

IN THIS NOTEBOOK:

- **ADVANCE HOUSE PRICE PREDICTION - DATA ANALYSIS**

- **ADVANCE HOUSE PRICE PREDICTION - FEATURE SELECTION**

- **DATASET:**

<https://www.kaggle.com/c/house-prices-advanced-regression-techniques/data>

DATA ANALYSIS - ADVANCE HOUSE PRICE PREDICTION

- Missing Values
 - Categorical missing values
 - Numerical missing valeues
- Relationships between independent and dependent feature(SalePrice)
- All The Numerical Variables
 - Temporal variables (Eg.)
 - Discrete features
 - Continuous features
 - Distribution of the Numerical Variables
- All Categorical Variables
 - Cardinality
 - Rare categorical variables
- Outliers
- Feature Scaling

Import Libraries

In [412...

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns

pd.pandas.set_option('display.max_columns', None)
```

Read data

```
In [413... dataset = pd.read_csv('train.csv')
```

```
In [414... dataset.shape
```

```
Out[414]: (1460, 81)
```

```
In [415... dataset.head()
```

```
Out[415]:
```

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	L
0	1	60	RL	65.0	8450	Pave	NaN	Reg		Lvl
1	2	20	RL	80.0	9600	Pave	NaN	Reg		Lvl
2	3	60	RL	68.0	11250	Pave	NaN	IR1		Lvl
3	4	70	RL	60.0	9550	Pave	NaN	IR1		Lvl
4	5	60	RL	84.0	14260	Pave	NaN	IR1		Lvl

```
In [416... dataset.dtypes
```

```
Out[416]: Id                int64
MSSubClass                int64
MSZoning                  object
LotFrontage              float64
LotArea                  int64
...
MoSold                   int64
YrSold                   int64
SaleType                 object
SaleCondition            object
SalePrice                int64
Length: 81, dtype: object
```

ALL MISSING VALUES

```
In [417... # check percentage of NaN values in each feature

# features_with_na = List of all features containing null values
features_with_na = []
for feature in dataset.columns:
    if dataset[feature].isnull().sum() > 1:
        features_with_na.append(feature)

# print the feature name along with percentage of missing value
for feature in features_with_na:
    print(feature, "- ", np.round(dataset[feature].isnull().mean(), 4), "% missing
```

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

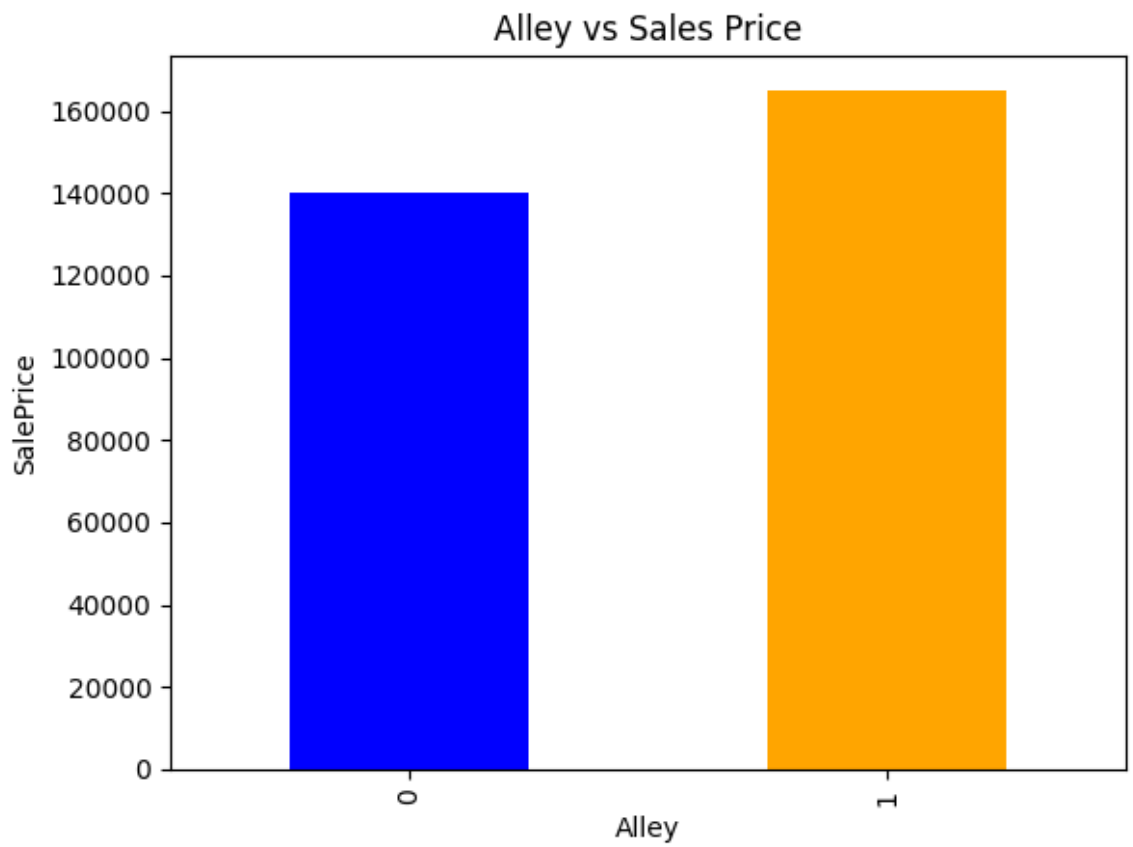
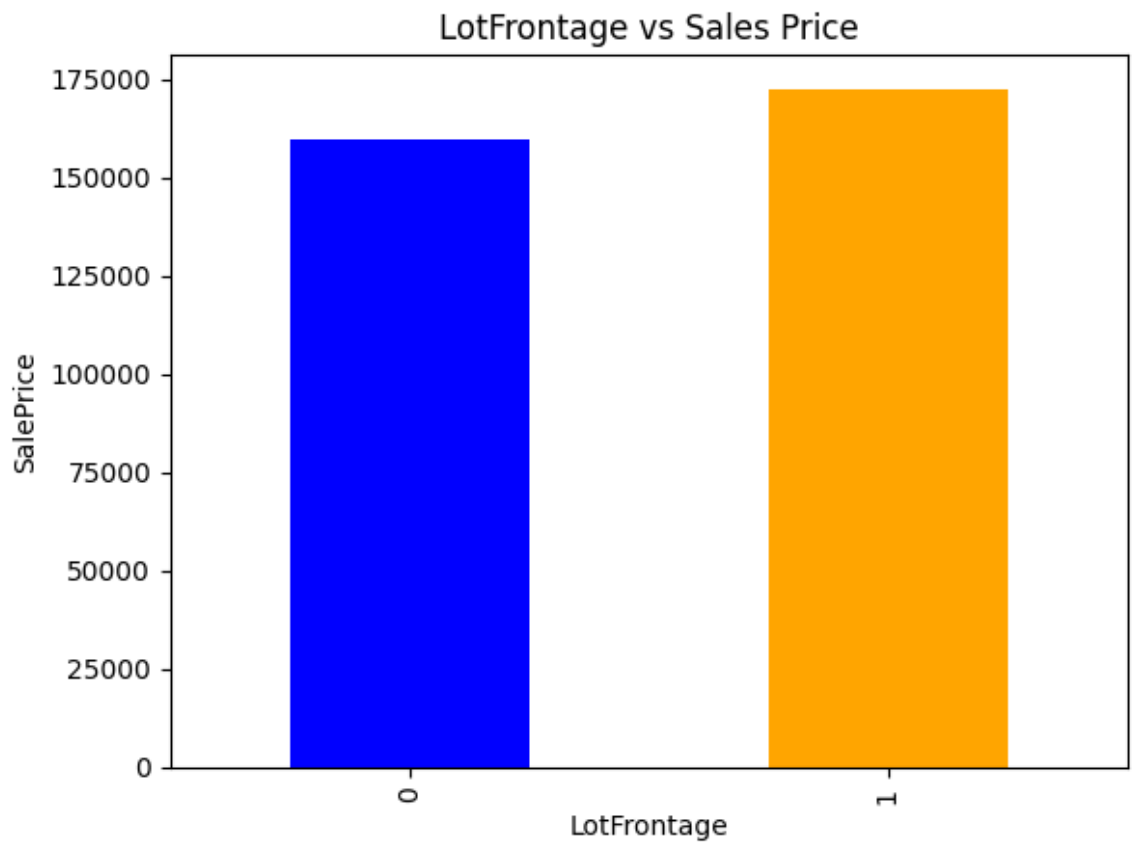
Since there are many missing values we need to find the relationship between the missing values and SalePrice

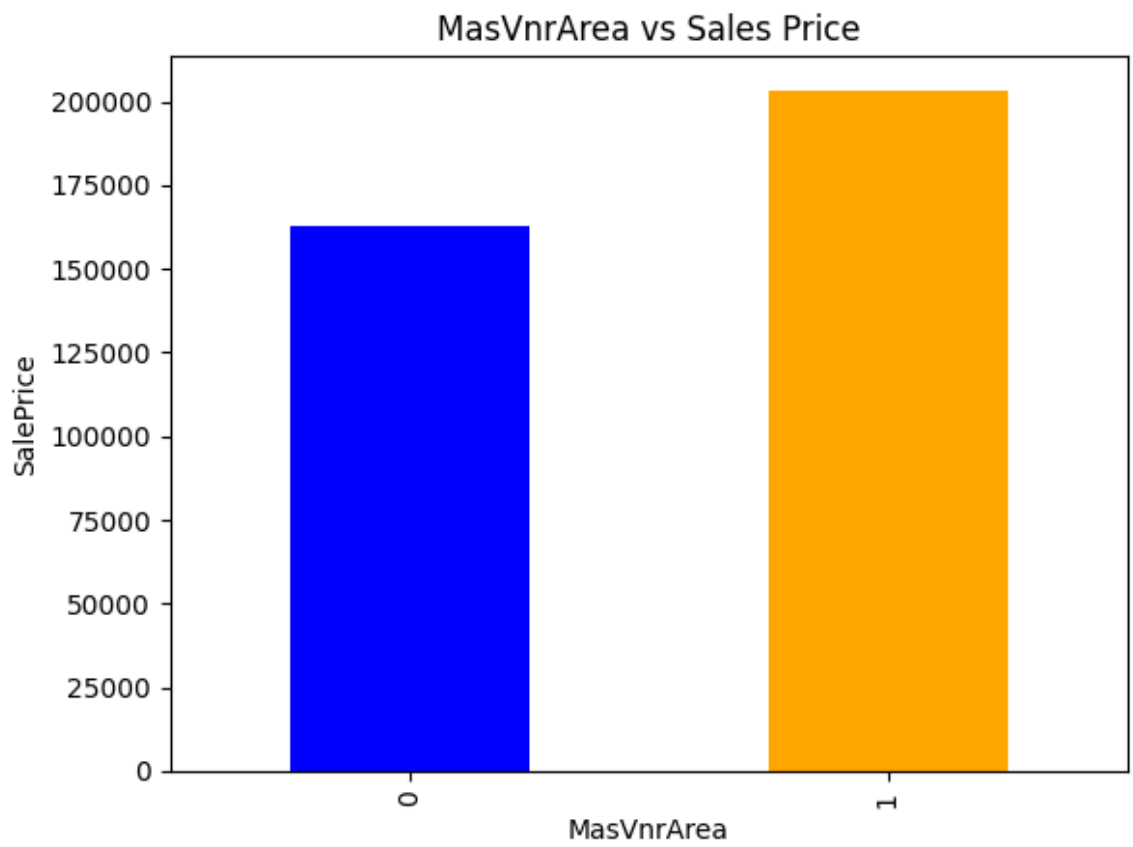
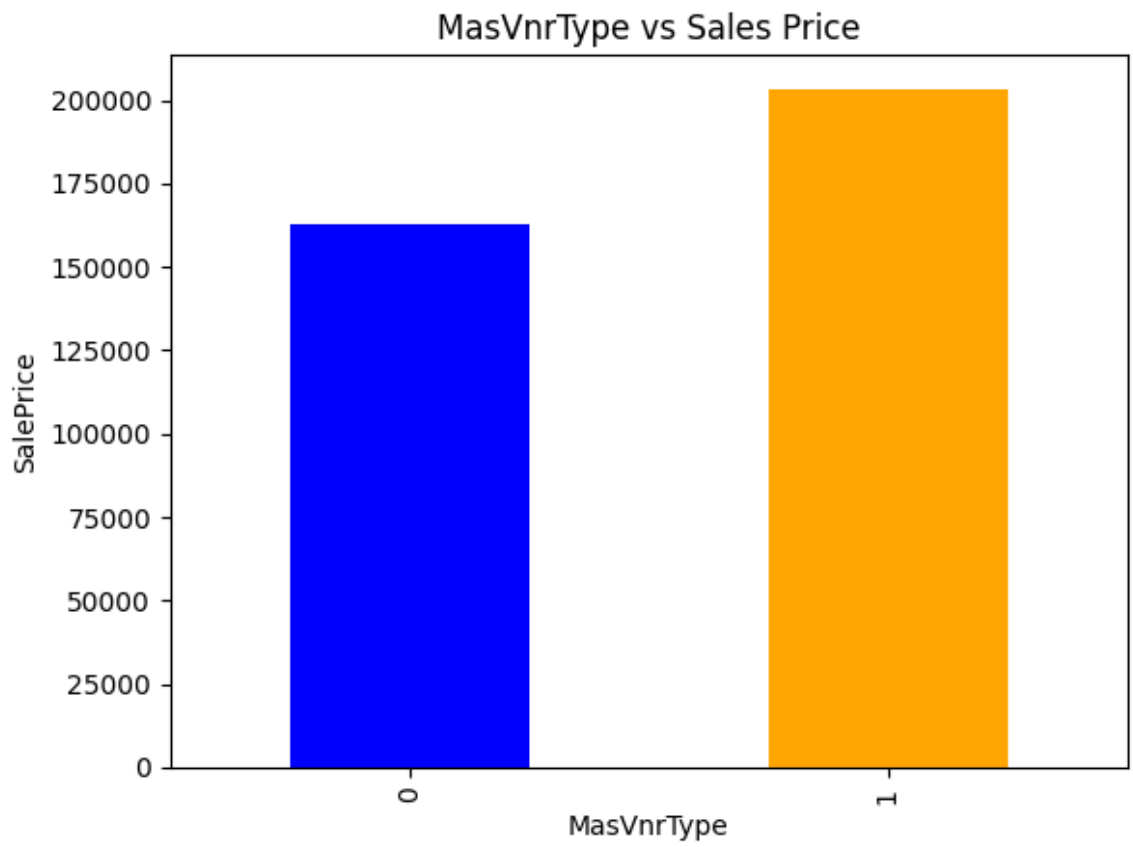
In [418...

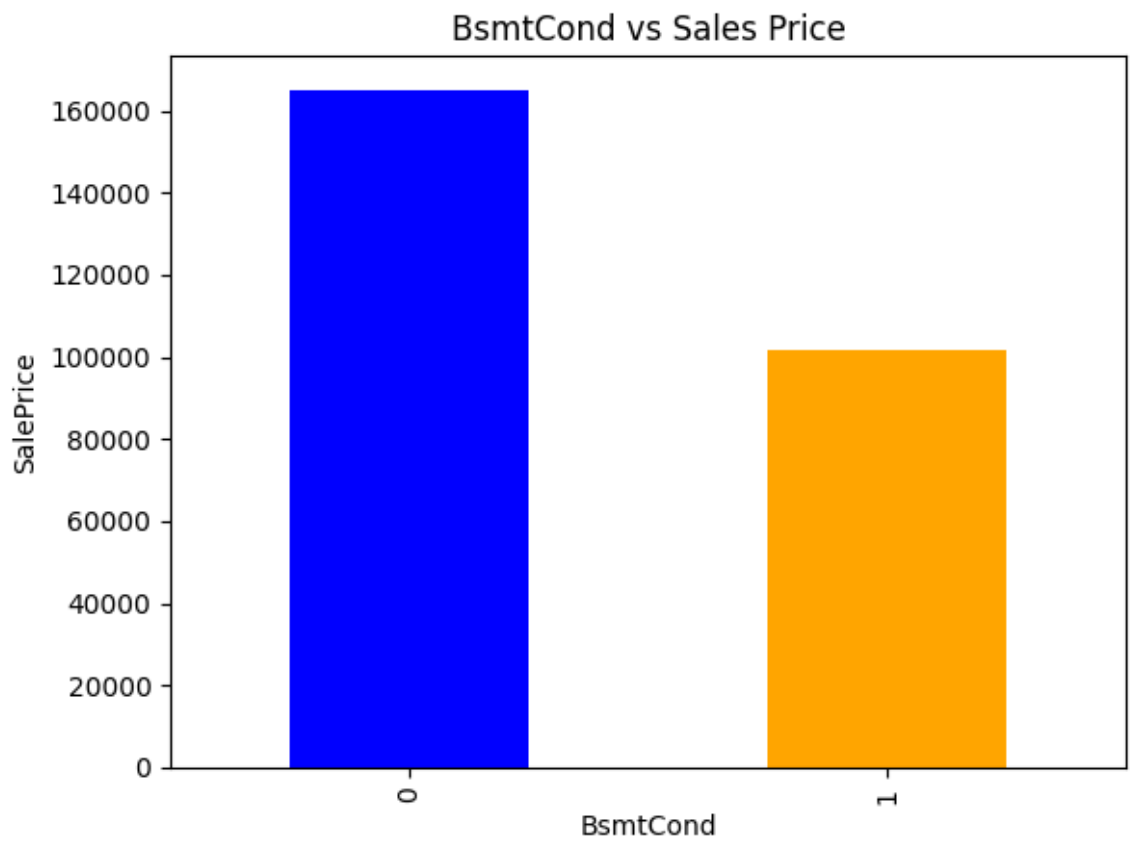
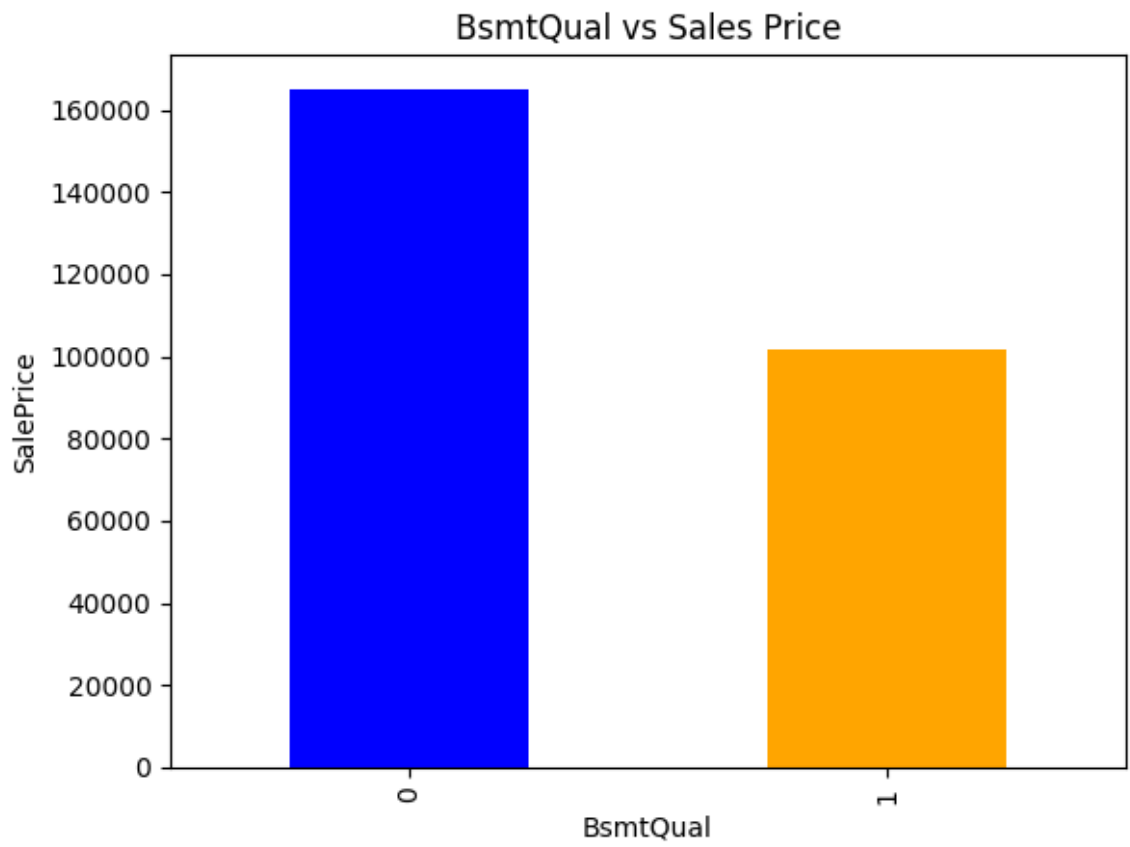
```
for feature in features_with_na:
    data = dataset.copy()

    # 1: missing value present; 0: otherwise; does this to the whole copy of dataset
    data[feature] = np.where(data[feature].isnull(), 1, 0)

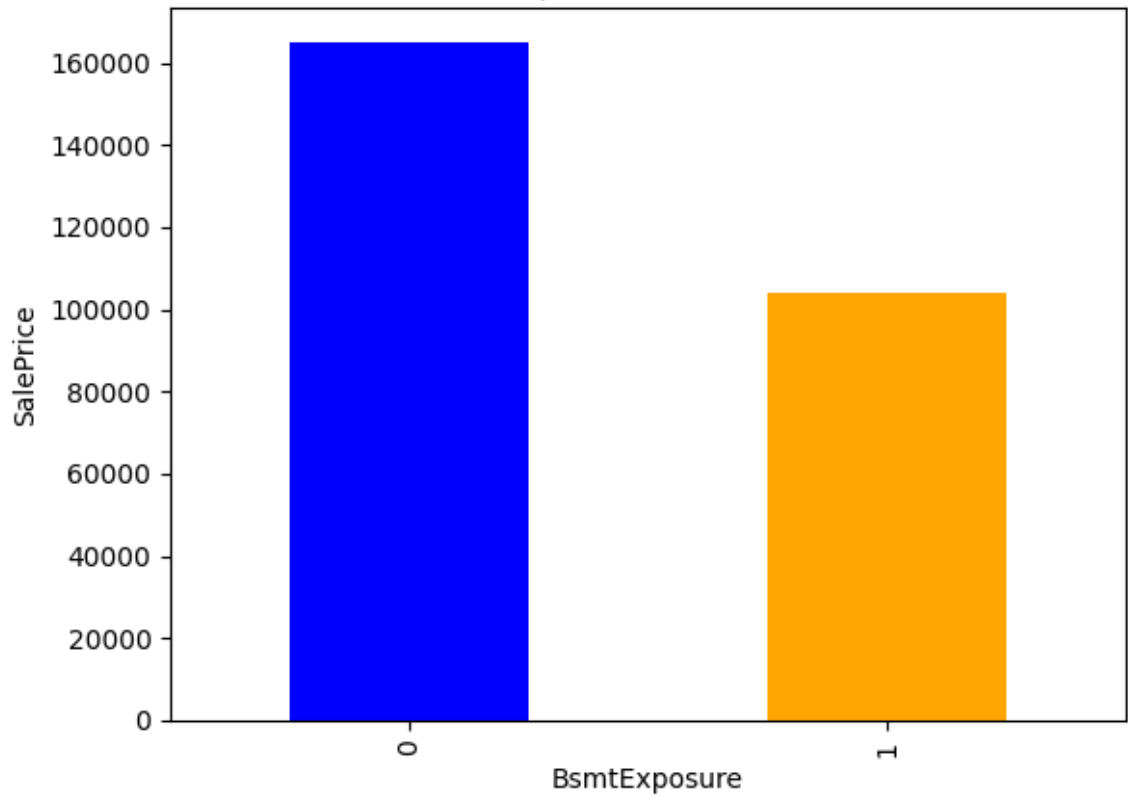
    median_sale_price_by_feature = data.groupby(feature)['SalePrice'].median()
    median_sale_price_by_feature.plot.bar(color=['blue', 'orange']) # blue = missing
    plt.xlabel(feature)
    plt.ylabel("SalePrice")
    plt.title(feature + " vs Sales Price")
    plt.show()
```



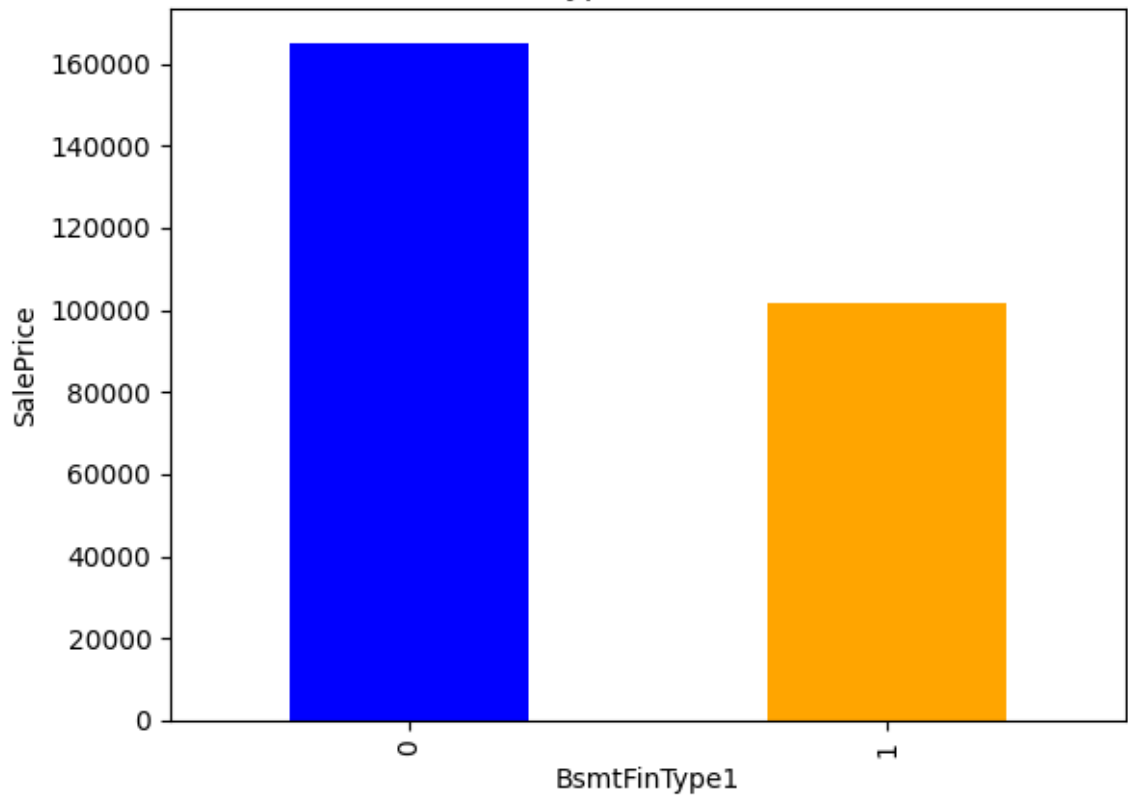




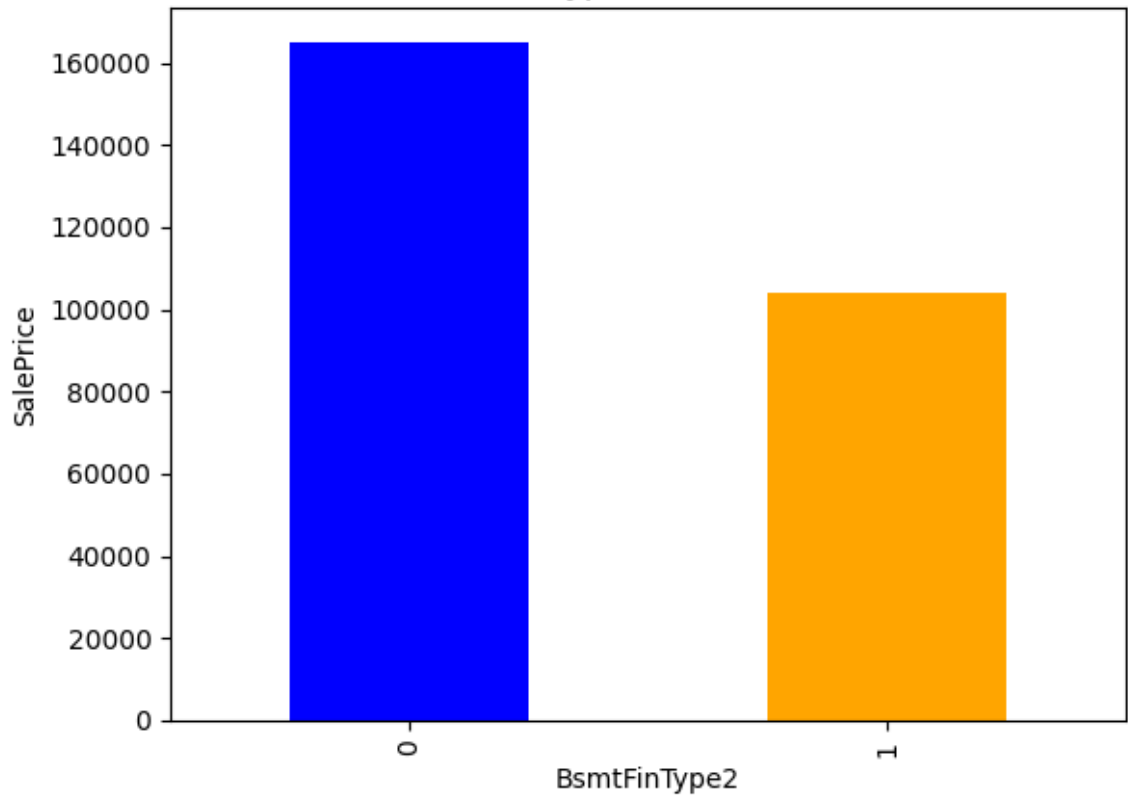
BsmtExposure vs Sales Price



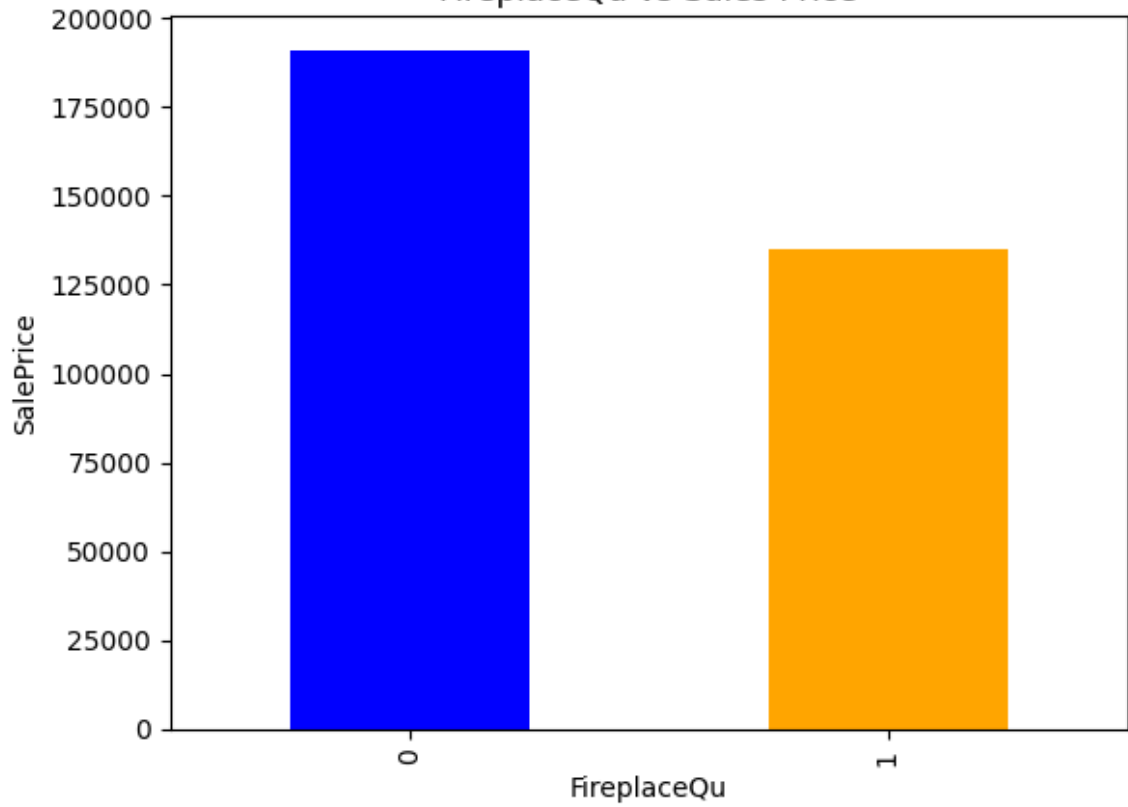
BsmtFinType1 vs Sales Price

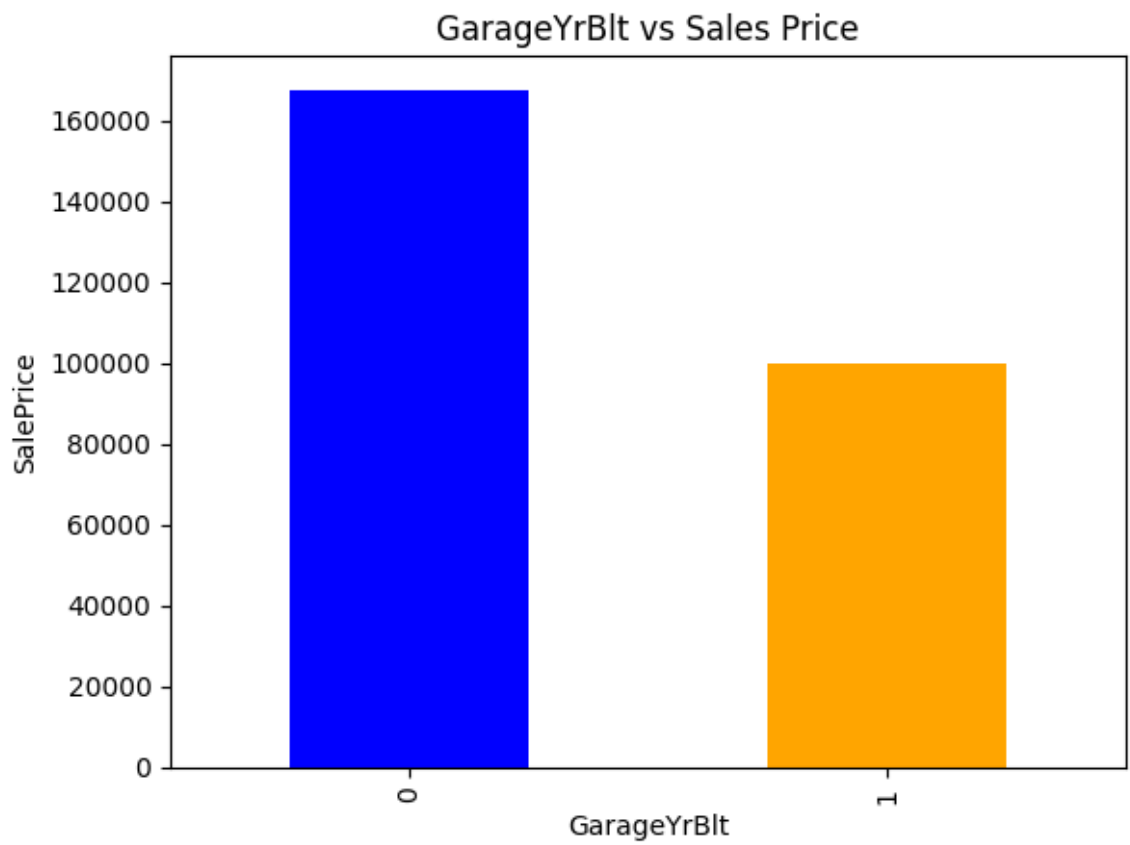
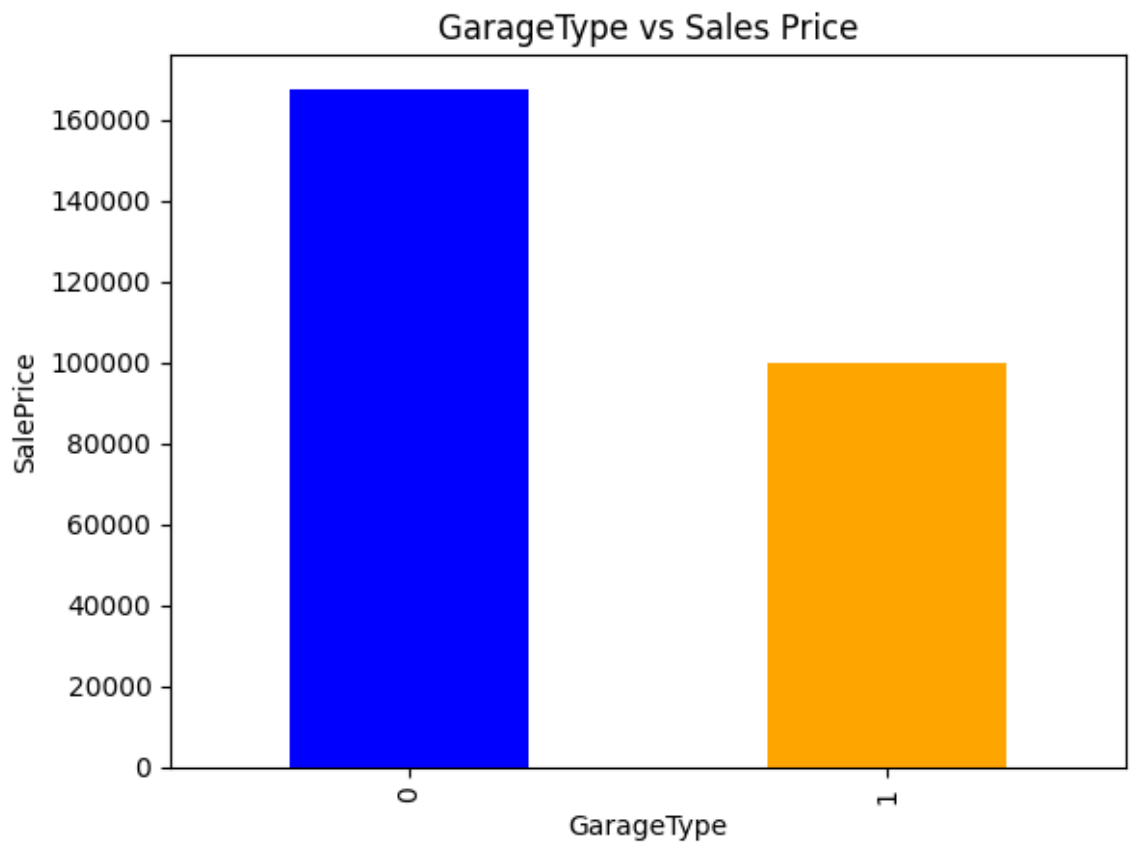


BsmtFinType2 vs Sales Price

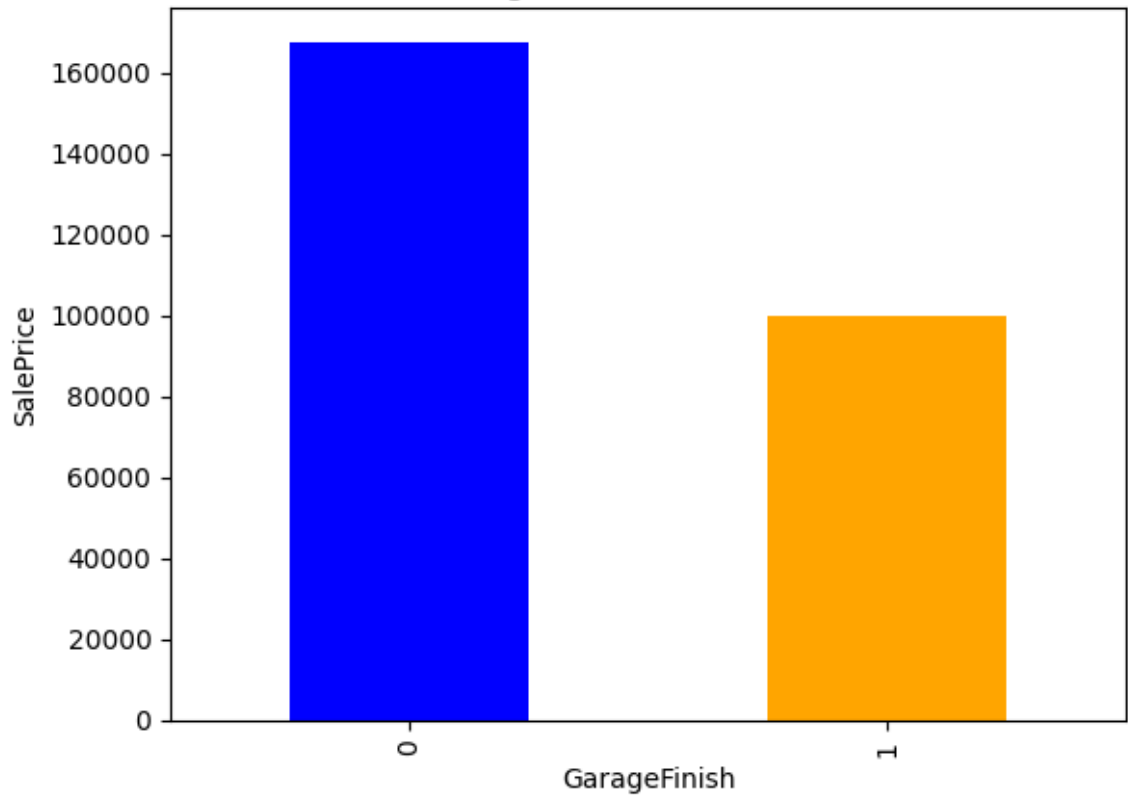


FireplaceQu vs Sales Price

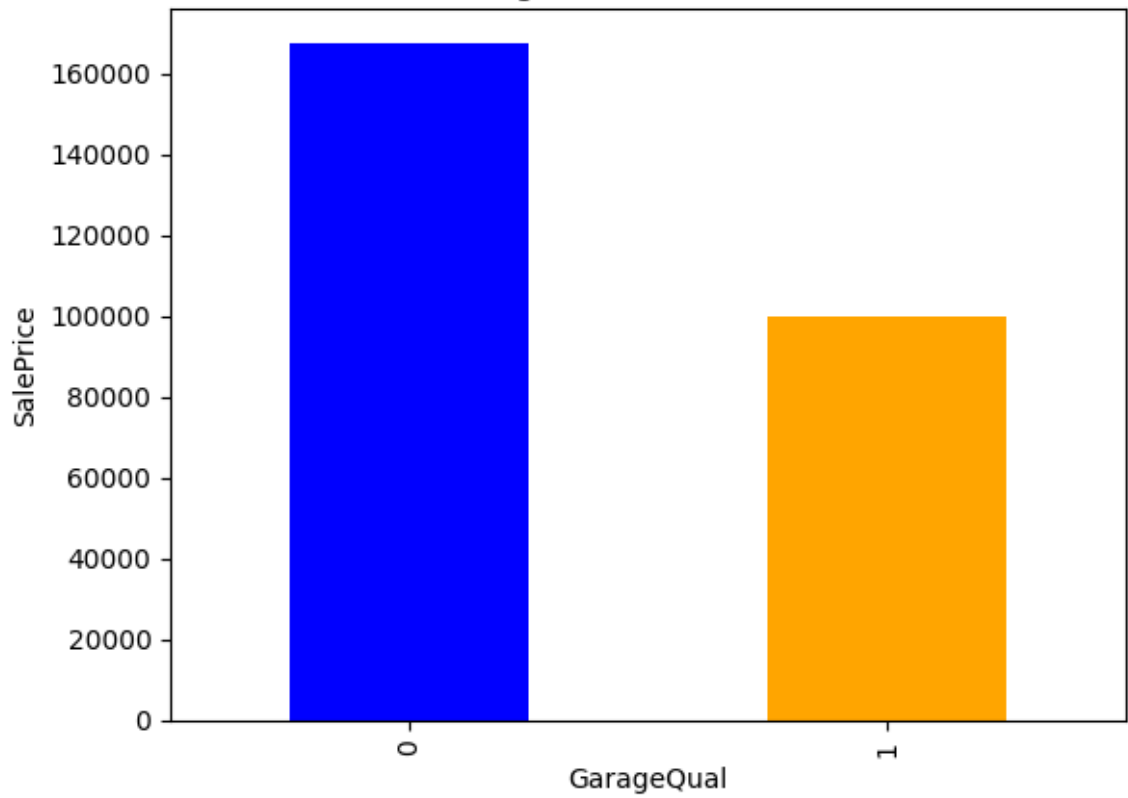


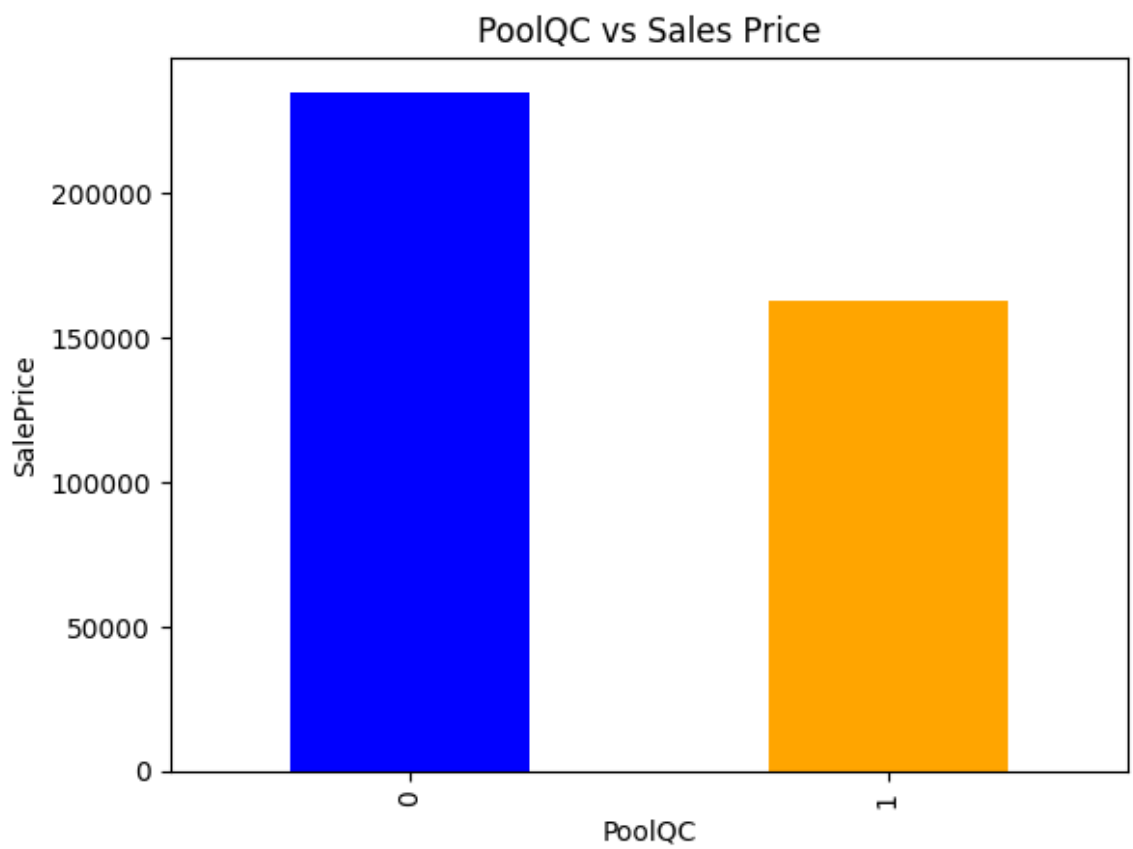
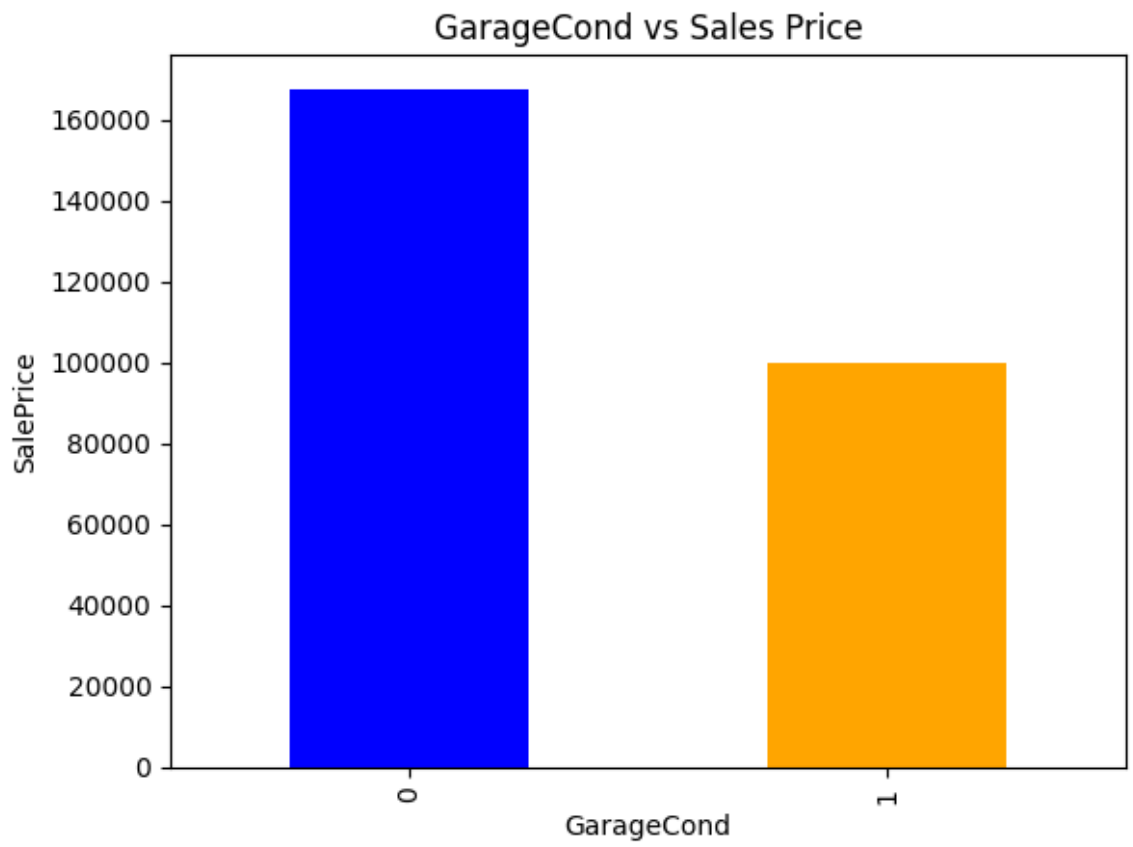


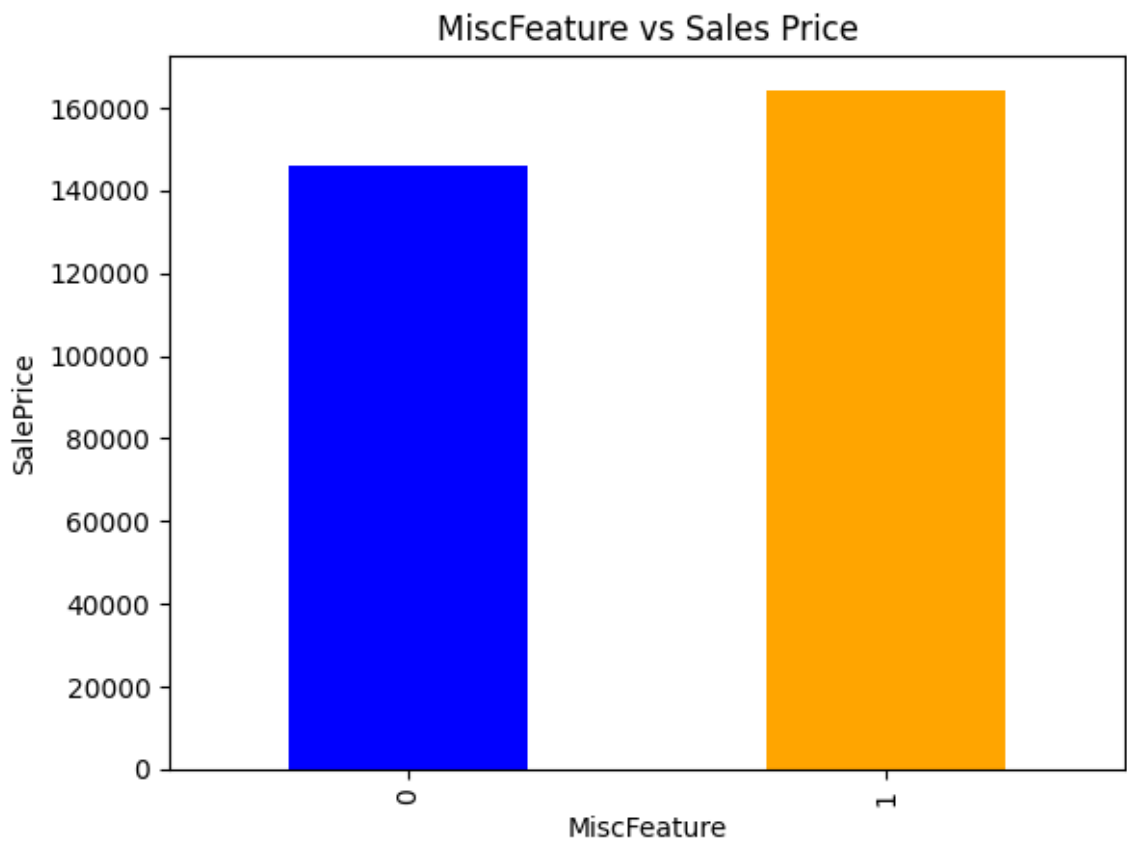
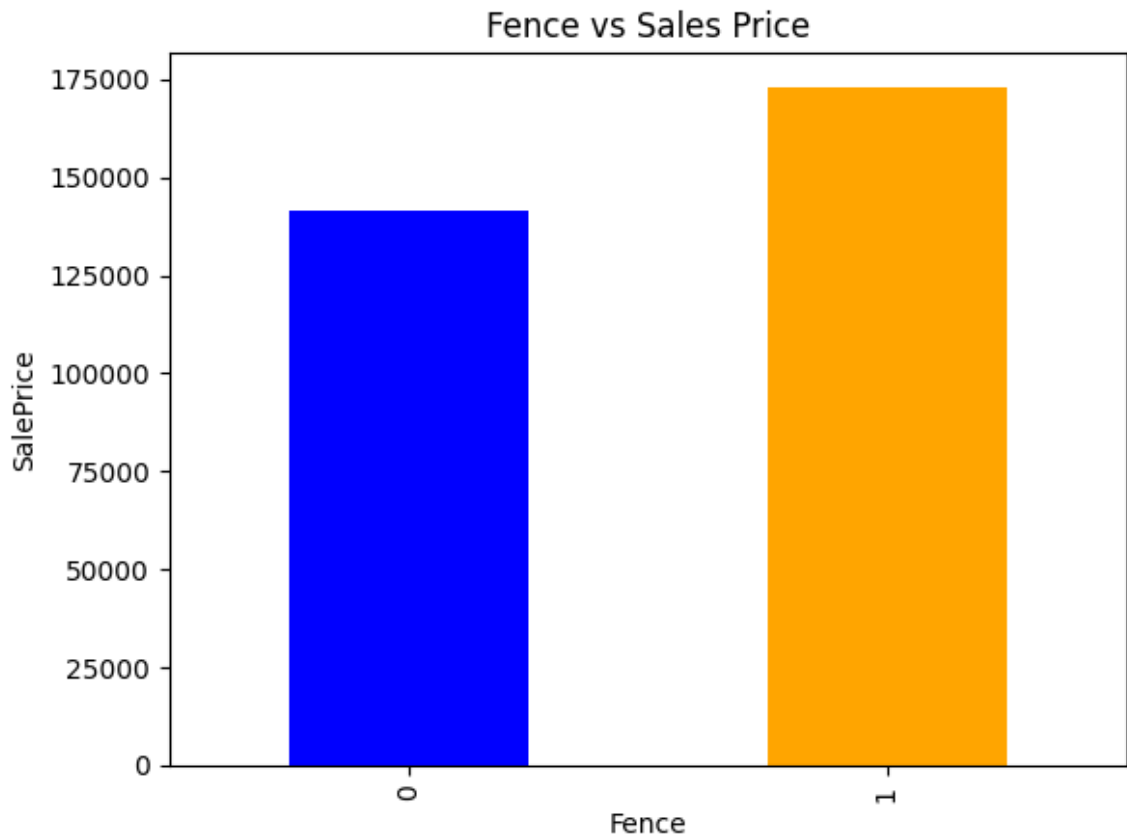
GarageFinish vs Sales Price



GarageQual vs Sales Price







Here the relation between the missing values and the sale price is clearly visible. So We need to replace these NaN values with something meaningful.

CATEGORICAL MISSING VALUES

In [419...

