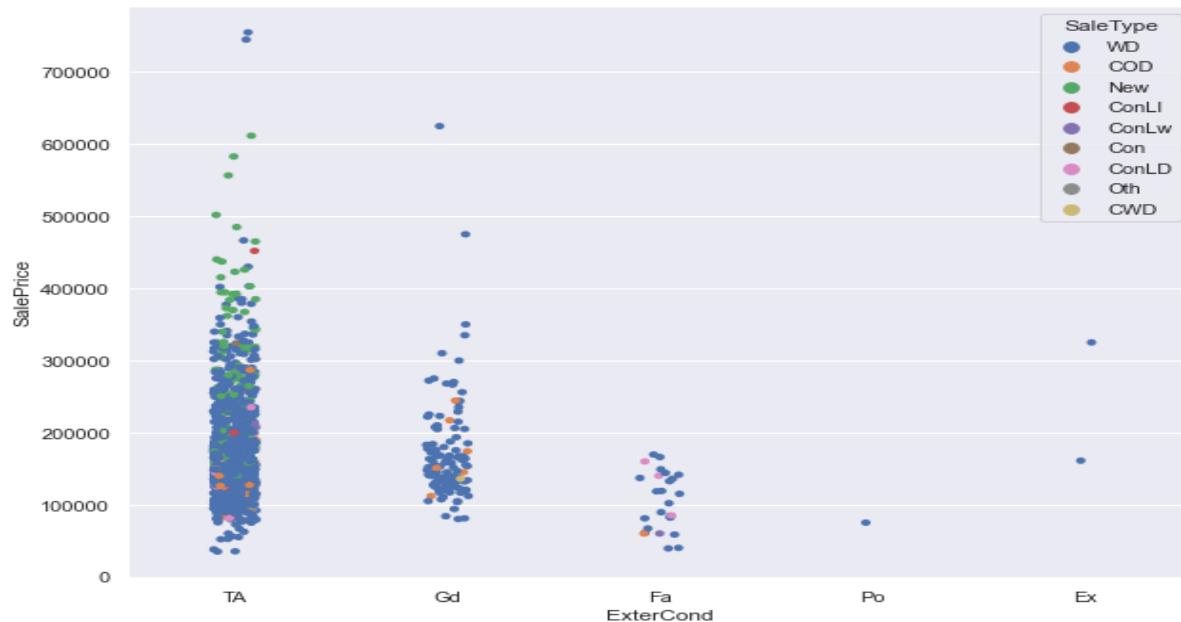
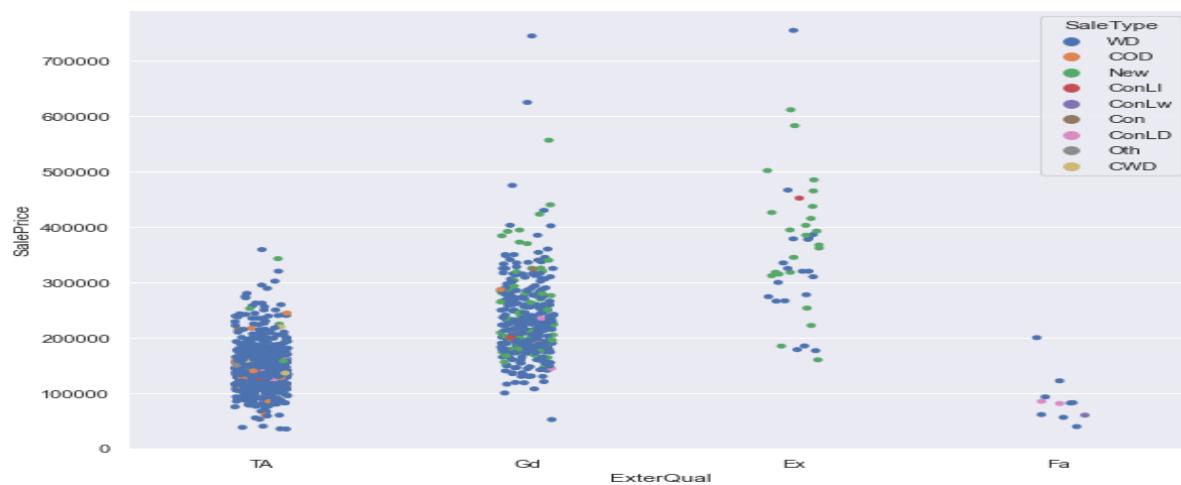
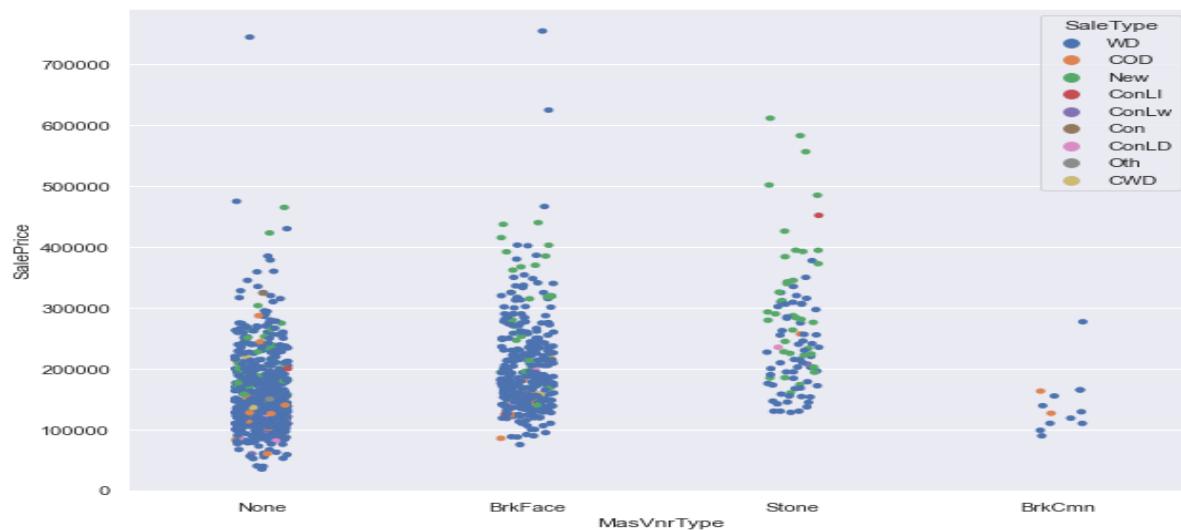
## HOUSE PRICE PREDICTION MODEL