```
# find nan/missing values in categorical features
cat_features_nan = []
for feature in dataset.columns:
    if dataset[feature].isnull().sum() > 1 and dataset[feature].dtypes == 'O':
        cat_features_nan.append(feature)

# print the feature name along with percentage of missing value
for feature in cat_features_nan:
    print(feature, "-", np.round(dataset[feature].isnull().mean(), 4), "% missing")

## Replace missing value with a new label
def replace_cat_feature(dataset, cat_features_nan):
    data = dataset.copy()
    data[cat_features_nan] = data[cat_features_nan].fillna('Missing')
    return data

dataset = replace_cat_feature(dataset, cat_features_nan)
dataset[cat_features_nan].isnull().sum()
```

```
Alley - 0.9377 % missing values
MasVnrType - 0.0055 % missing values
BsmtQual - 0.0253 % missing values
BsmtCond - 0.0253 % missing values
BsmtExposure - 0.026 % missing values
BsmtFinType1 - 0.0253 % missing values
BsmtFinType2 - 0.026 % missing values
FireplaceQu - 0.4726 % missing values
GarageType - 0.0555 % missing values
GarageFinish - 0.0555 % missing values
GarageQual - 0.0555 % missing values
GarageCond - 0.0555 % missing values
PoolQC - 0.9952 % missing values
Fence - 0.8075 % missing values
MiscFeature - 0.963 % missing values
```

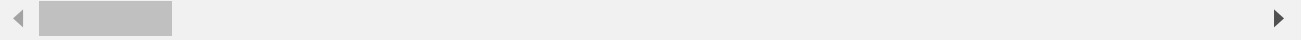
```
Out[419]: Alley      0
MasVnrType    0
BsmtQual      0
BsmtCond      0
BsmtExposure  0
BsmtFinType1  0
BsmtFinType2  0
FireplaceQu   0
GarageType    0
GarageFinish   0
GarageQual    0
GarageCond    0
PoolQC        0
Fence         0
MiscFeature    0
dtype: int64
```

In [420...

```
dataset.head()
```

Out[420]:

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour
0	1	60	RL	65.0	8450	Pave	Missing	Reg	Lvl
1	2	20	RL	80.0	9600	Pave	Missing	Reg	Lvl
2	3	60	RL	68.0	11250	Pave	Missing	IR1	Lvl
3	4	70	RL	60.0	9550	Pave	Missing	IR1	Lvl
4	5	60	RL	84.0	14260	Pave	Missing	IR1	Lvl



NUMERICAL MISSING VALUES

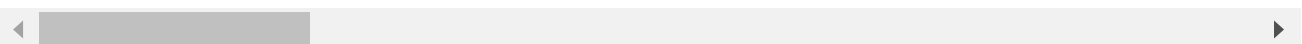
```
In [421... # find list of all numerical features
numerical_features = []
for feature in dataset.columns:
    if dataset[feature].dtypes != 'O':
        numerical_features.append(feature)
        print(feature)

dataset[numerical_features].head()
```

Id
 MSSubClass
 LotFrontage
 LotArea
 OverallQual
 OverallCond
 YearBuilt
 YearRemodAdd
 MasVnrArea
 BsmtFinSF1
 BsmtFinSF2
 BsmtUnfSF
 TotalBsmtSF
 1stFlrSF
 2ndFlrSF
 LowQualFinSF
 GrLivArea
 BsmtFullBath
 BsmtHalfBath
 FullBath
 HalfBath
 BedroomAbvGr
 KitchenAbvGr
 TotRmsAbvGrd
 Fireplaces
 GarageYrBlt
 GarageCars
 GarageArea
 WoodDeckSF
 OpenPorchSF
 EnclosedPorch
 3SsnPorch
 ScreenPorch
 PoolArea
 MiscVal
 MoSold
 YrSold
 SalePrice

Out[421]:

	Id	MSSubClass	LotFrontage	LotArea	OverallQual	OverallCond	YearBuilt	YearRemodAdd
0	1	60	65.0	8450	7	5	2003	2003
1	2	20	80.0	9600	6	8	1976	1976
2	3	60	68.0	11250	7	5	2001	2002
3	4	70	60.0	9550	7	5	1915	1970
4	5	60	84.0	14260	8	5	2000	2000



In [422...]

```

# find nan/missing values in numerical features
num_features_nan = []
for feature in dataset.columns:
    if feature in numerical_features and dataset[feature].isnull().sum() > 1:
        num_features_nan.append(feature)

# print the feature name along with percentage of missing value
  
```

```
for feature in num_features_nan:
    print(feature, "-", np.round(dataset[feature].isnull().mean(), 4), "% missing")
```

```
LotFrontage - 0.1774 % missing values
MasVnrArea - 0.0055 % missing values
GarageYrBlt - 0.0555 % missing values
```

In [423]...

```
# Replacing the numerical Missing Values

for feature in num_features_nan:
    # replace by using median since there could be outliers
    median_value = dataset[feature].median()

    # create a new feature to capture NaN values
    dataset[feature + '_NaN'] = np.where(dataset[feature].isnull(), 1, 0)
    dataset[feature].fillna(median_value, inplace = True)

dataset[num_features_nan].isnull().sum()
```

Out[423]: LotFrontage 0
MasVnrArea 0
GarageYrBlt 0
dtype: int64

TEMPORAL VARIABLES (Eg: Datetime Variables)

From the Dataset we have 4 year variables. We have to extract information from the datetime variables like number of years or number of days. One example in this specific scenario can be the difference between the year the house was built and the year the house was sold.

In [424]...

```
# find list of all features containing year data
year_feature = []
for feature in numerical_features:
    if 'Yr' in feature or 'Year' in feature:
        year_feature.append(feature)

year_feature
```

Out[424]: ['YearBuilt', 'YearRemodAdd', 'GarageYrBlt', 'YrSold']

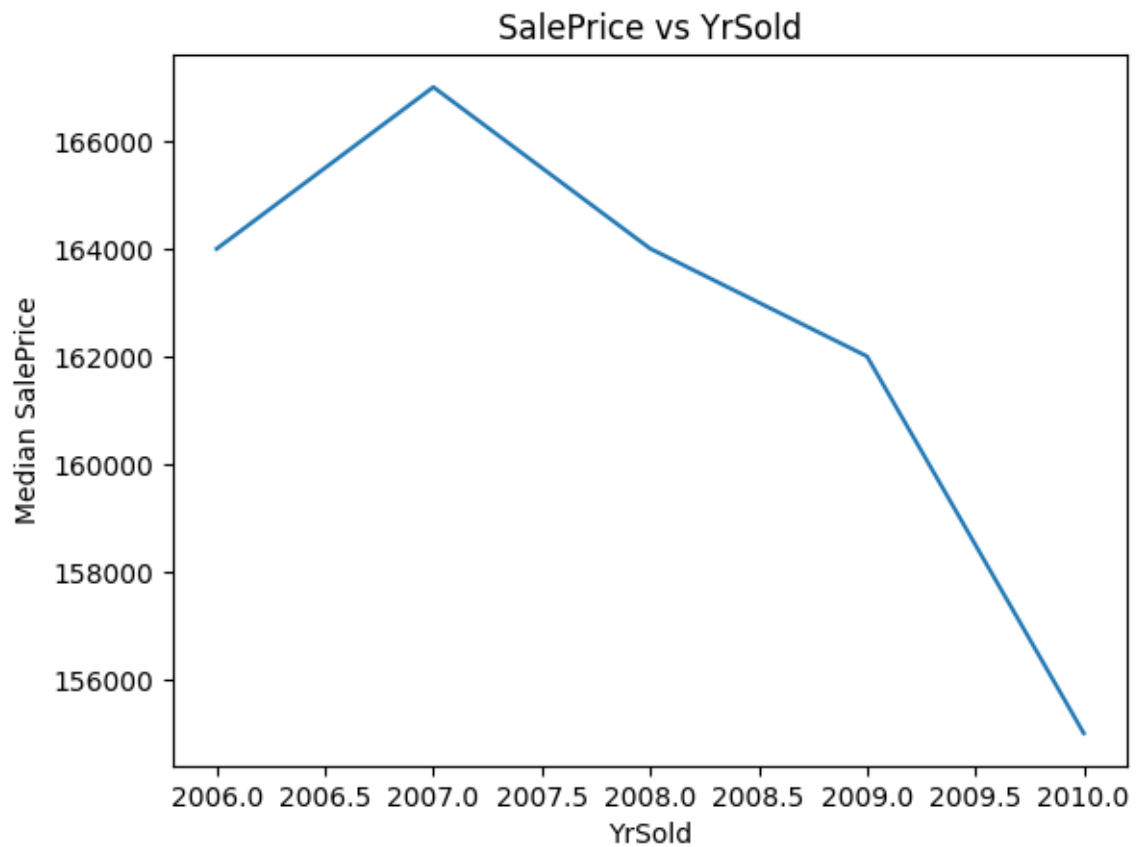
In [425]...

```
# check whether there is a relation between year the house is sold and the sales
print(dataset.groupby('YrSold')['SalePrice'].median().to_frame())

dataset.groupby('YrSold')['SalePrice'].median().plot()
plt.xlabel('YrSold')
plt.ylabel('Median SalePrice')
plt.title("SalePrice vs YrSold")
```

	SalePrice
YrSold	
2006	163995.0
2007	167000.0
2008	164000.0
2009	162000.0
2010	155000.0

Out[425]: Text(0.5, 1.0, 'SalePrice vs YrSold')

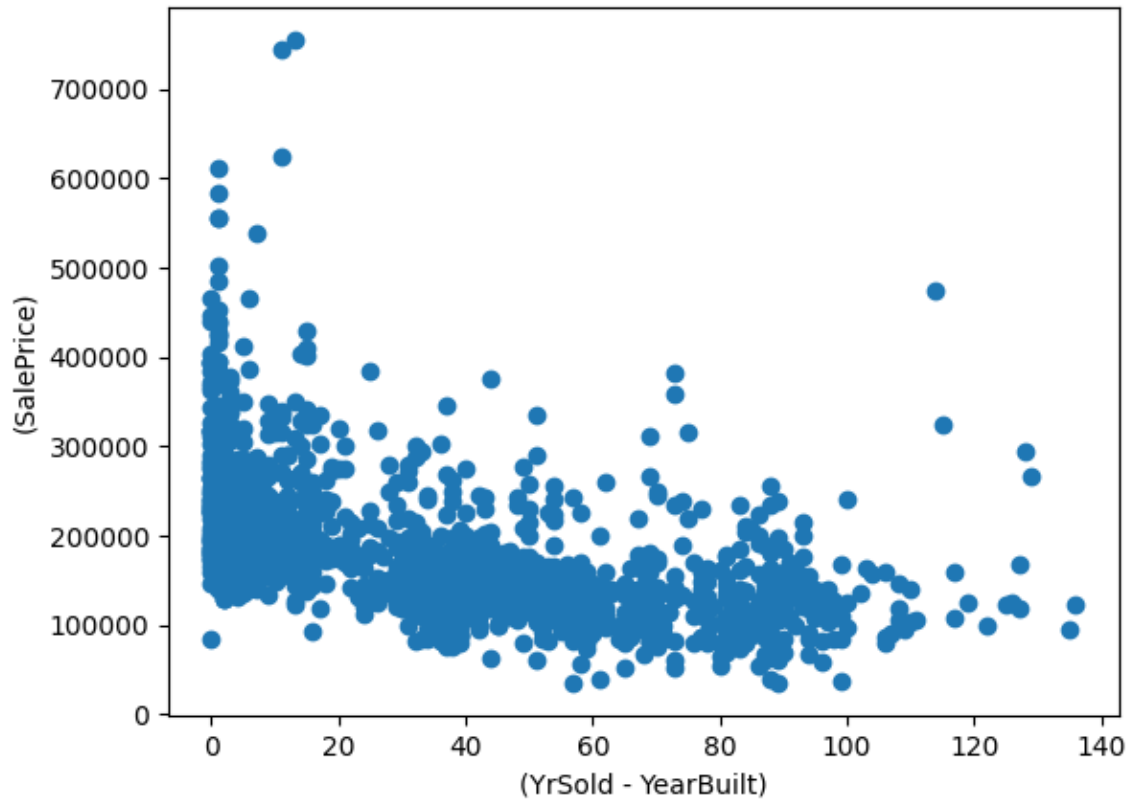


As YrSold increases, the price drops.

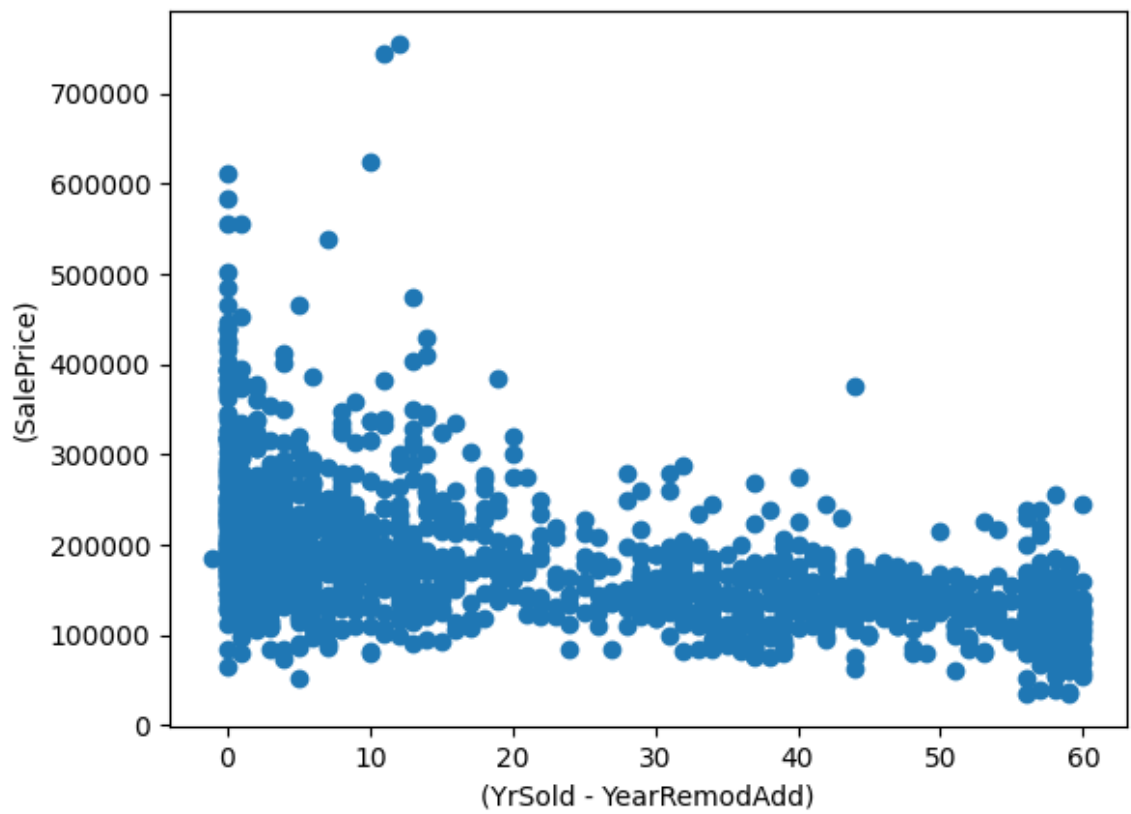
```
In [426... # compare the (['YrSold'] - rest of the features containing year data) with (SalePrice)
for feature in year_feature:
    if feature != 'YrSold':
        data = dataset.copy()
        data[feature] = data['YrSold'] - data[feature]

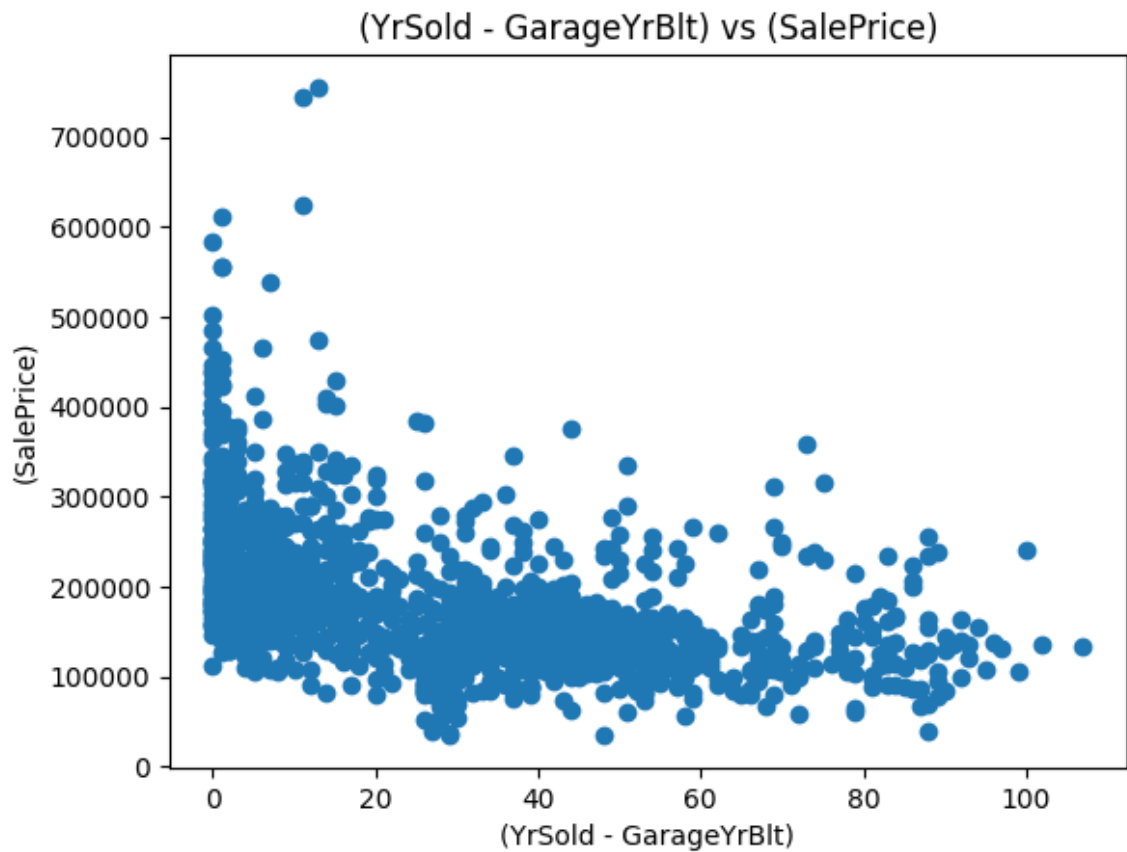
        plt.scatter(data[feature], data['SalePrice'])
        plt.xlabel('(YrSold - ' + feature + ')')
        plt.ylabel('(SalePrice)')
        plt.title('(YrSold - ' + feature + ') vs (SalePrice)')
        plt.show()
```

(YrSold - YearBuilt) vs (SalePrice)



(YrSold - YearRemodAdd) vs (SalePrice)





```
In [427... for feature in ['YearBuilt', 'YearRemodAdd', 'GarageYrBlt']:
             dataset[feature] = dataset['YrSold'] - dataset[feature] # gives number of yr
```

```
In [428... dataset[['YearBuilt', 'YearRemodAdd', 'GarageYrBlt']].head()
```

```
Out[428]:
```

	YearBuilt	YearRemodAdd	GarageYrBlt
0	5	5	5.0
1	31	31	31.0
2	7	6	7.0
3	91	36	8.0
4	8	8	8.0

DISCRETE VARIABLES

```
In [429... # find all discrete features
discrete_feature = []
for feature in numerical_features:
    if len(dataset[feature].unique()) < 25 and feature not in year_feature + ['Id']:
        discrete_feature.append(feature)

len(discrete_feature)
```

```
Out[429]: 17
```

```
In [430... discrete_feature
```

```
Out[430]: ['MSSubClass',
           'OverallQual',
           'OverallCond',
           'LowQualFinSF',
           'BsmtFullBath',
           'BsmtHalfBath',
           'FullBath',
           'HalfBath',
           'BedroomAbvGr',
           'KitchenAbvGr',
           'TotRmsAbvGrd',
           'Fireplaces',
           'GarageCars',
           '3SsnPorch',
           'PoolArea',
           'MiscVal',
           'MoSold']
```

```
In [431... dataset[discrete_feature].head()
```

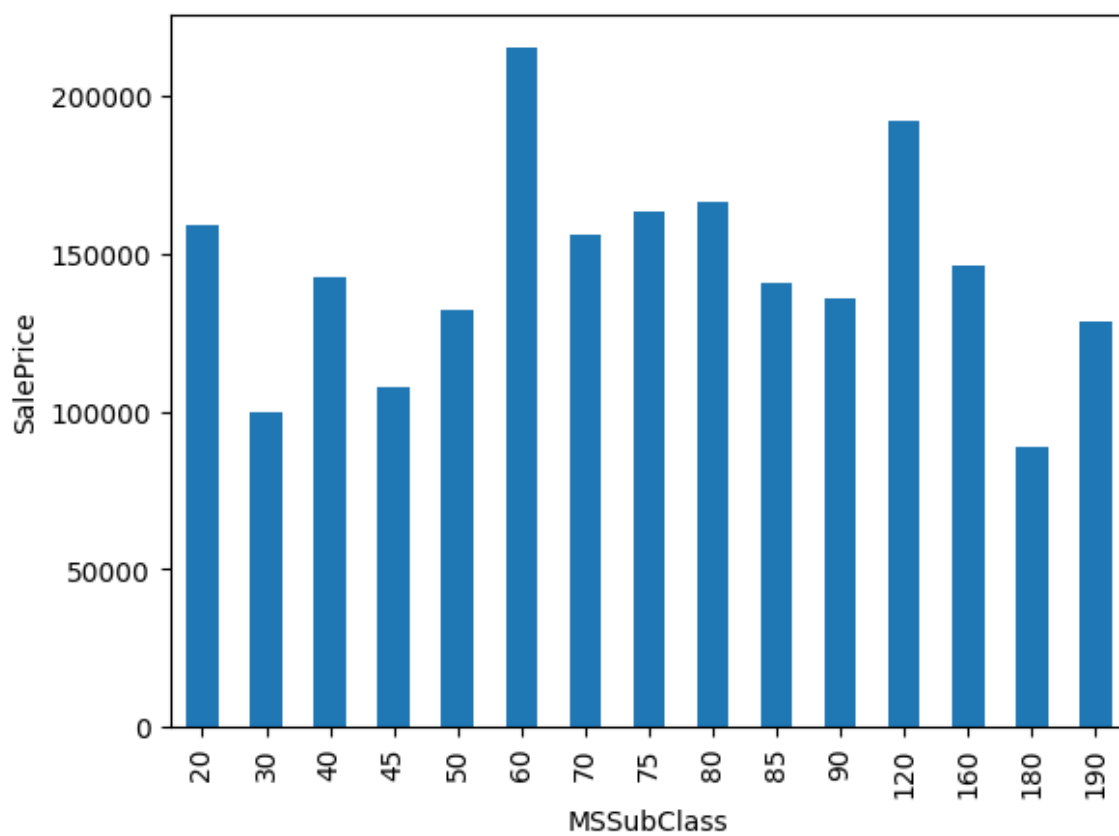
```
Out[431]:
```

	MSSubClass	OverallQual	OverallCond	LowQualFinSF	BsmtFullBath	BsmtHalfBath	FullBath
0	60	7	5	0	1	0	
1	20	6	8	0	0	1	
2	60	7	5	0	1	0	
3	70	7	5	0	1	0	
4	60	8	5	0	1	0	

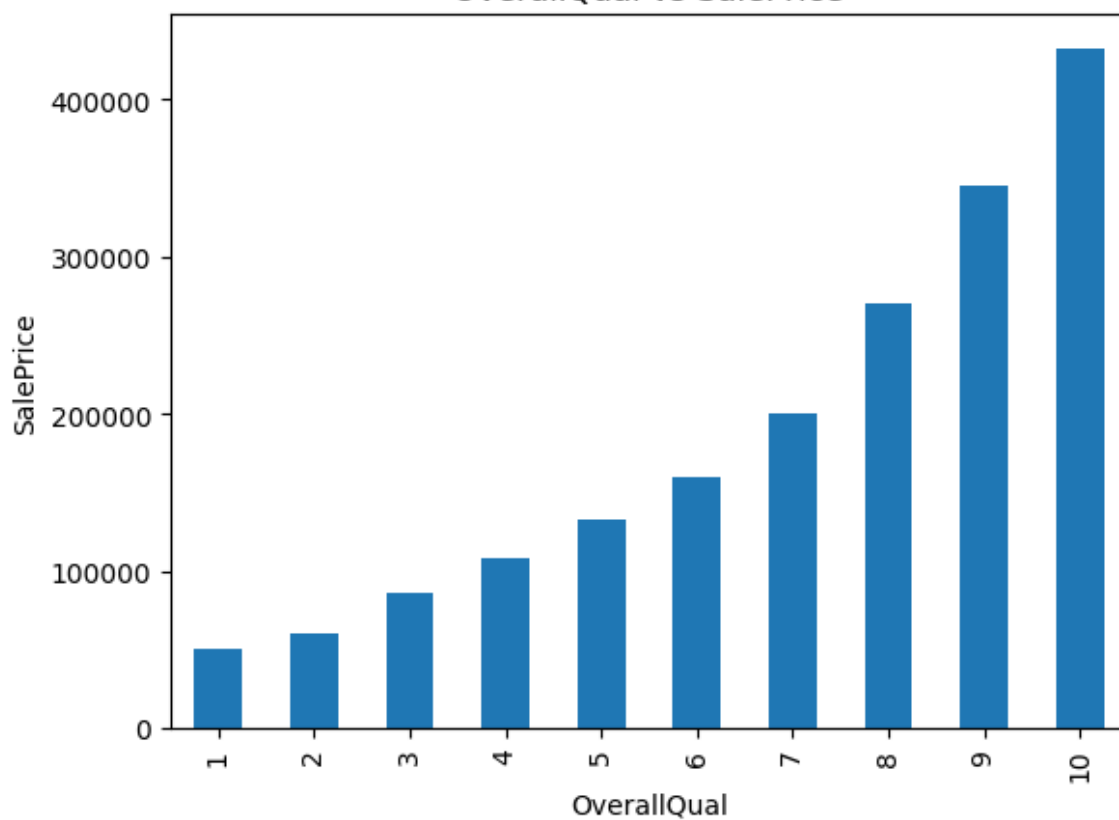
```
In [432... for feature in discrete_feature:
            data = dataset.copy()
            data.groupby(feature)['SalePrice'].median().plot.bar()

            plt.xlabel(feature)
            plt.ylabel('SalePrice')
            plt.title(feature + ' vs SalePrice')
            plt.show()
```

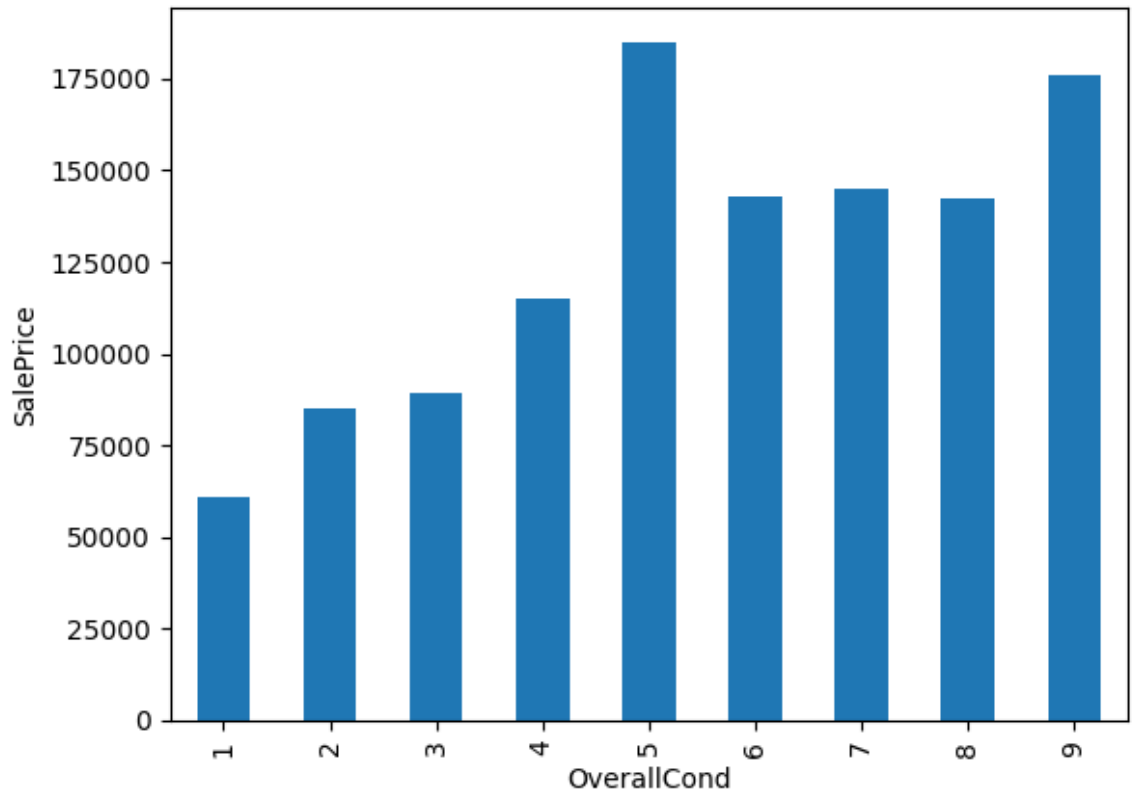
MSSubClass vs SalePrice



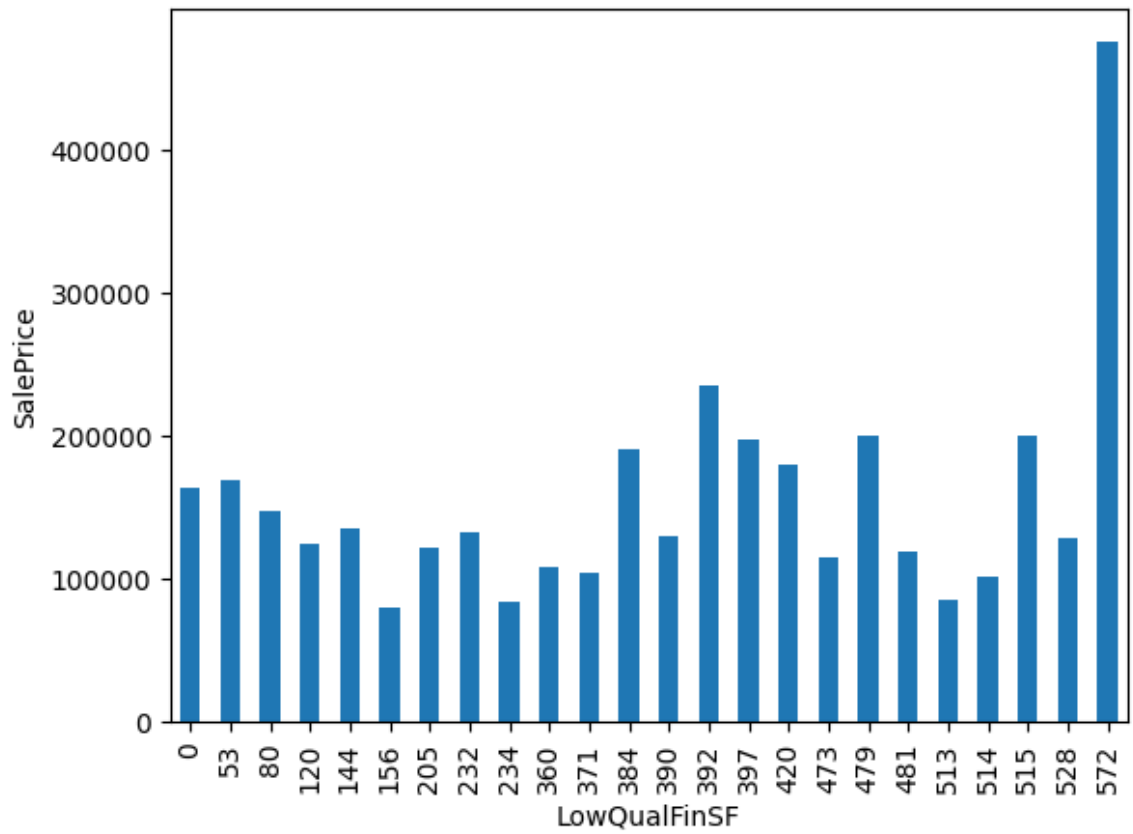
OverallQual vs SalePrice



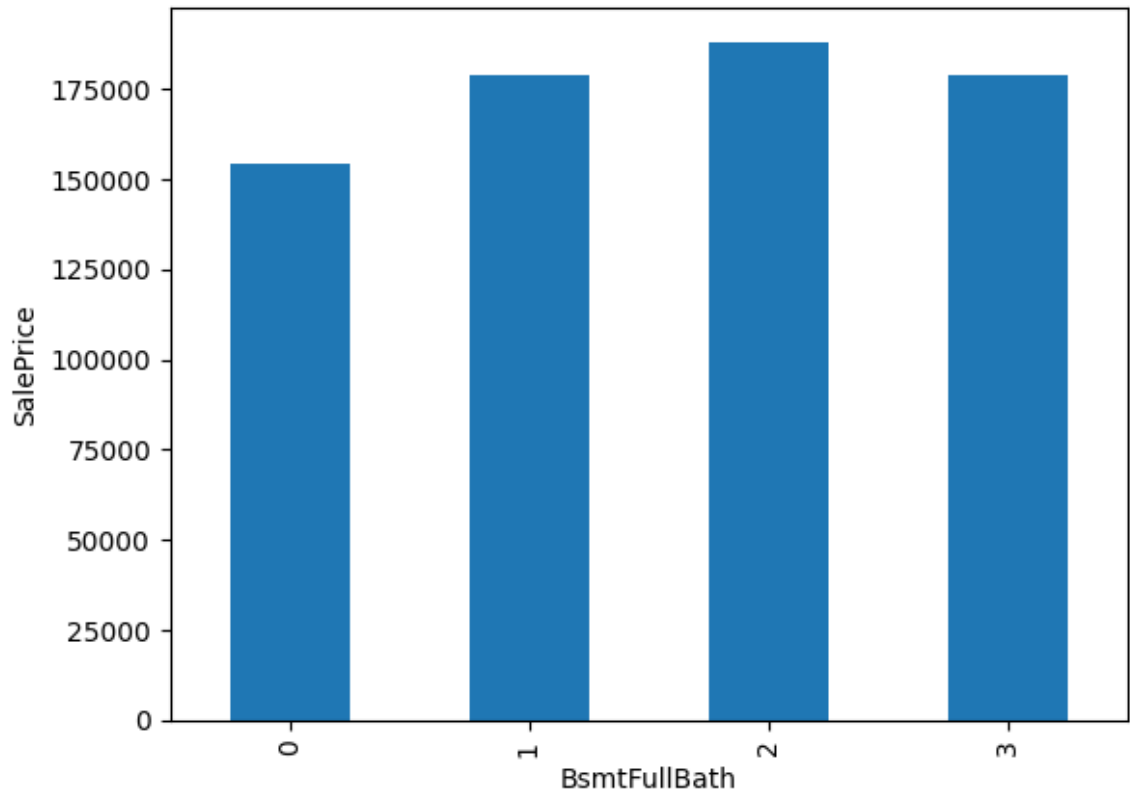
OverallCond vs SalePrice



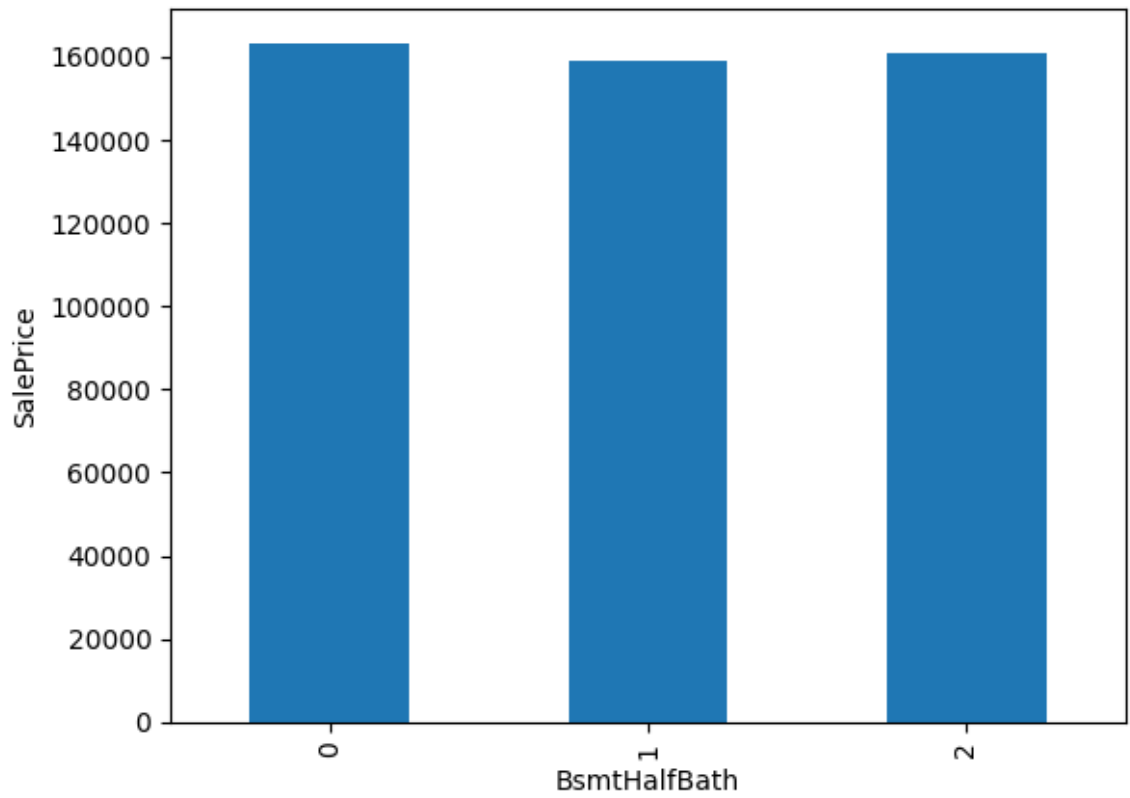
LowQualFinSF vs SalePrice

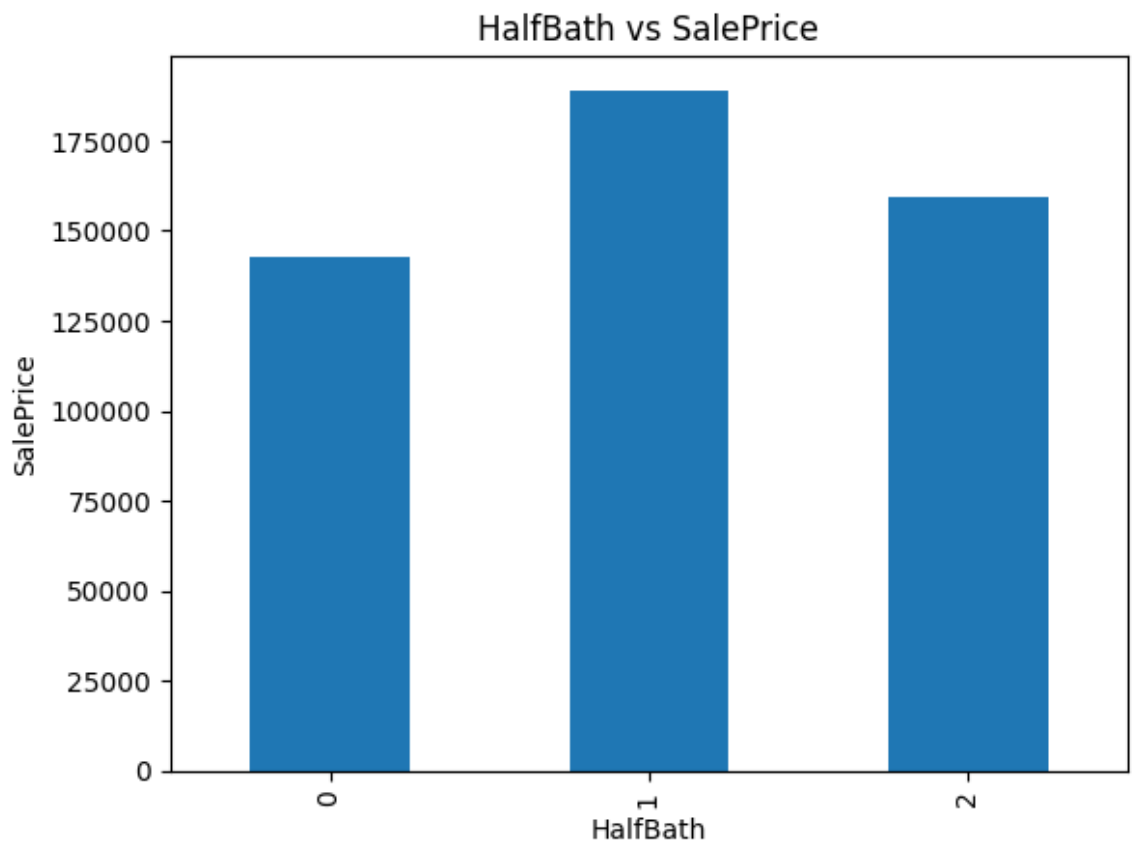
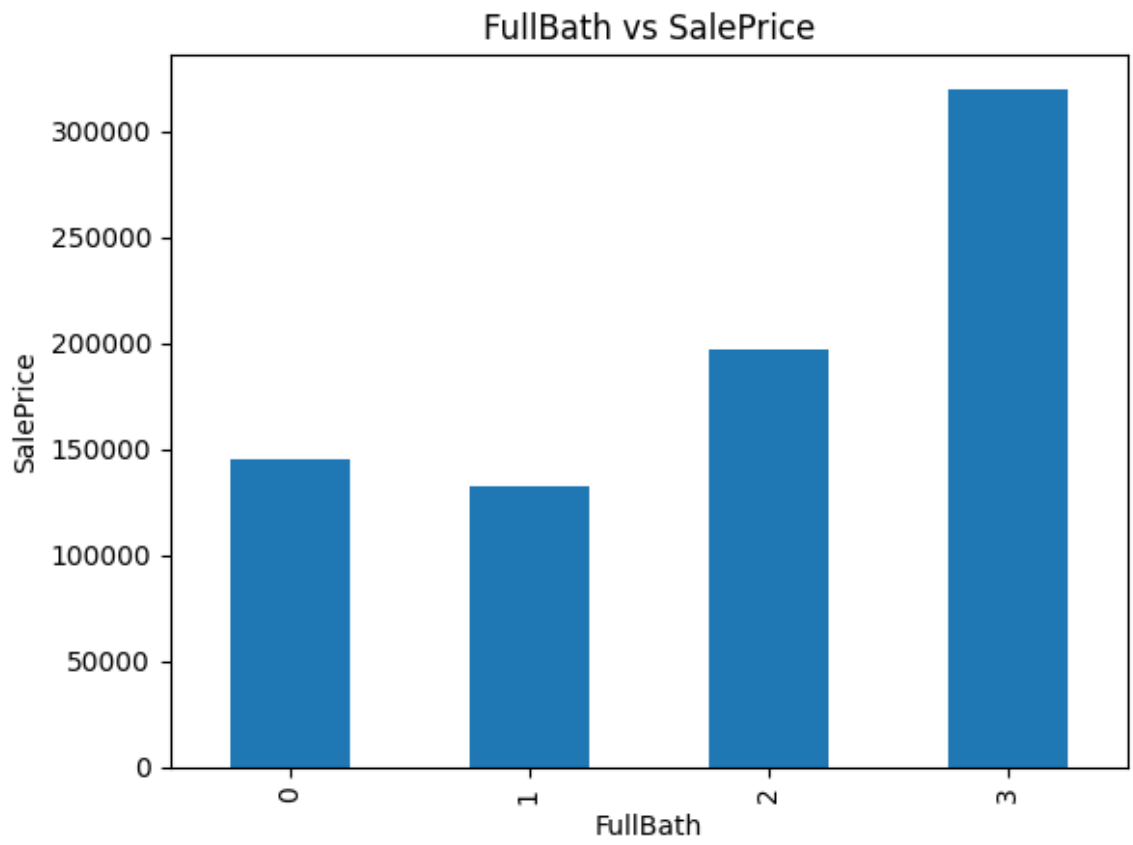


BsmtFullBath vs SalePrice

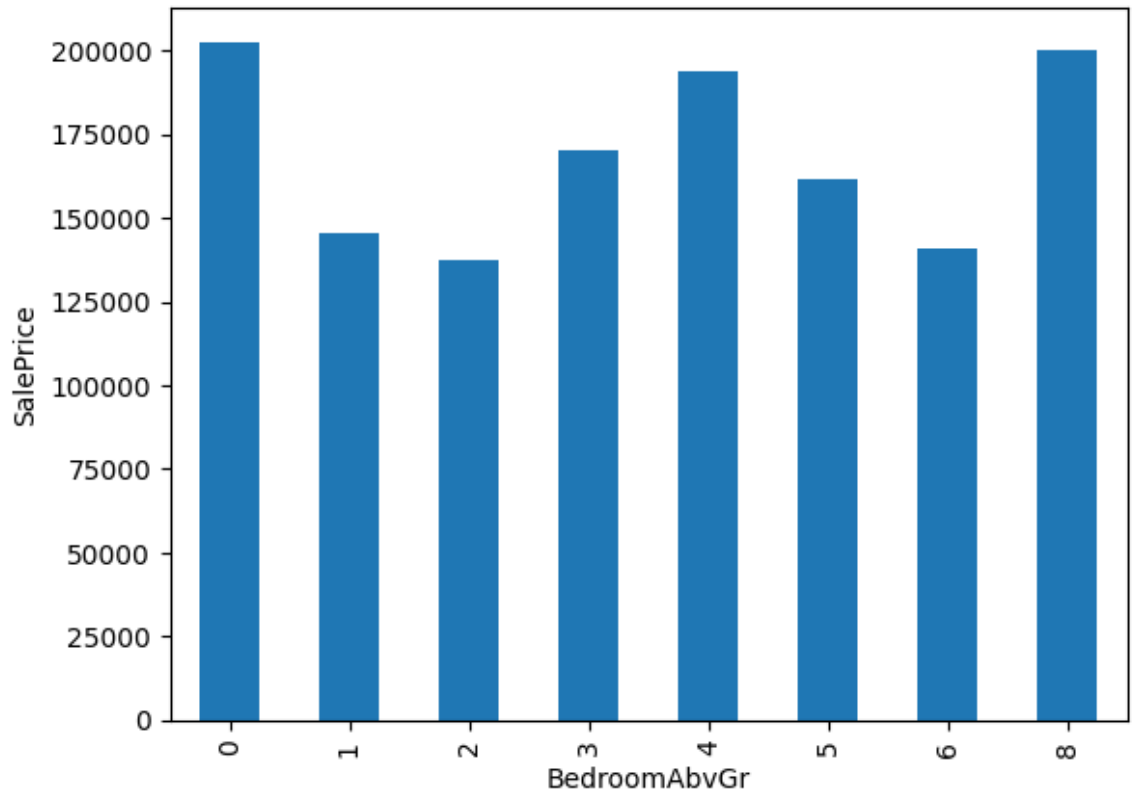


BsmtHalfBath vs SalePrice

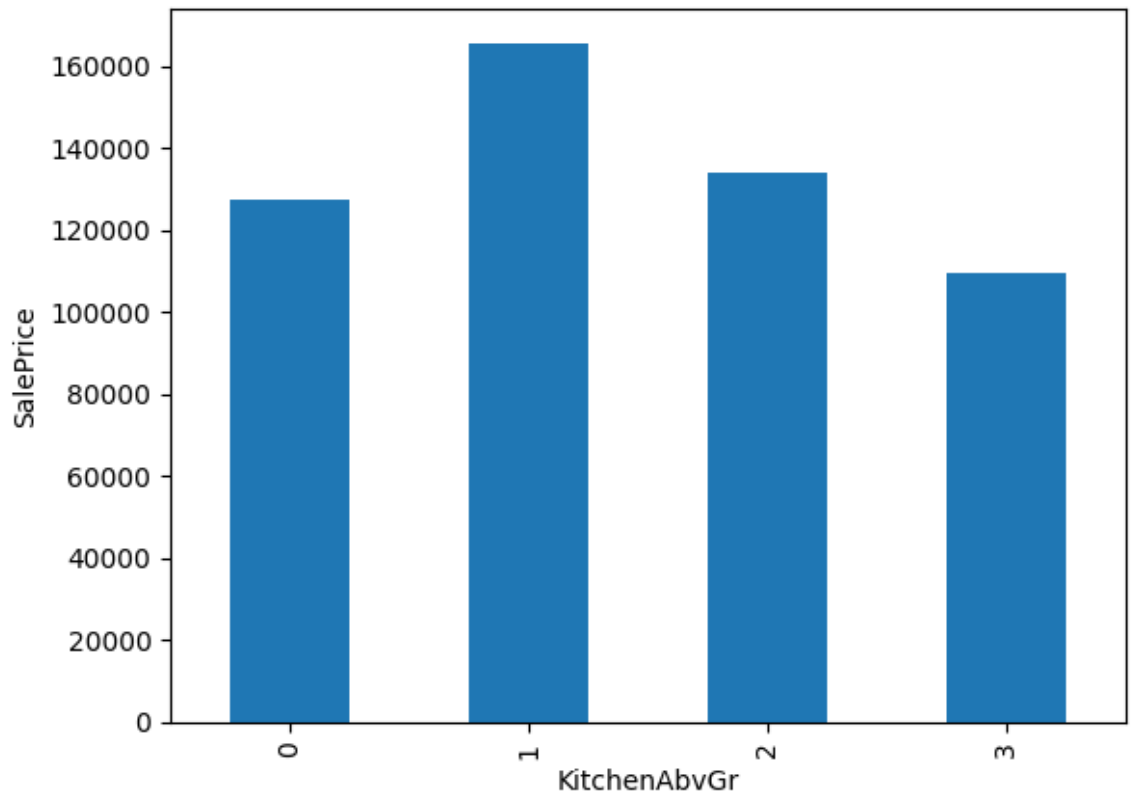




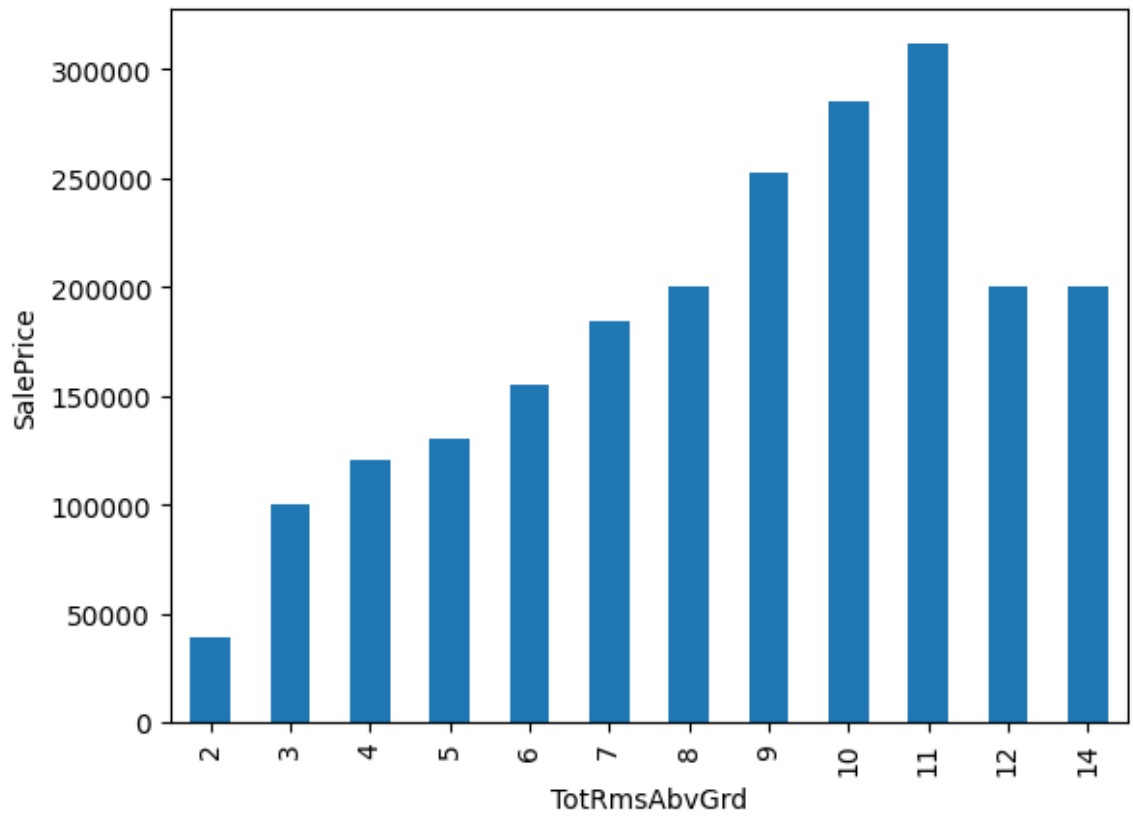
BedroomAbvGr vs SalePrice



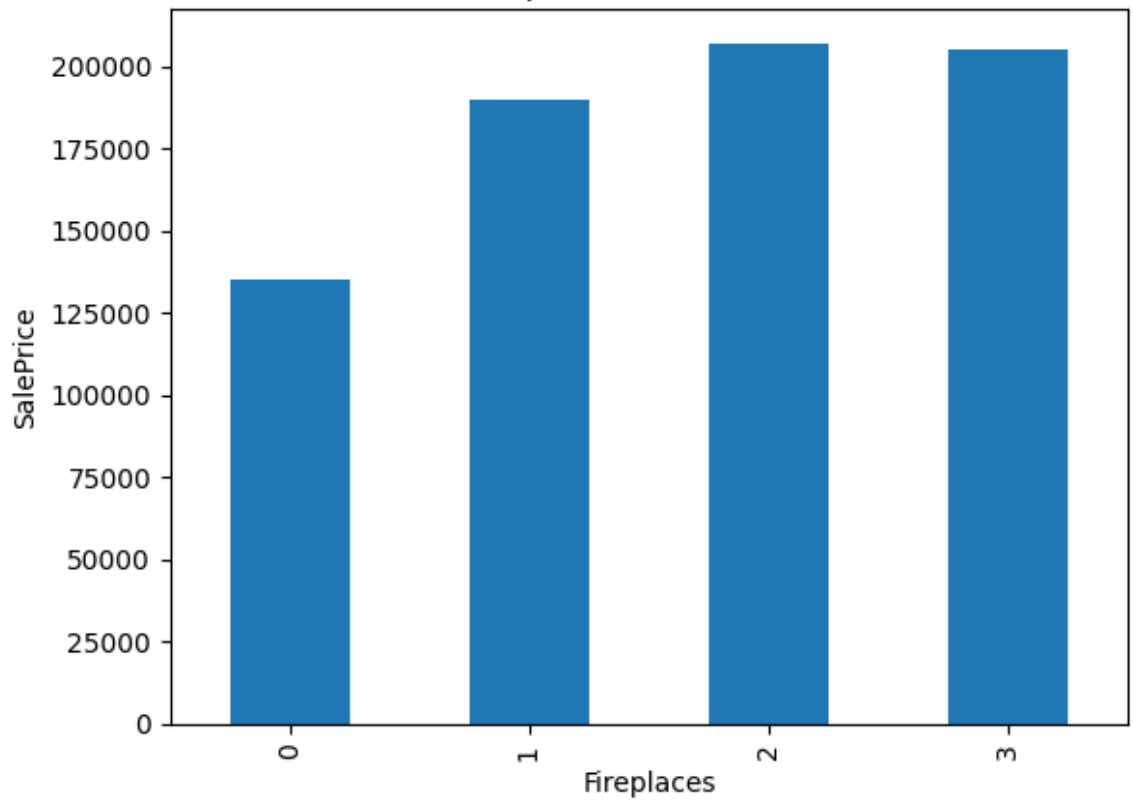
KitchenAbvGr vs SalePrice



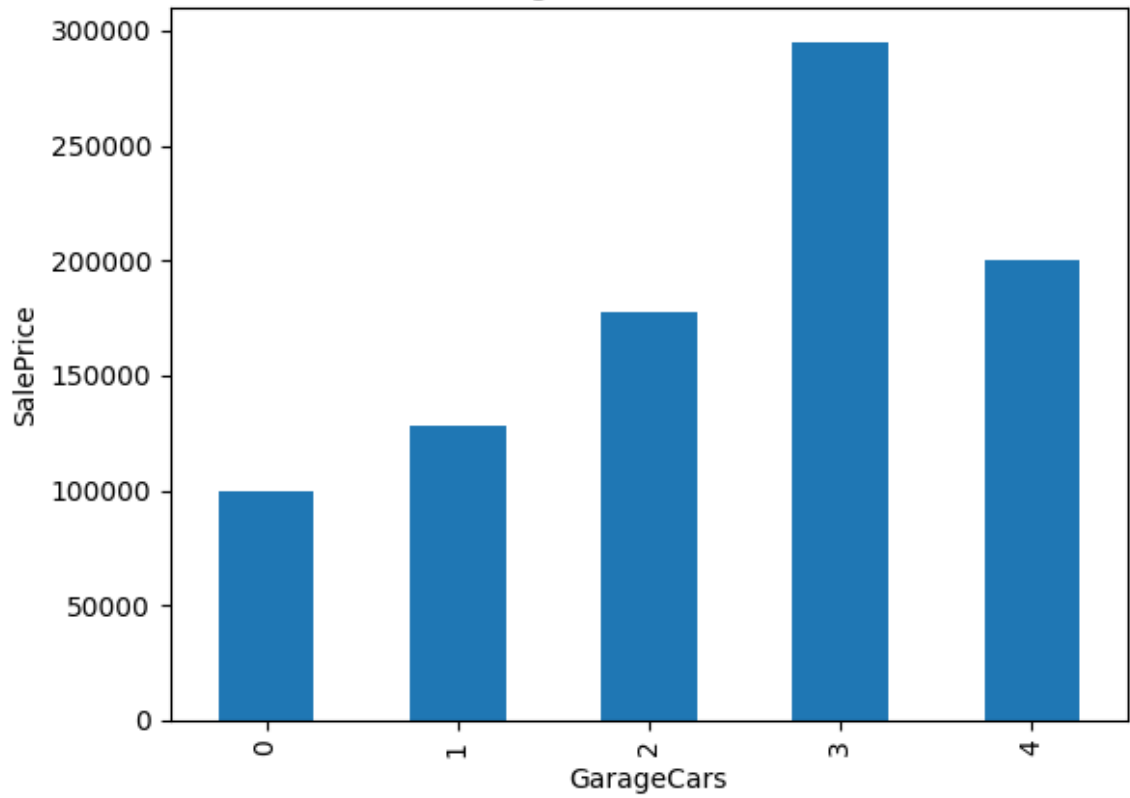
TotRmsAbvGrd vs SalePrice



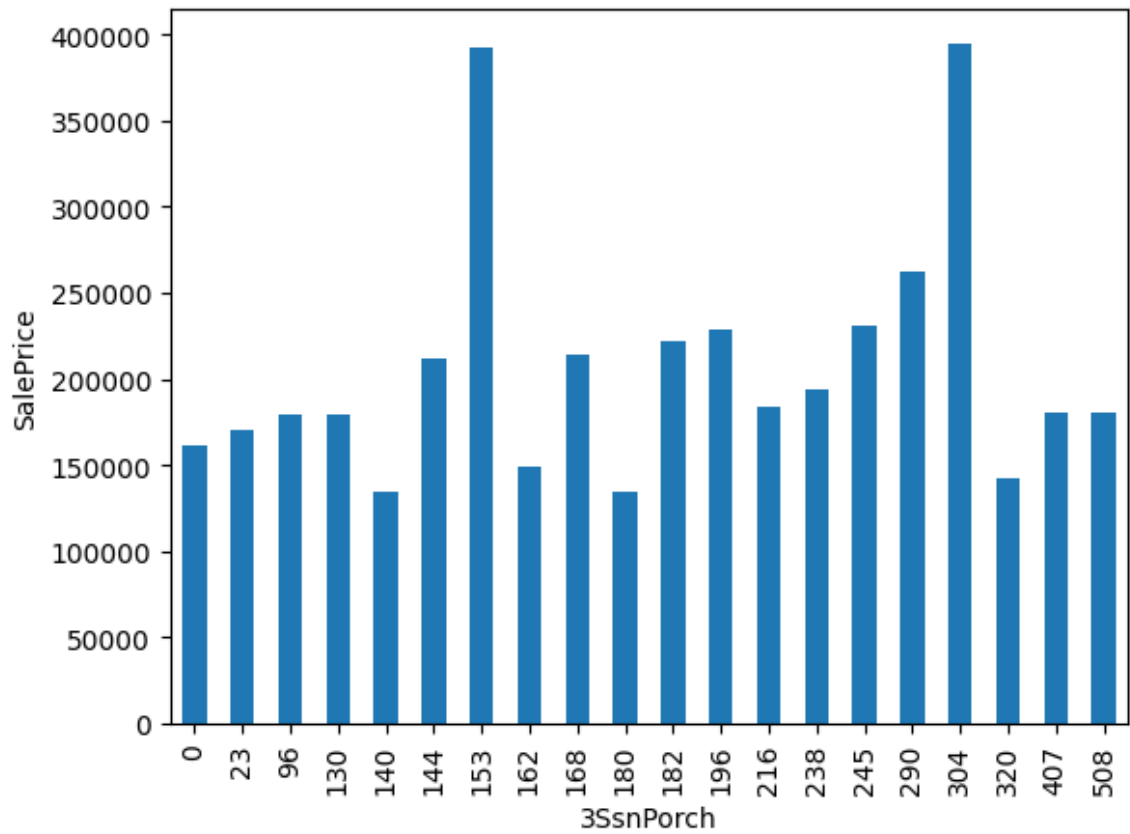
Fireplaces vs SalePrice



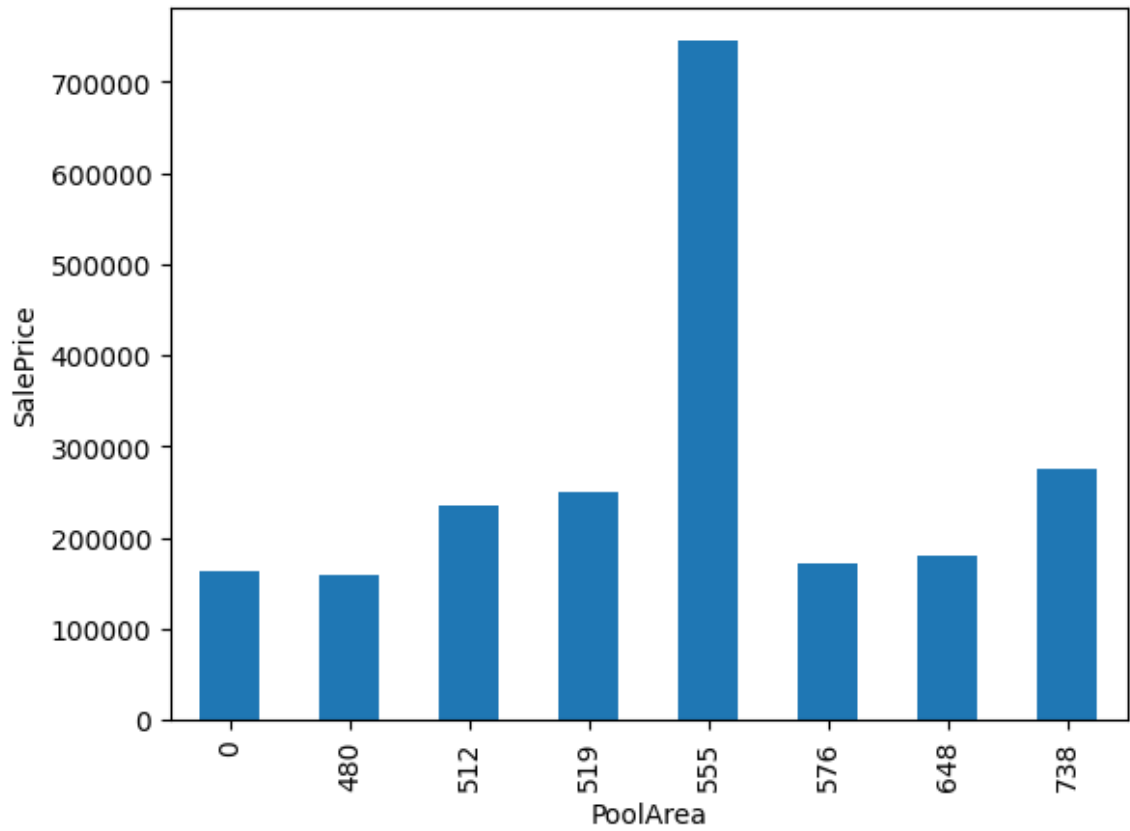
GarageCars vs SalePrice



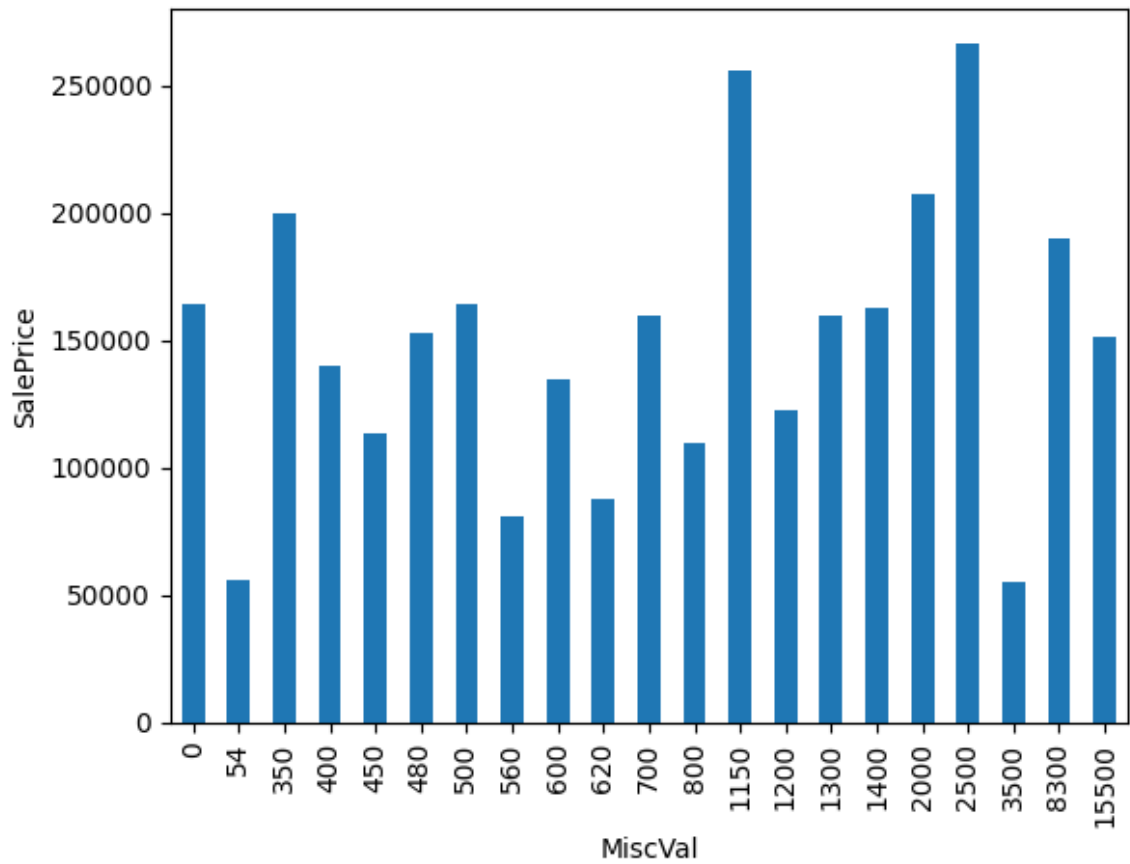
3SsnPorch vs SalePrice

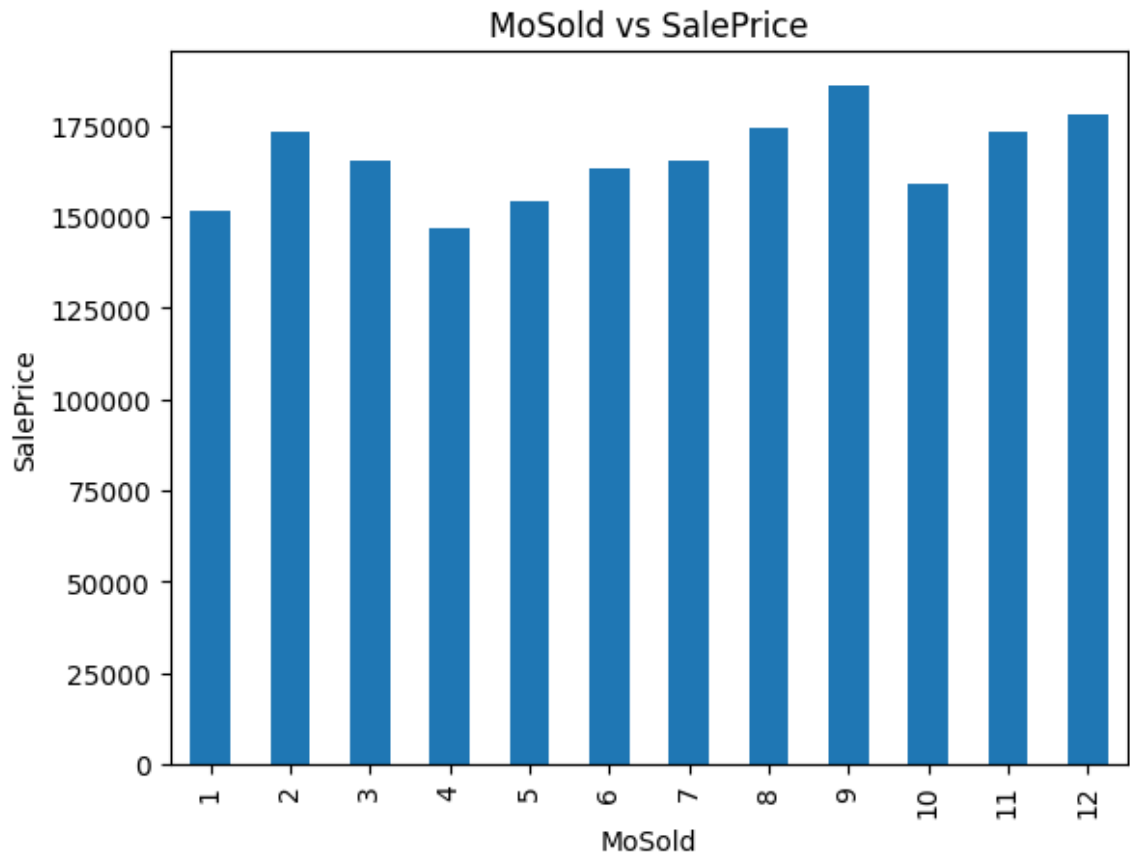


PoolArea vs SalePrice



MiscVal vs SalePrice





CONTINUOUS VARIABLES