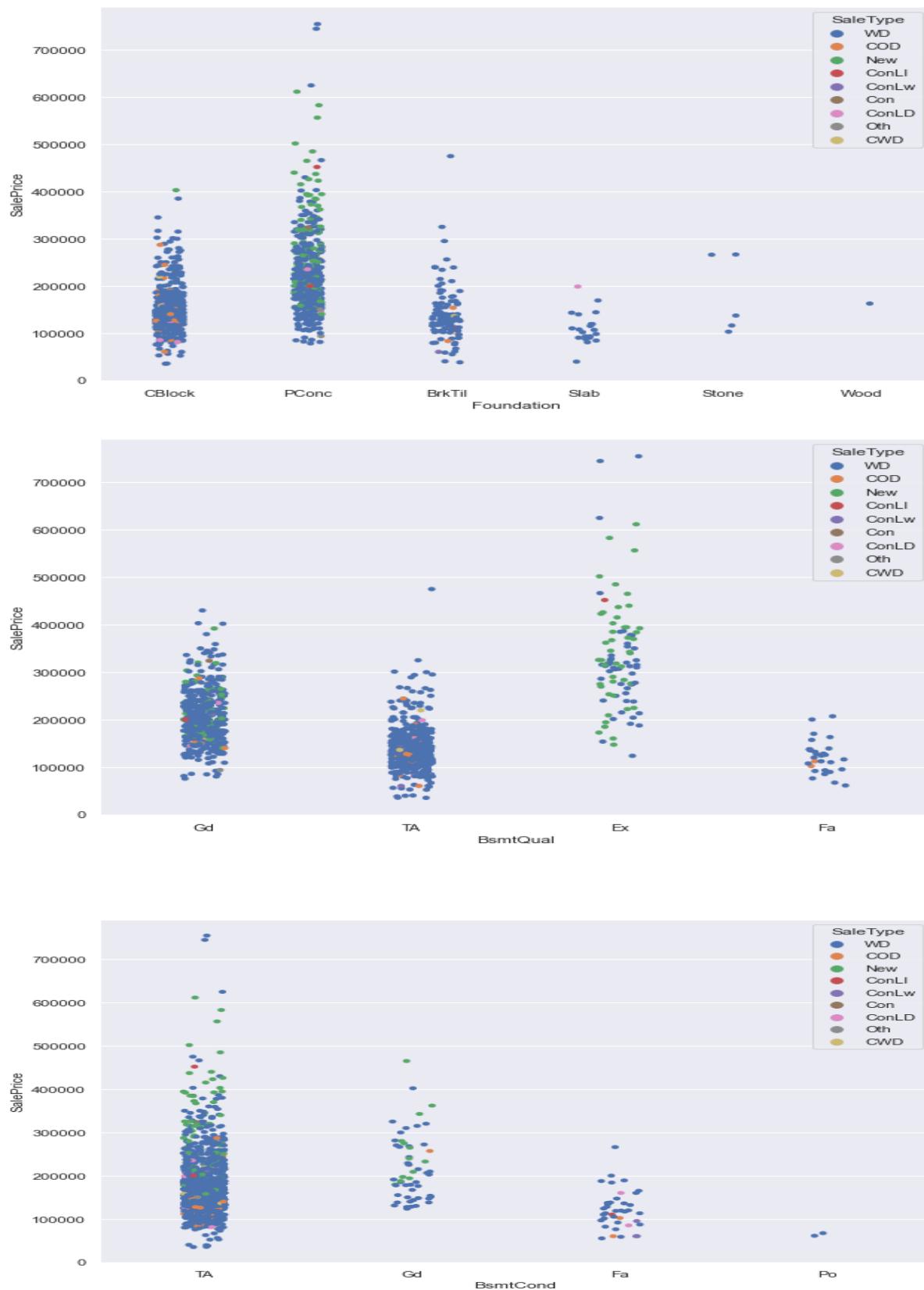
## HOUSE PRICE PREDICTION MODEL



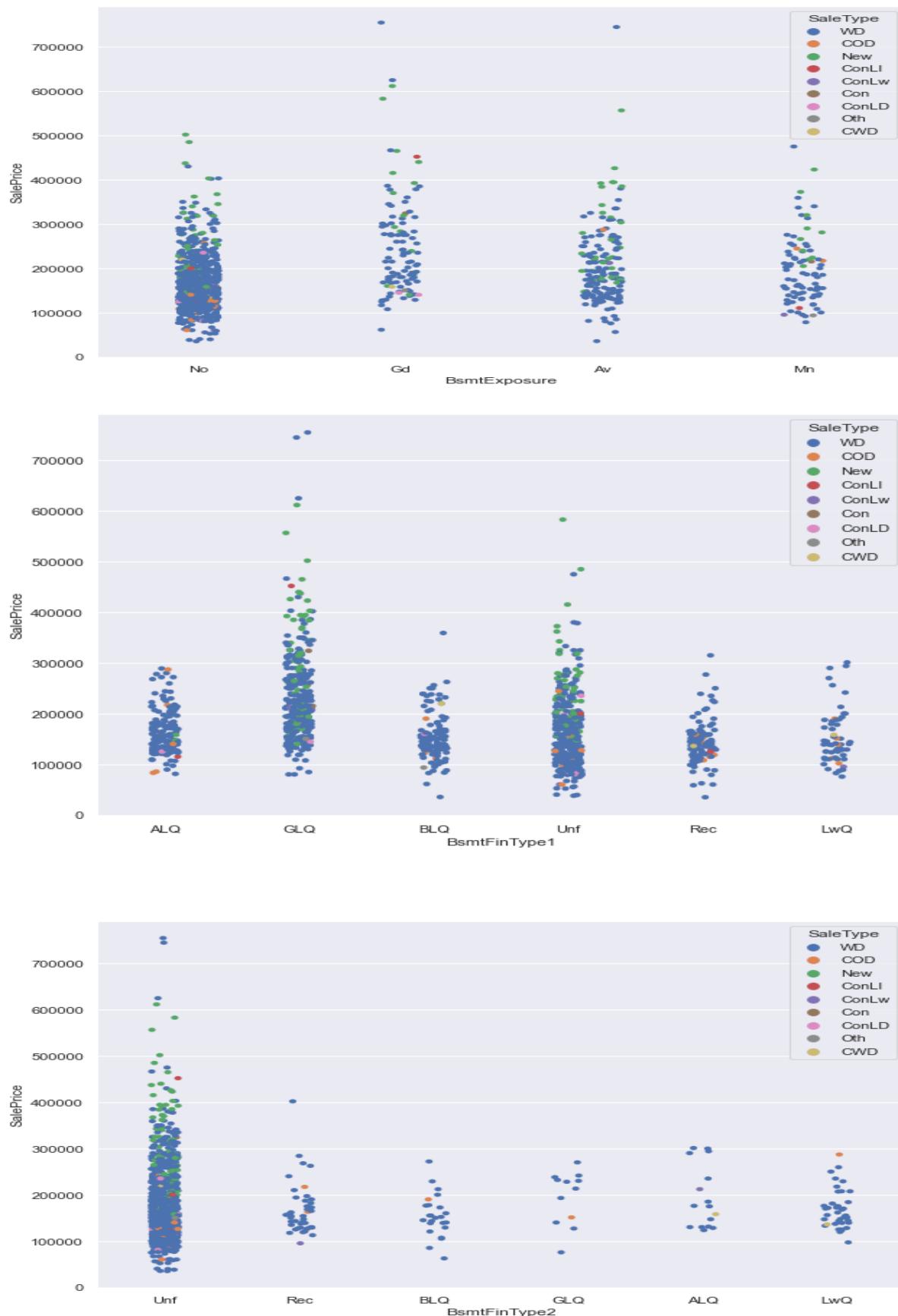
## HOUSE PRICE PREDICTION MODEL



## HOUSE PRICE PREDICTION MODEL



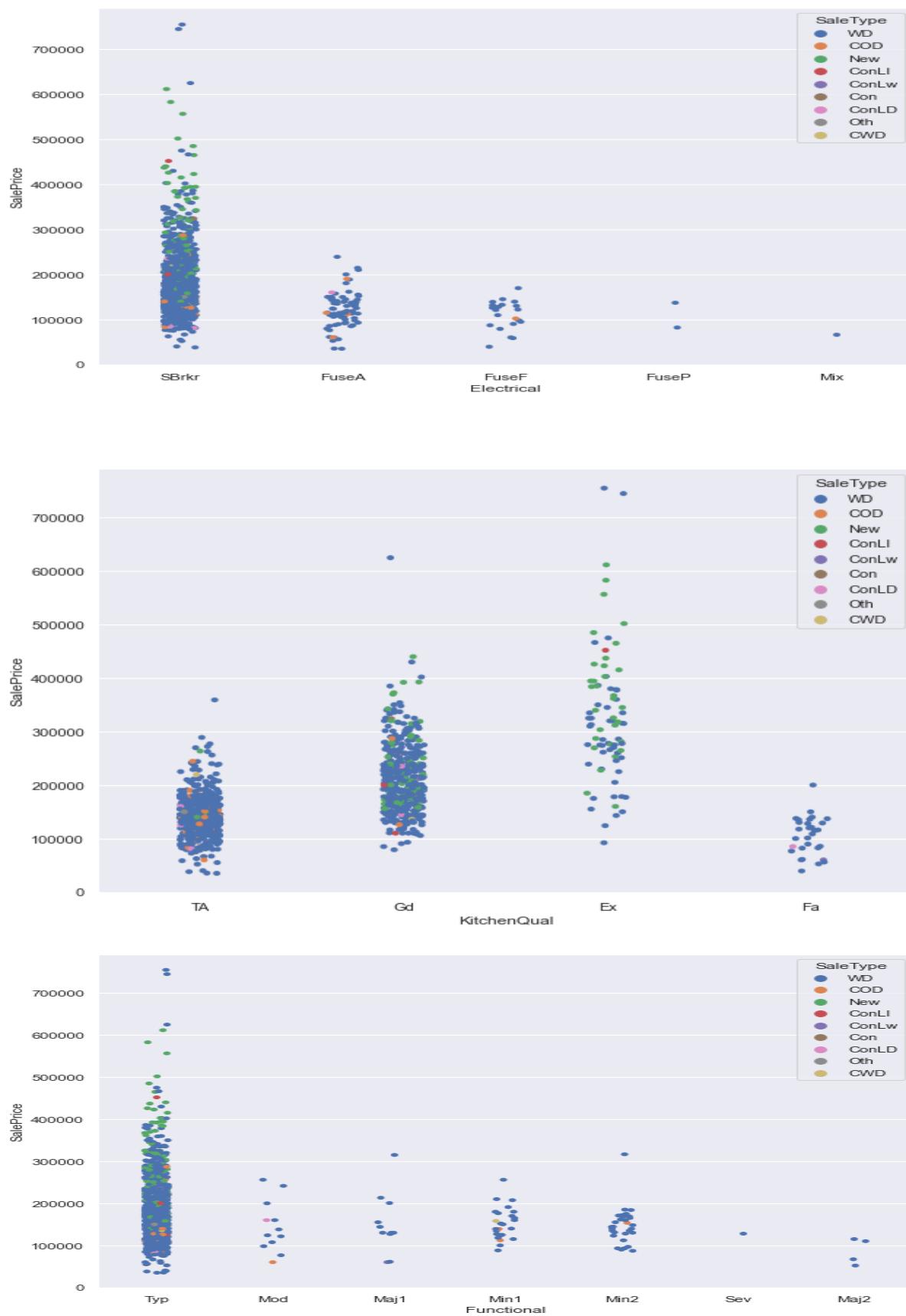
## HOUSE PRICE PREDICTION MODEL



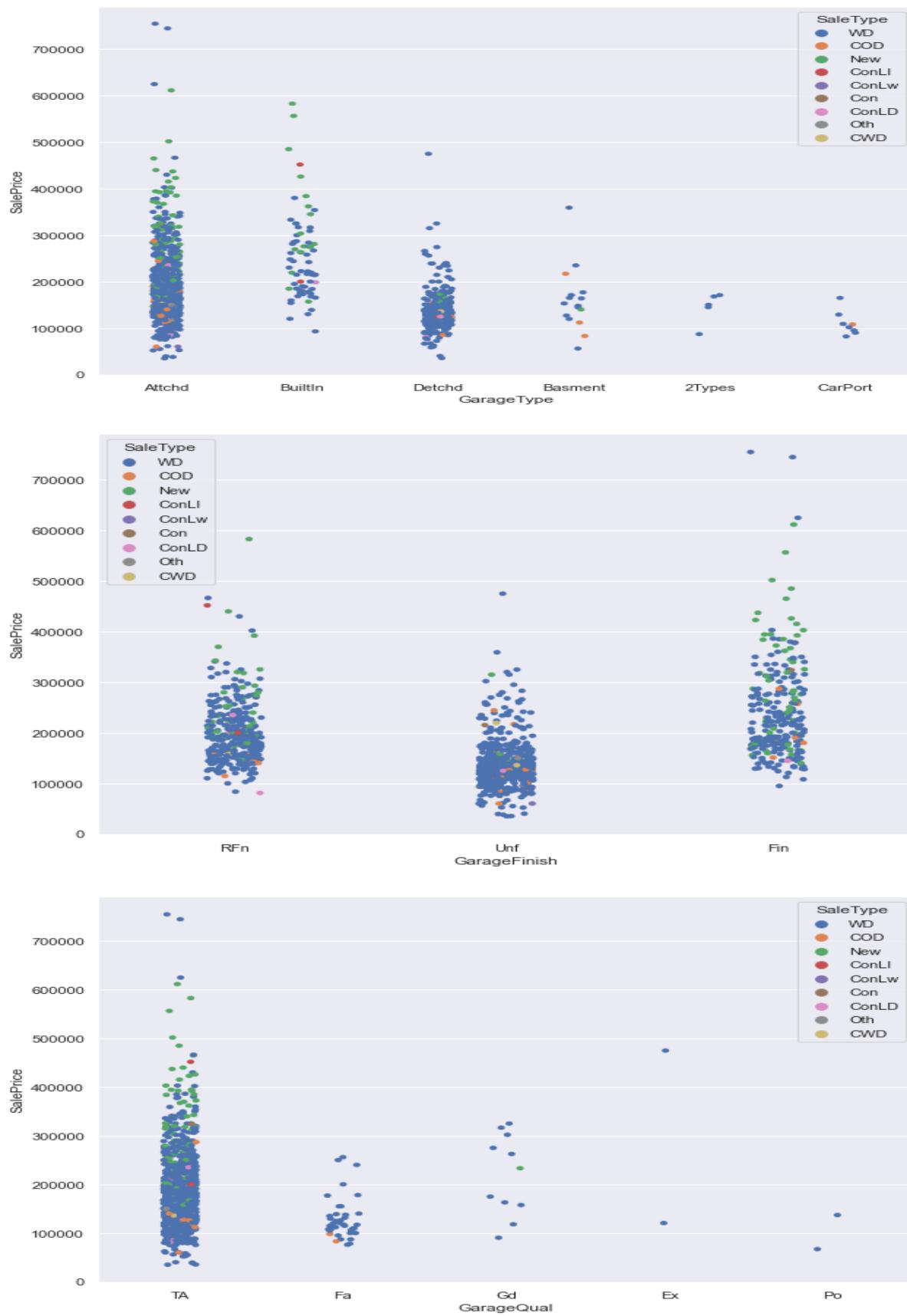
## HOUSE PRICE PREDICTION MODEL



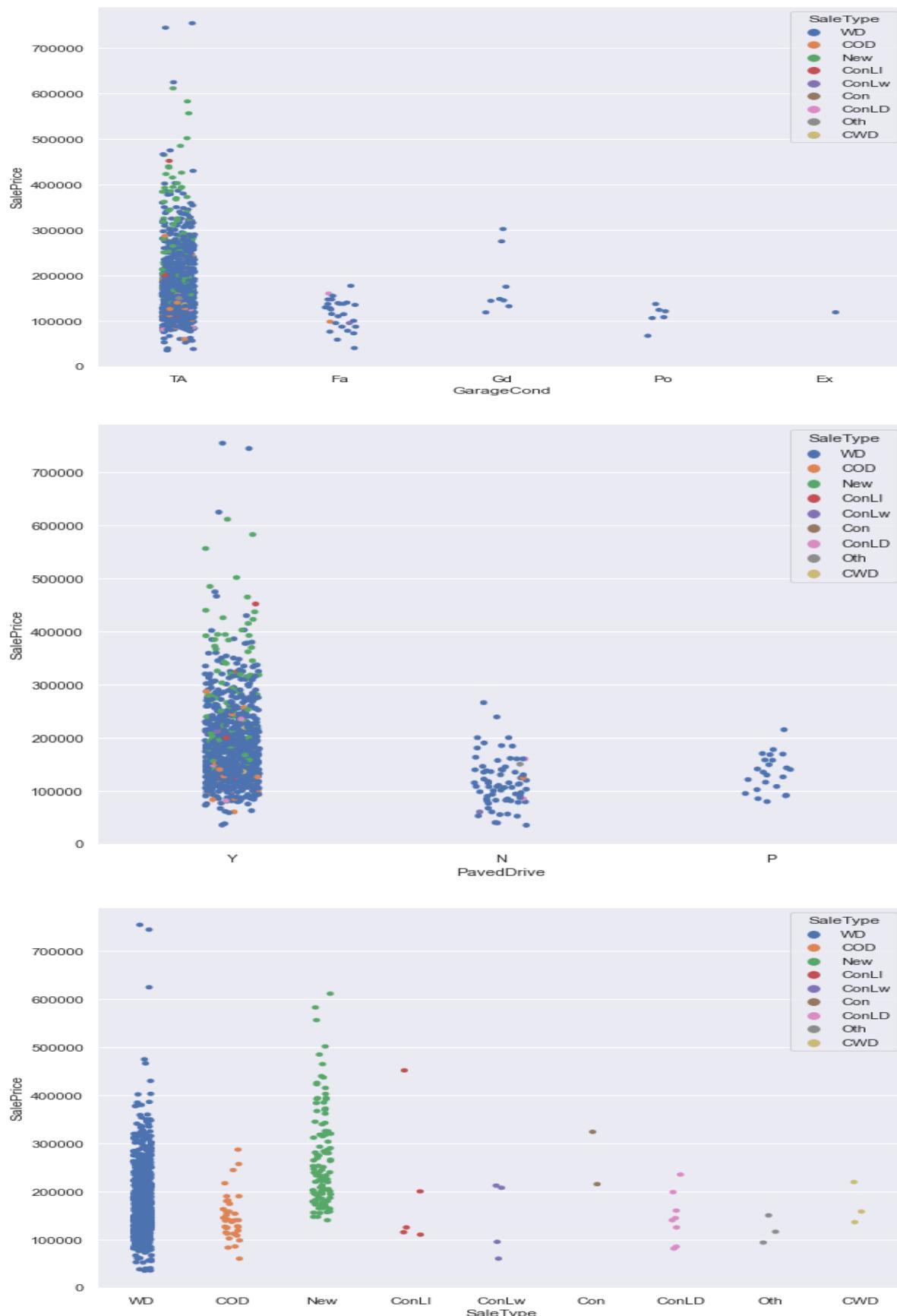
## HOUSE PRICE PREDICTION MODEL



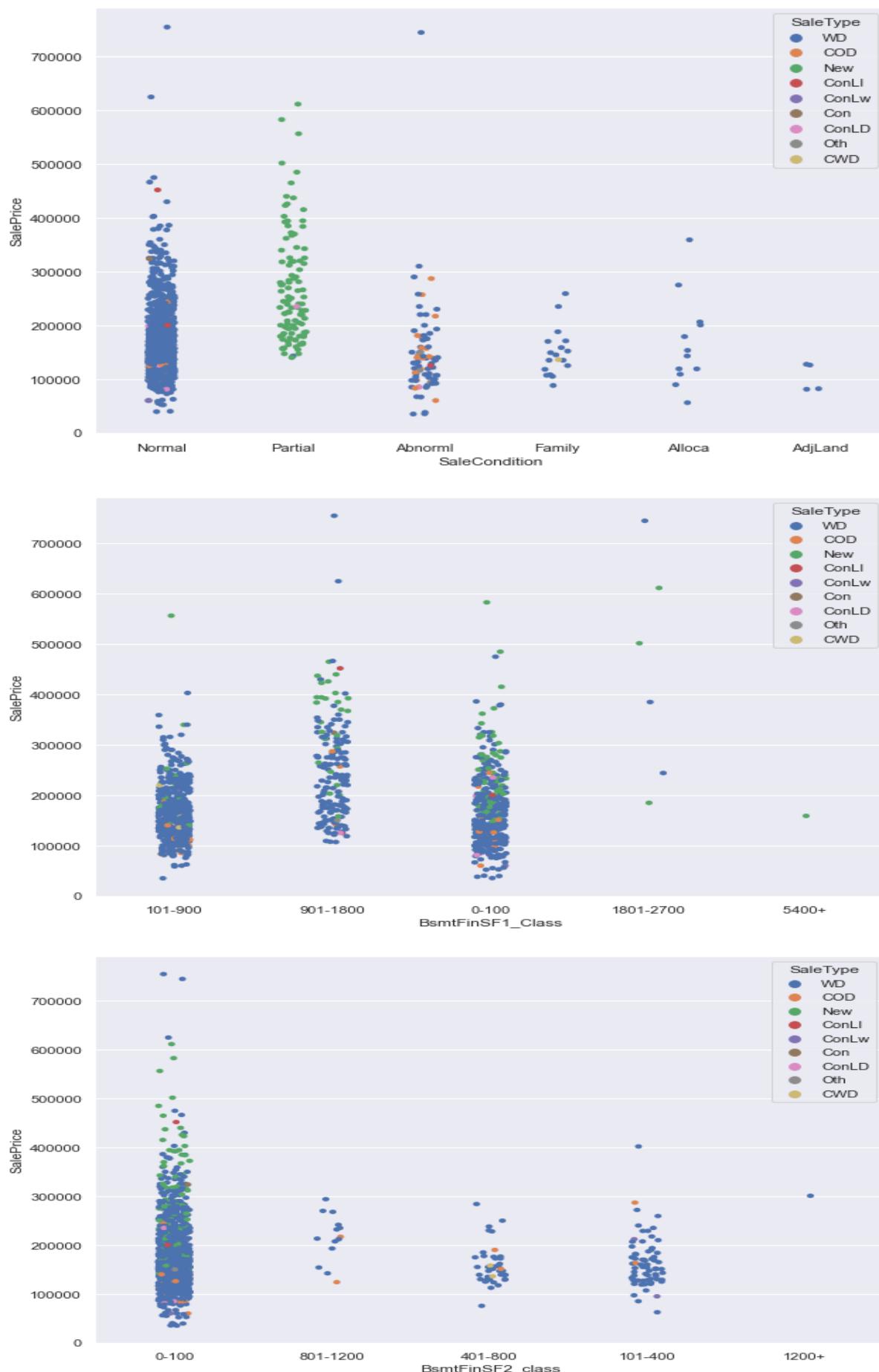
## HOUSE PRICE PREDICTION MODEL



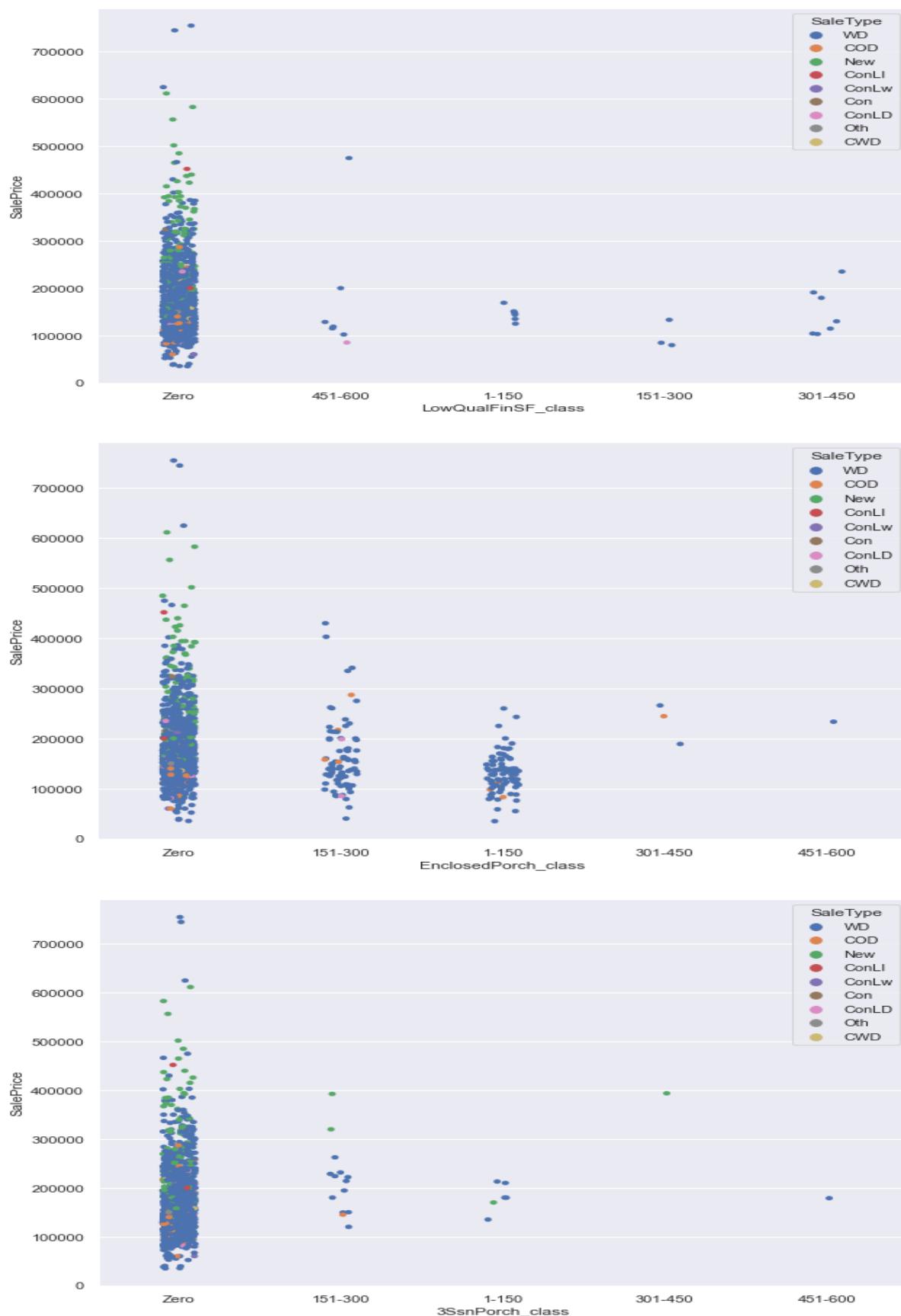
## HOUSE PRICE PREDICTION MODEL



## HOUSE PRICE PREDICTION MODEL



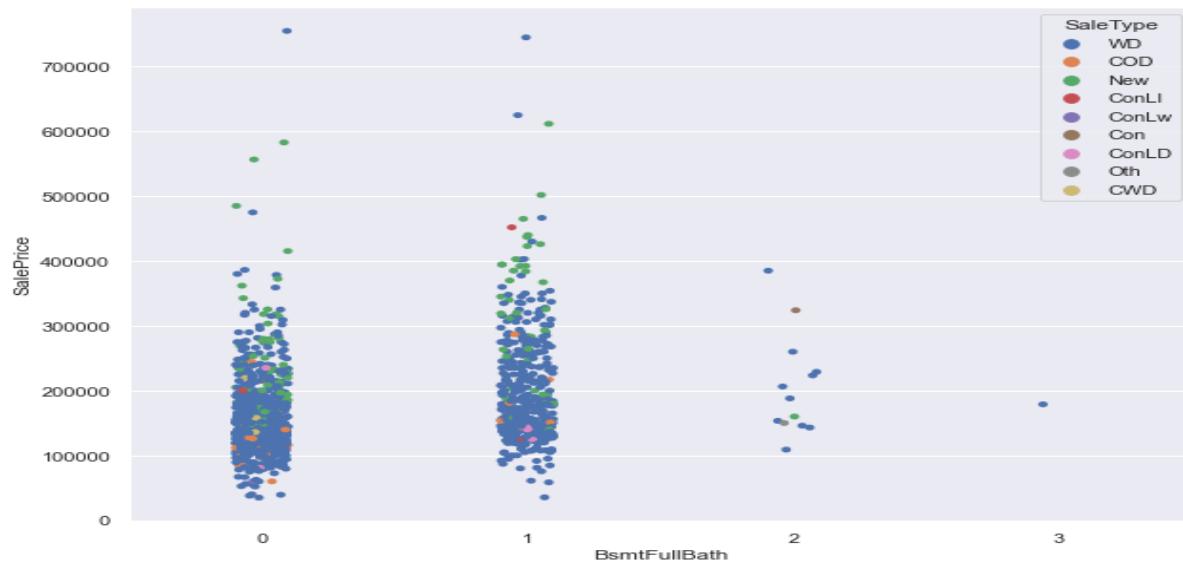
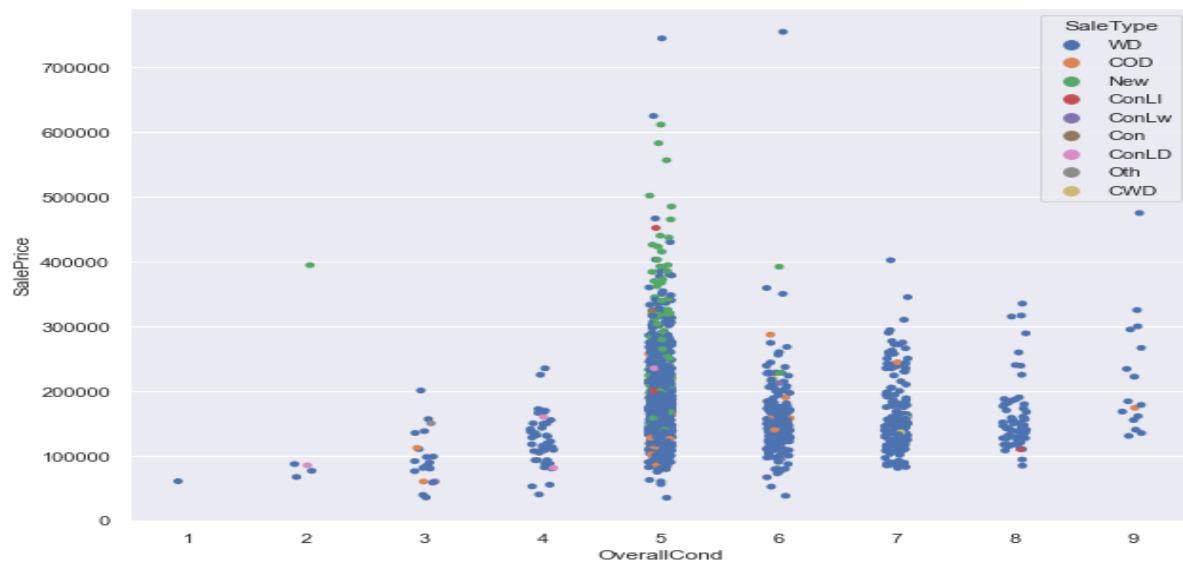
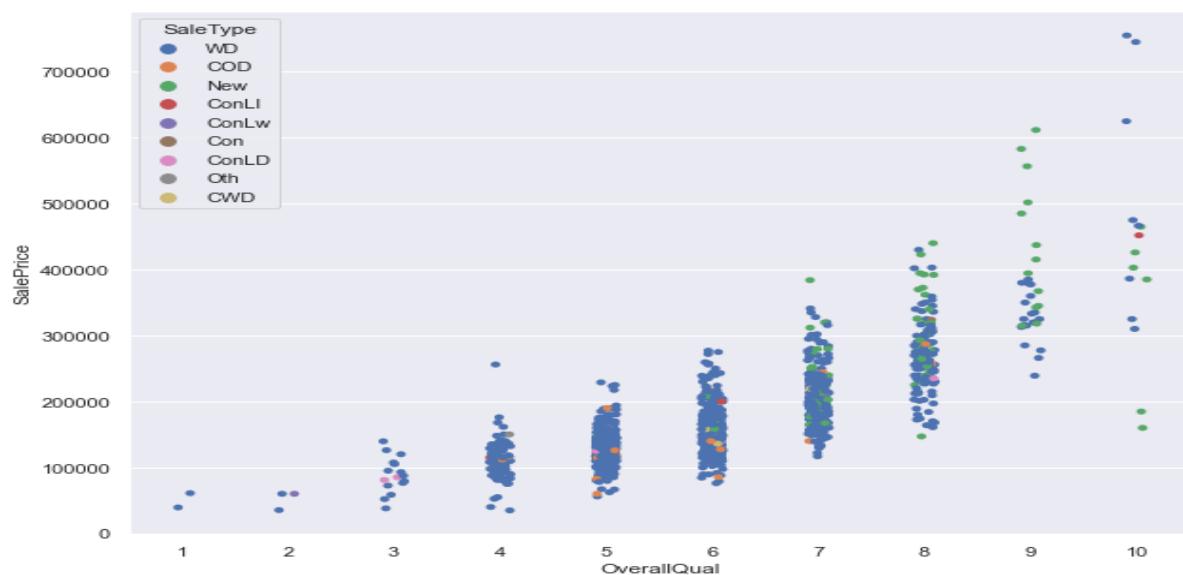
## HOUSE PRICE PREDICTION MODEL