```
In [433... # find all continuous features
continuous_feature = []
for feature in numerical_features:
    if feature not in discrete_feature + year_feature + ['Id']:
        continuous_feature.append(feature)

len(continuous_feature)
```

Out[433]: 16

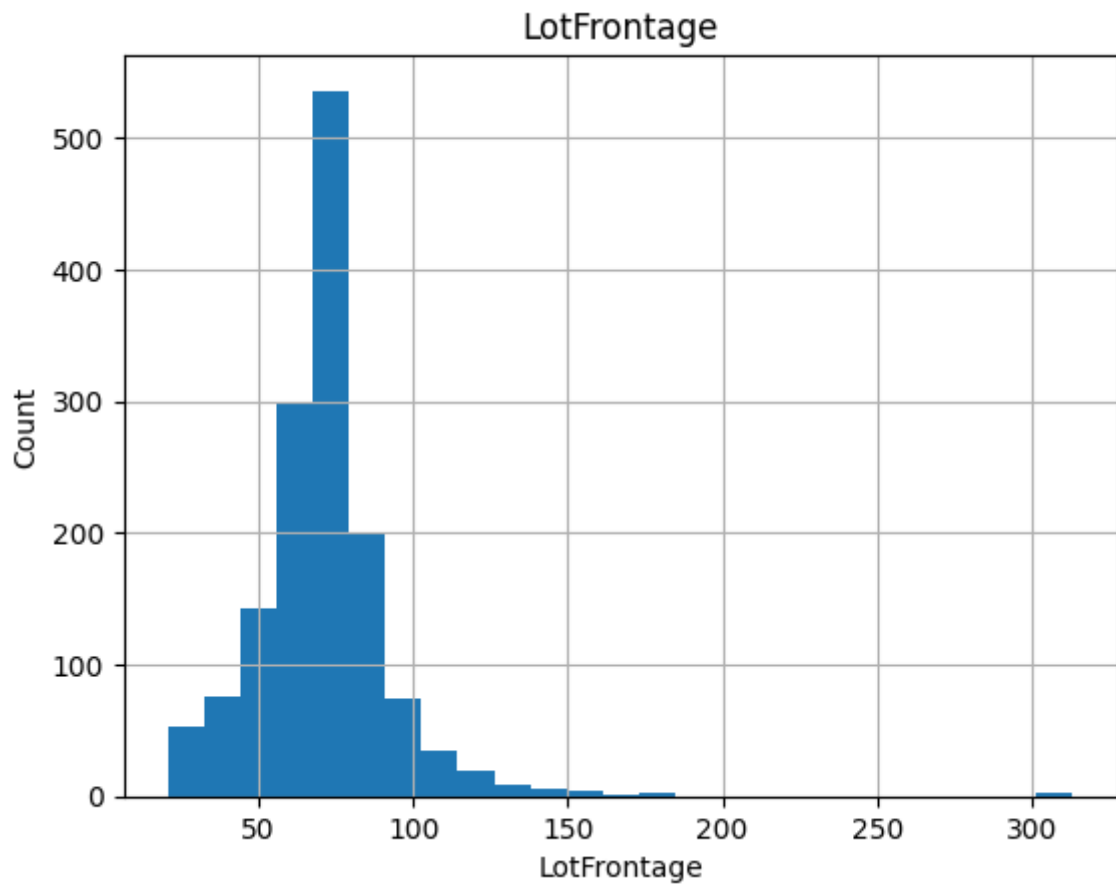
```
In [434... continuous_feature
```

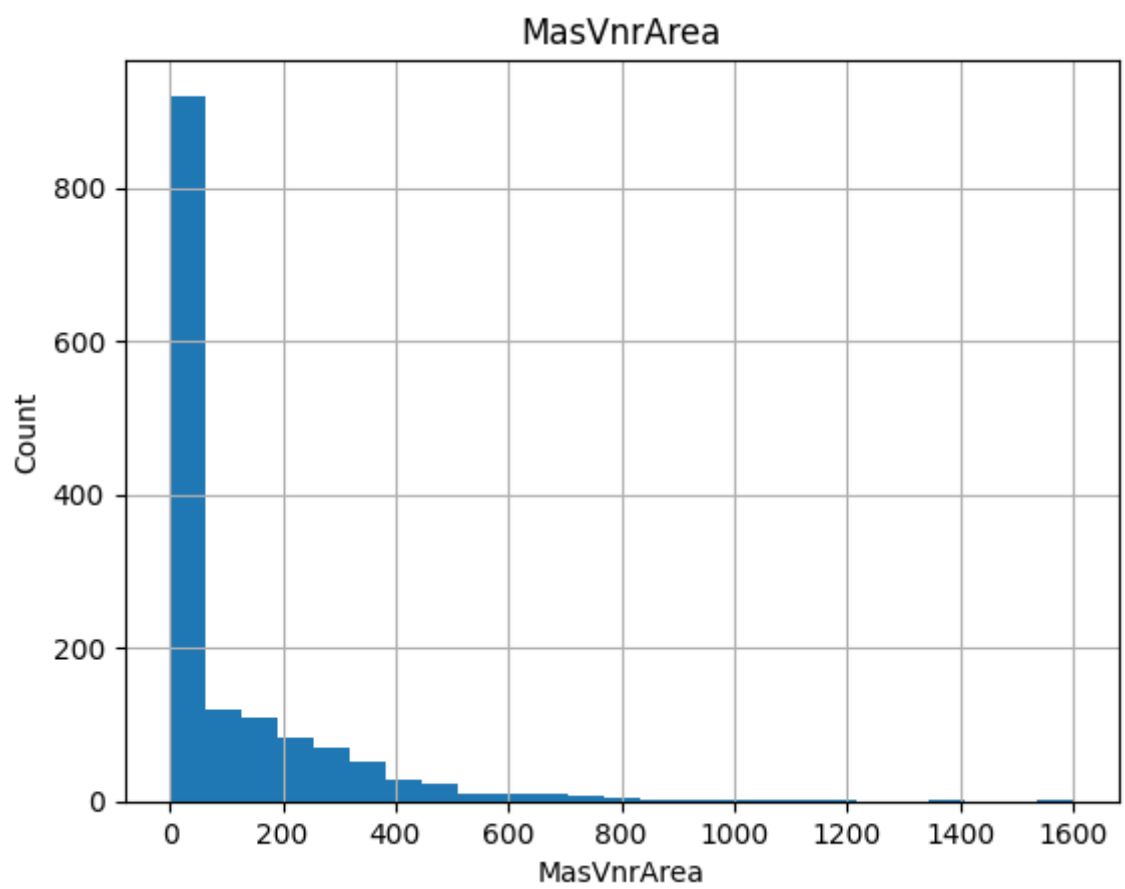
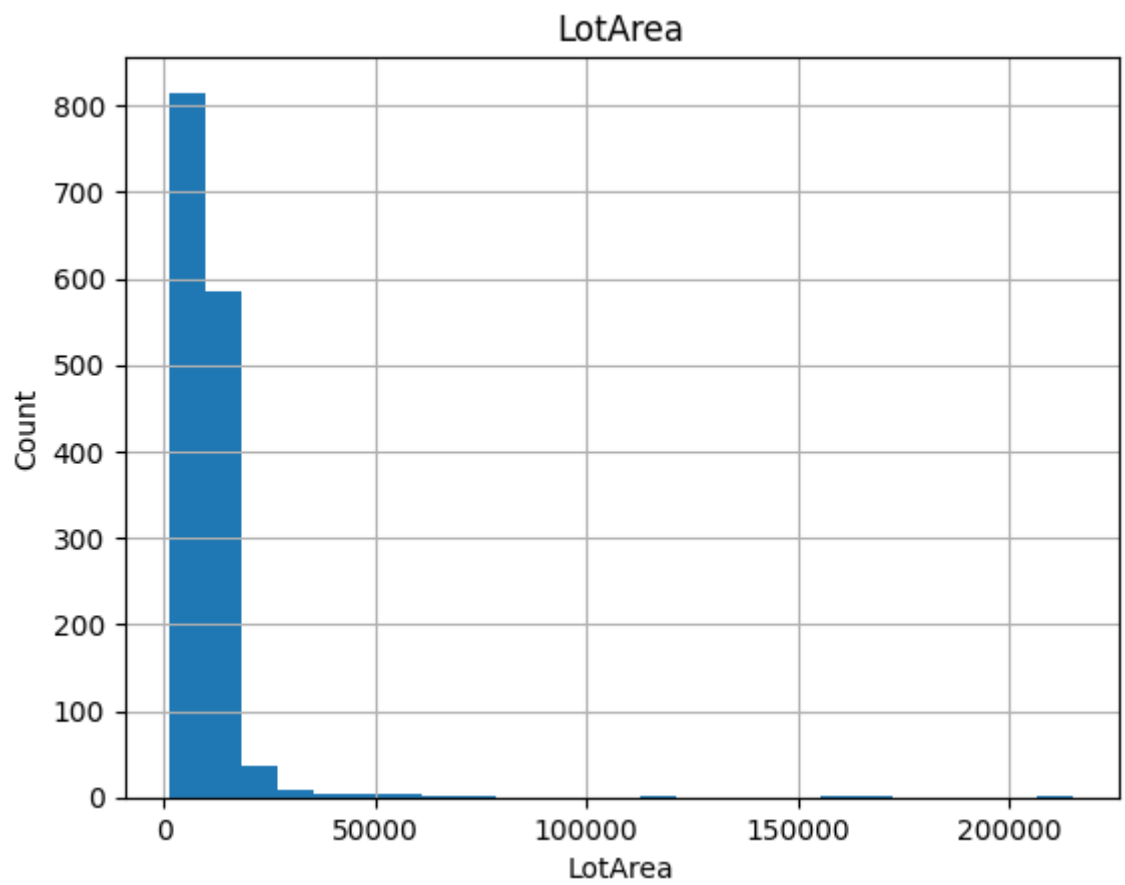
```
Out[434]: ['LotFrontage',
'LotArea',
'MasVnrArea',
'BsmtFinSF1',
'BsmtFinSF2',
'BsmtUnfSF',
'TotalBsmtSF',
'1stFlrSF',
'2ndFlrSF',
'GrLivArea',
'GarageArea',
'WoodDeckSF',
'OpenPorchSF',
'EnclosedPorch',
'ScreenPorch',
'SalePrice']
```

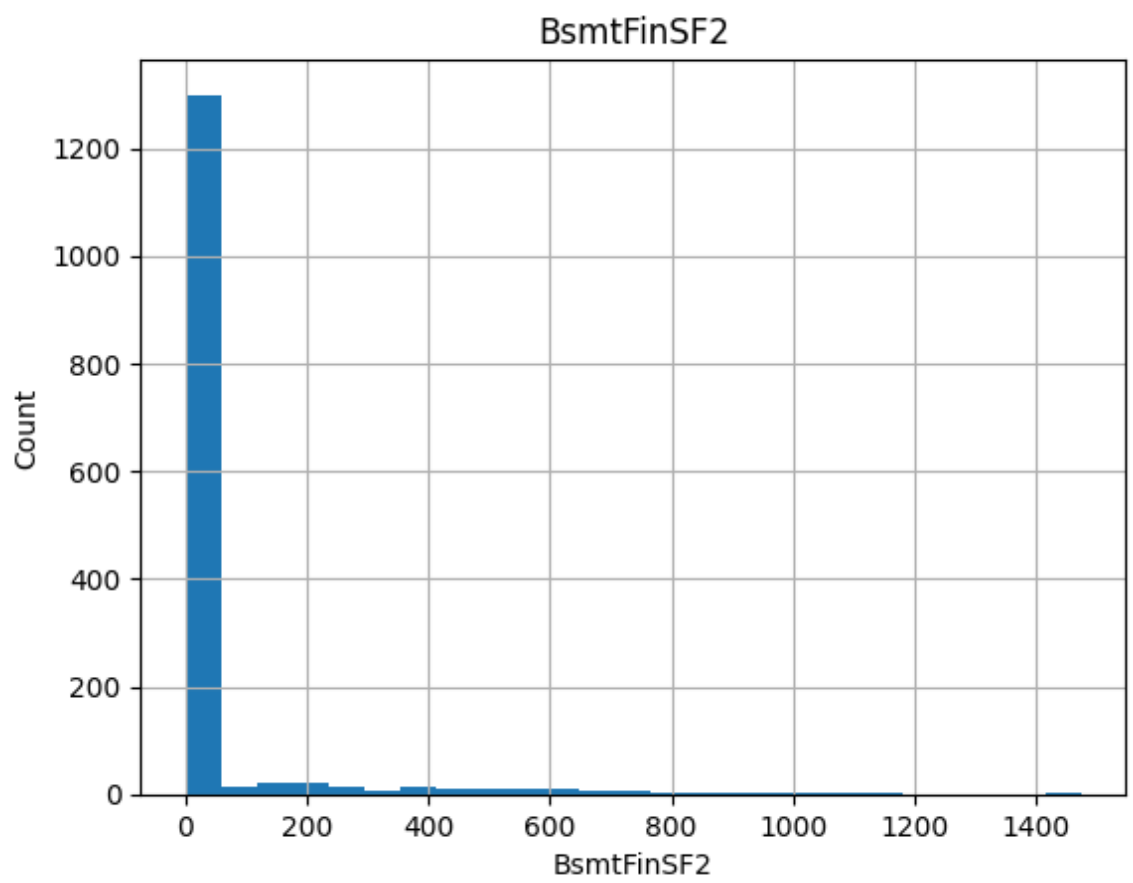
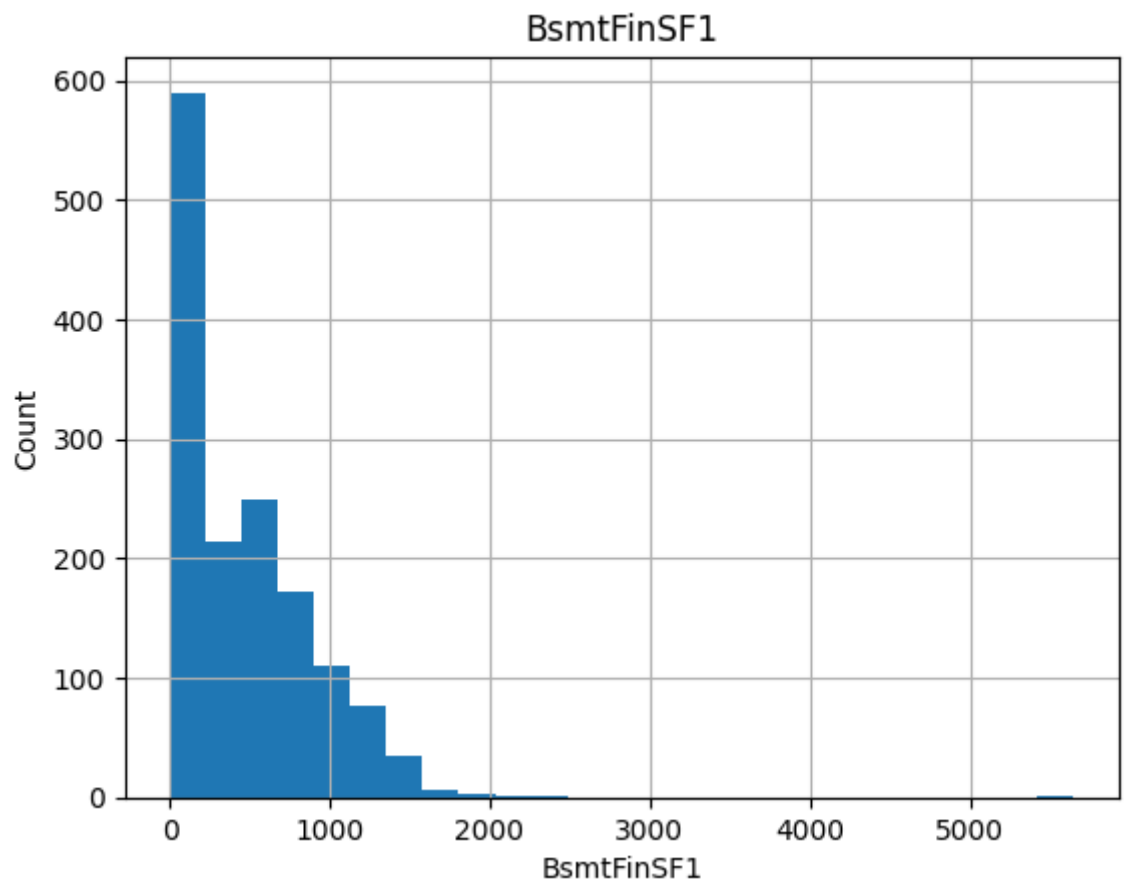
In [435... *# analyse the continuous values by creating histograms to understand the distrib*

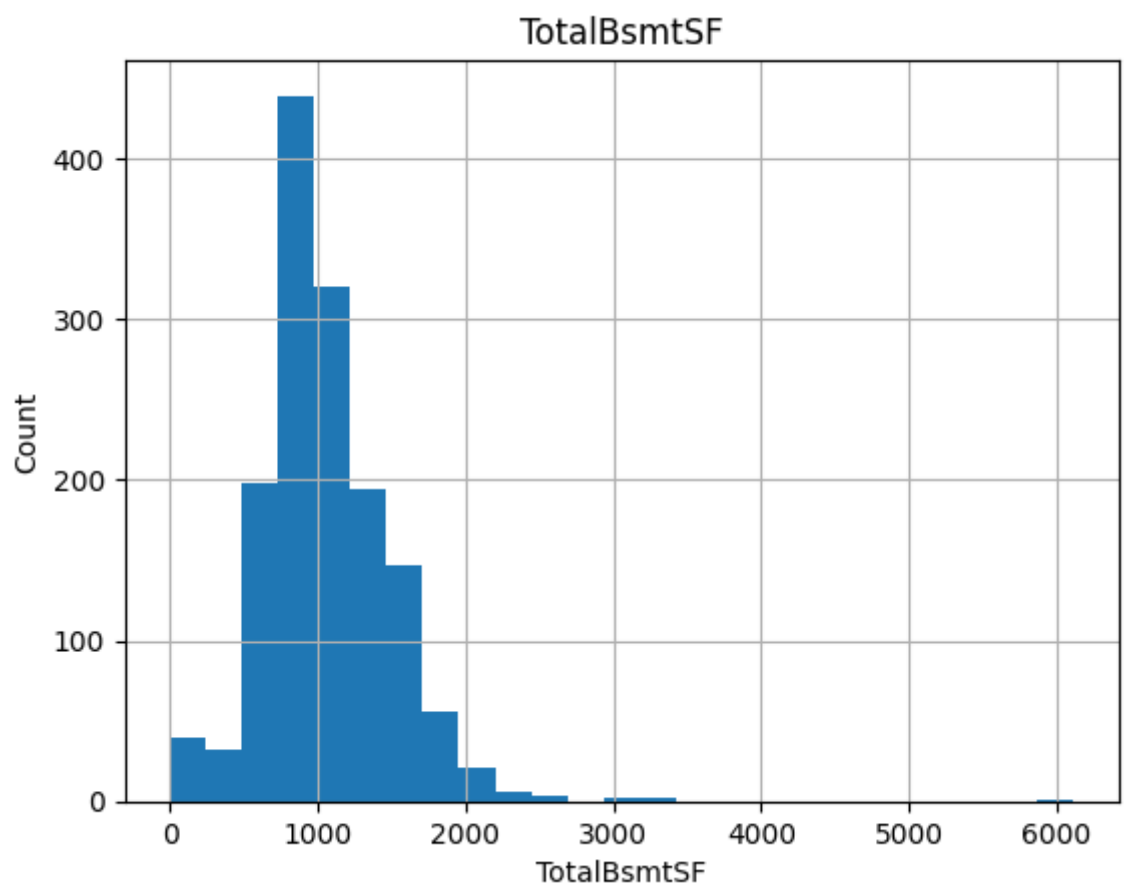
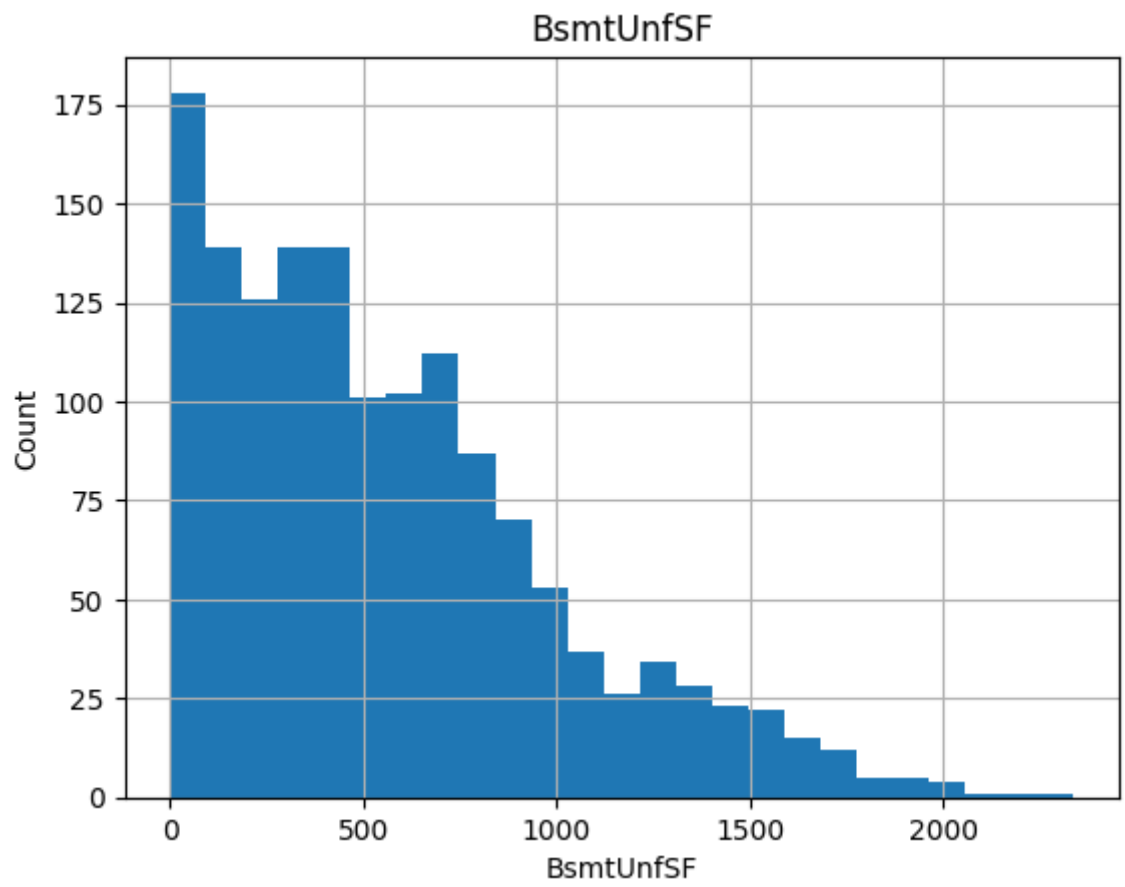
```
for feature in continuous_feature:
    data=dataset.copy()
    data[feature].hist(bins = 25)