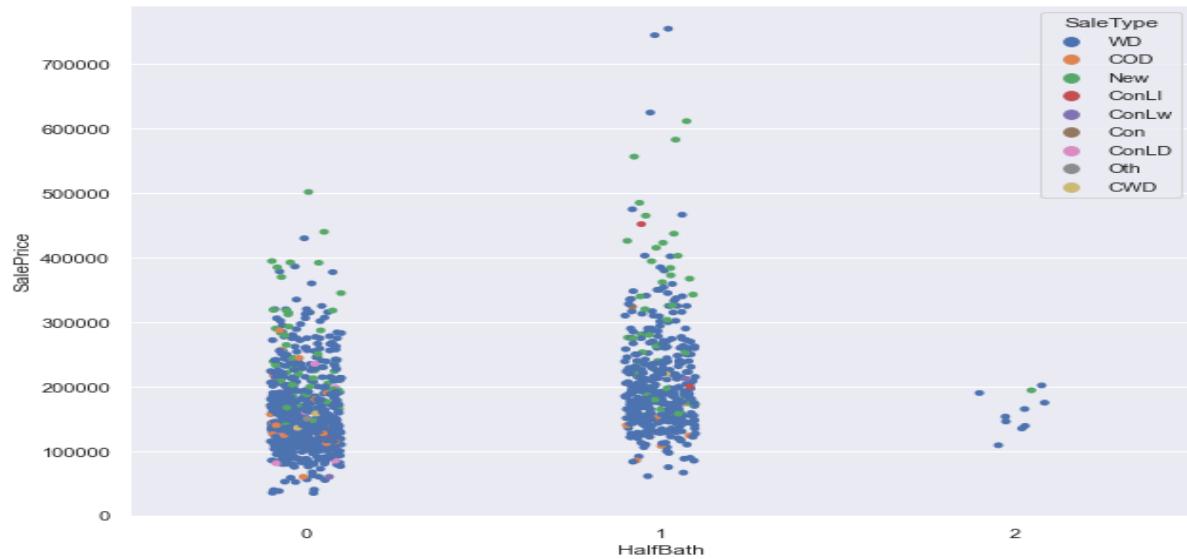
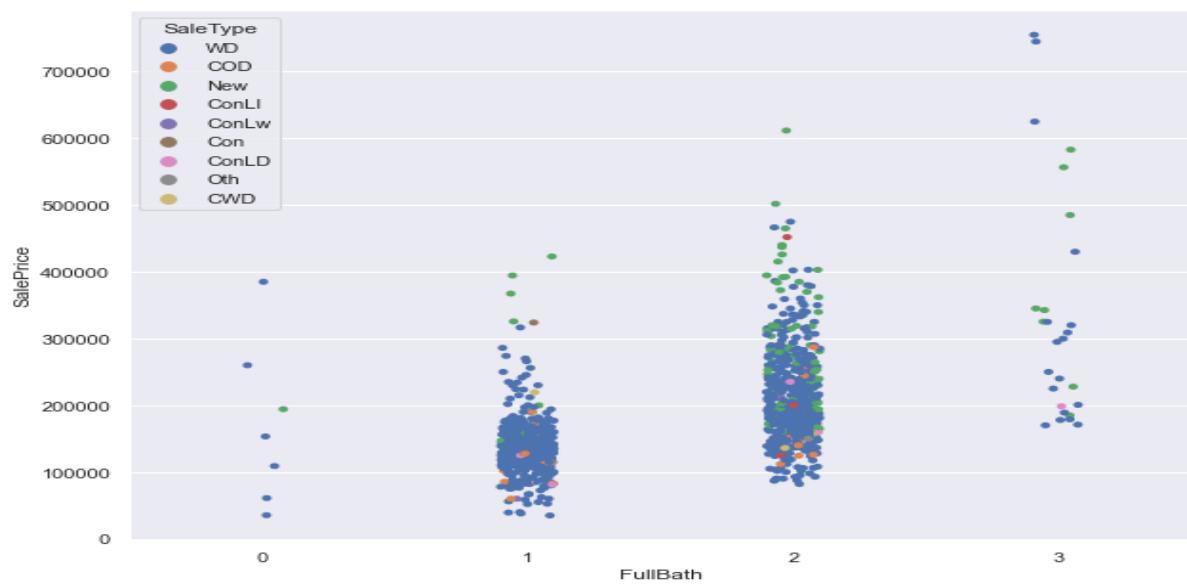
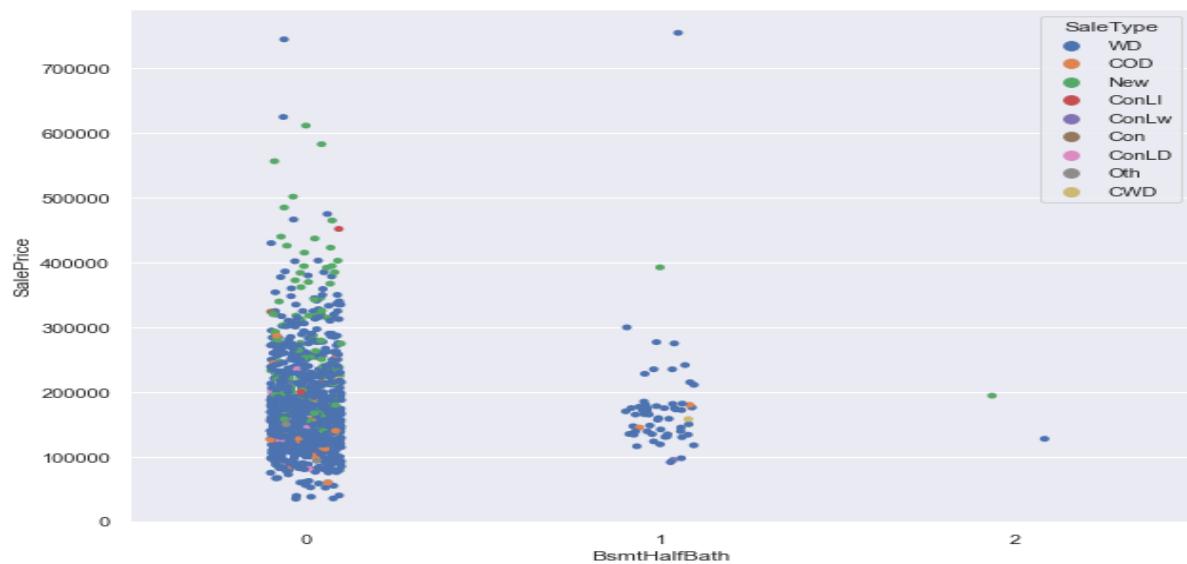
## HOUSE PRICE PREDICTION MODEL



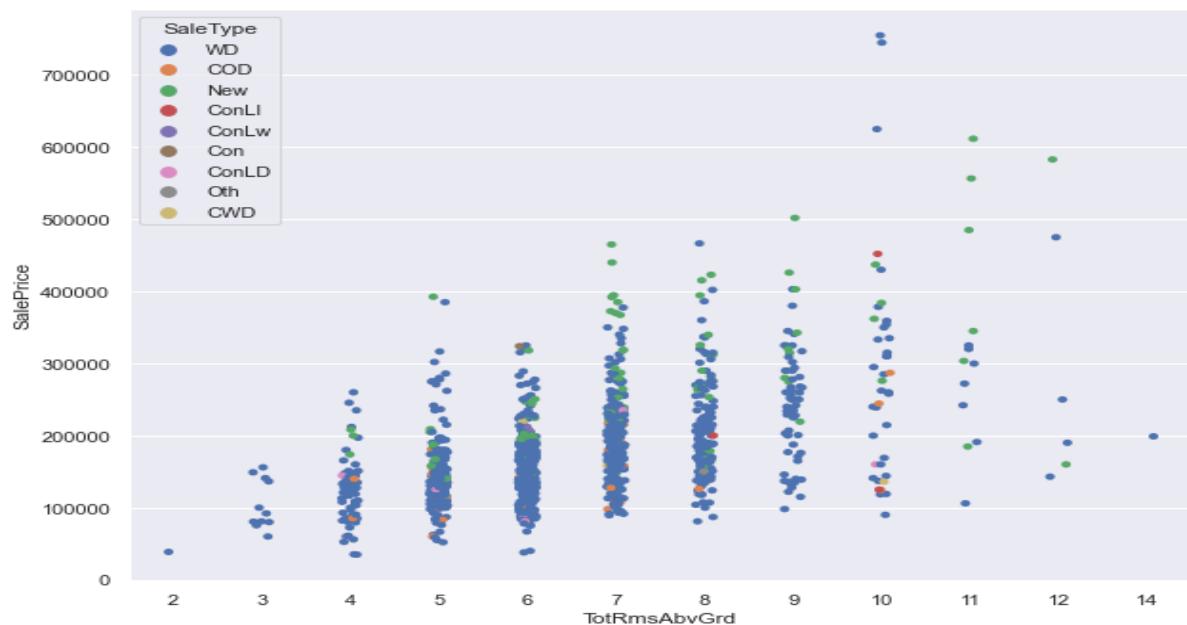
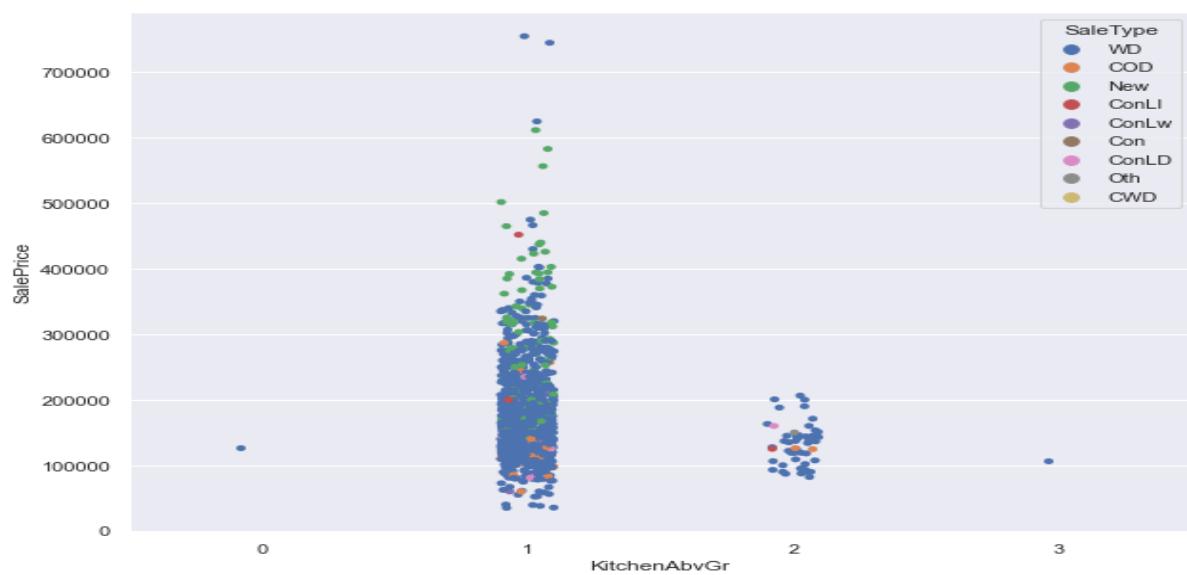
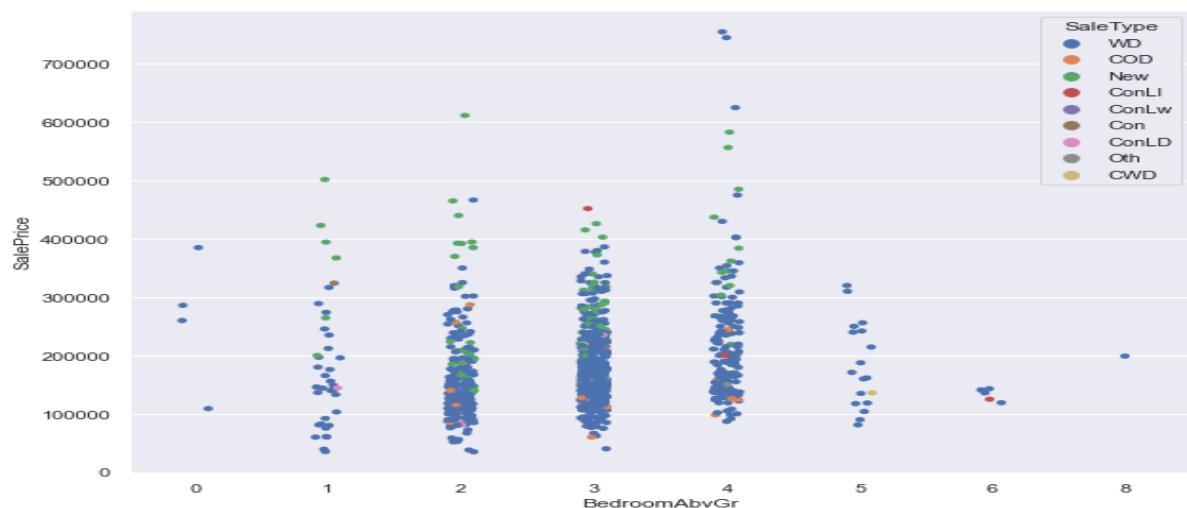
## HOUSE PRICE PREDICTION MODEL



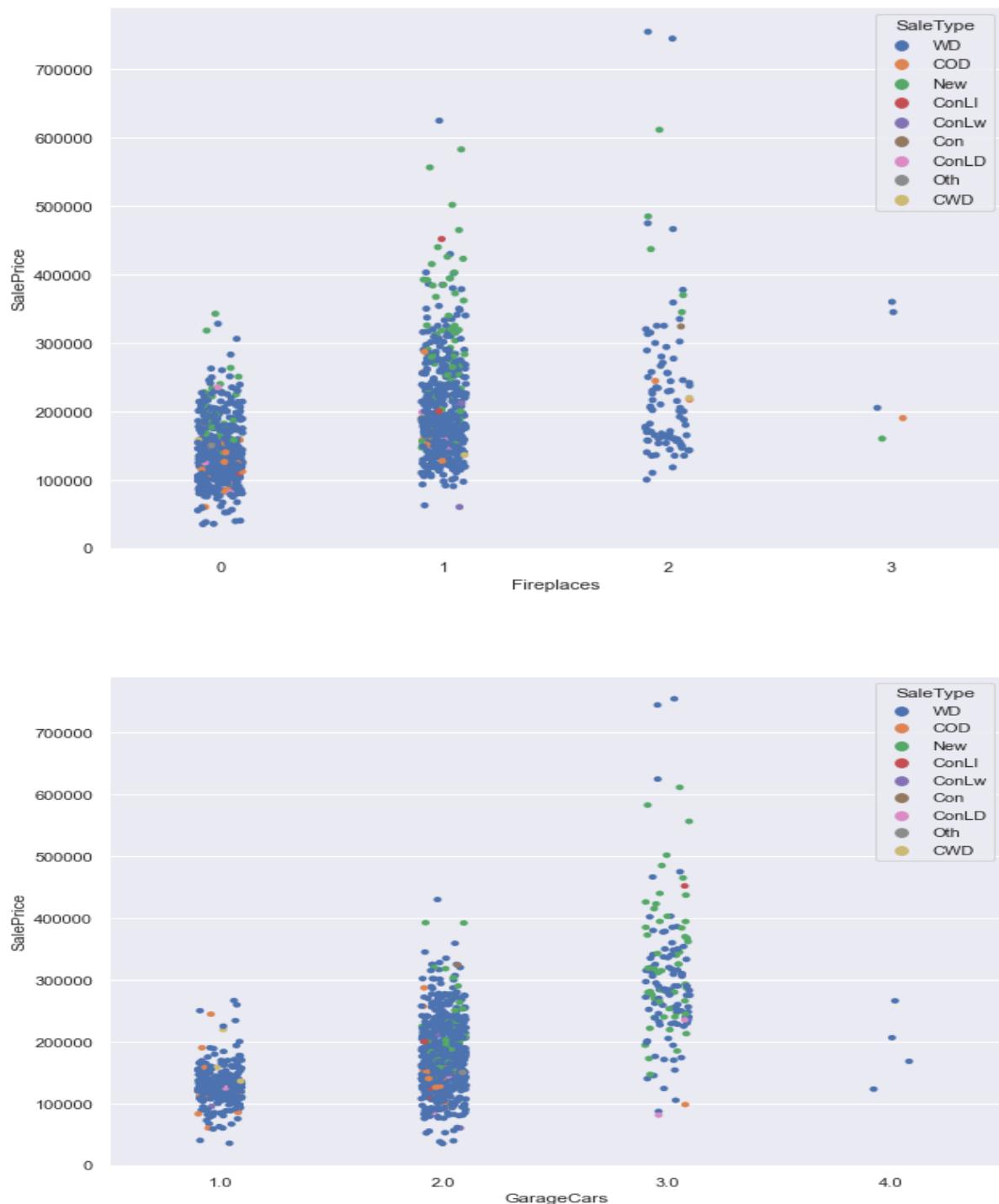
## HOUSE PRICE PREDICTION MODEL



## HOUSE PRICE PREDICTION MODEL



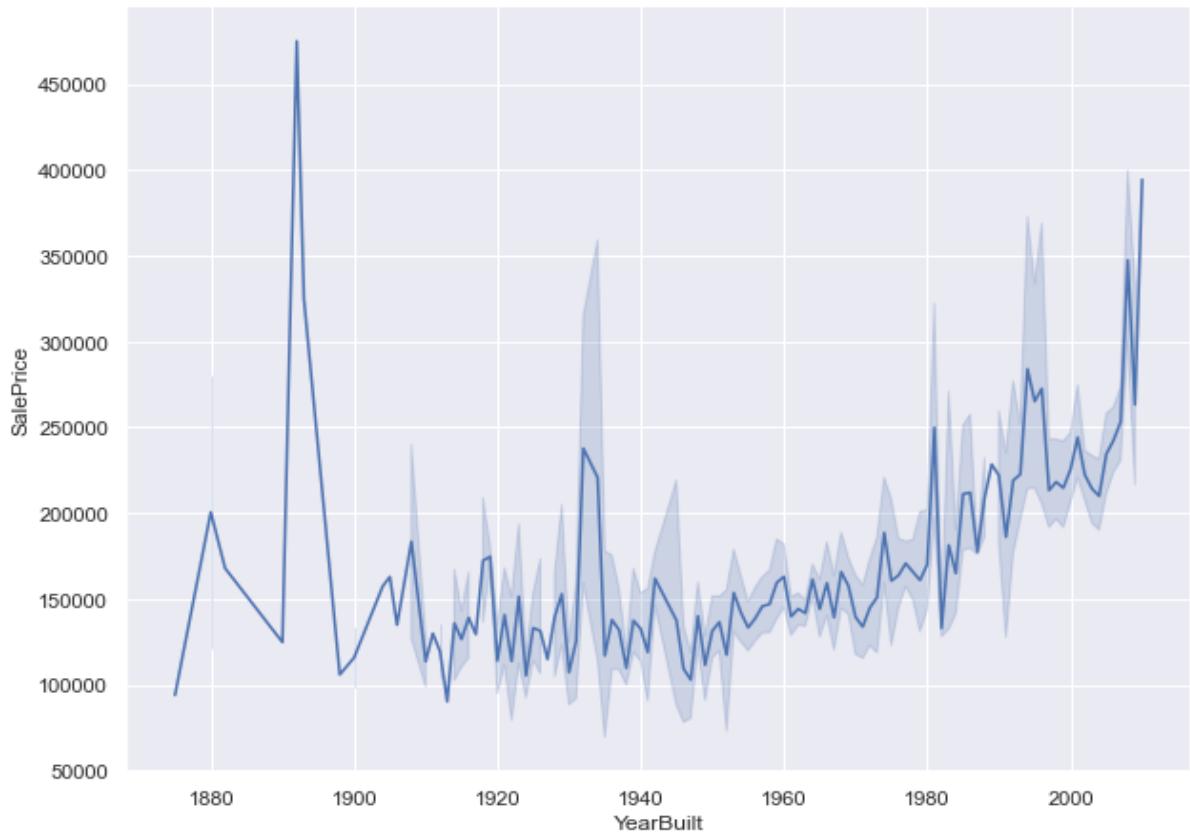
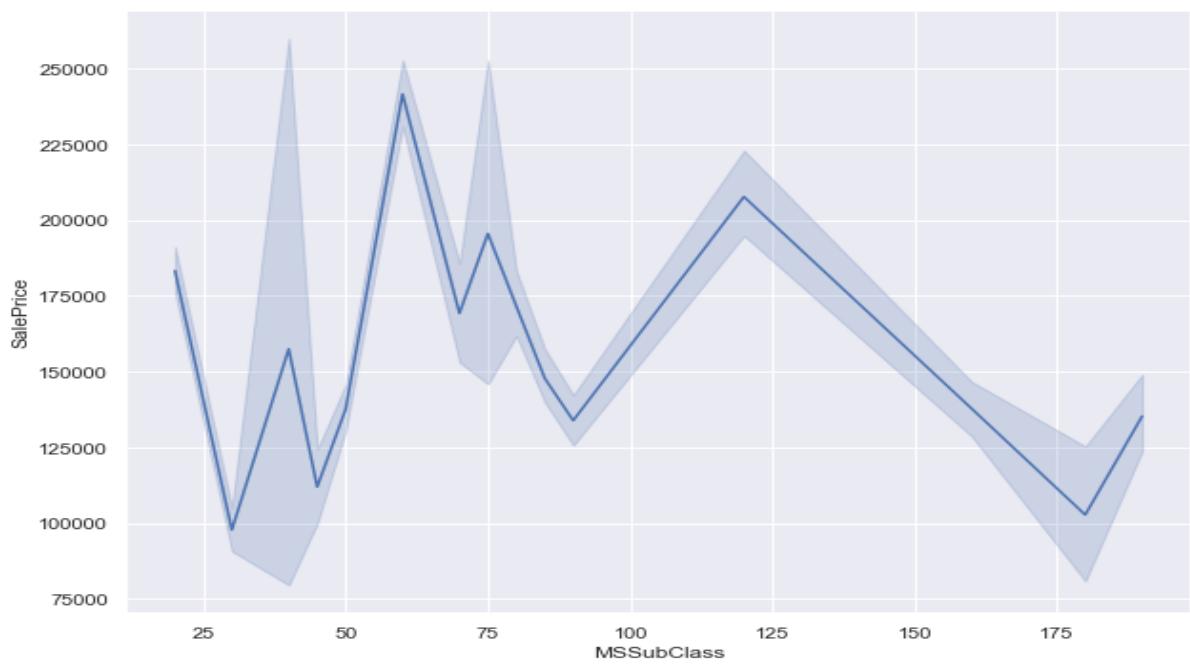
## HOUSE PRICE PREDICTION MODEL



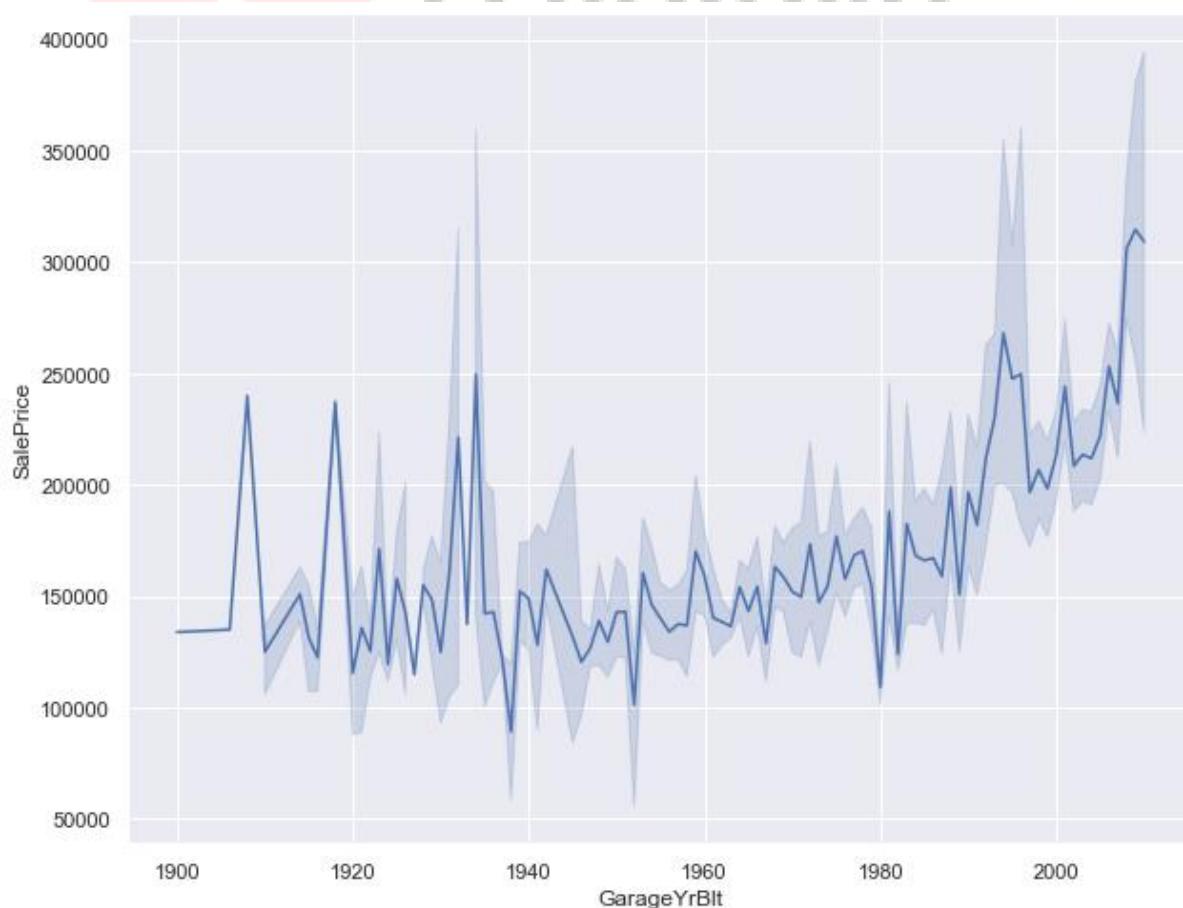
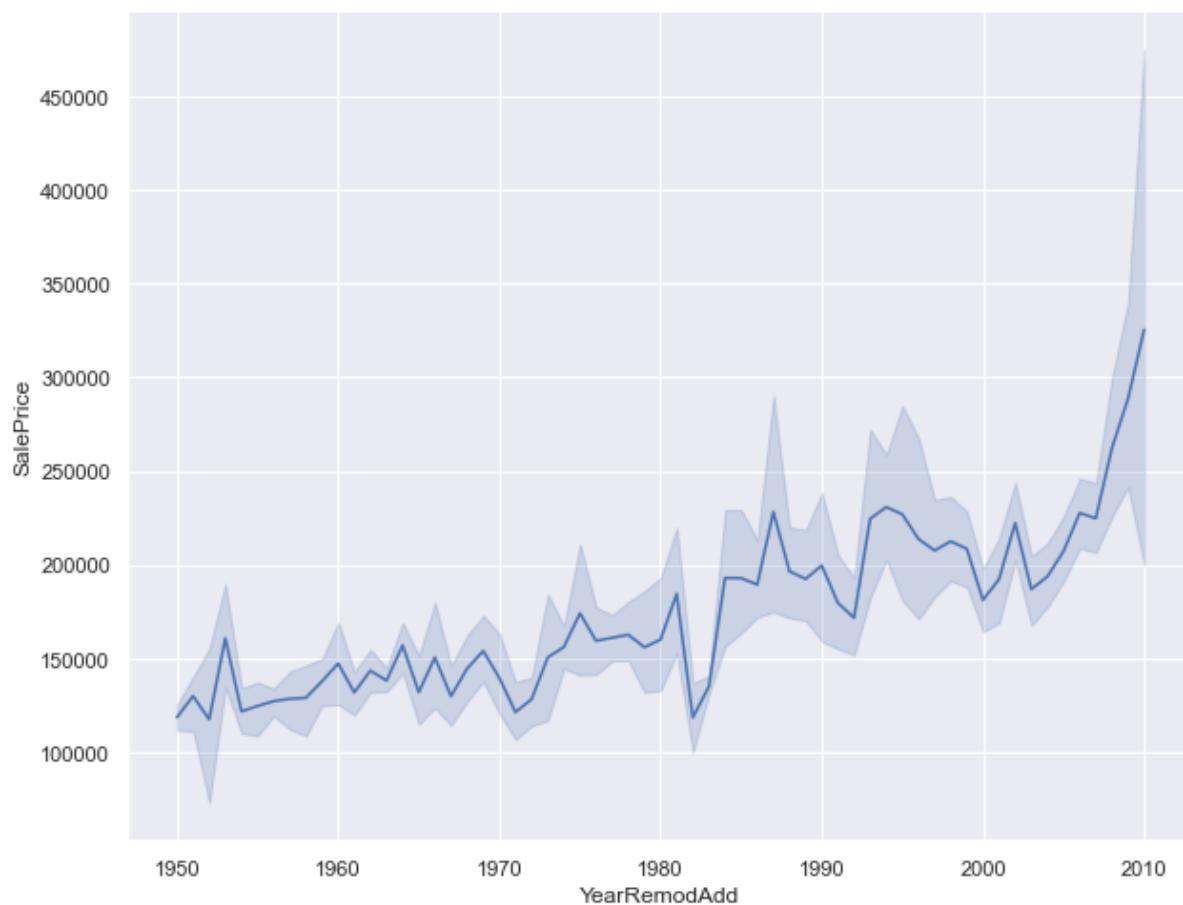
- Keys: Scatter plot is shown with sale price of property, by which can estimate about how column is impacting sale price