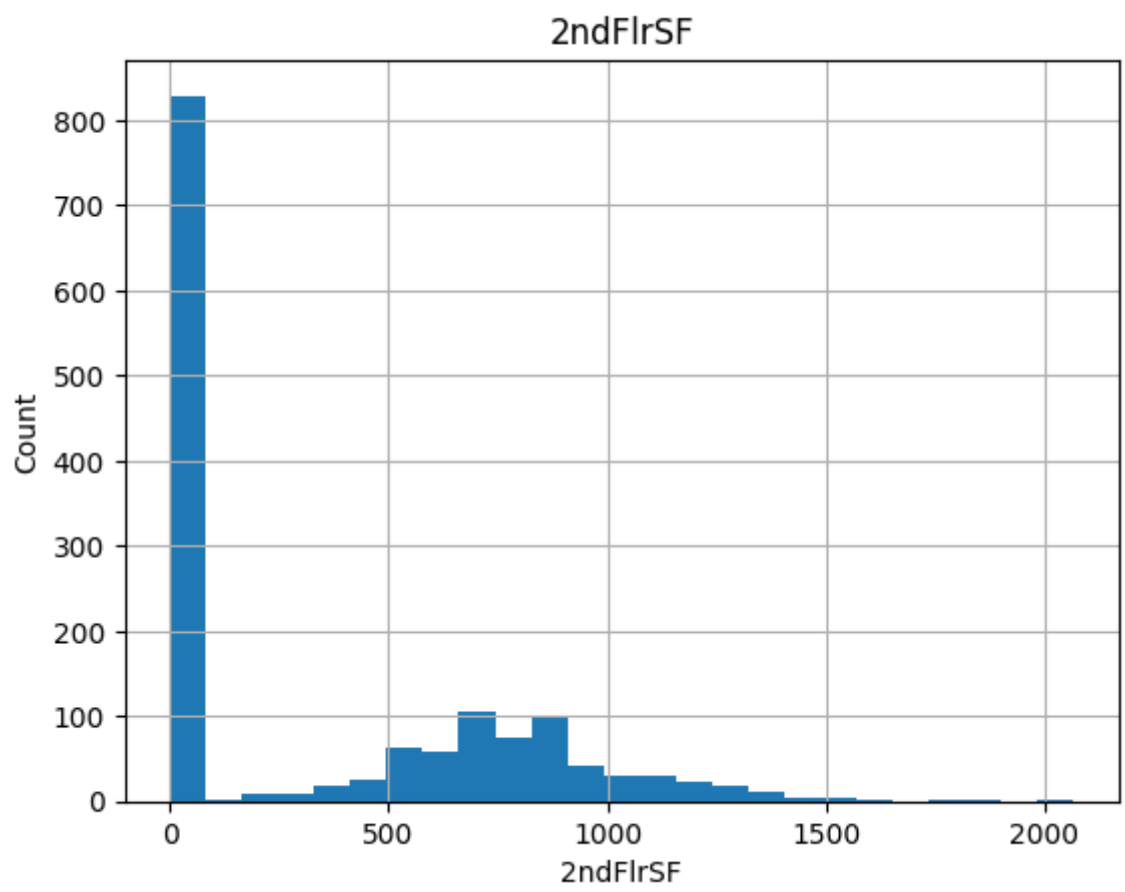
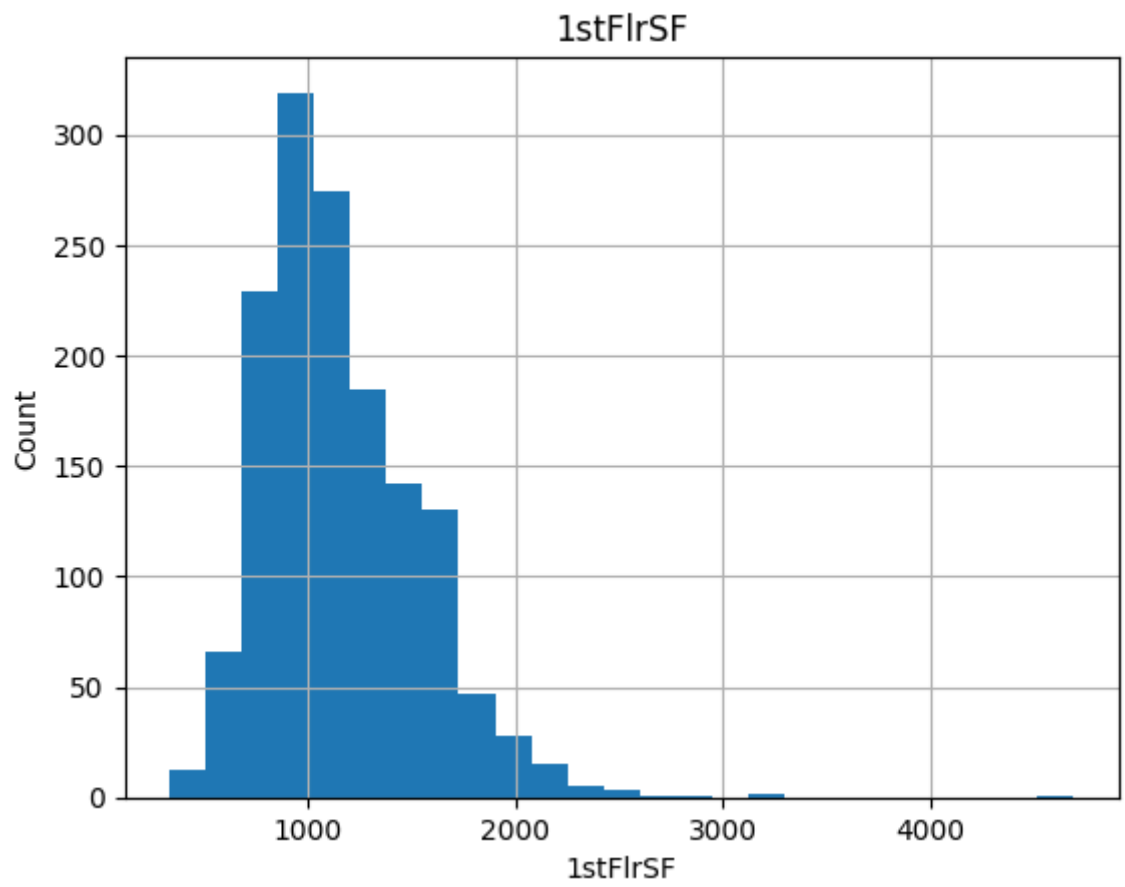
    plt.xlabel(feature)
    plt.ylabel("Count")
    plt.title(feature)
    plt.show()
```

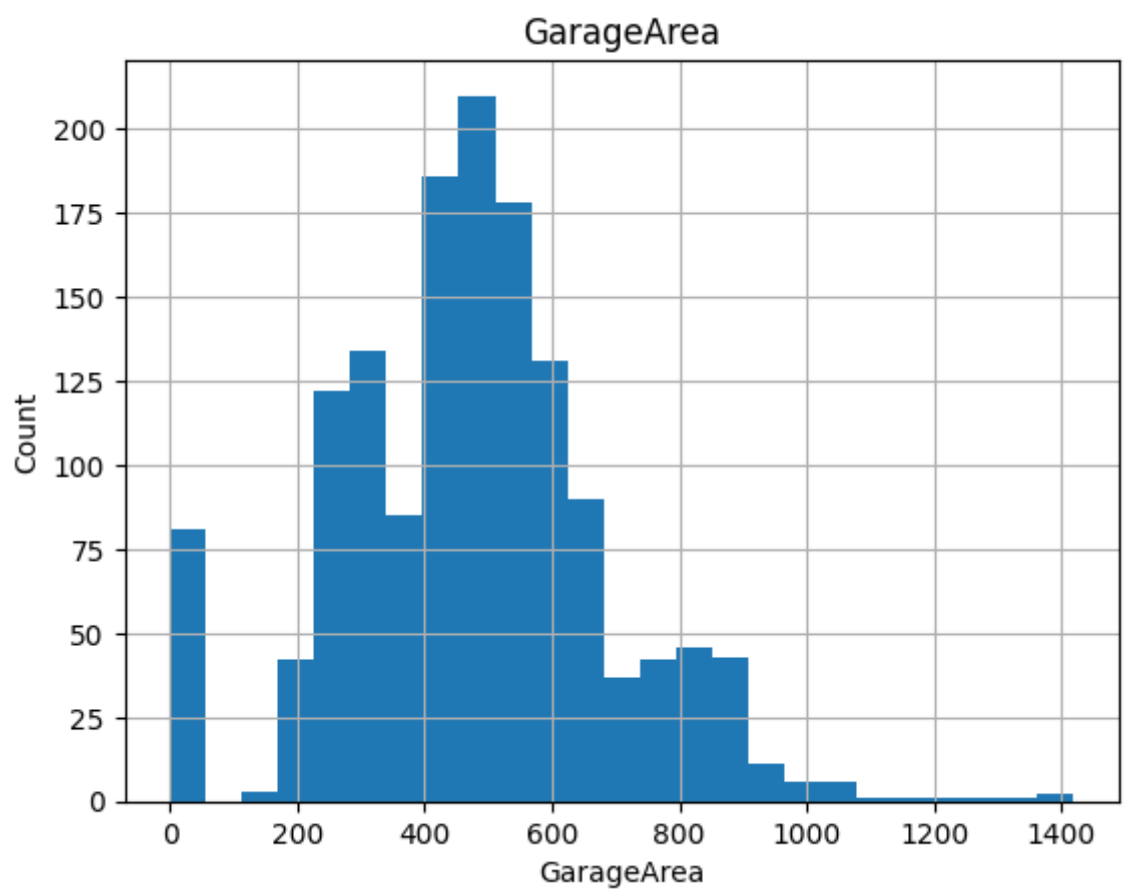
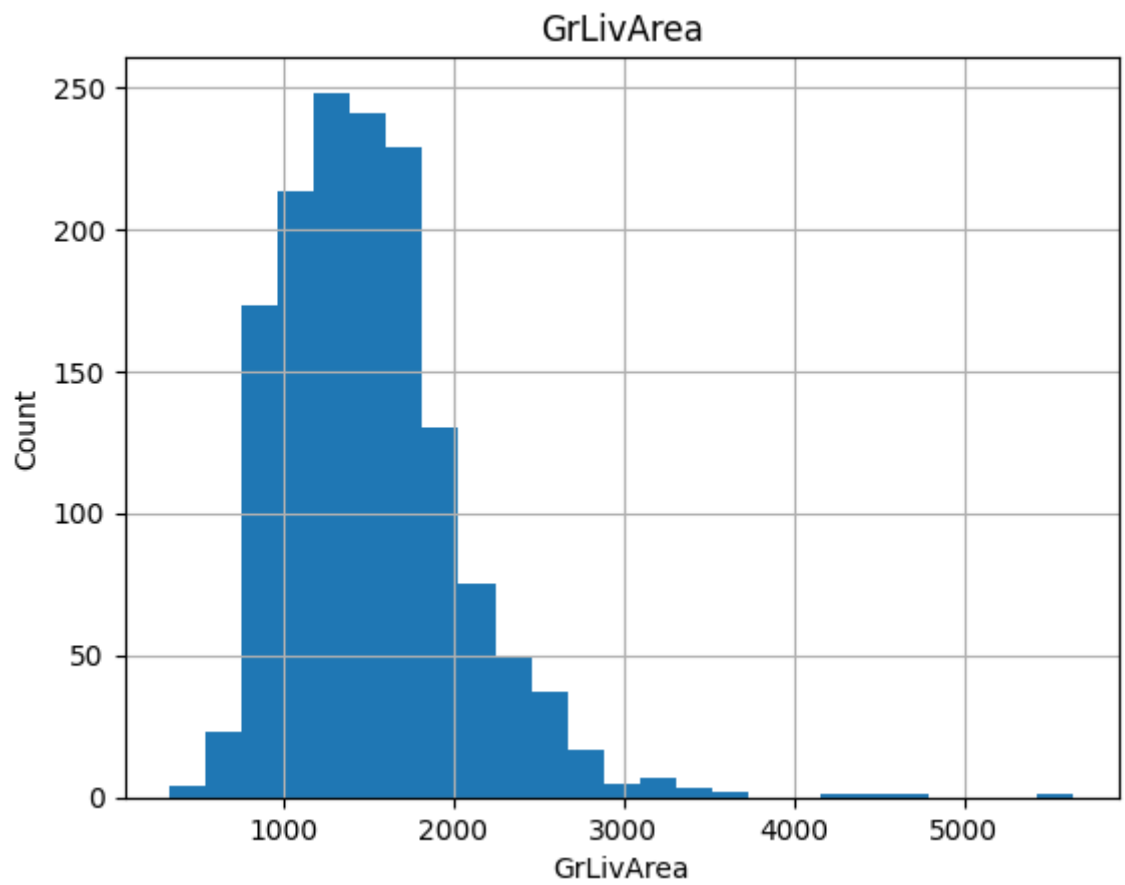


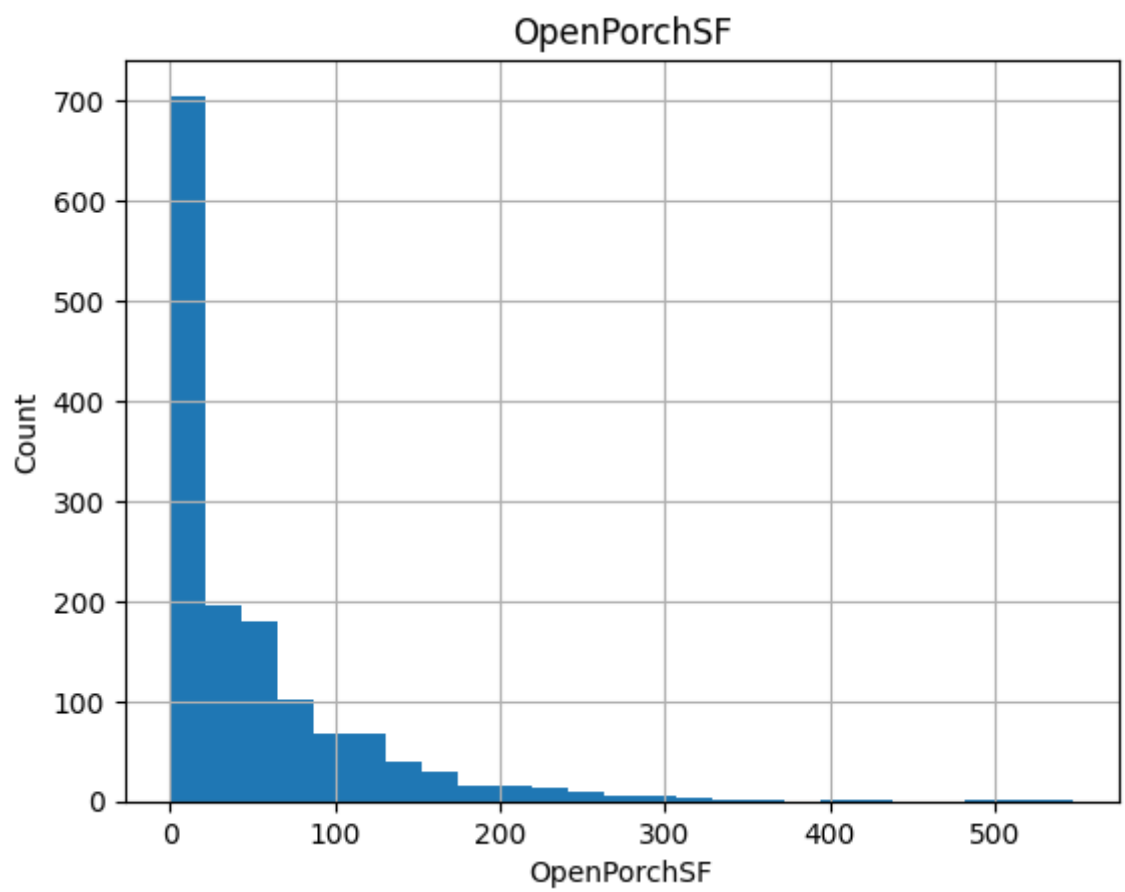
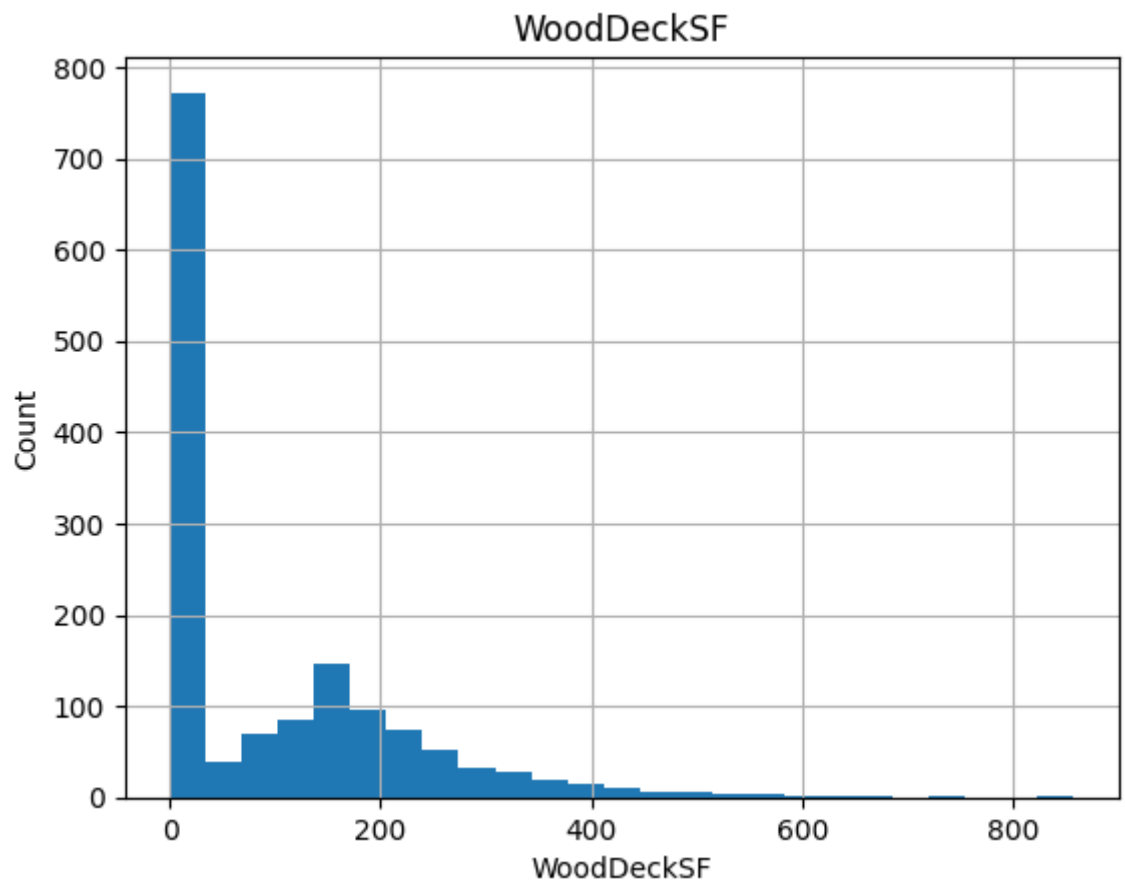


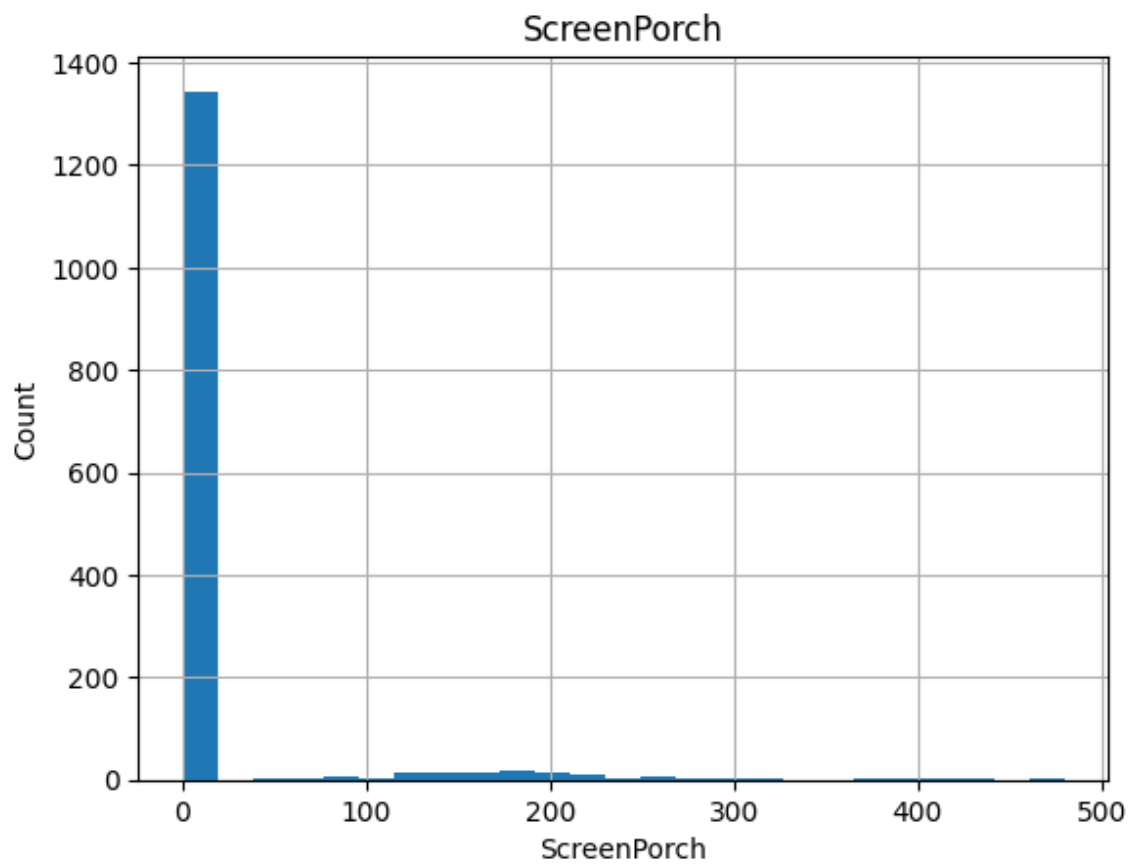
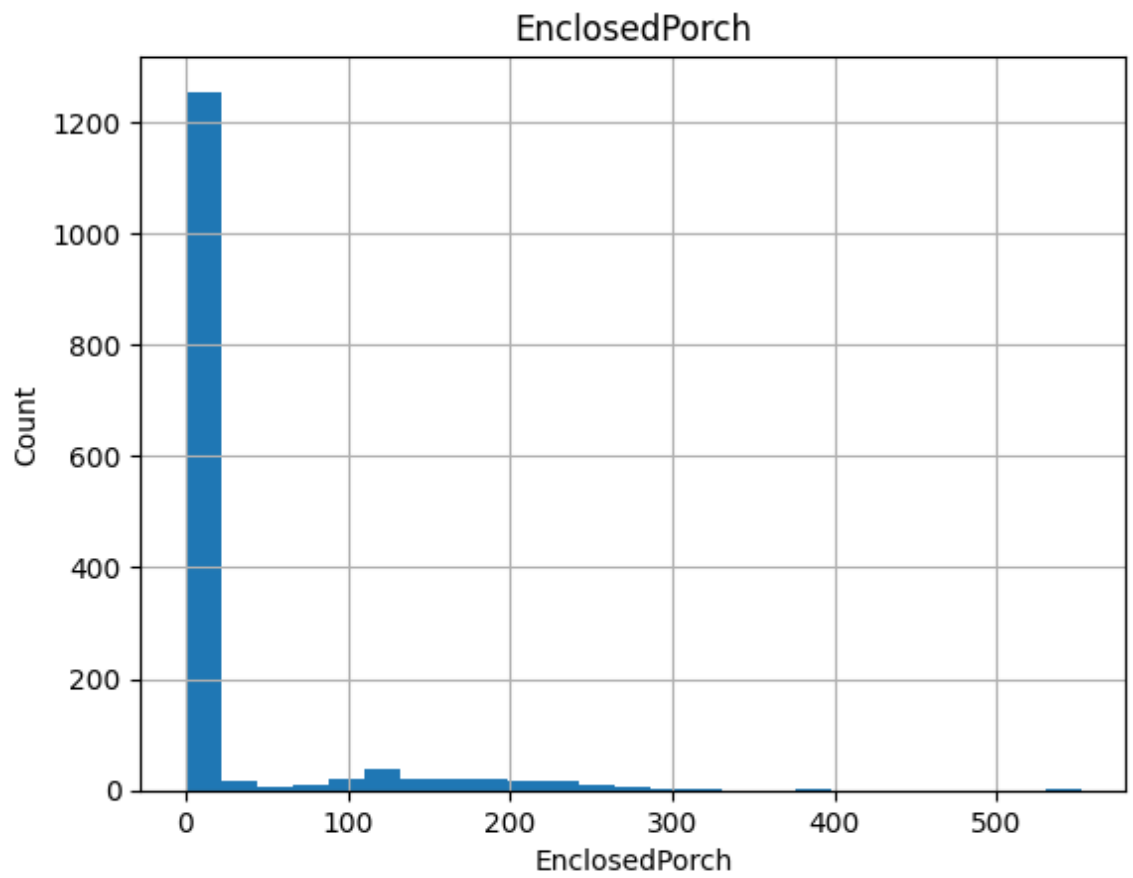


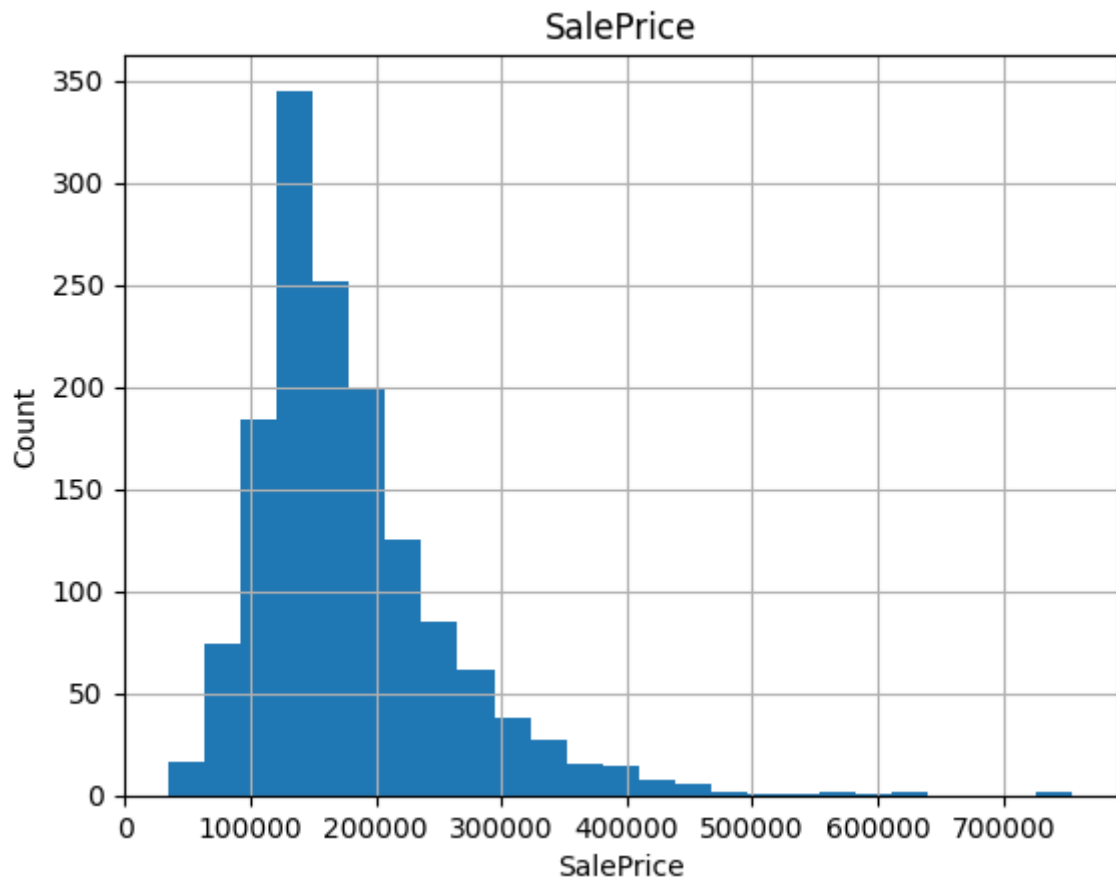










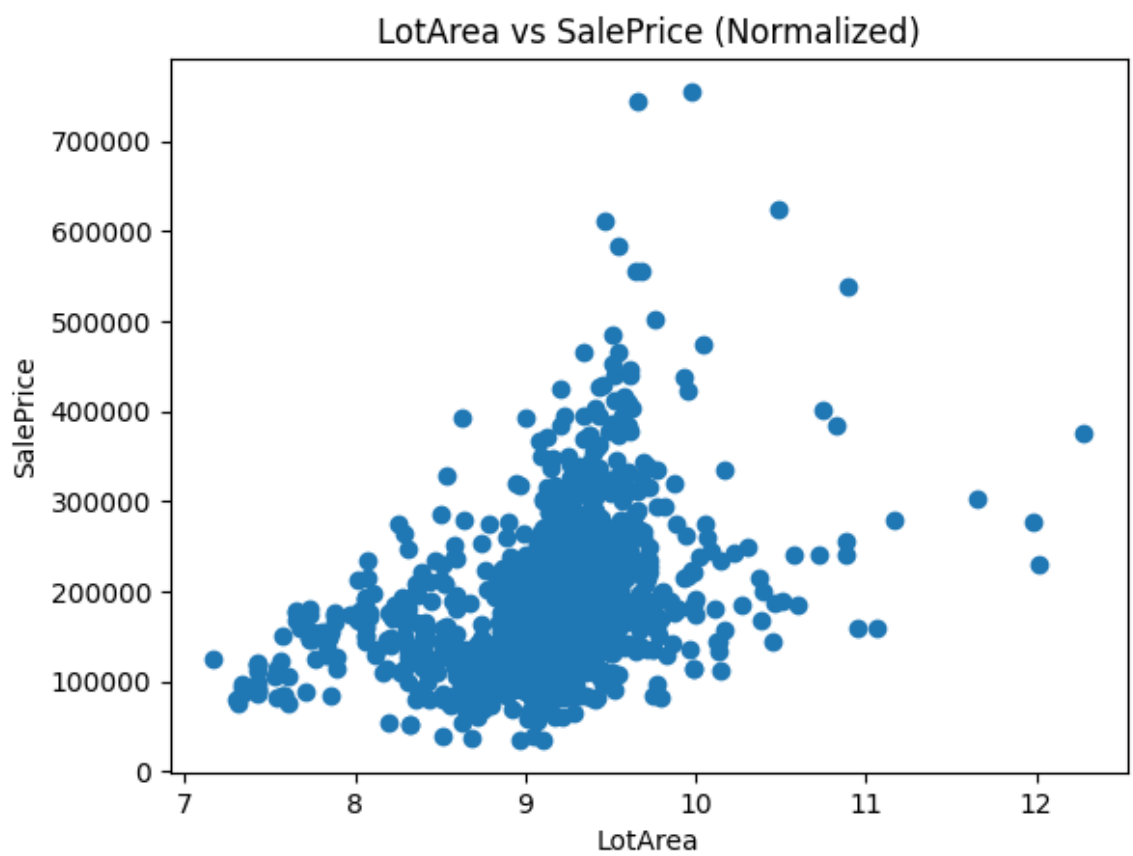
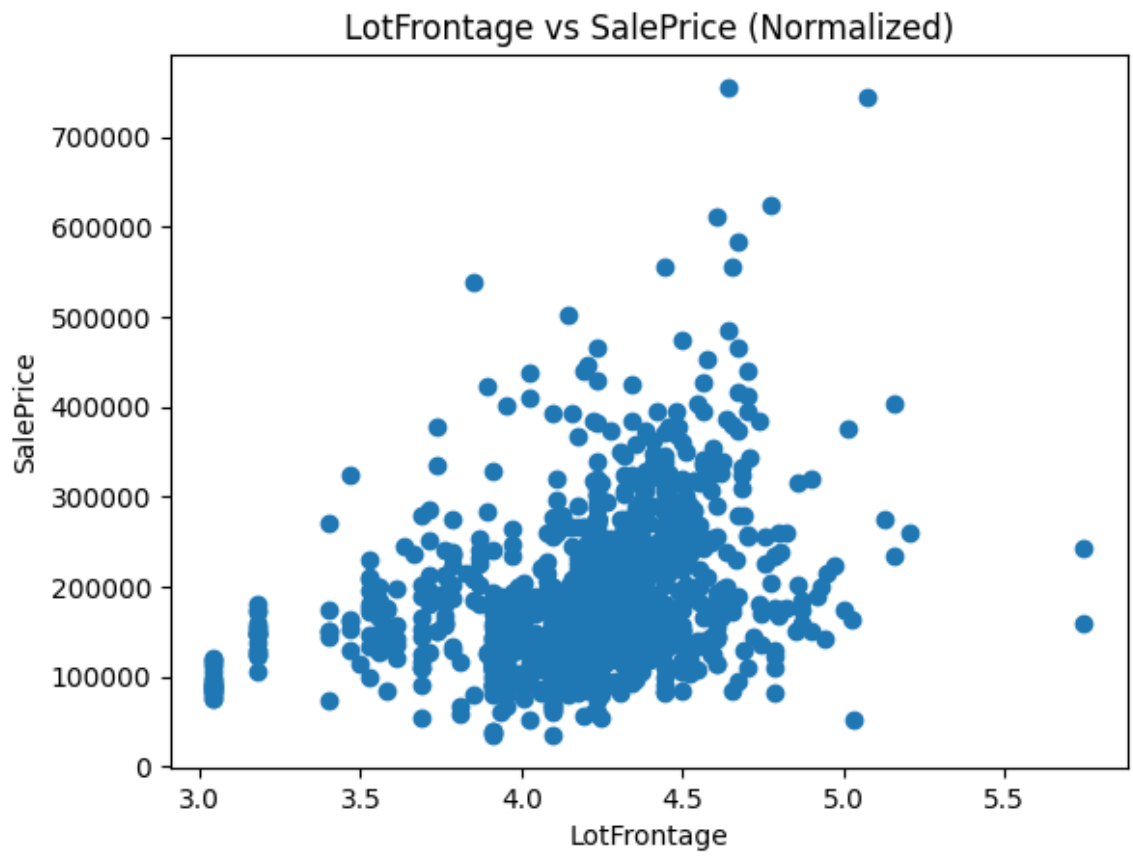


In [436...

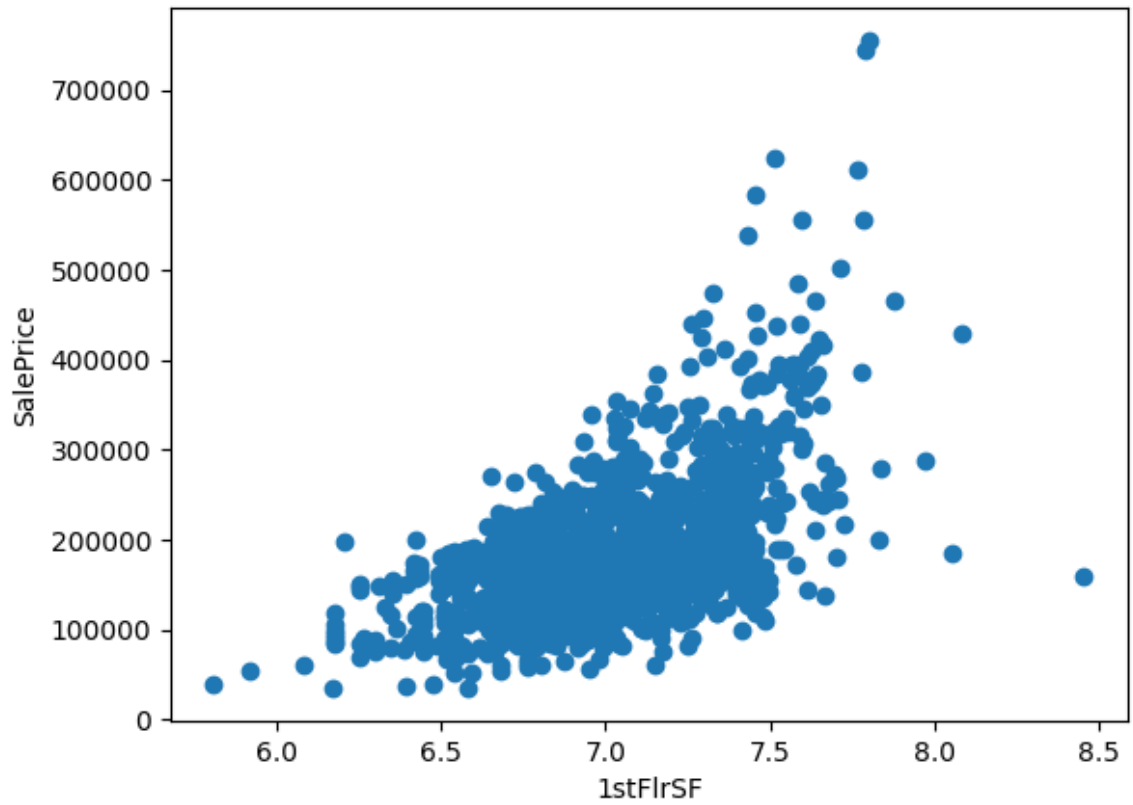
```
# use logarithmic transformation to normalize the above skewed data
skewed_num_features=['LotFrontage', 'LotArea', '1stFlrSF', 'GrLivArea', 'SalePrice']

for feature in skewed_num_features:
    if 0 in data[feature].unique(): # because log(0) is not defined
        pass
    else:
        dataset[feature]=np.log(dataset[feature])

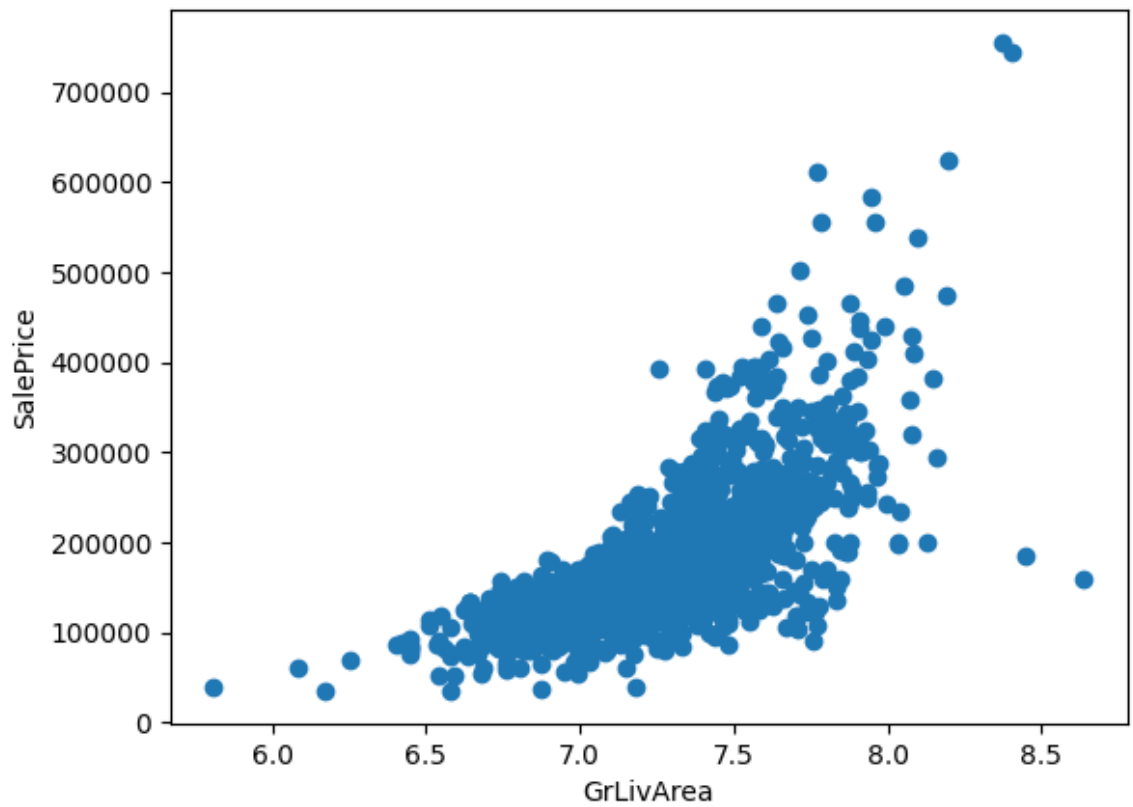
        plt.scatter(dataset[feature], dataset['SalePrice'])
        plt.xlabel(feature)
        plt.ylabel('SalePrice')
        plt.title(feature + ' vs SalePrice (Normalized)')
        plt.show()
```

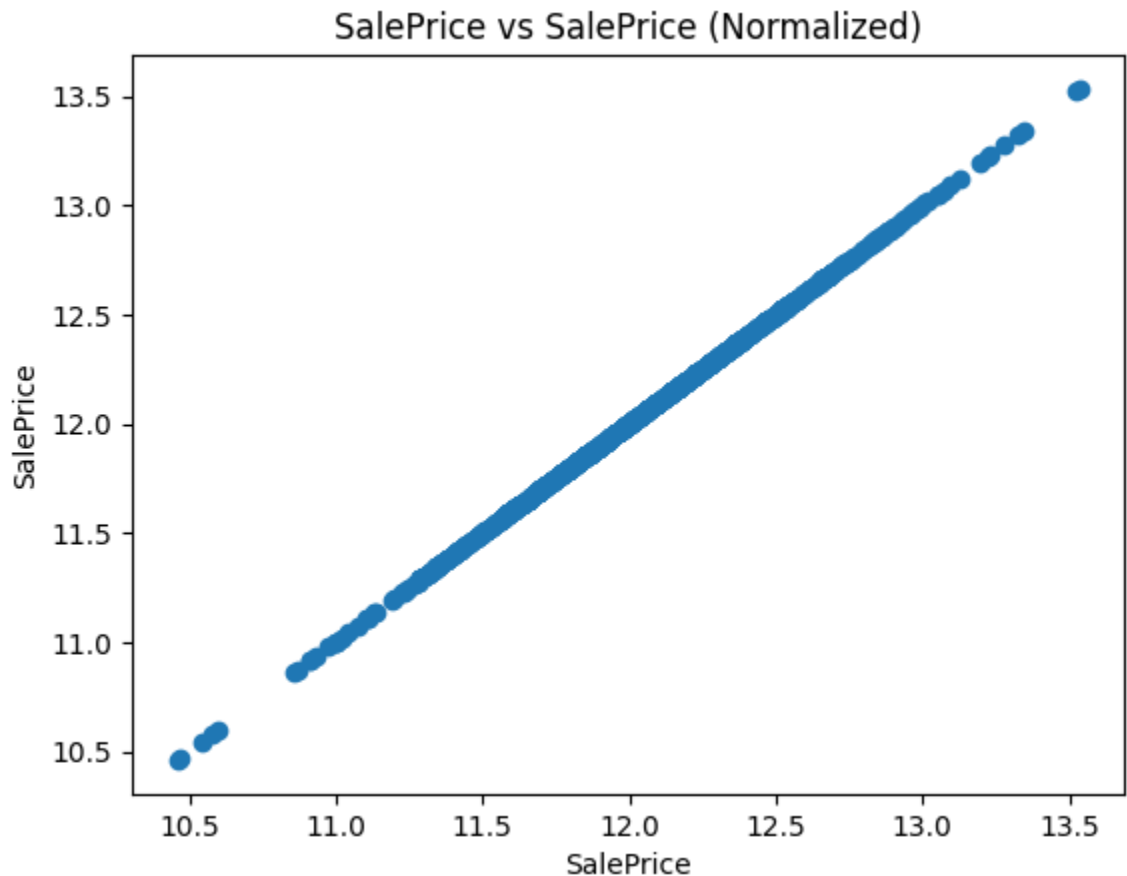


1stFlrSF vs SalePrice (Normalized)



GrLivArea vs SalePrice (Normalized)





In [437]: `dataset.head()`

Out[437]:

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour
0	1	60	RL	4.174387	9.041922	Pave	Missing	Reg	Lvl
1	2	20	RL	4.382027	9.169518	Pave	Missing	Reg	Lvl
2	3	60	RL	4.219508	9.328123	Pave	Missing	IR1	Lvl
3	4	70	RL	4.094345	9.164296	Pave	Missing	IR1	Lvl
4	5	60	RL	4.430817	9.565214	Pave	Missing	IR1	Lvl

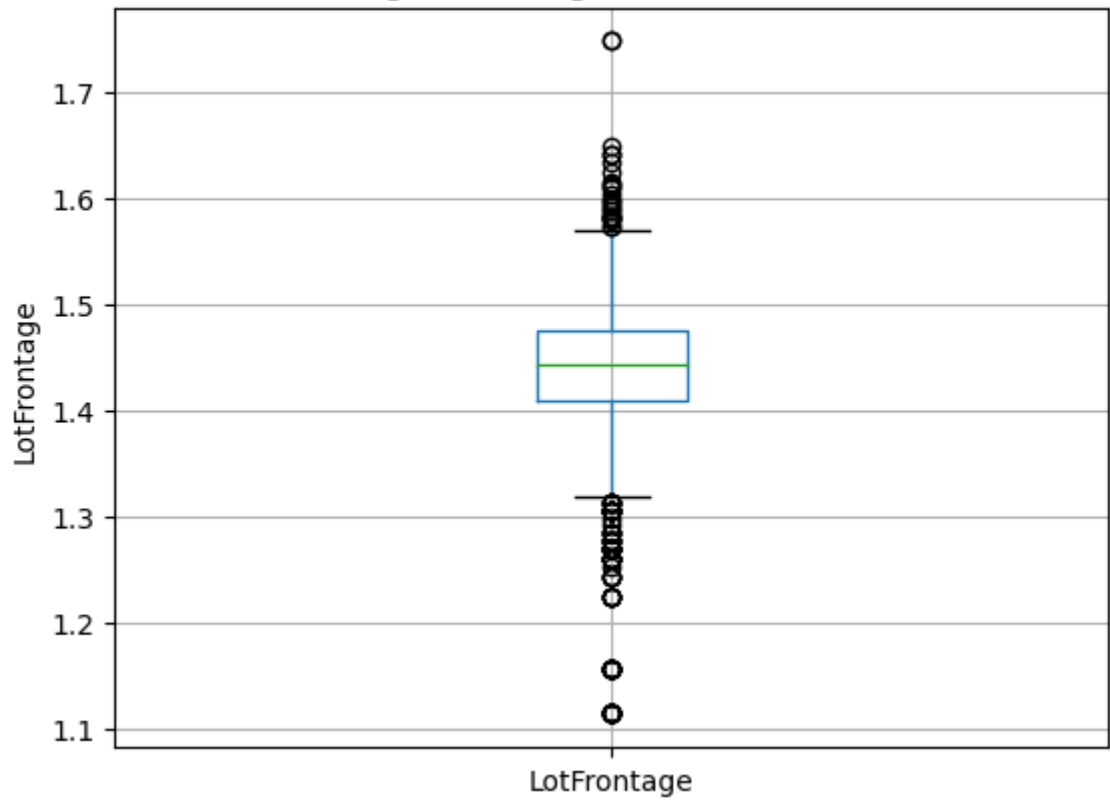
OUTLIERS

In [438]:

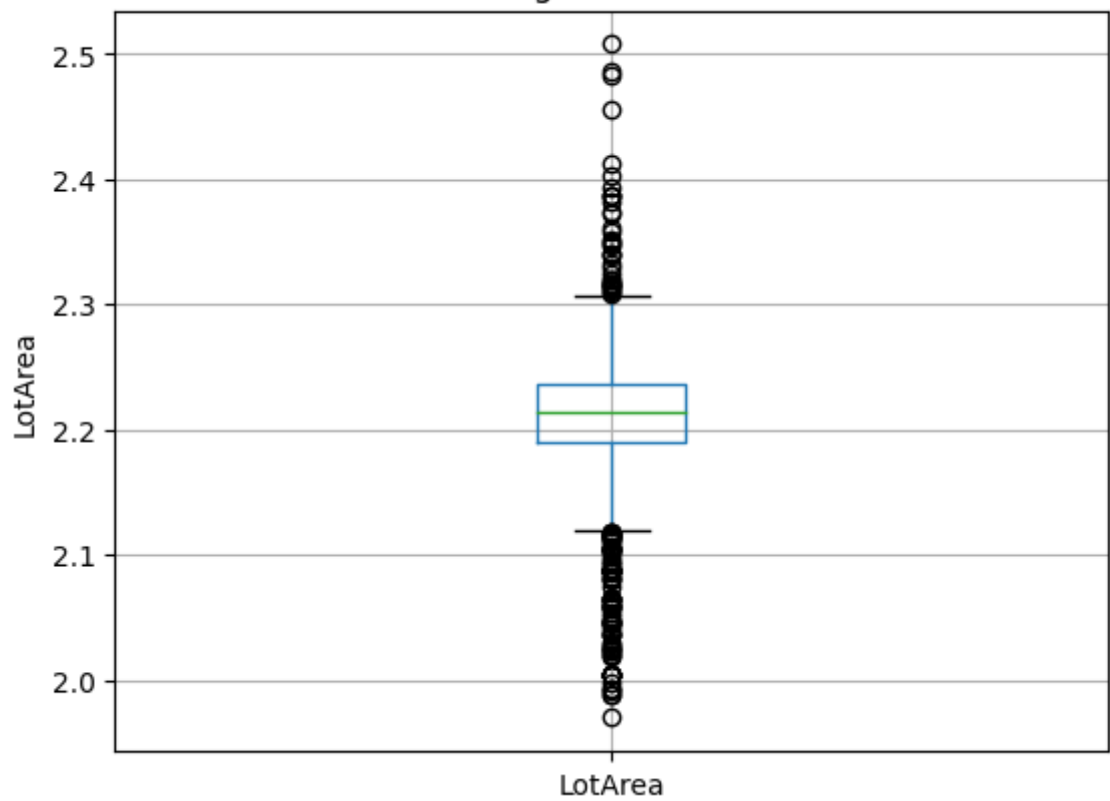
```
# use boxplots
for feature in continuous_feature:
    data=dataset.copy()
    if 0 in data[feature].unique():
        pass
    else:
        data[feature]=np.log(data[feature])