## HOUSE PRICE PREDICTION MODEL

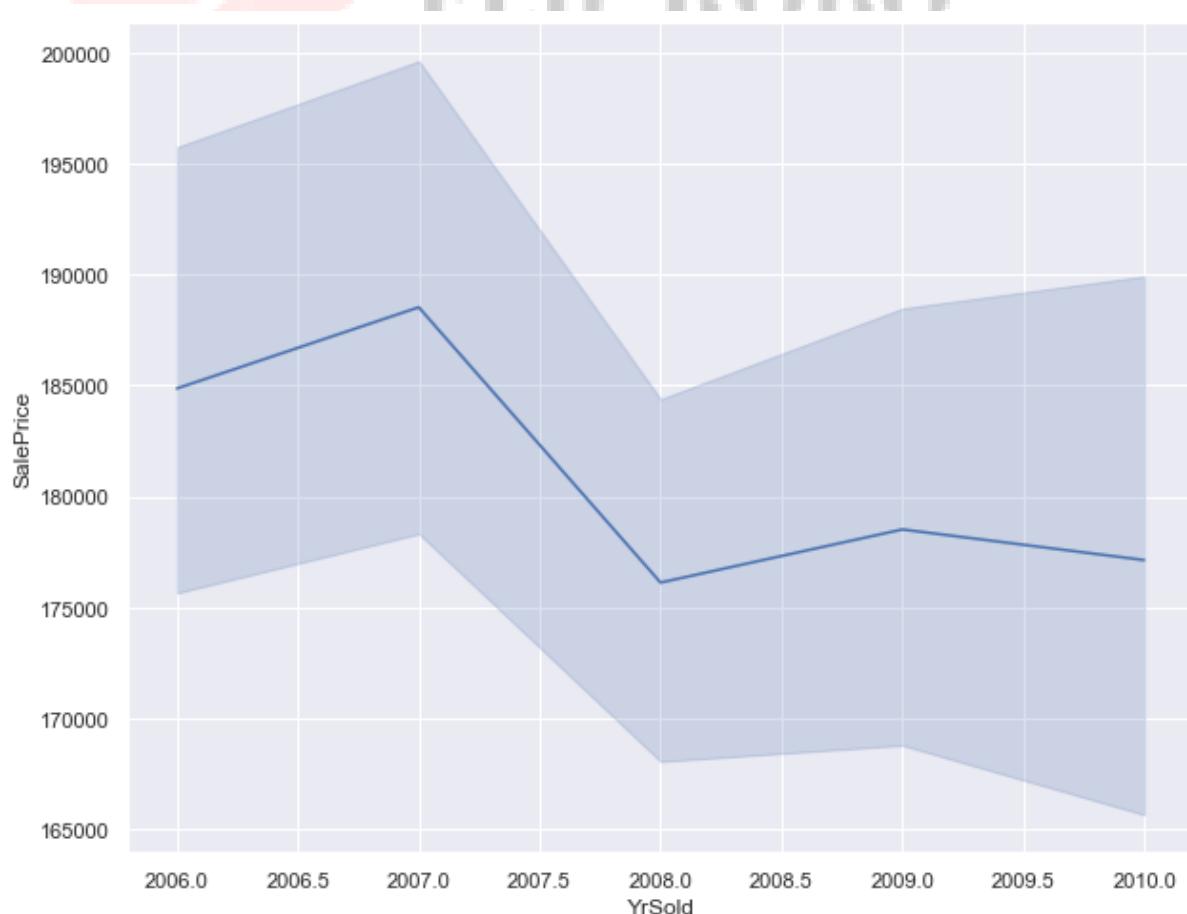
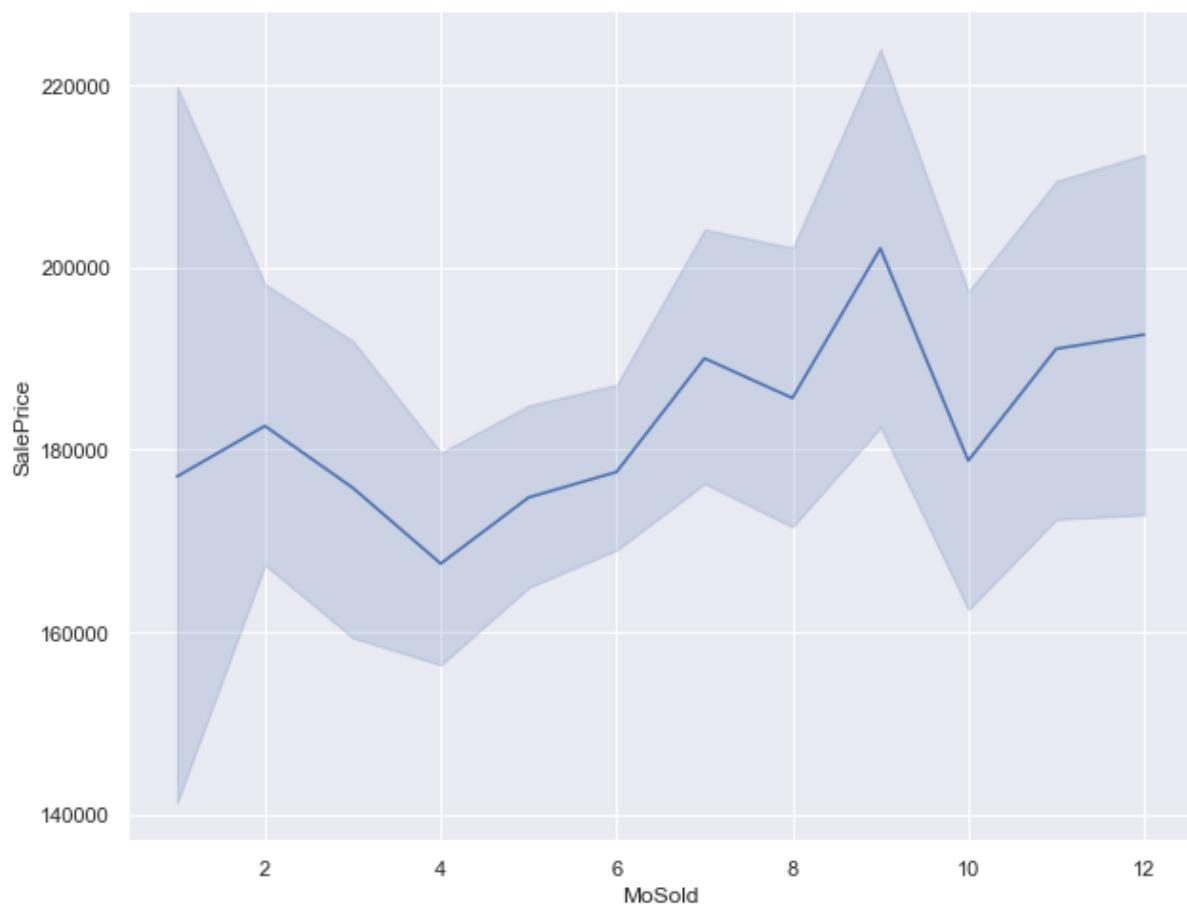
➤ **Line Plot:**



## HOUSE PRICE PREDICTION MODEL



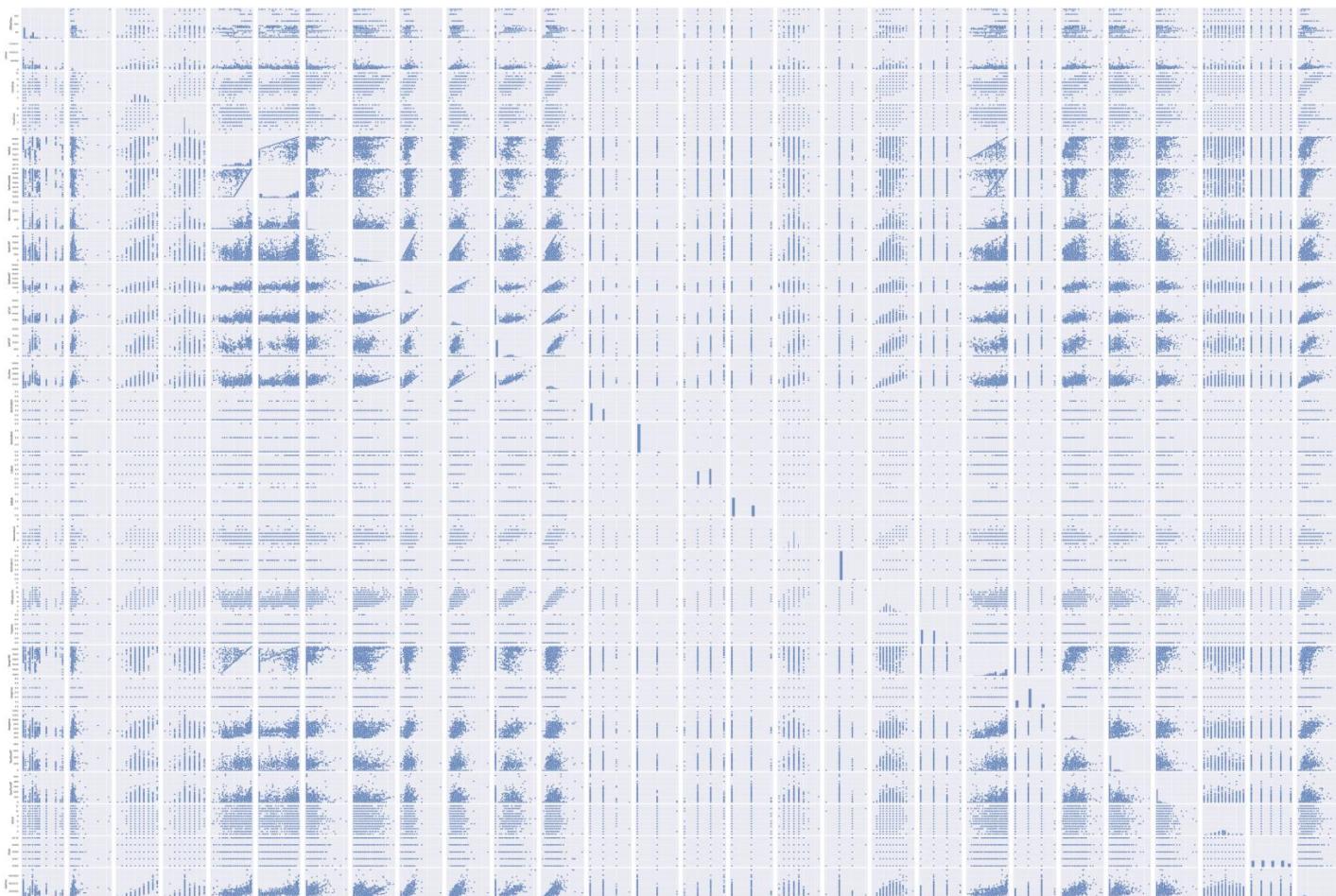
## HOUSE PRICE PREDICTION MODEL



**Key Points:**

- Scatter plot is shown with sale price of property, by which can estimate about how column is impacting sale price
- Increase in yearRemodAdd and GarageYrbit, both are leading to increase sale price.
- Sale price is high as older "Yrsold" is.
- One can see the impact of these columns on saleprice

## PAIR PLOT OF DATASET



➤ **Encoding of dataset:**

```
# Copy of train_df and test_df
encoded_train = train_df.copy()
encoded_test = test_df.copy()

from sklearn.preprocessing import OrdinalEncoder
oe = OrdinalEncoder() # created instance of Ordinal Encoder
```

```
object_type = list(encoded_train.select_dtypes(include = 'object').keys())
len(object_type)
```

45

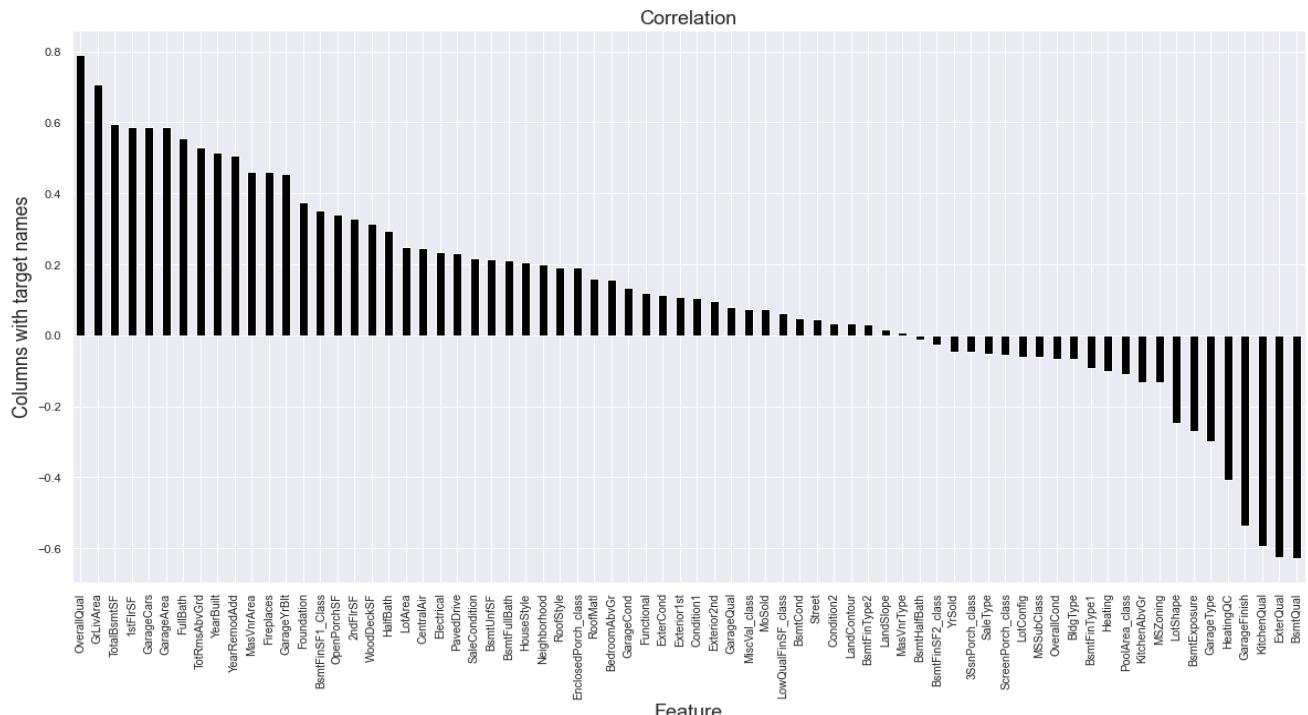
```
oe.fit(encoded_train[object_type]) # fit from train dataset
encoded_test[object_type] = oe.transform(encoded_test[object_type])

encoded_train[object_type] = oe.transform(encoded_train[object_type])
```

We have to encode feature column, therefore I have used ordinal encoding.