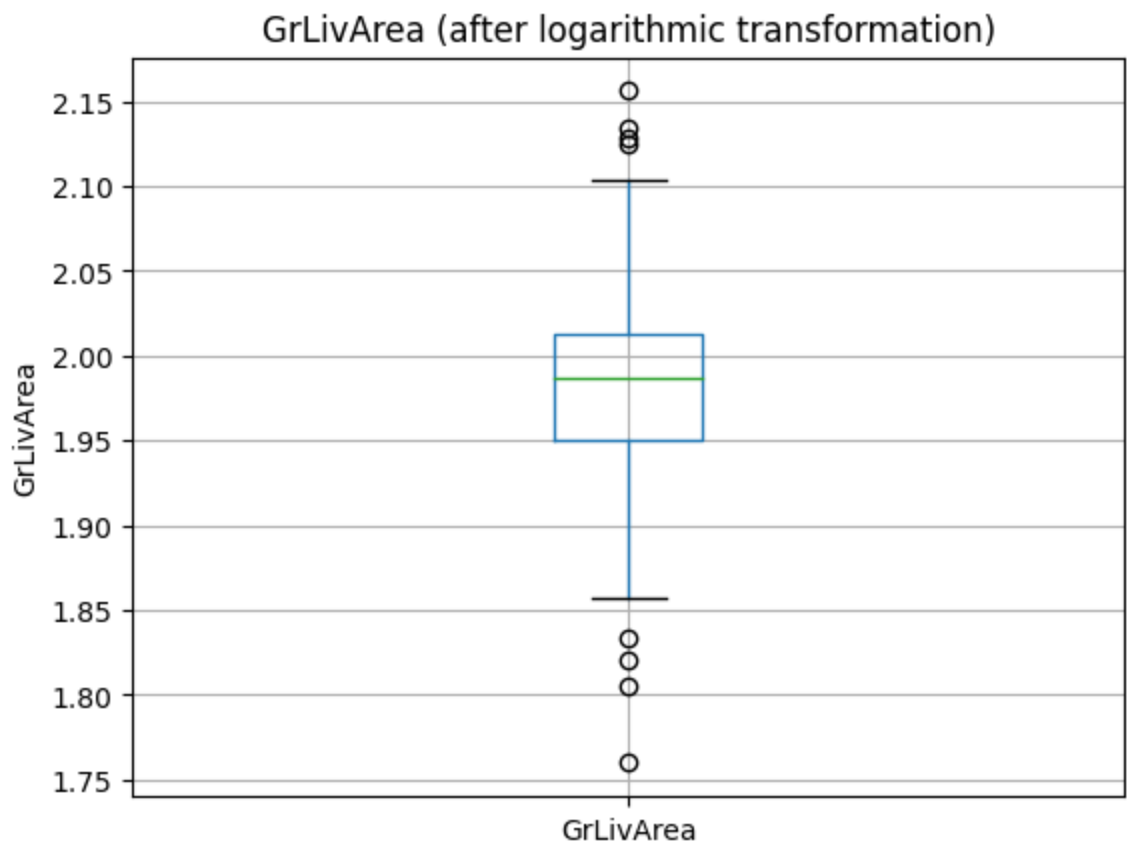
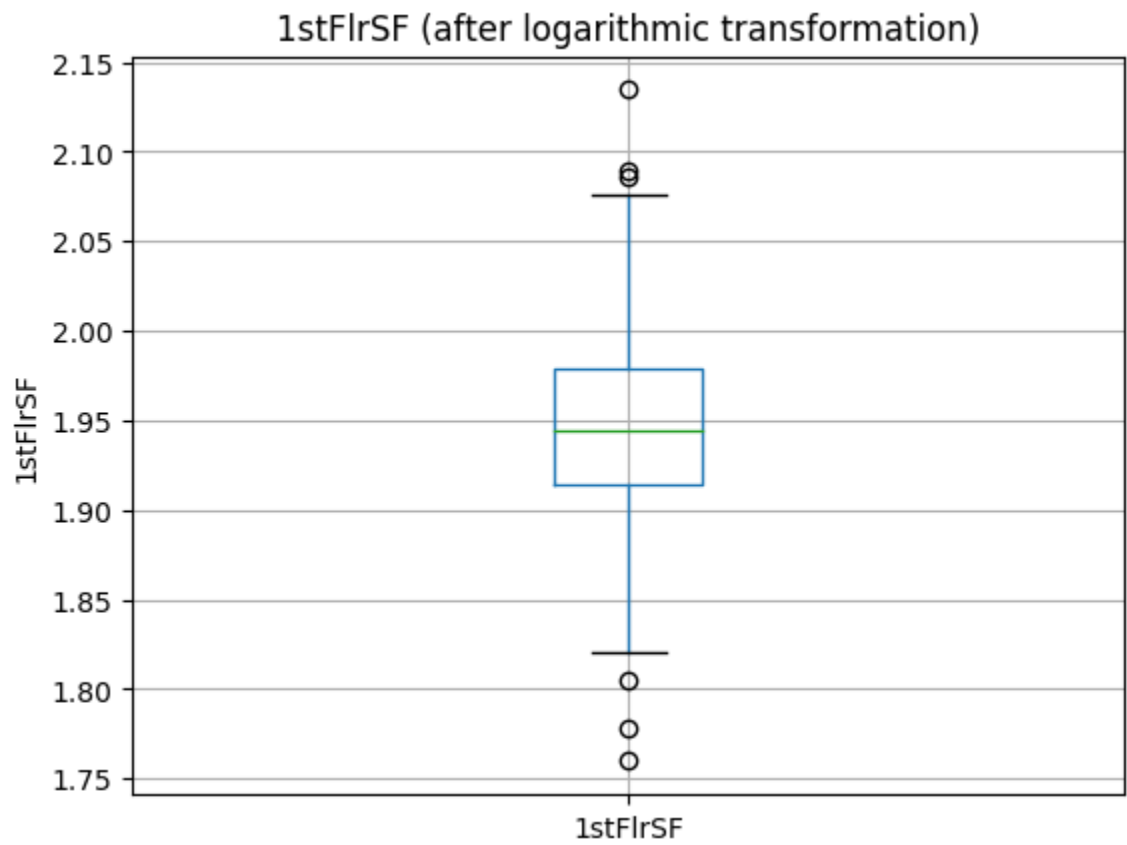
    data.boxplot(column = feature)
    plt.ylabel(feature)
    plt.title(feature + " (after logarithmic transformation)")
    plt.show()
```

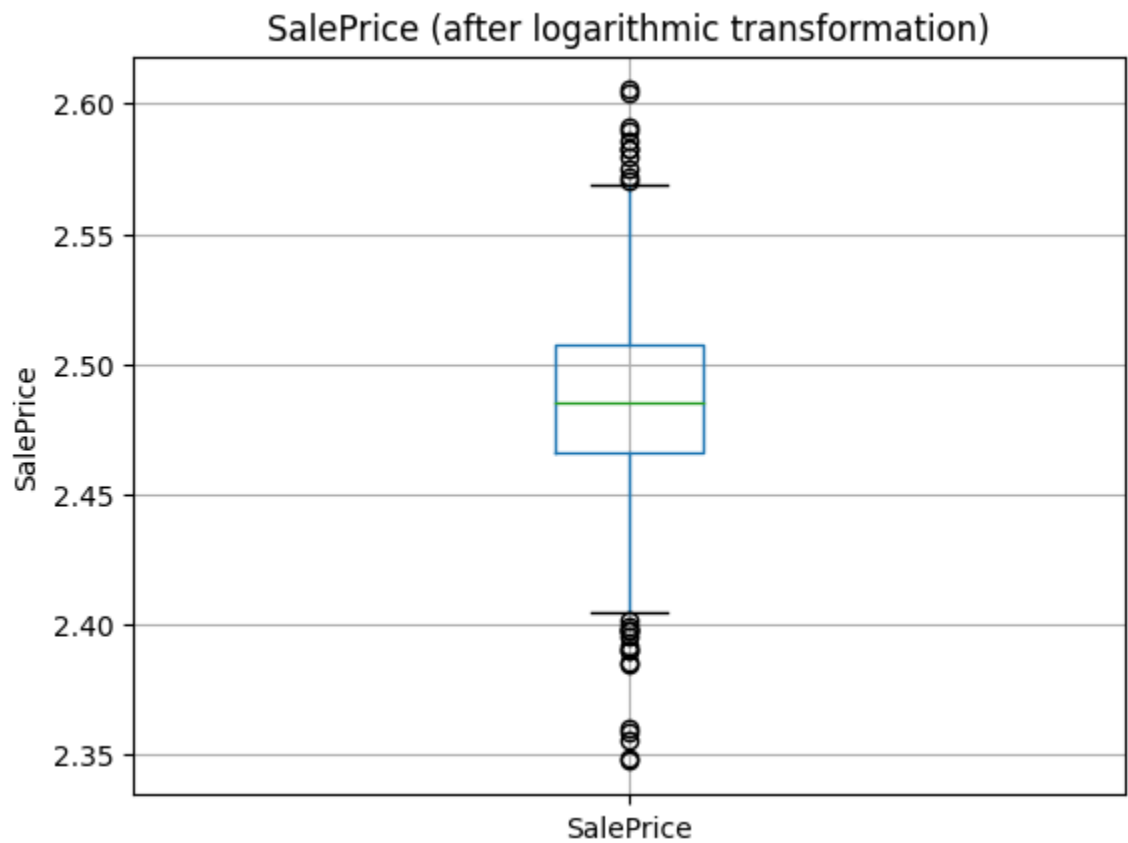
LotFrontage (after logarithmic transformation)



LotArea (after logarithmic transformation)







CATEGORICAL VARIABLES

In [439...

```
# find all categorical features
categorical_features = []
for feature in dataset.columns:
    if data[feature].dtypes == 'O':
        categorical_features.append(feature)

categorical_features
```

```
Out[439]: ['MSZoning',
           'Street',
           'Alley',
           'LotShape',
           'LandContour',
           'Utilities',
           'LotConfig',
           'LandSlope',
           'Neighborhood',
           'Condition1',
           'Condition2',
           'BldgType',
           'HouseStyle',
           'RoofStyle',
           'RoofMatl',
           'Exterior1st',
           'Exterior2nd',
           'MasVnrType',
           'ExterQual',
           'ExterCond',
           'Foundation',
           'BsmtQual',
           'BsmtCond',
           'BsmtExposure',
           'BsmtFinType1',
           'BsmtFinType2',
           'Heating',
           'HeatingQC',
           'CentralAir',
           'Electrical',
           'KitchenQual',
           'Functional',
           'FireplaceQu',
           'GarageType',
           'GarageFinish',
           'GarageQual',
           'GarageCond',
           'PavedDrive',
           'PoolQC',
           'Fence',
           'MiscFeature',
           'SaleType',
           'SaleCondition']
```

```
In [440]: dataset[categorical_features].head()
```

```
Out[440]:
```

	MSZoning	Street	Alley	LotShape	LandContour	Utilities	LotConfig	LandSlope	Neight
0	RL	Pave	Missing	Reg		Lvl	AllPub	Inside	Gtl
1	RL	Pave	Missing	Reg		Lvl	AllPub	FR2	Gtl
2	RL	Pave	Missing	IR1		Lvl	AllPub	Inside	Gtl
3	RL	Pave	Missing	IR1		Lvl	AllPub	Corner	Gtl
4	RL	Pave	Missing	IR1		Lvl	AllPub	FR2	Gtl

CARDINALITY OF CATEGORICAL VARIABLES

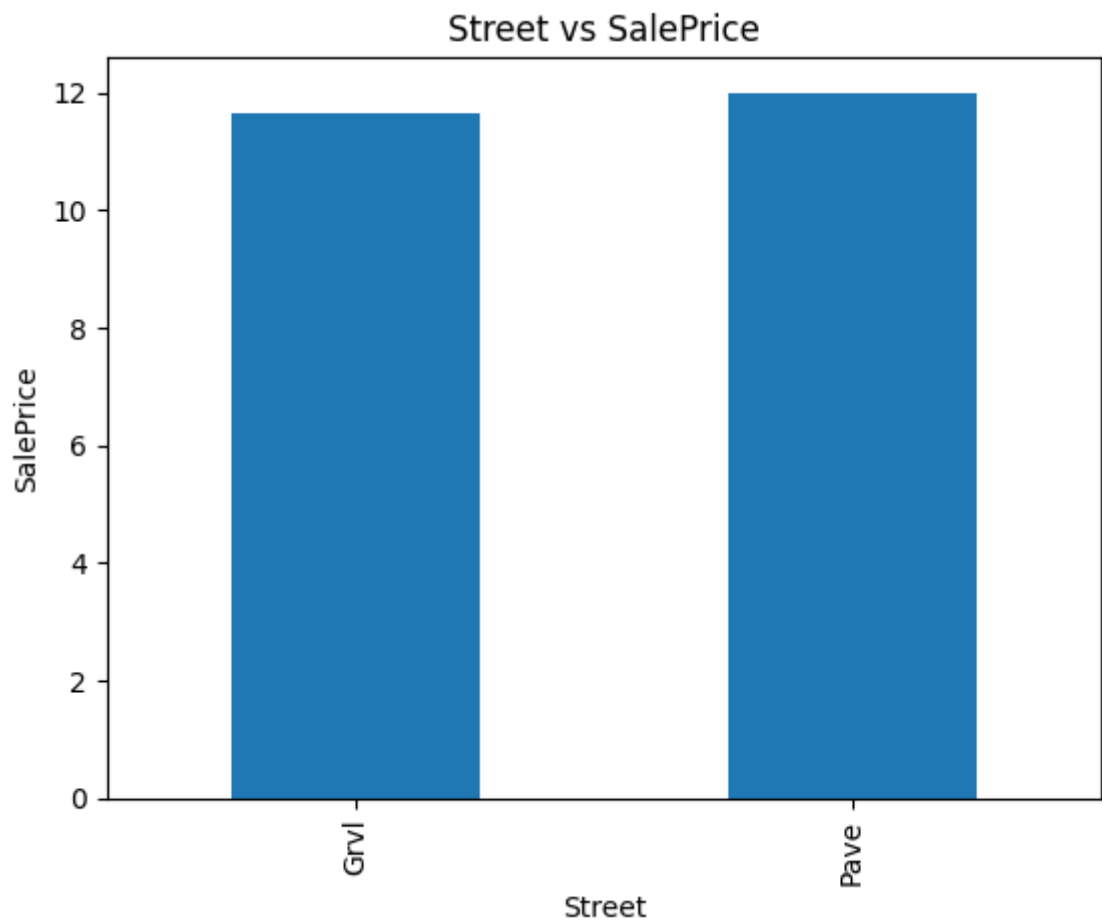
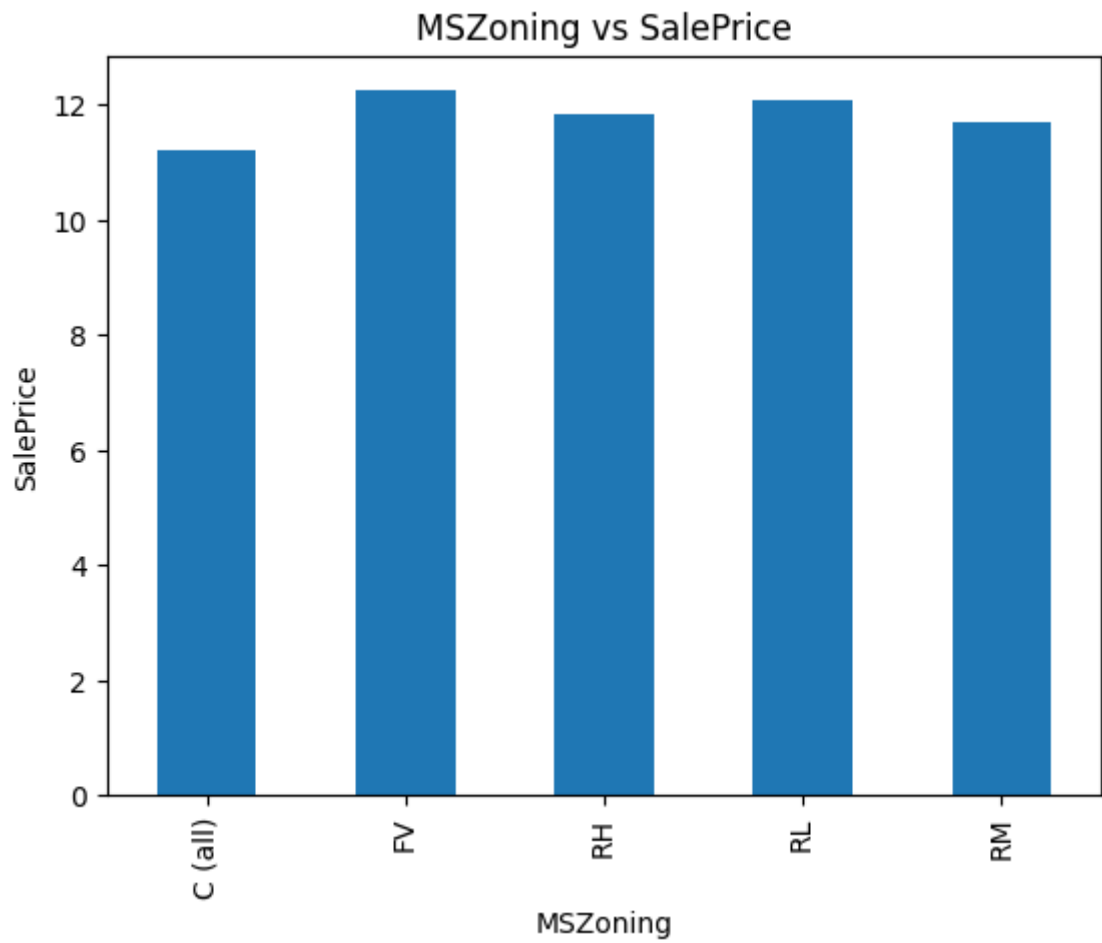
```
In [441... for feature in categorical_features:
    print('Feature: {}, Number of categories: {}'.format(feature, len(dataset[fe

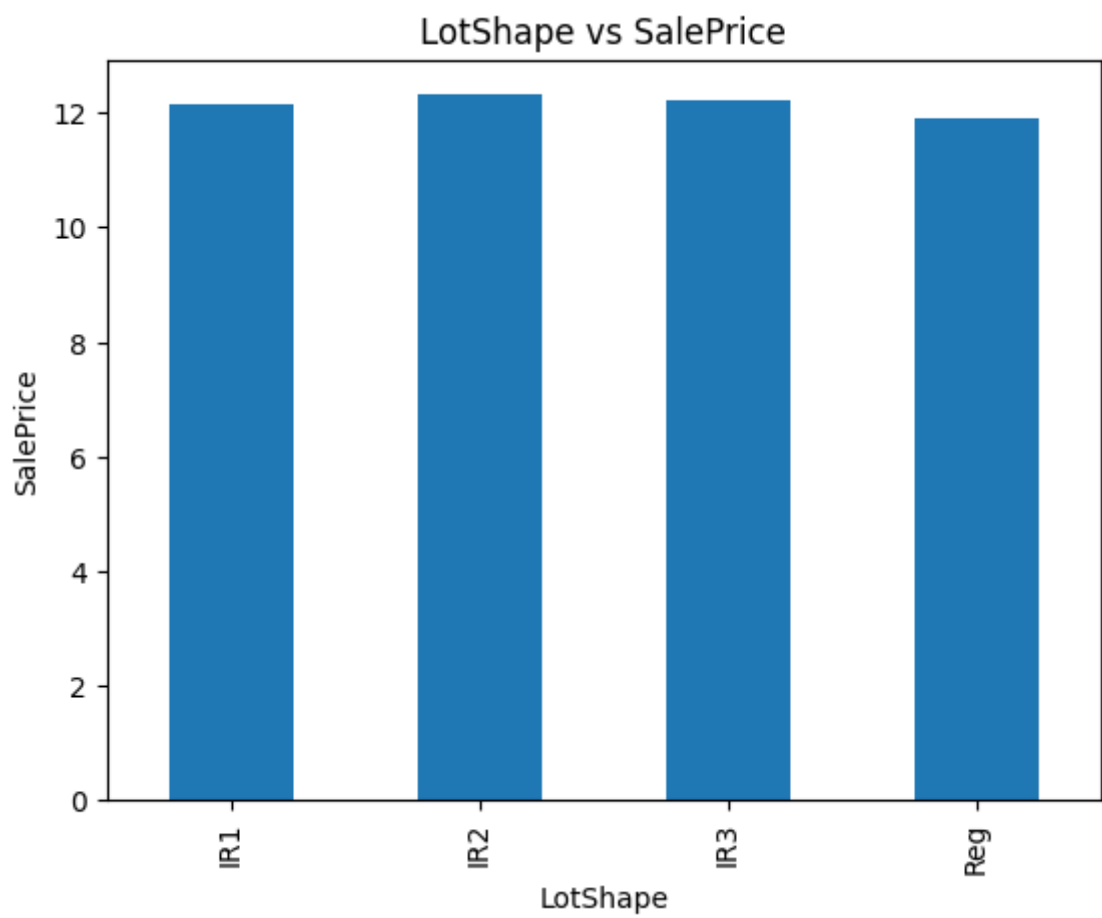
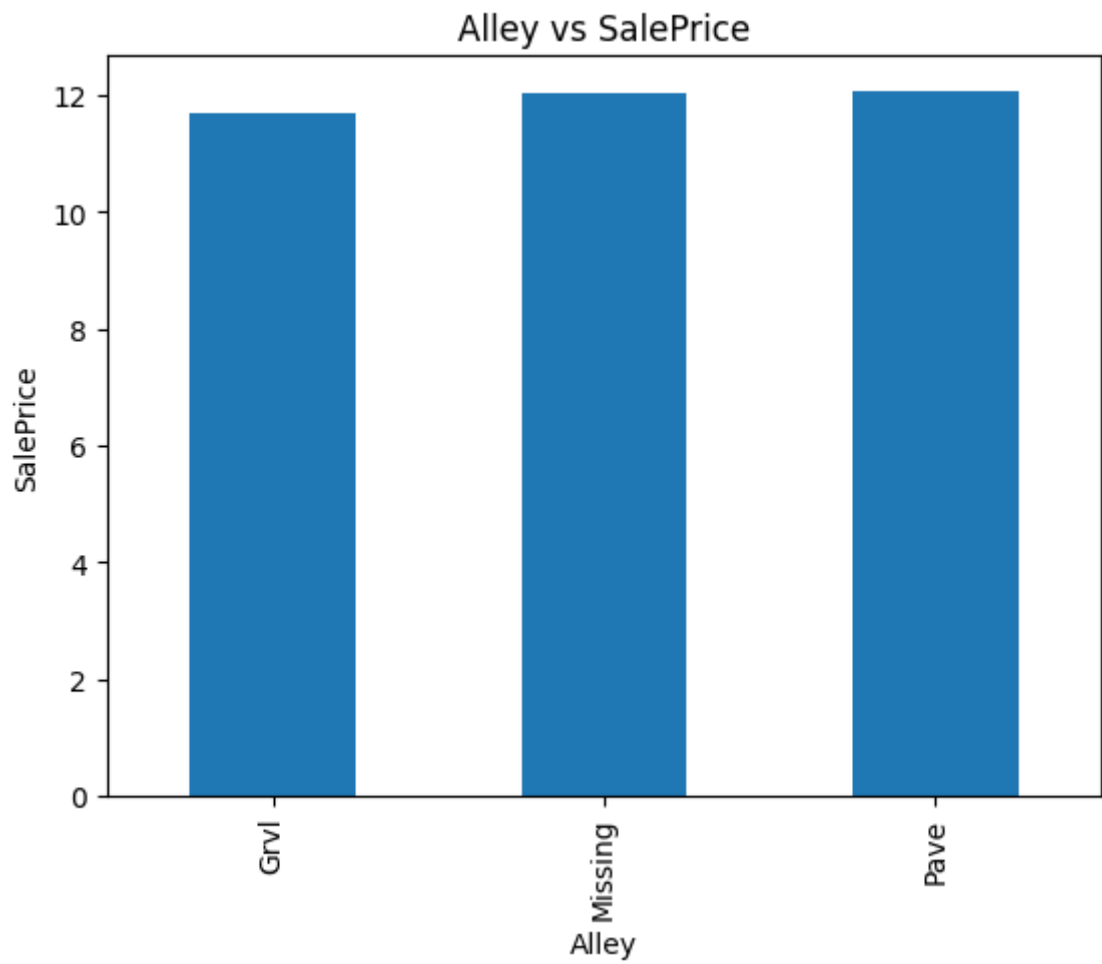
Feature: MSZoning, Number of categories: 5
Feature: Street, Number of categories: 2
Feature: Alley, Number of categories: 3
Feature: LotShape, Number of categories: 4
Feature: LandContour, Number of categories: 4
Feature: Utilities, Number of categories: 2
Feature: LotConfig, Number of categories: 5
Feature: LandSlope, Number of categories: 3
Feature: Neighborhood, Number of categories: 25
Feature: Condition1, Number of categories: 9
Feature: Condition2, Number of categories: 8
Feature: BldgType, Number of categories: 5
Feature: HouseStyle, Number of categories: 8
Feature: RoofStyle, Number of categories: 6
Feature: RoofMatl, Number of categories: 8
Feature: Exterior1st, Number of categories: 15
Feature: Exterior2nd, Number of categories: 16
Feature: MasVnrType, Number of categories: 5
Feature: ExterQual, Number of categories: 4
Feature: ExterCond, Number of categories: 5
Feature: Foundation, Number of categories: 6
Feature: BsmtQual, Number of categories: 5
Feature: BsmtCond, Number of categories: 5
Feature: BsmtExposure, Number of categories: 5
Feature: BsmtFinType1, Number of categories: 7
Feature: BsmtFinType2, Number of categories: 7
Feature: Heating, Number of categories: 6
Feature: HeatingQC, Number of categories: 5
Feature: CentralAir, Number of categories: 2
Feature: Electrical, Number of categories: 5
Feature: KitchenQual, Number of categories: 4
Feature: Functional, Number of categories: 7
Feature: FireplaceQu, Number of categories: 6
Feature: GarageType, Number of categories: 7
Feature: GarageFinish, Number of categories: 4
Feature: GarageQual, Number of categories: 6
Feature: GarageCond, Number of categories: 6
Feature: PavedDrive, Number of categories: 3
Feature: PoolQC, Number of categories: 4
Feature: Fence, Number of categories: 5
Feature: MiscFeature, Number of categories: 5
Feature: SaleType, Number of categories: 9
Feature: SaleCondition, Number of categories: 6
```

```
In [442... # find the relationship between categorical variables and SalePrice
for feature in categorical_features:
    data = dataset.copy()
    data.groupby(feature)['SalePrice'].median().plot.bar()

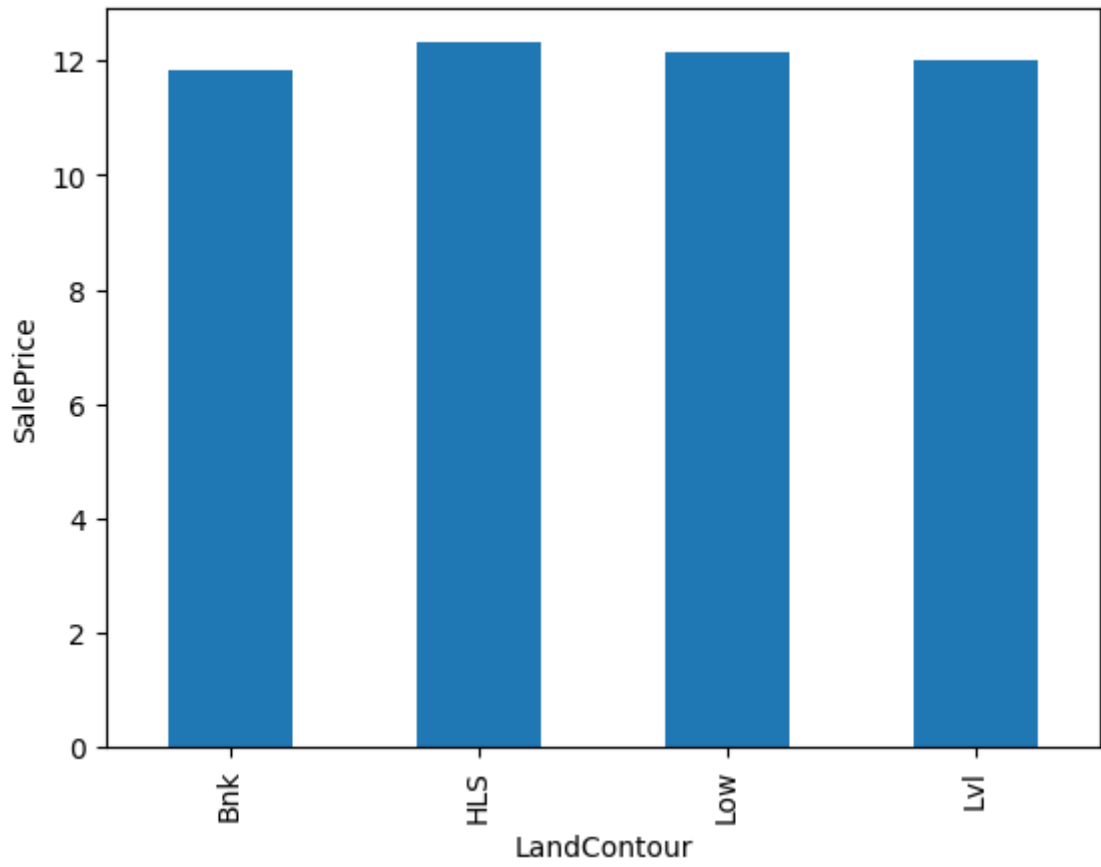
    plt.xlabel(feature)
    plt.ylabel('SalePrice')
```

```
plt.title(feature + ' vs SalePrice')  
plt.show()
```

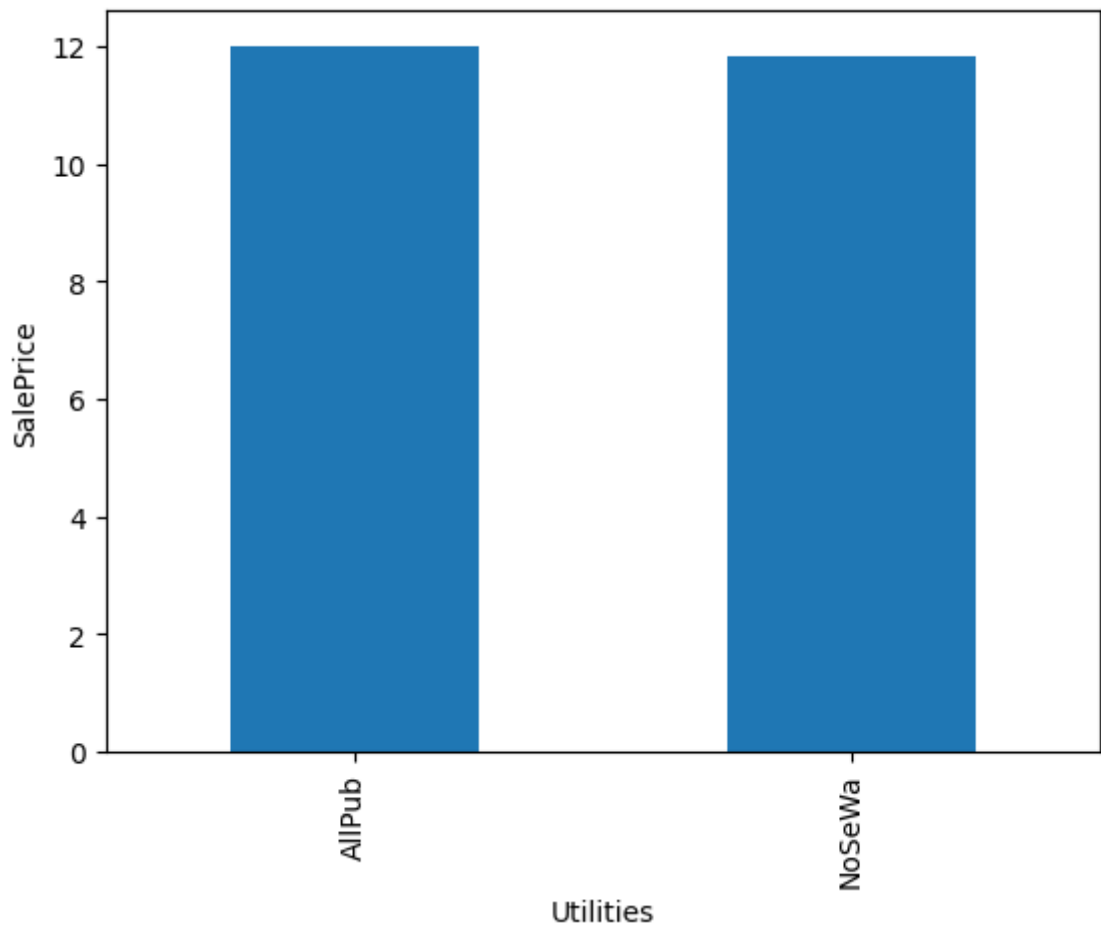


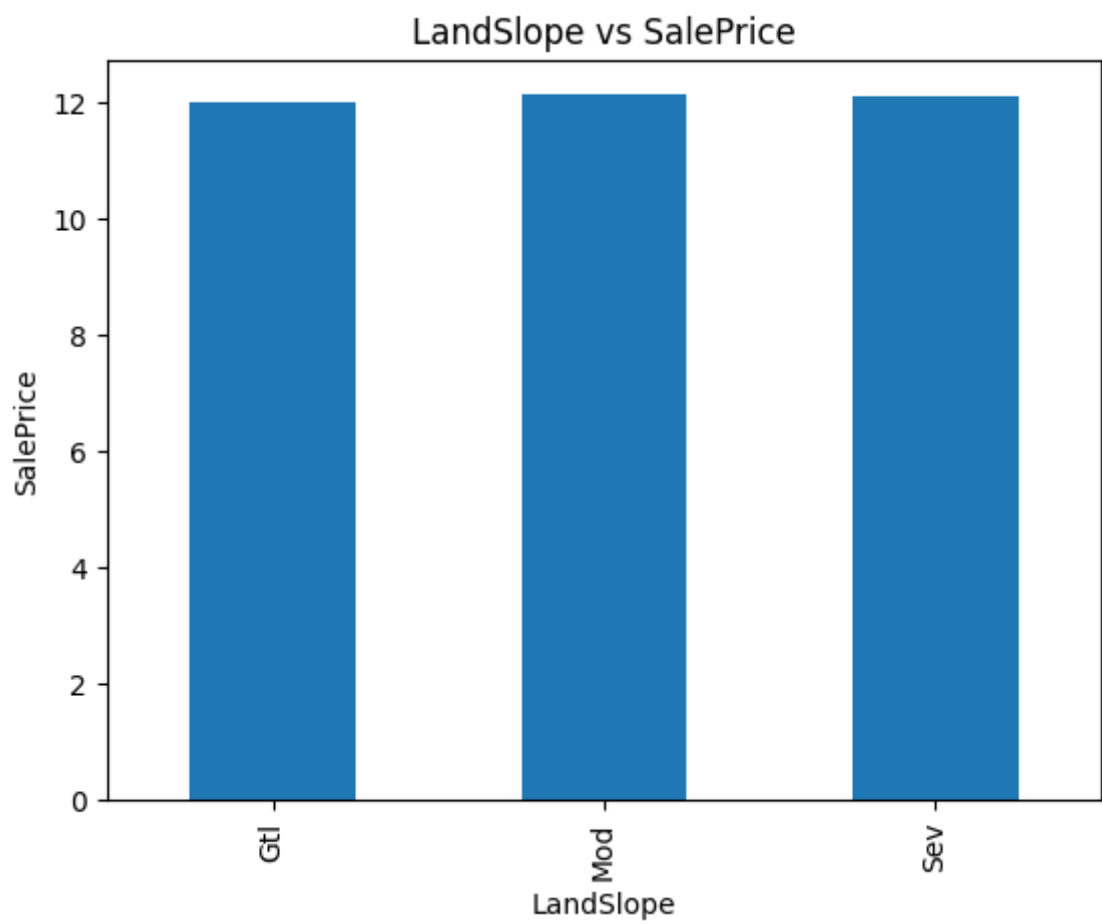
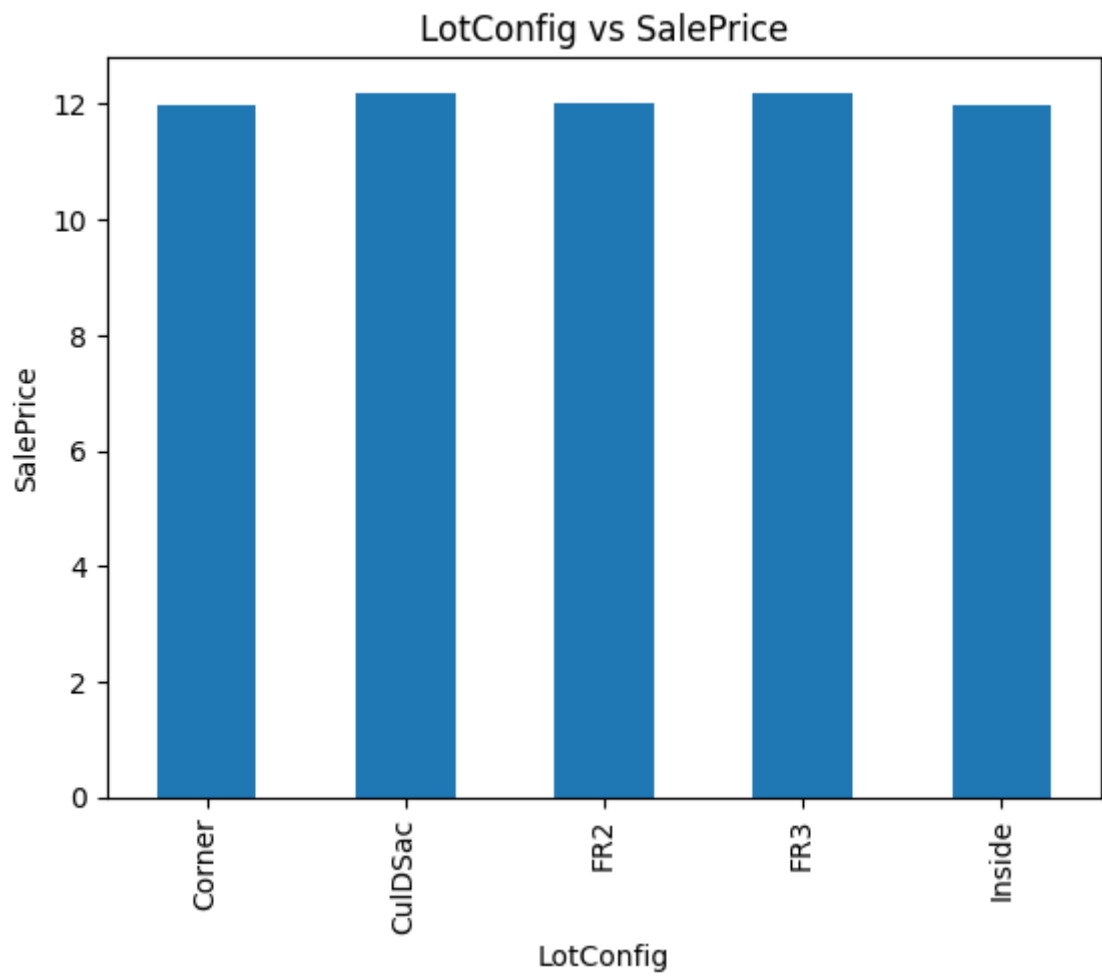


LandContour vs SalePrice

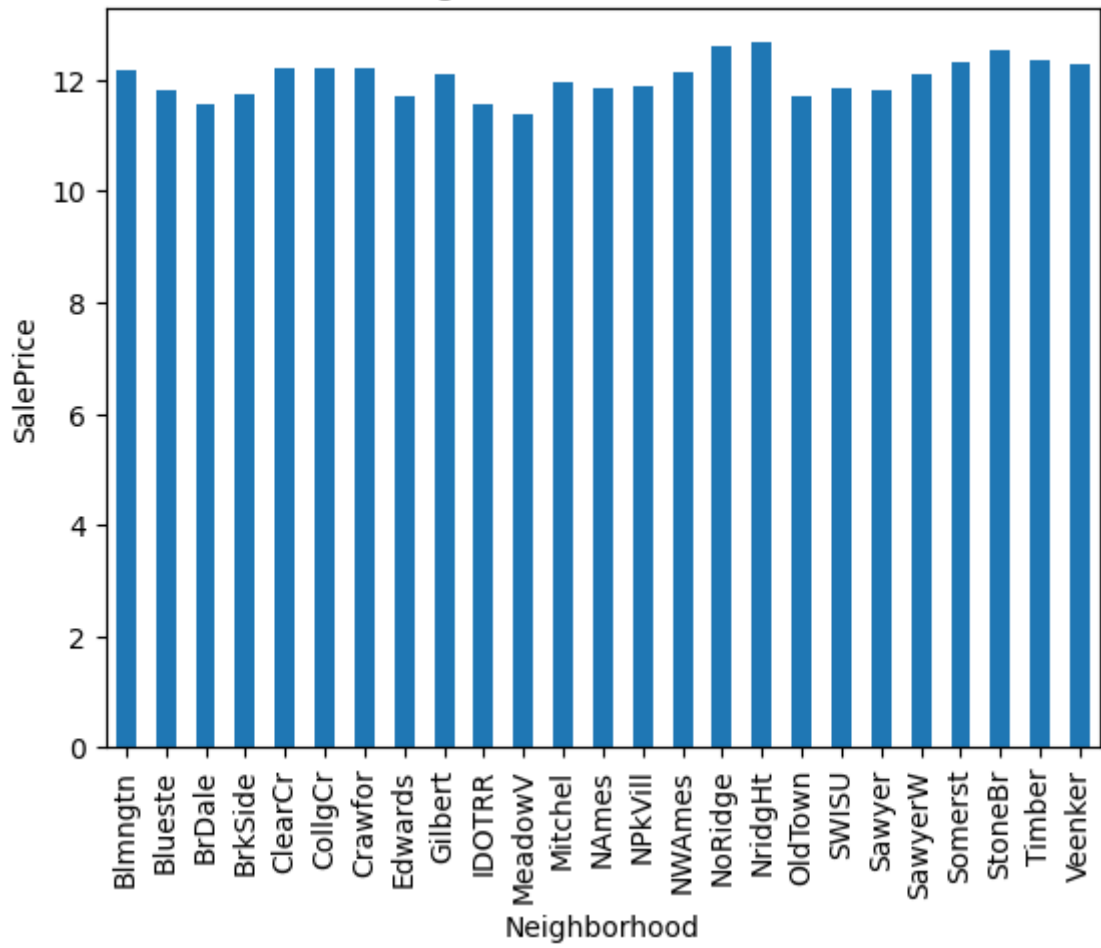


Utilities vs SalePrice

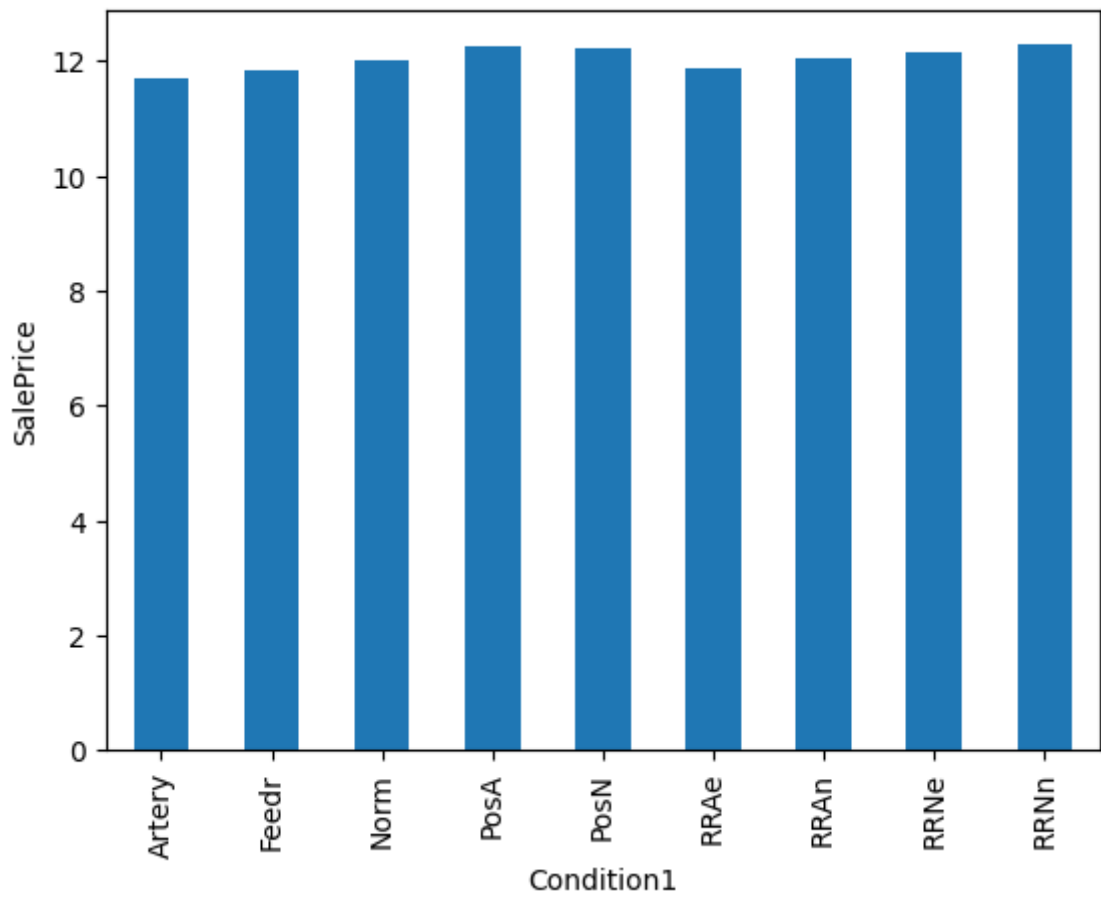


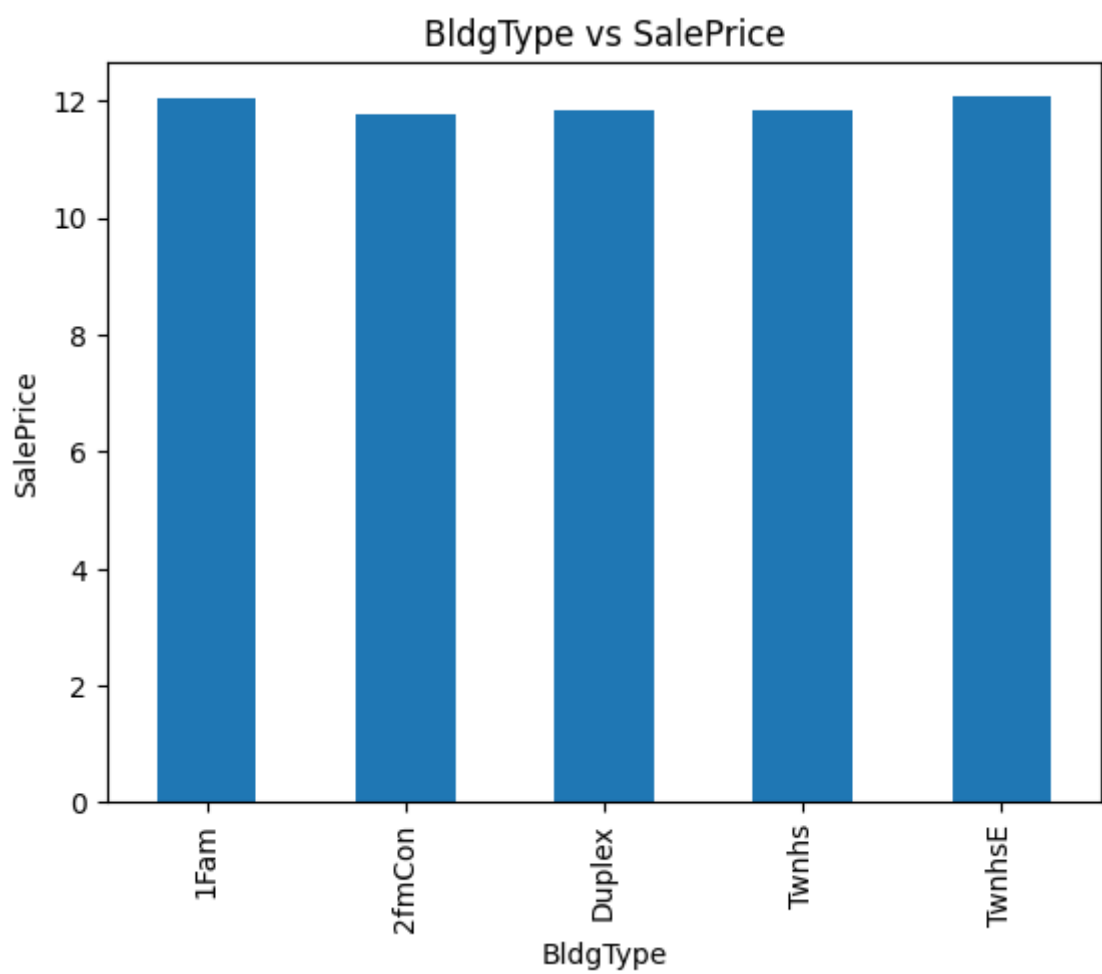
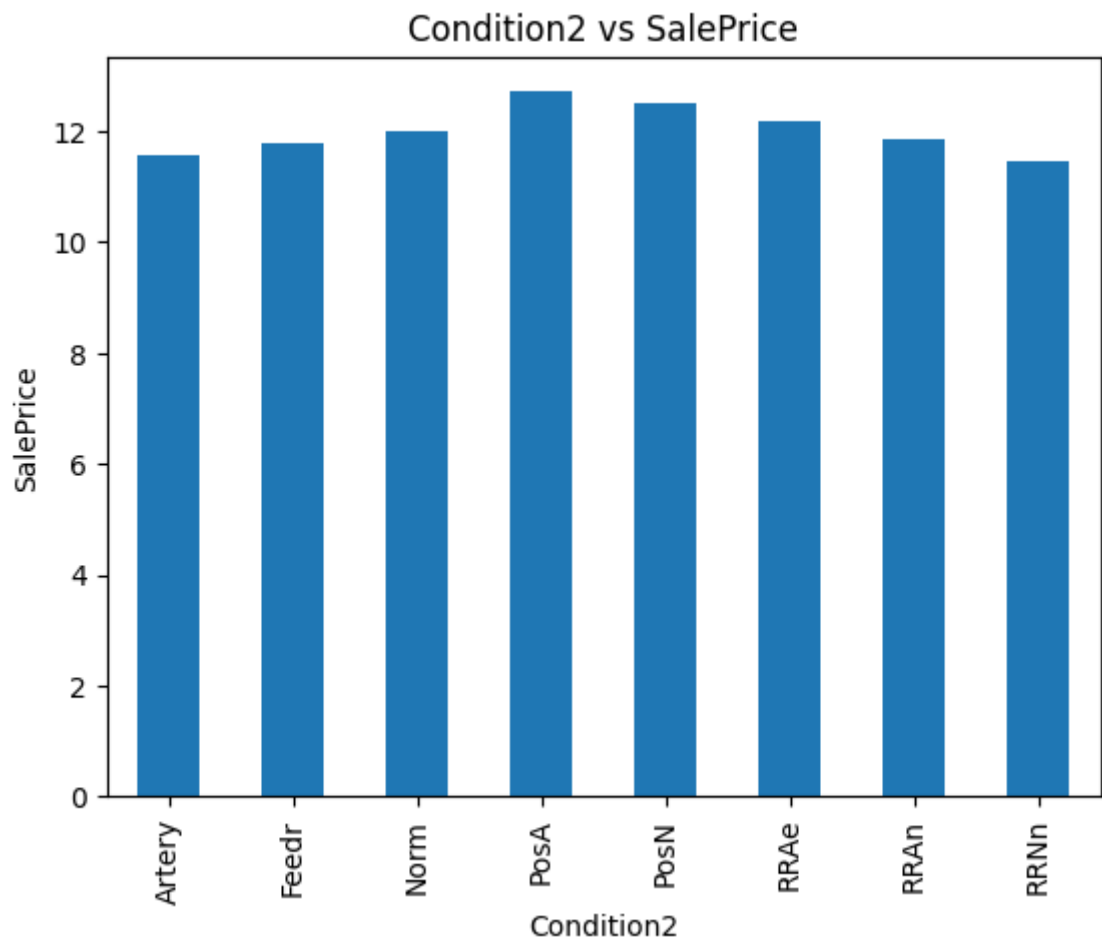


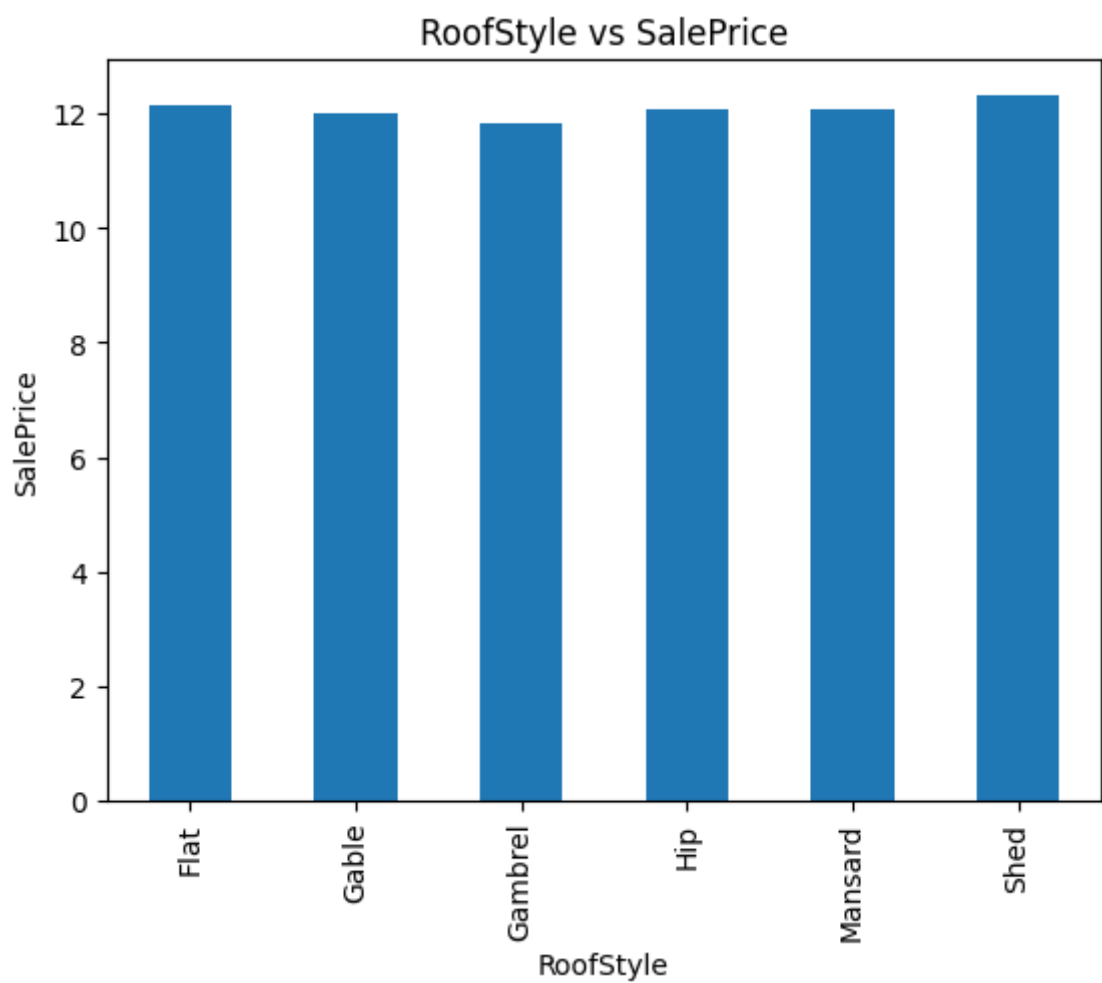
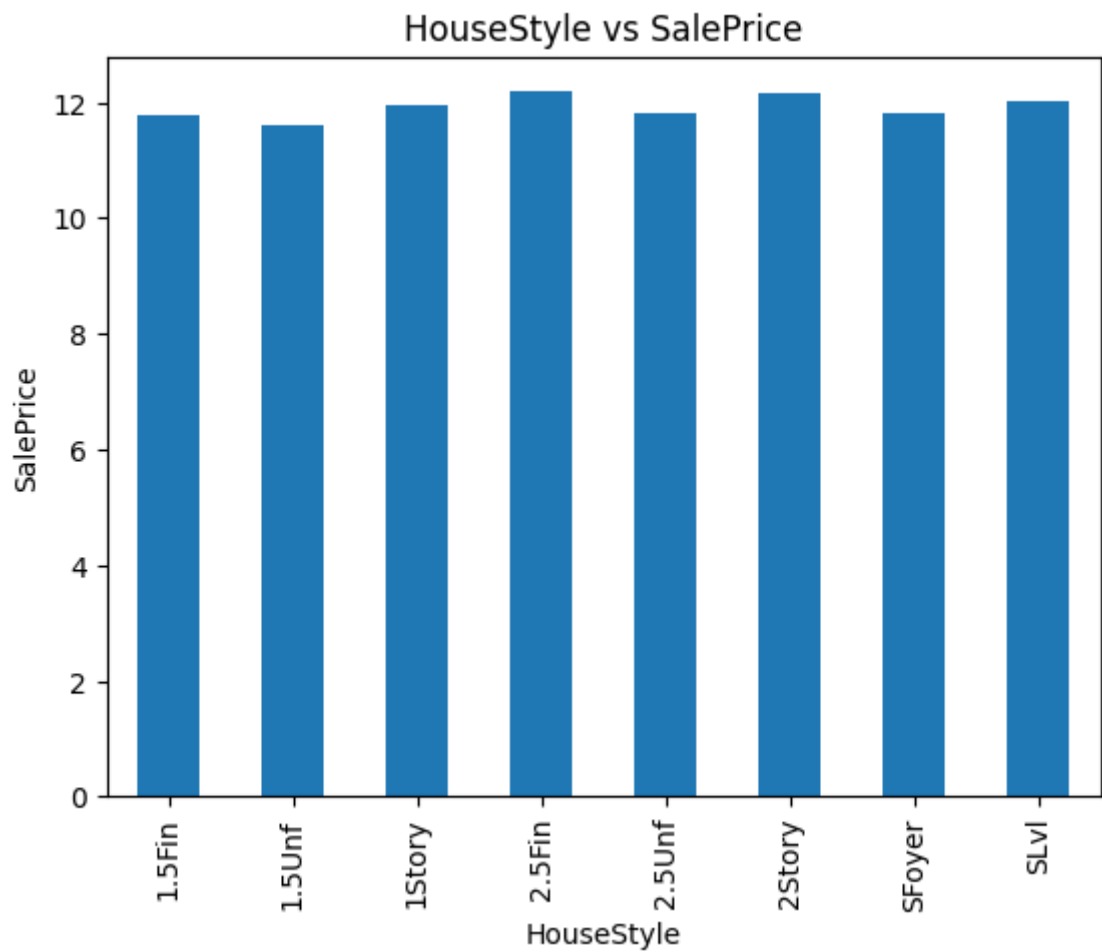
Neighborhood vs SalePrice



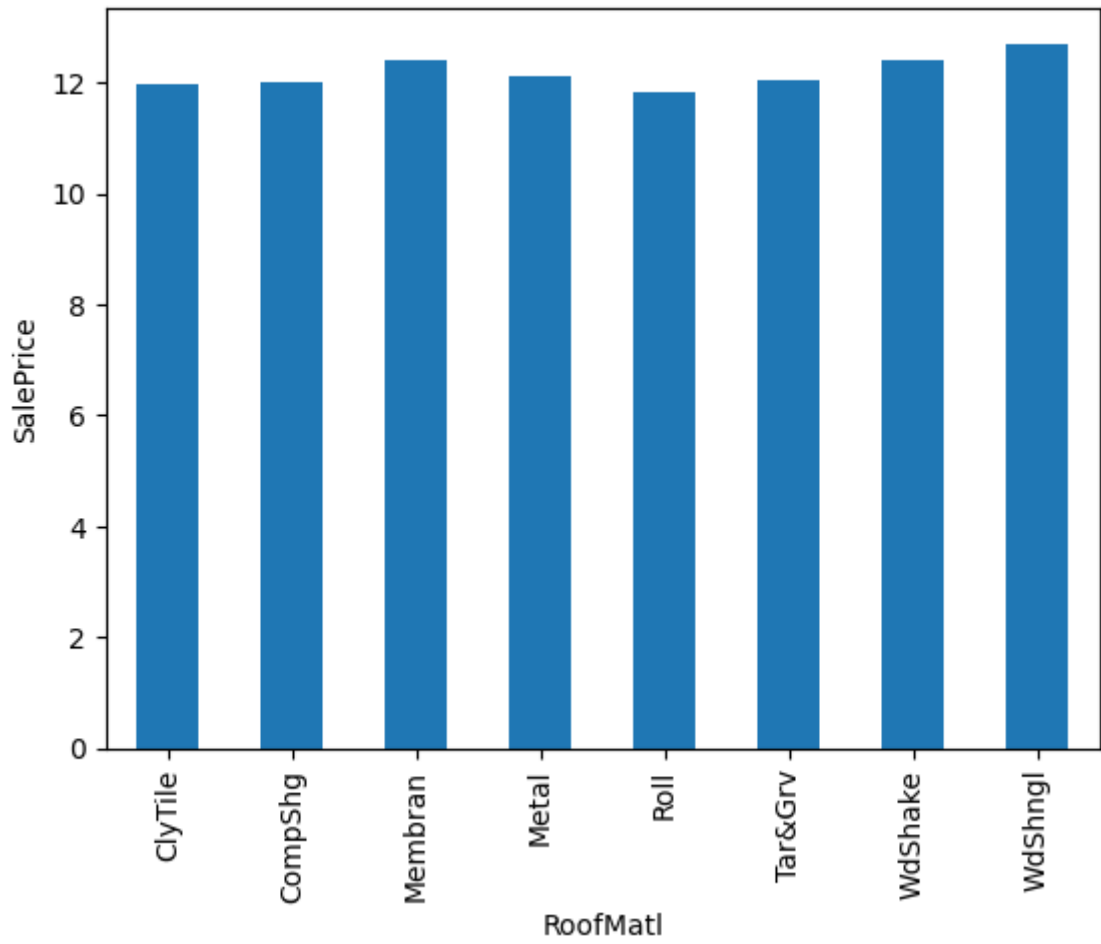
Condition1 vs SalePrice



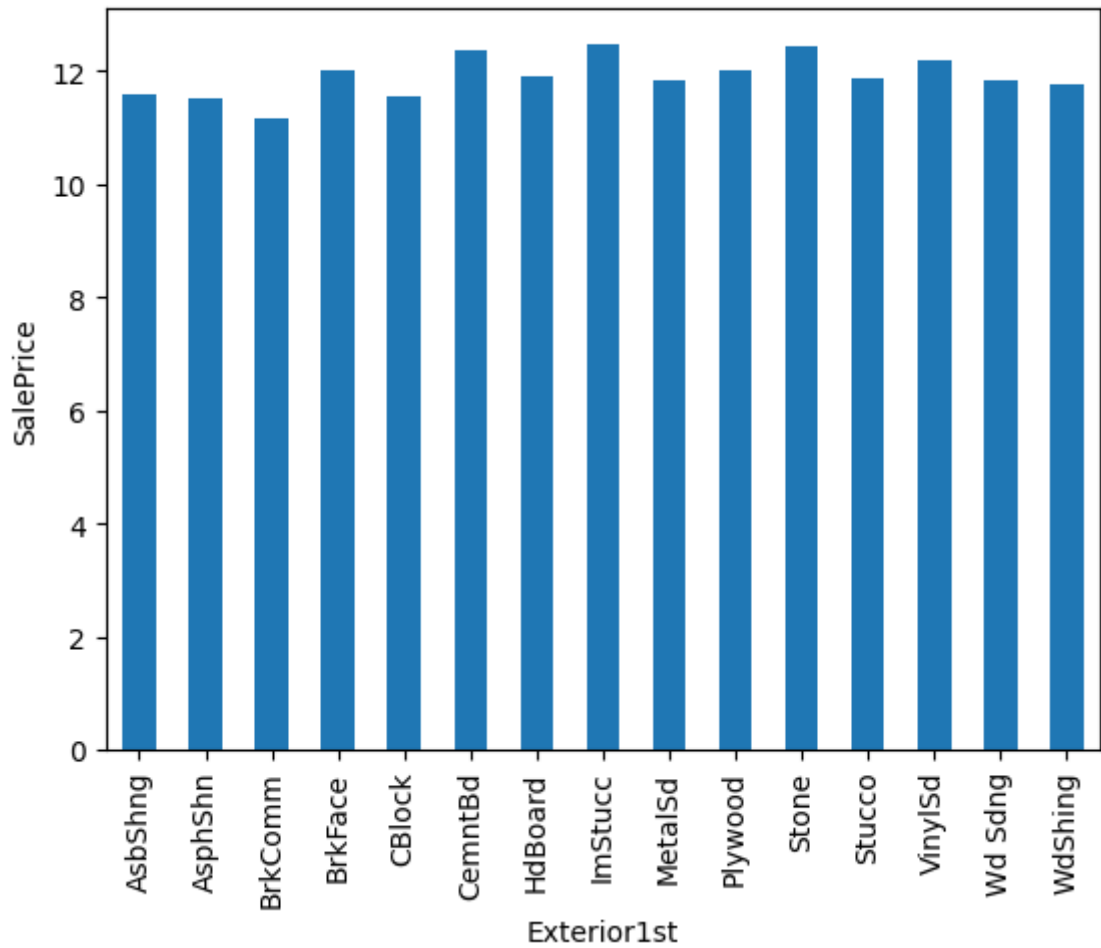


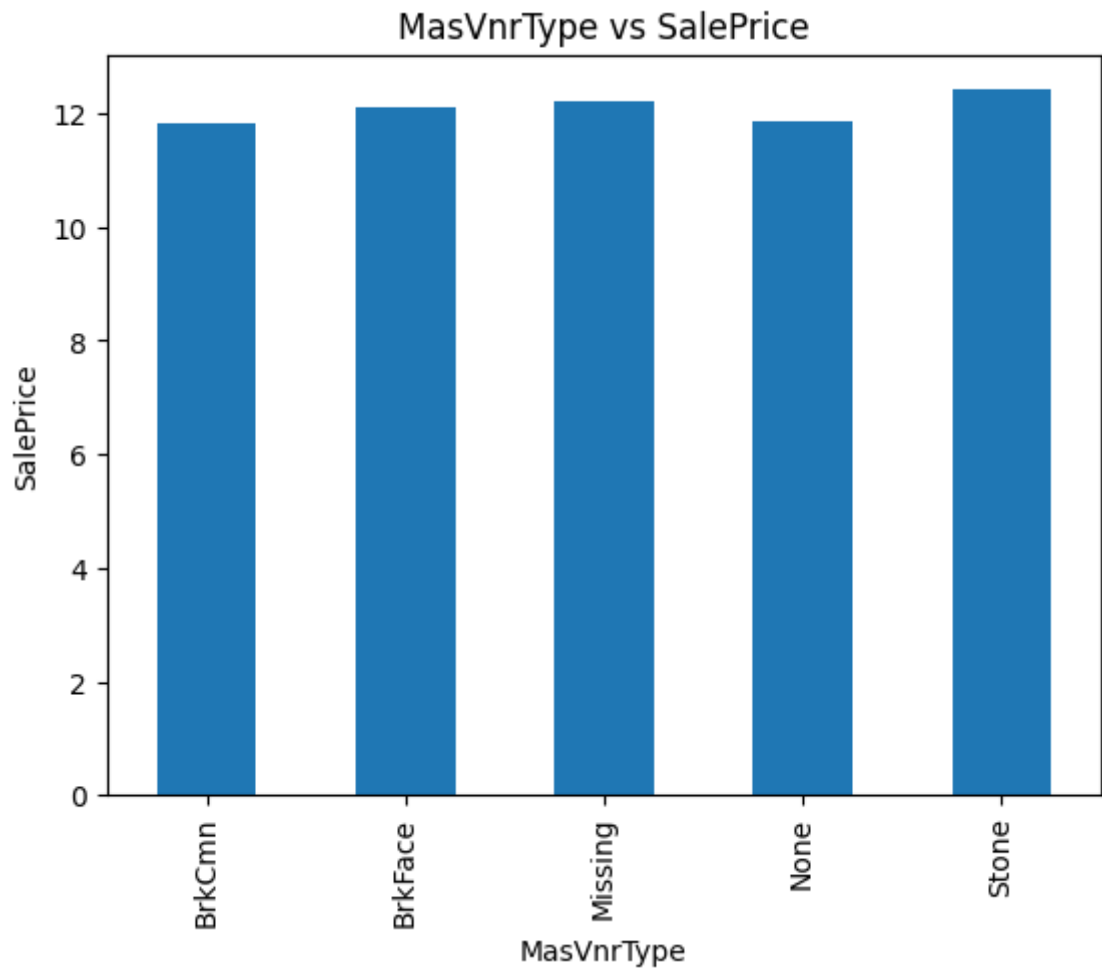
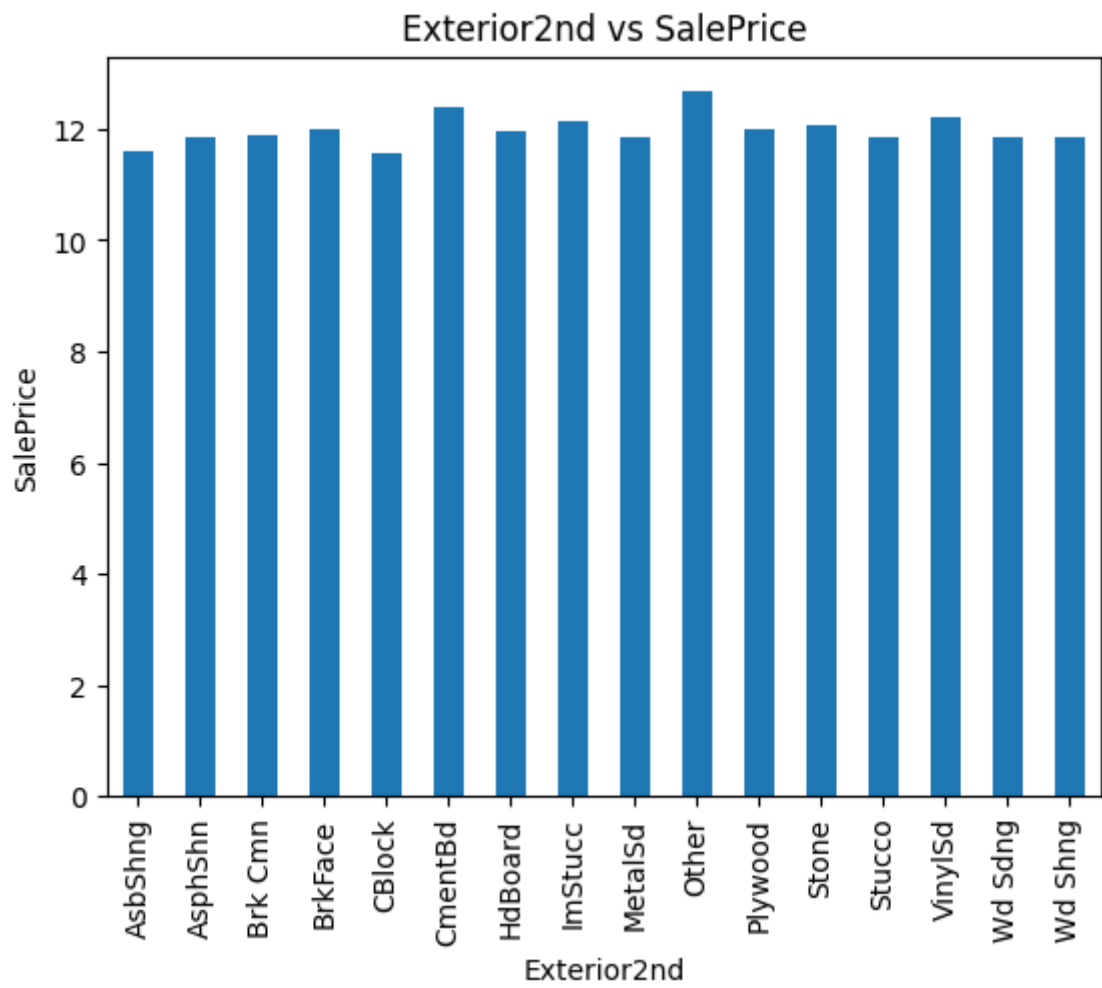