I have also checked model performance by applying OneHotEncoding but ordinal encoding is giving best result as compare to OneHotEncoding.

➤ **Data Inputs- Logic- Output Relationships**



“OverallQual” and “GLivArea” both are positively correlated with target variable.

“BsmtQual” and “ExterQual” both are negatively correlated with target variable

### ➤ State the set of assumptions

#### OUTLIERS:

Outliers are not removed from the dataset because it was giving high loss of data; therefore I have worked on existing dataset without removing outliers

##### 1. Try zscore technique

```
: from scipy.stats import zscore  
  
: z = np.abs(zscore(encoded_train) )  
df_z = encoded_train[(z < 3).all(axis = 1)]  
df_z.shape  
  
# (465, 73)  
  
: (465, 73)  
  
: (encoded_train.shape[0] - df_z.shape[0] ) / encoded_train.shape[0]*100  
# 60.18835616438356  
  
: 60.18835616438356
```

##### 2. IQR Technique

```
: Q1 = encoded_train.quantile(0.25)  
Q3 = encoded_train.quantile(0.75)  
IQR = Q3 - Q1  
  
: df_IQR = encoded_train[((encoded_train < (Q1 - 1.5*IQR)) | (encoded_train > (Q3 + 1.5*IQR))).any(axis = 1)]  
df_IQR.shape  
# (84, 73)  
  
: (84, 73)  
  
: (encoded_train.shape[0] - df_IQR.shape[0] ) / encoded_train.shape[0] * 100  
# This method is giving 60% data loss, therefore we can not implement this method to the dataset  
: 92.8082191780822
```

```
# As both method of removing outliers, IQR and zscore is giving high loss data, therefore, I choose to work with existing data without removing outliers.
```

## Separating dataset into x1 and y1 form

```
x = encoded_train.drop(columns= ['SalePrice'])
y = encoded_train['SalePrice']

print('shape of x', x.shape)
print('shape of y', y.shape)
print('shape of test dataset: ', encoded_test.shape)

shape of x (1168, 72)
shape of y (1168,)
shape of test dataset: (292, 72)
```

Separating Dataset into x and y form, where x is feature and y is target variable

## Treatment of Skewness of columns:

```
from sklearn.preprocessing import PowerTransformer
pt = PowerTransformer(method='yeo-johnson')
from scipy import stats

pt.fit(x[num_col])
encoded_train[num_col] = pt.transform(encoded_train[num_col])

encoded_test[num_col] = pt.transform(encoded_test[num_col])

encoded_train[num_col].skew()
```

|              | Skewness Value |
|--------------|----------------|
| MSSubClass   | 0.064007       |
| LotArea      | 0.032509       |
| YearBuilt    | -0.126641      |
| YearRemodAdd | -0.225131      |
| MasVnrArea   | 0.439526       |
| BsmtUnfSF    | -0.284390      |
| TotalBsmtSF  | 0.286779       |
| 1stFlrSF     | -0.002391      |
| 2ndFlrSF     | 0.280208       |
| GrLivArea    | -0.000054      |
| GarageYrBlt  | -0.136293      |
| GarageArea   | 0.000291       |
| WoodDeckSF   | 0.113026       |
| OpenPorchSF  | -0.002749      |
| MoSold       | -0.035838      |
| YrSold       | 0.112893       |
|              | dtype: float64 |

I have applied many skewness treatment on this dataset for removing skewness. After applying many skewness treatment this applied treatment is working best to remove skewness.

## Multicollinearity

Multicollinearity is removed using variance inflation factor,

Here cal\_vif is a my created function to calculate vif value of column, this function is created using variance inflation factor imported from statsmodels library

```
from statsmodels.stats.outliers_influence import variance_inflation_factor

# function to calculate VIF
def cal_vif(data):
    vif = pd.DataFrame()
    vif['Columns Name'] = data.columns # columns name
    vif['VIF'] = [variance_inflation_factor(data.values, i) for i in range(data.shape[1])]
    return (vif)
```

| Columns Name |              | VIF          | Columns Name |              | VIF          | Columns Name |             | VIF         |
|--------------|--------------|--------------|--------------|--------------|--------------|--------------|-------------|-------------|
| 0            | MSSubClass   | 3.362308     | 0            | MSSubClass   | 3.334524     | 0            | MSSubClass  | 3.329888    |
| 1            | LotArea      | 2.861014     | 1            | LotArea      | 2.835146     | 1            | LotArea     | 2.832445    |
| 2            | YearBuilt    | 13909.329062 | 2            | YearBuilt    | 8265.180342  | 2            | YearBuilt   | 6226.374354 |
| 3            | YearRemodAdd | 16549.199285 | 3            | YearRemodAdd | 15535.781100 | 3            | MasVnrArea  | 1.743527    |
| 4            | MasVnrArea   | 1.765547     | 4            | MasVnrArea   | 1.761914     | 4            | BsmtUnfSF   | 3.305136    |
| 5            | BsmtUnfSF    | 3.366975     | 5            | BsmtUnfSF    | 3.323654     | 5            | TotalBsmtSF | 24.854914   |
| 6            | TotalBsmtSF  | 25.011037    | 6            | TotalBsmtSF  | 24.872871    | 6            | 1stFlrSF    | 651.897880  |
| 7            | 1stFlrSF     | 673.603256   | 7            | 1stFlrSF     | 651.901830   | 7            | 2ndFlrSF    | 131.568237  |
| 8            | 2ndFlrSF     | 135.612987   | 8            | 2ndFlrSF     | 131.568789   | 8            | GrLivArea   | 1065.169049 |
| 9            | GrLivArea    | 1093.190135  | 9            | GrLivArea    | 1065.505588  | 9            | GarageArea  | 14.152404   |
| 10           | GarageYrBlt  | 25164.114865 | 10           | GarageArea   | 14.339255    | 10           | WoodDeckSF  | 1.813654    |
| 11           | GarageArea   | 17.948075    | 11           | WoodDeckSF   | 1.822872     | 11           | OpenPorchSF | 1.818040    |
| 12           | WoodDeckSF   | 1.844064     | 12           | OpenPorchSF  | 1.829156     | 12           | MoSold      | 6.642341    |
| 13           | OpenPorchSF  | 1.834544     | 13           | MoSold       | 6.642457     | 13           | YrSold      | 6014.344955 |
| 14           | MoSold       | 6.651594     | 14           | YrSold       | 12182.377295 |              |             |             |
| 15           | YrSold       | 14756.152411 |              |              |              |              |             |             |

**Step 1:** As we can see, GarageYrBlt column is showing high vif value, therefore I have remove it.

**Step 2:** YearRemodYr column is showing high vif, therefore this column is also removed

**Step 3:** Now, maximum column is in an acceptable range, and for those column which are showing high vif, we cannot delete much columns because it will lead to high loss of data.

### Scaling:

Scaling of column, this step apply just before applying Machine learning of dataset, by followed this I have applied this Standard Scaling technique to the dataset

```

: from sklearn.preprocessing import StandardScaler
ss = StandardScaler() # Instance of Standard Scaler

: # Scaling Training dataset
x[x.columns] = ss.fit_transform(x[x.columns])
x.head()

:   MSSubClass MSZoning LotArea Street LotShape LandContour LotConfig LandSlope Neighborhood Condition1 Condition2 BldgType HouseStyle
0   1.508301 -0.021646 -0.620616 0.058621 -1.373107 0.318473 0.606420 -0.226126 0.142224 -0.037339 -0.023979 2.985495 -0.549930
1   -0.877042 -0.021646 0.600903 0.058621 -1.373107 0.318473 0.606420 3.295414 -0.024227 -0.037339 -0.023979 -0.403288 -0.549930
2    0.077095 -0.021646 -0.063075 0.058621 -1.373107 0.318473 -1.220661 -0.226126 0.475125 -0.037339 -0.023979 -0.403288 1.030838
3   -0.877042 -0.021646 0.141424 0.058621 -1.373107 0.318473 0.606420 -0.226126 0.308675 -0.037339 -0.023979 -0.403288 -0.549930
4   -0.877042 -0.021646 0.686902 0.058621 -1.373107 0.318473 -0.611634 -0.226126 0.308675 -0.037339 -0.023979 -0.403288 -0.549930

```

```

print(encoded_test.shape)
print(x.shape)

(292, 70)
(1168, 70)

y.shape

(1168,)

```

## ➤ Model/s Development and Evaluation

In Machine Learning, I have applied 3 base models of regression problem, after this I have also applied boosting techniques over this dataset and then selected one of them based on their performance.

After choosing best model, Hyper Parameter tuning is applied over it to increase the accuracy of model.

In final, finalize the final model for this dataset based on its performance.

In this machine learning process, I have created a function named “ML\_Model” for smooth functioning of machine learning,

This function will give Training and testing data accuracy along with Mean Squared error and Mean Absolute Error and cv score on different-different cross fold values.

## # Libraries used in process of Machine Learning:

### Machine learning

```
from sklearn.model_selection import train_test_split
from sklearn.metrics import r2_score, mean_squared_error, mean_absolute_error
from sklearn.model_selection import GridSearchCV
from sklearn.model_selection import cross_val_score
from sklearn.tree import DecisionTreeRegressor
from sklearn.linear_model import LinearRegression
from sklearn.svm import SVR
from sklearn.neighbors import KNeighborsRegressor
from sklearn.ensemble import RandomForestRegressor
from sklearn.ensemble import GradientBoostingRegressor
from sklearn.ensemble import AdaBoostRegressor
import xgboost
from xgboost import XGBRegressor
```

# creating instances for ML model and storing them into a list of models.

After storing instances into model list, apply ML\_Model function with parameter model and dataset (x, y)

```
: lr = LinearRegression()
dtr = DecisionTreeRegressor()
knn = KNeighborsRegressor()

models = [lr, dtr, knn]
ML_Model(models, x, y)
```

# Here, ML\_Model is my created function which will give performance of the model based on various performing factor accuracy, MSE, MAE and cv score.

### # Linear Regression Model Performance :

```
LinearRegression() is giving best accuracy 0.8546463998110927 on random state of 37
Training accuracy is : 0.8326217719611156
Testing accuracy is : 0.8569051712981776
-----
Mean Squared Error: 719212202.5057057
Mean Absolute Error: 20269.256900642973
-----
Cross value score
cv score 0.692381605033251 at 2 cross fold
cv score 0.7412202227321302 at 3 cross fold
cv score 0.756604352733875 at 4 cross fold
cv score 0.7449757887406958 at 5 cross fold
cv score 0.7585327316079106 at 6 cross fold
cv score 0.7293223736940927 at 7 cross fold
```

---

- This algorithm is not giving over fitted model, but it is having high difference between cv and test data accuracy

### # DecisionTreeRegressor Model Performance:

```
DecisionTreeRegressor() is giving best accuracy 0.8247455396270994 on random state of 12
Training accuracy is : 1.0
Testing accuracy is : 0.7861177776969679
-----
Mean Squared Error: 1633946569.980057
Mean Absolute Error: 25167.05128205128
-----
Cross value score
cv score 0.5002159106685531 at 2 cross fold
cv score 0.6858012782406403 at 3 cross fold
cv score 0.653282831632314 at 4 cross fold
cv score 0.7187454884073861 at 5 cross fold
cv score 0.7326260069099656 at 6 cross fold
cv score 0.6979664901626247 at 7 cross fold
```

---

- By applying this model, getting overfitted model, as its traing accuracy is higher than testing accuracy
- Training Accuracy > Testing Accuracy

### # KNeighborsRegressor Model Performance

```
KNeighborsRegressor()  is giving best accuracy 0.7213838331827729 on random state of 34
Training accuracy is : 0.8383771779034621
Testing accuracy is : 0.7855729216403198
-----
Mean Squared Error: 1028584891.2813675
Mean Absolute Error: 21261.302564102563
-----
Cross value score
cv score 0.7238722970882102 at 2 cross fold
cv score 0.7378697405780926 at 3 cross fold
cv score 0.7296852724791547 at 4 cross fold
cv score 0.7319152856571005 at 5 cross fold
cv score 0.7303304025735947 at 6 cross fold
cv score 0.727269195081118 at 7 cross fold
-----
```

- It is having less accuracy and also having high difference between cv score value and accuracy value

### # Applied Boosting Techniques:

```
from sklearn.ensemble import RandomForestRegressor
from sklearn.ensemble import GradientBoostingRegressor
from sklearn.ensemble import AdaBoostRegressor
rfr = RandomForestRegressor()
gbr = GradientBoostingRegressor()
abr = AdaBoostRegressor()

models = [rfr, gbr, abr]
ML_Model(models, x, y)
```

- Bagging and Boosting techniques are applied on the dataset because normal algorithm was not giving enough result or good result

### # RandomForestRegressor Model Performance:

```
RandomForestRegressor()  is giving best accuracy 0.8602367294100308 on random state of 19
Training accuracy is :  0.9739880325466944
Testing accuracy is : 0.8829821893882084
-----
Mean Squared Error:  698728252.8761048
Mean Absolute Error:  17877.57965811966
-----
Cross value score
cv score 0.8086910312176738 at 2 cross fold
cv score 0.8588107218647809 at 3 cross fold
cv score 0.8390309744118734 at 4 cross fold
cv score 0.8414786265655934 at 5 cross fold
cv score 0.8517135755216628 at 6 cross fold
cv score 0.828347543734054 at 7 cross fold
```

---

- Having difference in accuracy of training and testing dataset, not much but having,

### # AdaBoostRegressor Model Performance:

```
AdaBoostRegressor()  is giving best accuracy 0.8078649544014966 on random state of 39
Training accuracy is :  0.8849071806084483
Testing accuracy is : 0.8230374575266938
-----
Mean Squared Error:  968630008.4141606
Mean Absolute Error:  22007.300167520945
-----
Cross value score
cv score 0.7634368146856891 at 2 cross fold
cv score 0.7935140877758146 at 3 cross fold
cv score 0.7845199169851335 at 4 cross fold
cv score 0.792675764114988 at 5 cross fold
cv score 0.7828340585920563 at 6 cross fold
cv score 0.7842546755765206 at 7 cross fold
```

---

- ML\_Model function is created for smooth operation of machine learning
- Giving Good results

### # XGBoostRegressor Model Performance:

```

Training accuracy is : 0.9999741751479383
Testing accuracy is : 0.8839930186140819
-----
Mean Squared Error: 718522764.5732626
Mean Absolute Error: 17632.69255252849
-----
Cross value score
cv score 0.8353491059844554 at 2 cross fold
cv score 0.8511448028030243 at 3 cross fold
cv score 0.8206761353698495 at 4 cross fold
cv score 0.8388832744206691 at 5 cross fold
cv score 0.8594095268629013 at 6 cross fold
cv score 0.8313952730129088 at 7 cross fold
-----
```

- This model is giving good result but giving over fitted model

### # GradientBoostingRegressor Model Performance:

```

GradientBoostingRegressor() is giving best accuracy 0.8917976009767302 on random state of 37
Training accuracy is : 0.9726752237962414
Testing accuracy is : 0.9049629969315385
-----
Mean Squared Error: 477667662.18253416
Mean Absolute Error: 16008.743091677125
-----
Cross value score
cv score 0.8388940013932447 at 2 cross fold
cv score 0.8750609850293148 at 3 cross fold
cv score 0.8513533945379466 at 4 cross fold
cv score 0.8611798083967706 at 5 cross fold
cv score 0.8711118848754861 at 6 cross fold
cv score 0.8503287678925674 at 7 cross fold
-----
```

- This model is giving good result as it not having much difference in accuracy of training and testing dataset
- And also not having high cv value and accuracy difference

| Models | Training Accuracy | Testing Accuracy | CV Score | Difference |
|--------|-------------------|------------------|----------|------------|
|--------|-------------------|------------------|----------|------------|

| #                            |          |          |          |       |
|------------------------------|----------|----------|----------|-------|
| RandomForestRegressor        | 0.973988 | 0.882982 | 0.85881  | 0.024 |
| # GradientBoosting Regressor | 0.972675 | 0.904962 | 0.87506  | 0.030 |
| # AdaBoostRegressor          | 0.884907 | 0.823037 | 0.793514 | 0.030 |
| # XGBRegressor               | 0.999974 | 0.883993 | 0.859409 | 0.025 |

- We can observed that least difference between testing accuracy and cv, found in xgboost but it is giving overfitted model and
- GradientBoost and AdaBoostRegressor is both are giving same cv and accuracy difference, and we select GradientBoosting because it is showing more accuracy as compare to adaboost

### Ensemble Techniques:

Based on model performance, we can see, GradientBoosting is performing best as compare to other ML Models, Its cv value is also very close to accuracy of model.

Hence I had applied hyper parameter tuning for this Model to increase the accuracy of model

```
: from sklearn.model_selection import GridSearchCV
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.30, random_state = 37)

: parameter = {'alpha' : [0.9, 0.09, 0.1],
               'learning_rate' : [0.1, 0.01],
               'max_depth' : [ 3, 4, 5],
               'min_samples_leaf' : [1, 2, 3],
               'min_samples_split': [2,3,4],
               'n_estimators': [100, 50, 10]}

: gcv = GridSearchCV(estimator = GradientBoostingRegressor(), param_grid = parameter, cv = 7)
gcv.fit(x_train, y_train)

: GridSearchCV(cv=7, estimator=GradientBoostingRegressor(),
              param_grid={'alpha': [0.9, 0.09, 0.1],
                          'learning_rate': [0.1, 0.01], 'max_depth': [3, 4, 5],
                          'min_samples_leaf': [1, 2, 3],
                          'min_samples_split': [2, 3, 4],
                          'n_estimators': [100, 50, 10]})
```

## HOUSE PRICE PREDICTION MODEL

```
: gcv.best_params_
: {'alpha': 0.09,
 'learning_rate': 0.1,
 'max_depth': 3,
 'min_samples_leaf': 1,
 'min_samples_split': 3,
 'n_estimators': 100}

: # 'parameter' these are the parameters which have to be applied on given model and then have to choose best
# one.
# GridSearchCV will help in that

: gbr = GradientBoostingRegressor(alpha = 0.9, learning_rate = 0.1, max_depth = 3, min_samples_leaf=1,
                                 min_samples_split=3 , n_estimators = 100)

models = [gbr]
ML_Model(models, x, y)
```

```
GradientBoostingRegressor(min_samples_split=3) is giving best accuracy 0.8925363168251541 on random state 0
Training accuracy is : 0.9709969915980627
Testing accuracy is : 0.9012624616581131
-----
Mean Squared Error:  558346098.2162528
Mean Absolute Error:  16042.284810874853
-----
Cross value score
cv score 0.8349190207390833 at 2 cross fold
cv score 0.8766740857081136 at 3 cross fold
cv score 0.8531874797993466 at 4 cross fold
cv score 0.8616931484565242 at 5 cross fold
cv score 0.8710084859163643 at 6 cross fold
cv score 0.8507762906963834 at 7 cross fold
-----
```

FLIP FLOP

After applying hyper parameter tuning, we can see, cv and accuracy has little decreased, which is good for model

### Final Model:

```
: x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.30, random_state = 6)

final_model = GradientBoostingRegressor(alpha = 0.9, learning_rate = 0.1, max_depth = 3, min_samples_leaf=1,
                                         min_samples_split=3 , n_estimators = 100)
final_model.fit(x_train, y_train)
final_pred = final_model.predict(x_test)
final_pred_train = final_model.predict(x_train)

train_accuracy = r2_score(y_train, final_pred_train )
test_accuracy = r2_score(y_test, final_pred )
print('Training accuracy: ', train_accuracy)
print('Testing accuracy: ', test_accuracy)
print('_____')
print('Mean squared error: ', mean_squared_error(y_test, final_pred ) )
print('Mean absolute error: ', mean_absolute_error(y_test, final_pred ) )

cv_score = cross_val_score(final_model, x, y, cv = 3 ).mean()
print('cv score', cv_score , 'at', i, 'cross fold')
```

## HOUSE PRICE PREDICTION MODEL

Training accuracy: 0.9709969915980627

Testing accuracy: 0.9006411168987434

---

Mean squared error: 561859710.4438533

Mean absolute error: 16020.68811082674

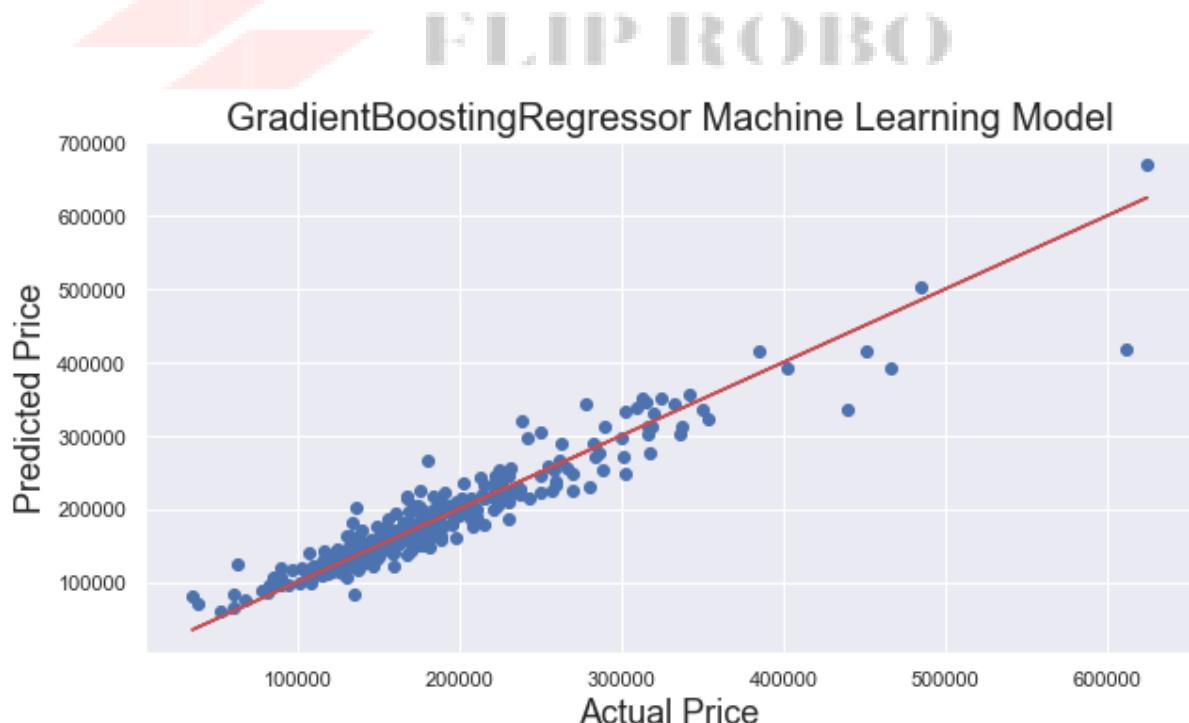
cv score 0.8773090209451722 at MiscVal\_class cross fold

### Graph for Model Performance:

```
: plt.figure(figsize = (10, 5))
plt.scatter(x = y_test, y = final_pred, color = 'b')
plt.plot(y_test, y_test, color = 'r')
plt.xlabel('Actual Price', fontsize= 18 )
plt.ylabel('Predicted Price', fontsize = 18)
plt.title('GradientBoostingRegressor Machine Learning Model', fontsize = 20)
```

In below mentioned chart, we can see how model is performing,

As our model is having 90% accuracy based on this, its chart is visible to us



## Deploying & Loading Model:

### Deploy Model

```
: import pickle  
filename = 'house_price_predictor.pkl'  
pickle.dump(final_model, open(filename, 'wb'))
```

```
:
```

### Loading model

```
: load_model = pickle.load(open('house_price_predictor.pkl', 'rb'))  
result = load_model.score(x_test, y_test)  
print(result)
```

0.9006411168987434

## Conclusion:

```
: predicted = np.array(load_model.predict(x_test))  
original = np.array(y_test)  
# convert columns in to np.array
```

```
: # print(predicted.shape)  
# print(original.shape)  
# print(x_test.shape)  
# print(y_test.shape)
```

```
conclusion = pd.DataFrame({'Actual Price': original, 'Predicted Price': predicted},  
                           index = range(len(original)))
```

```
conclusion['Predicted Price']= conclusion['Predicted Price'].apply(lambda x: round(x, 1))
```

```
conclusion.sample(10)
```

|     | Actual Price | Predicted Price |
|-----|--------------|-----------------|
| 348 | 278000       | 343380.4        |
| 70  | 320000       | 329704.3        |
| 286 | 625000       | 670691.5        |
| 317 | 175900       | 225094.2        |
| 62  | 197900       | 209751.9        |
| 196 | 302000       | 332550.9        |
| 229 | 130000       | 137226.3        |
| 142 | 128950       | 126746.8        |
| 135 | 117000       | 128221.3        |
| 312 | 196000       | 179084.5        |

## ➤ **Hardware and Software Requirements and Tools Used**

All used libraries:

```
: from sklearn.model_selection import train_test_split
from sklearn.metrics import r2_score, mean_squared_error, mean_absolute_error
from sklearn.model_selection import GridSearchCV
from sklearn.model_selection import cross_val_score
from sklearn.tree import DecisionTreeRegressor
from sklearn.linear_model import LinearRegression
from sklearn.svm import SVR
from sklearn.neighbors import KNeighborsRegressor
from sklearn.ensemble import RandomForestRegressor
from sklearn.ensemble import GradientBoostingRegressor
from sklearn.ensemble import AdaBoostRegressor
import xgboost
from xgboost import XGBRegressor
```

**Pandas:** This library used for dataframe operations

**Numpy:** This library gives statistical computation for smooth functioning

**Matplotlib:** Used for visualization

**Seaborn:** This library is also used for visualization

**Sklearn:** This library having so many machine learning module and we can import them from this library

**Pickle:** This is used for deploying the model

**Imblearn:** This library is import to get SMOTE technique for balance the data

**Scipy:** It is import to perform outlier removing technique using zscore

**Warning:** To avoid unwanted warning shows in the output

I am giving this requirement and tool used, based on my laptop configuration.

**Operating System:**      **Window 11**

**RAM:**                    **8 GB**

**Processor:**

**i5 10th Generation**

**Software:**

**Jupyter Notebook,**

➤ **Observations from the whole problem.**

- i) Dataset was having null value.
- ii) Impute values according to column relation with each other
- iii) Multicollinearity was not present in the dataset
- iv) Created group of range for few columns
- v) Skewness is removing by using PowerTransformer

➤ **Learning Outcomes of the Study in respect of Data Science**

My learnings: - the power of visualization is helpful for the understanding the data and graphical representation, its help me to understand that what data is trying to say,

Data cleaning is one of the most important step to remove missing value or null value. Perform mean or median imputation based on its effect on the column

Various algorithms I have used in this dataset and to get best result and save that model. The best algorithm is GradientBoostingClassifier.

The challenge I faced while working on this project, one record was wrong updated in the Exterior1st and Exterior2nd columns which was created issue when I was encoding the data for further machine learning process.

➤ **Limitations of this work and Scope for Future Work**

Need to train this same model with data in Lacs or more than that, then only it would be become best model and it can perform in best way.