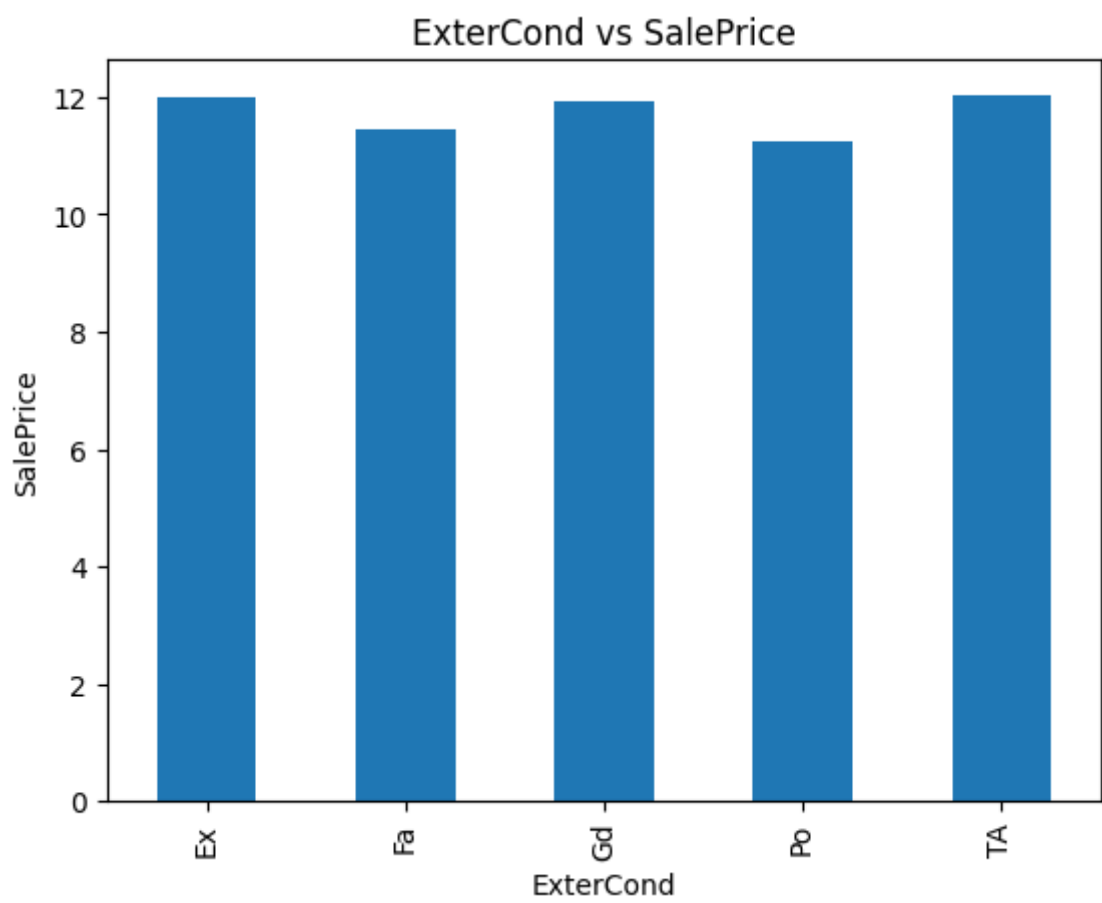
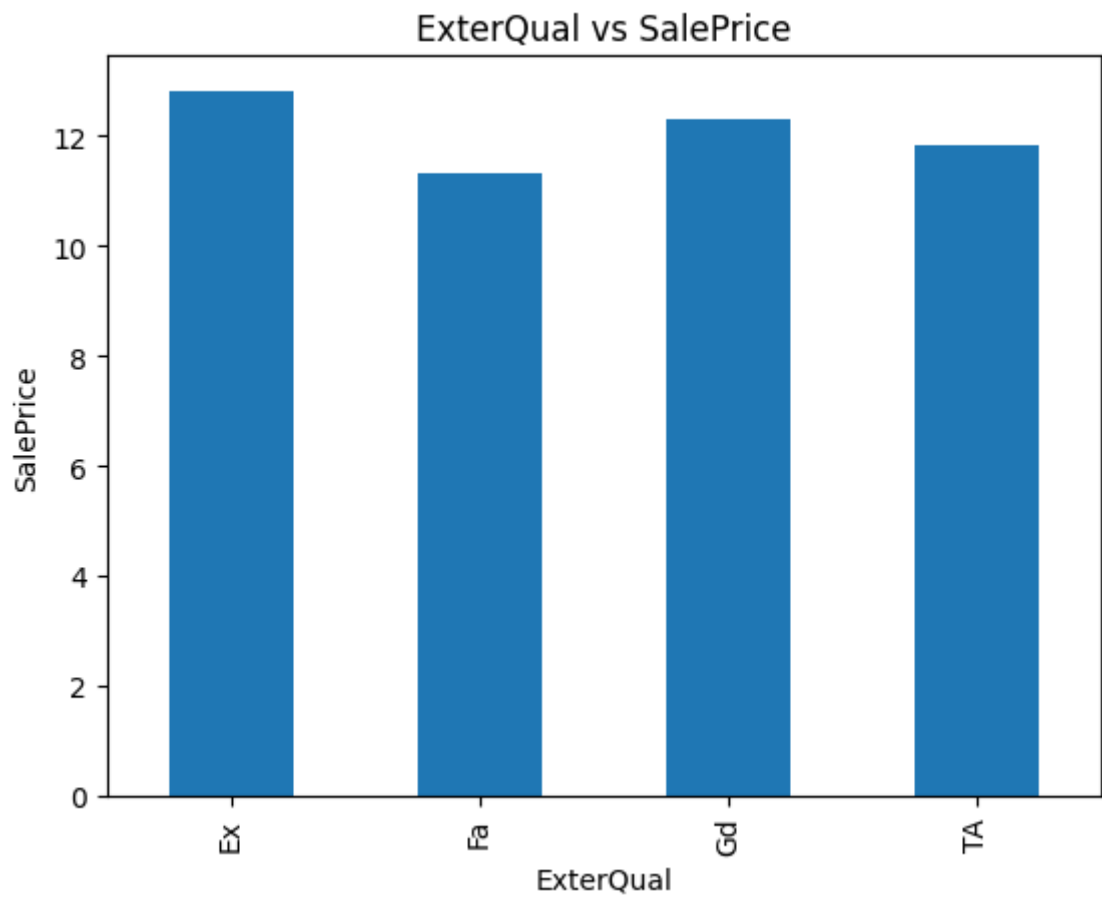
RoofMatl vs SalePrice

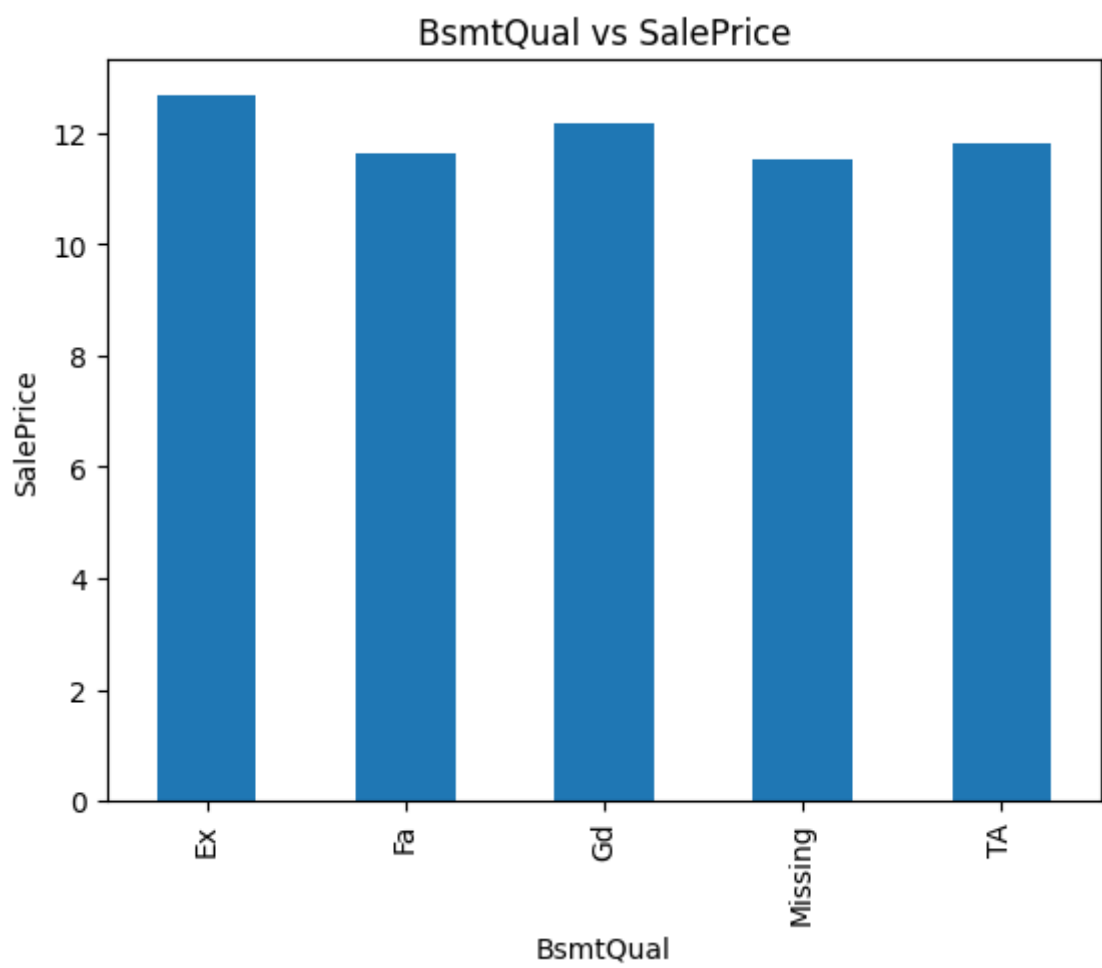
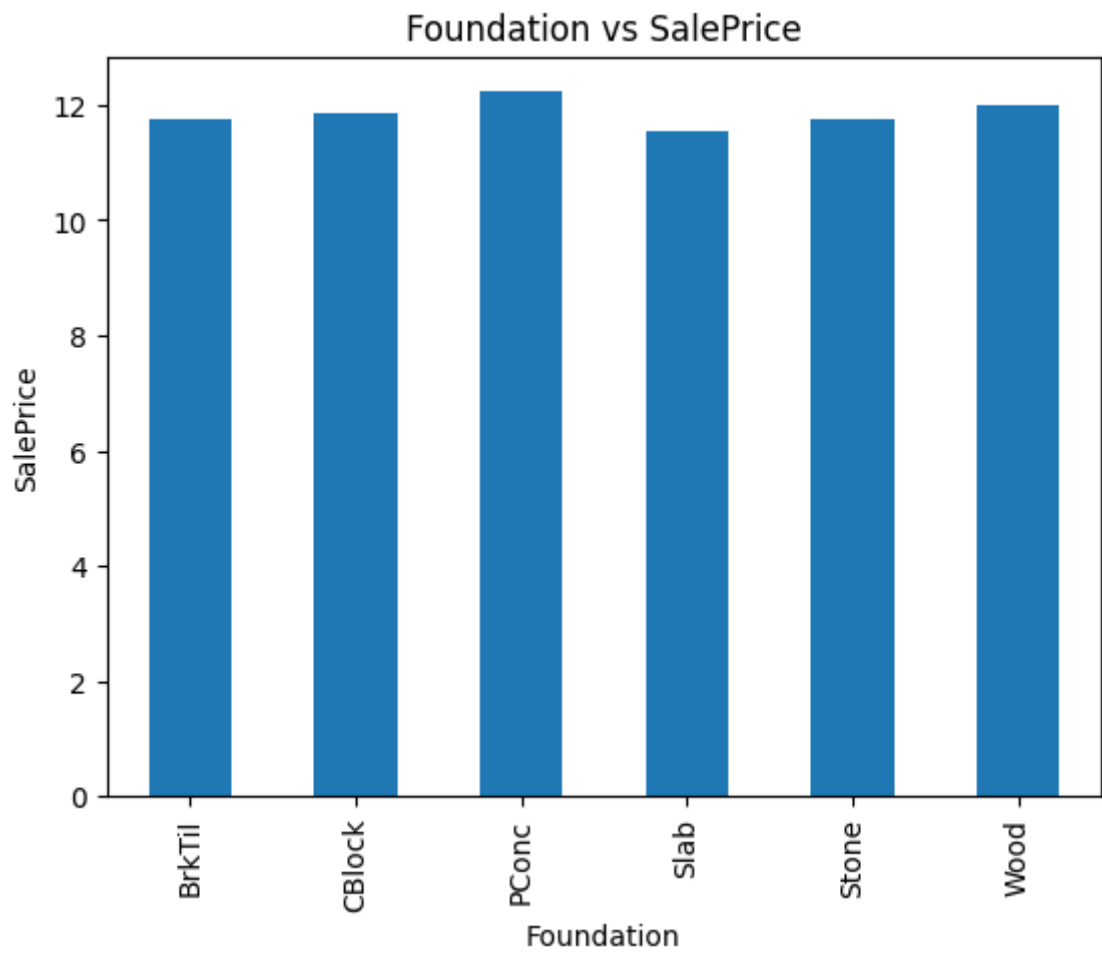


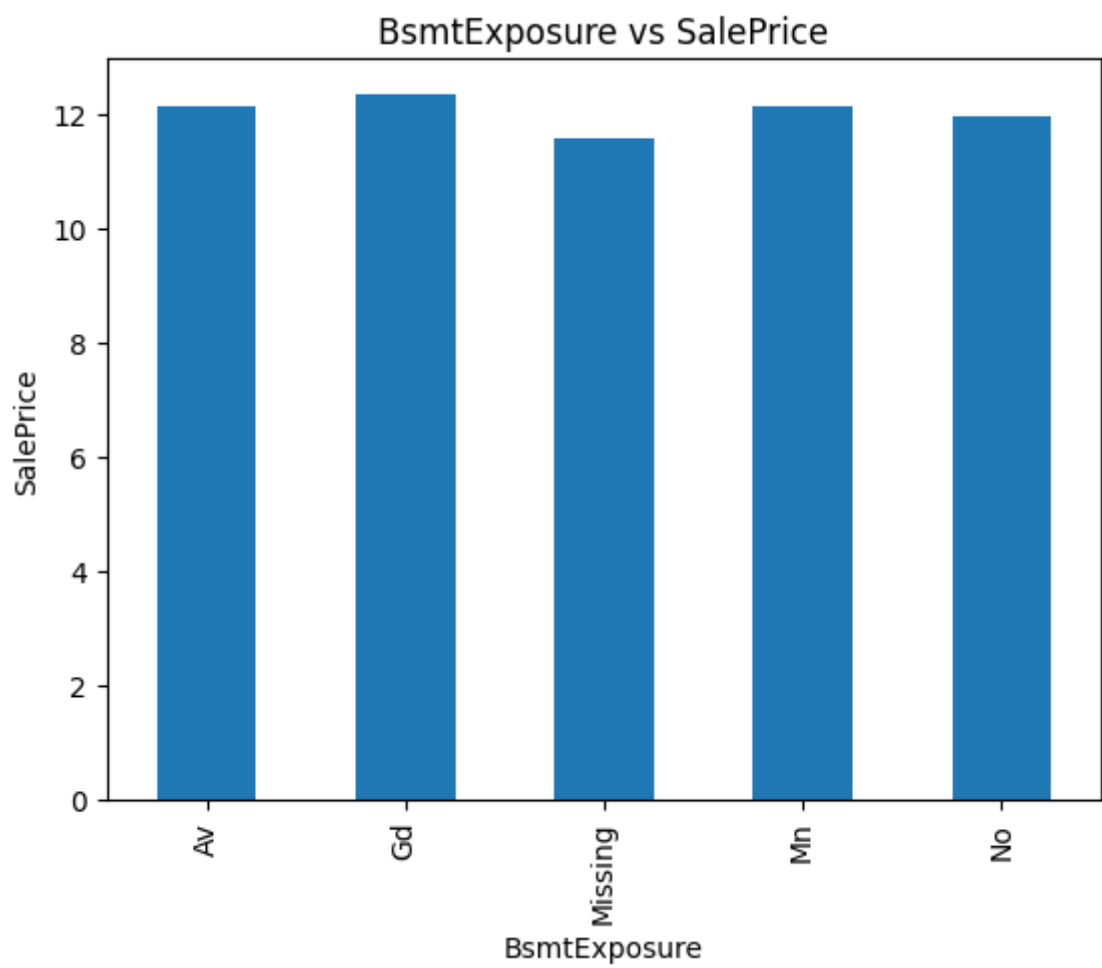
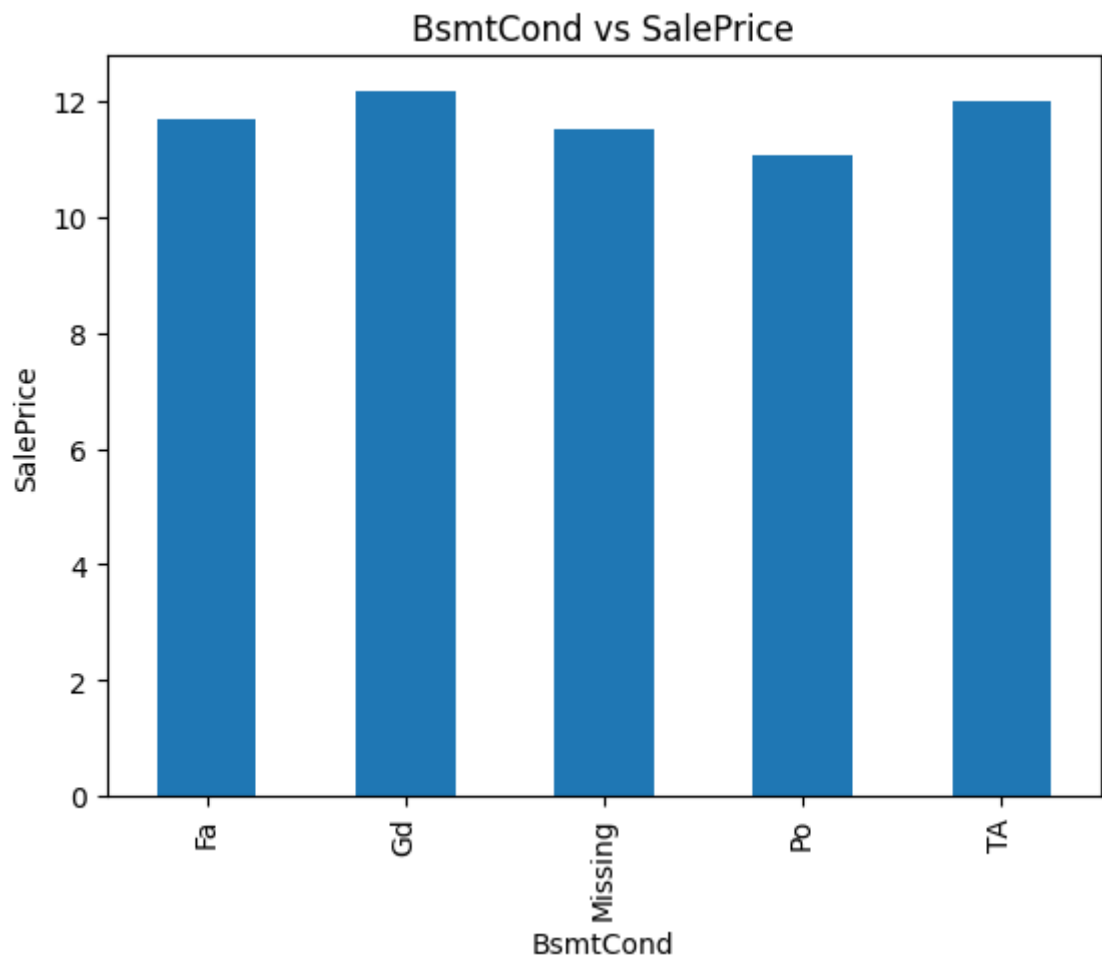
Exterior1st vs SalePrice

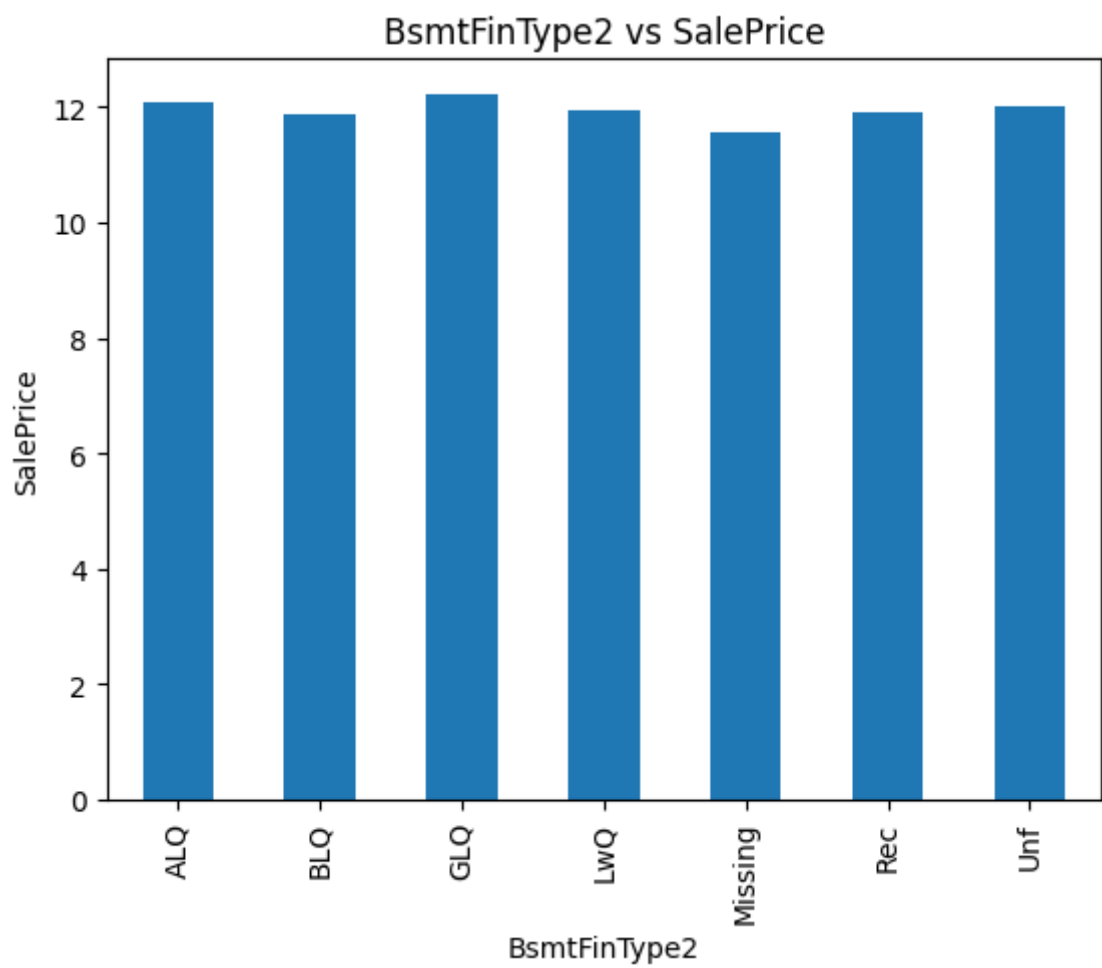
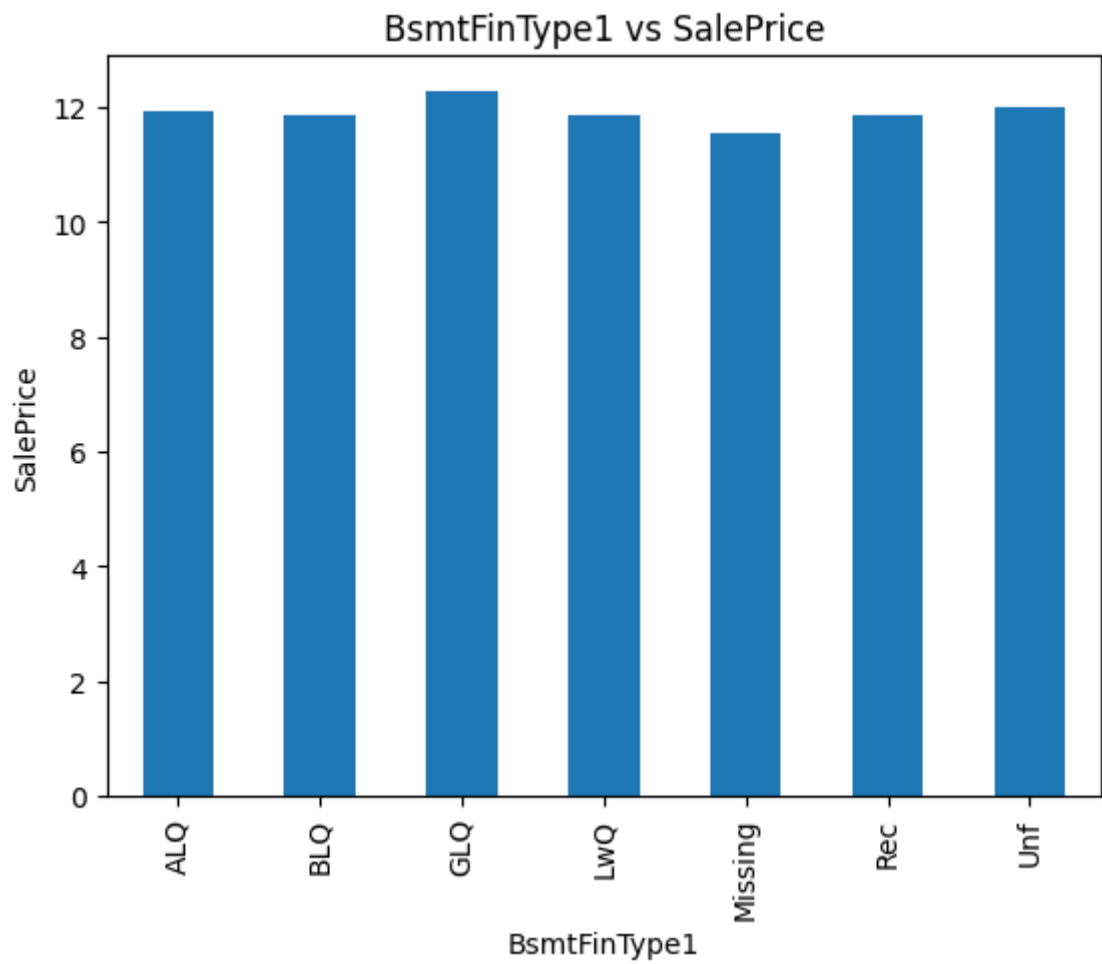


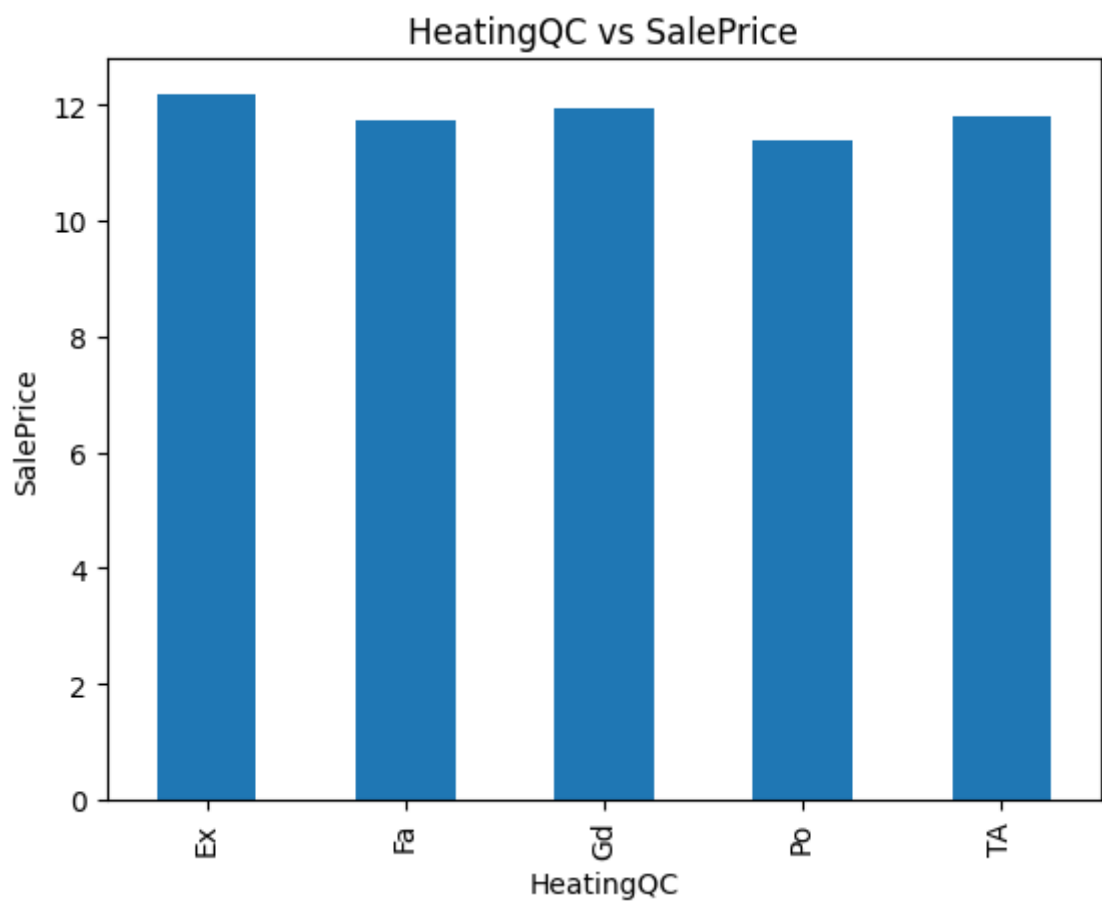
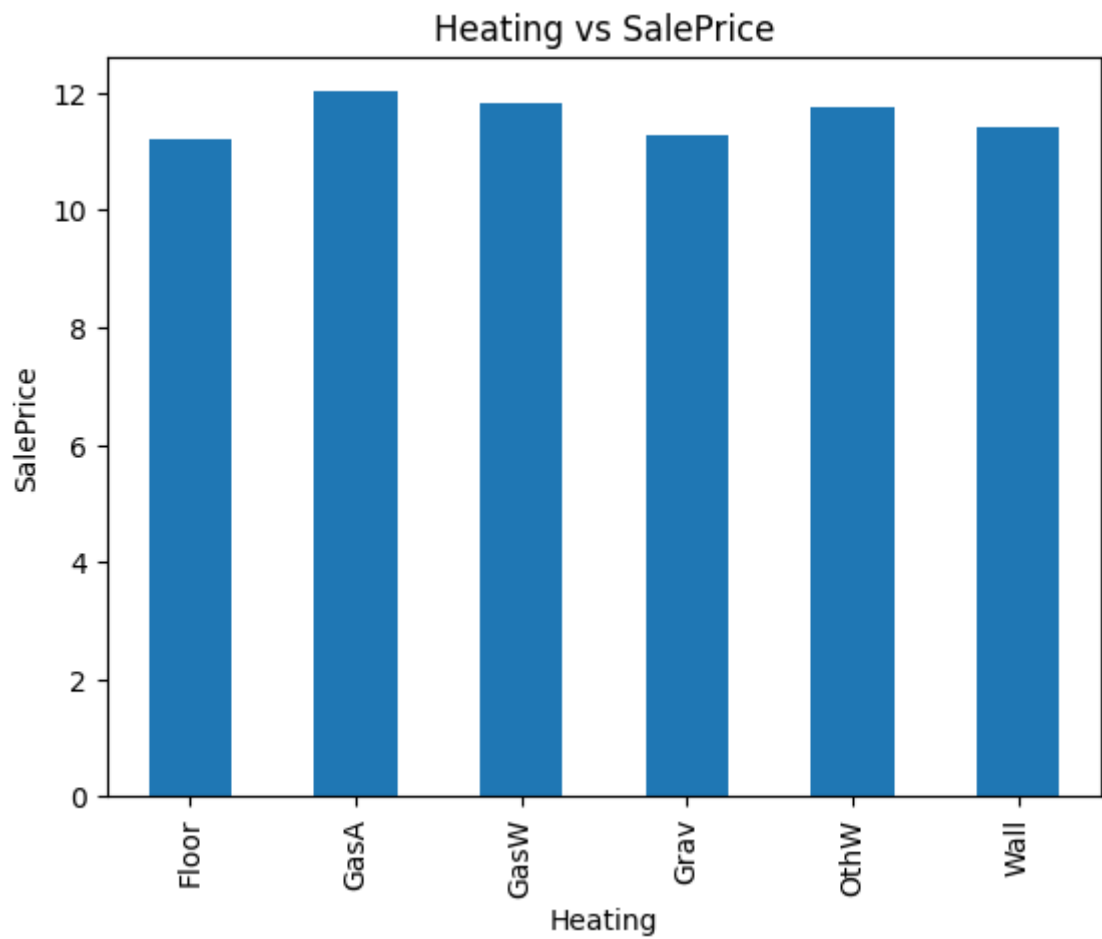


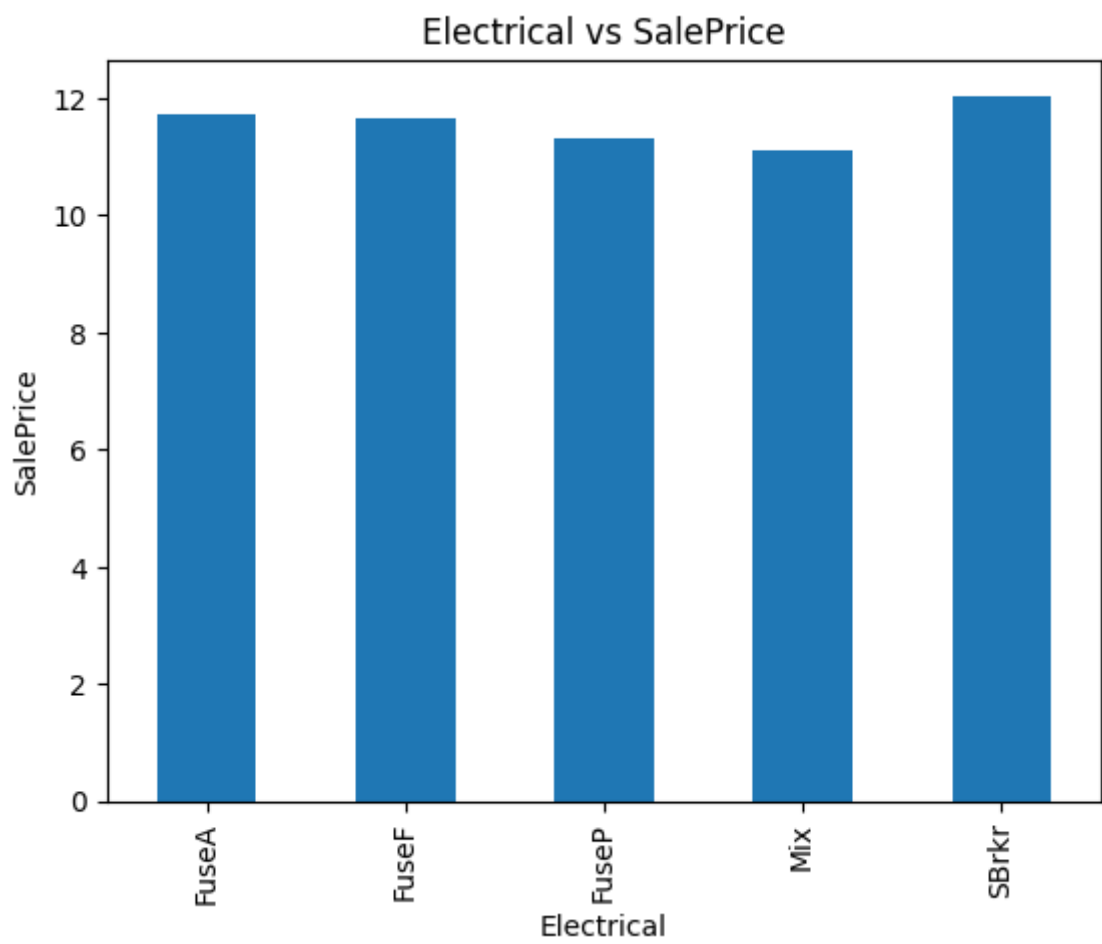
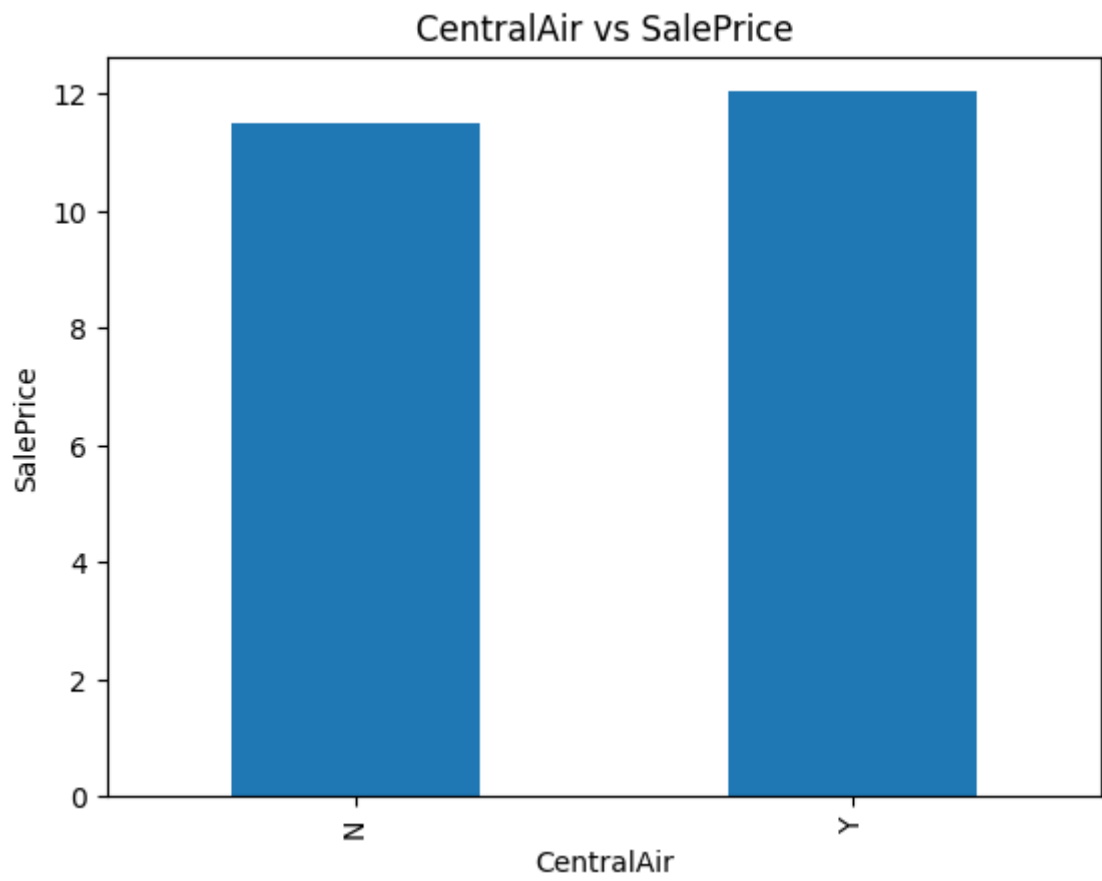




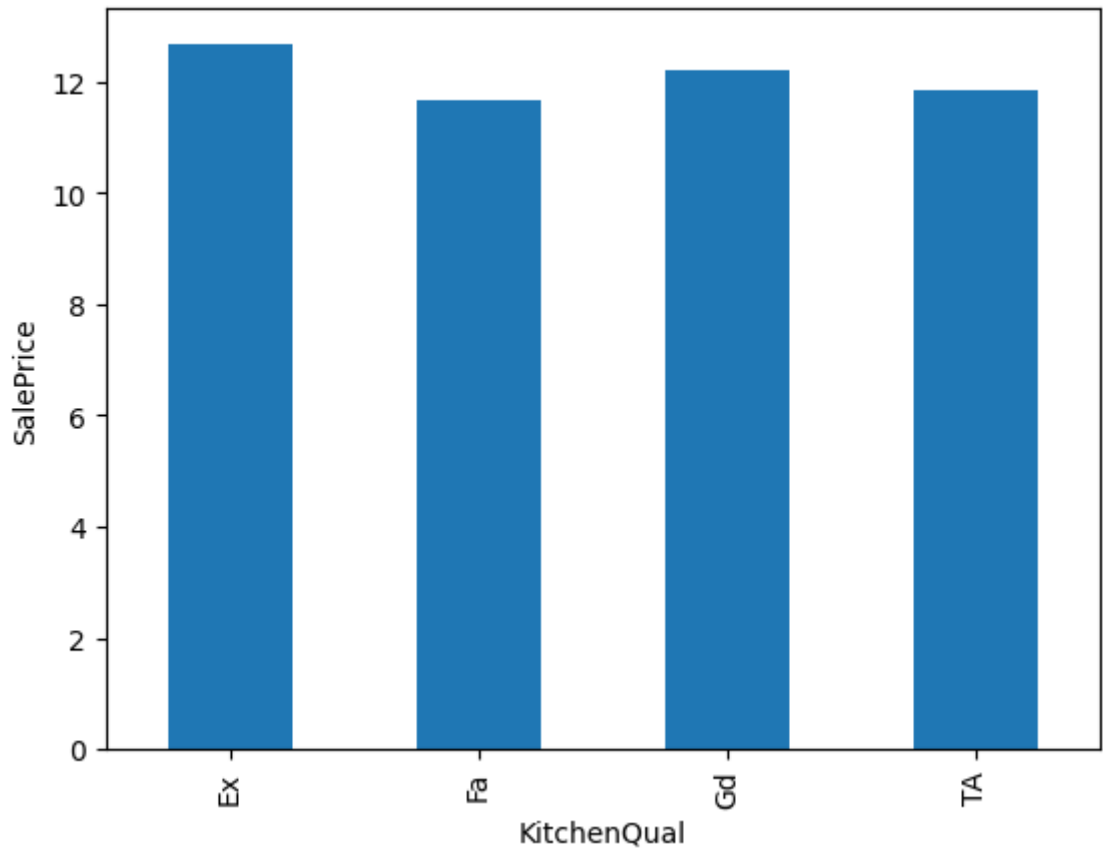




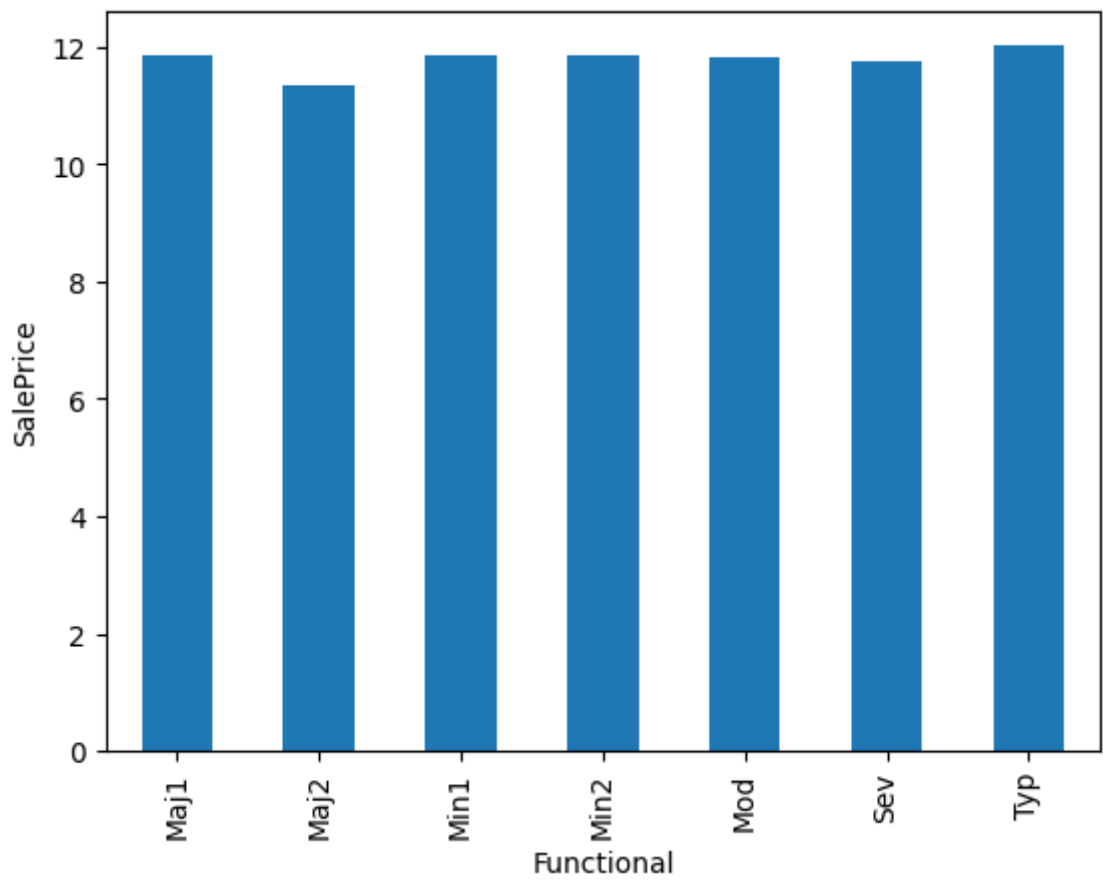




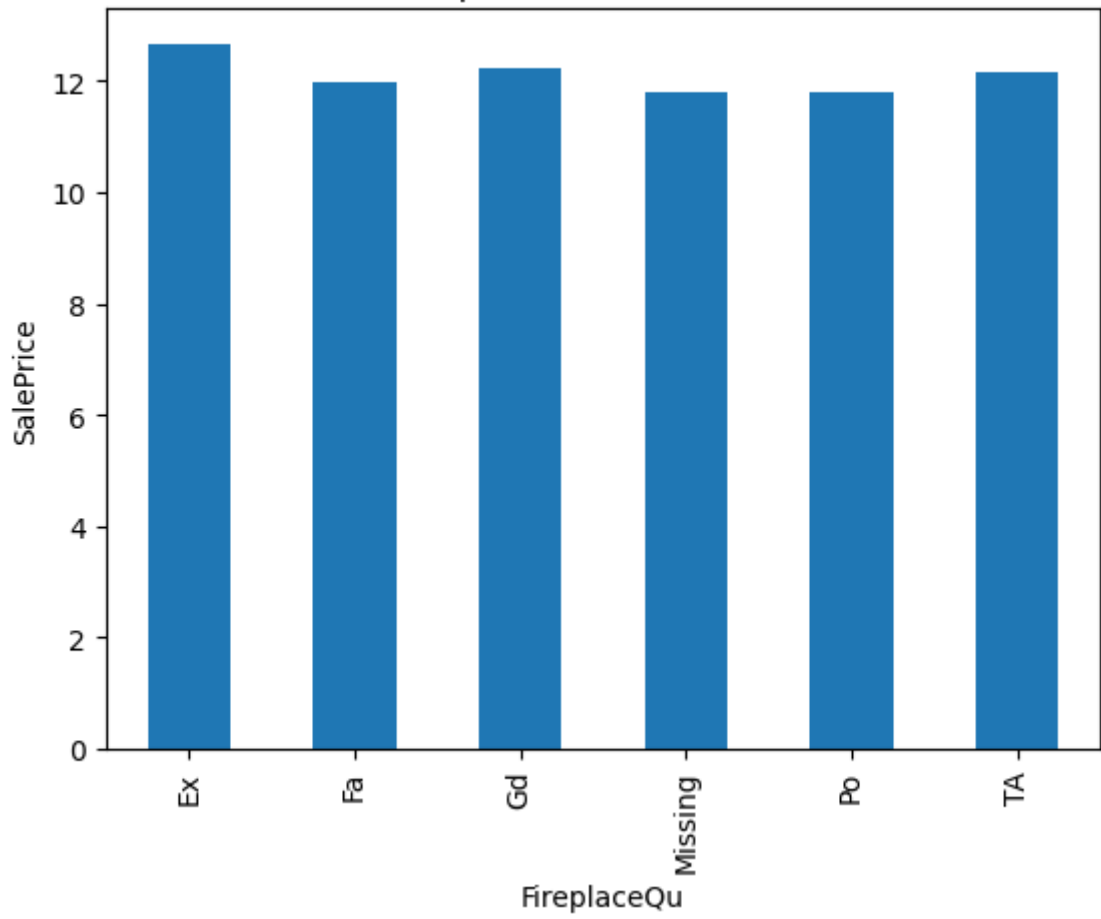
KitchenQual vs SalePrice



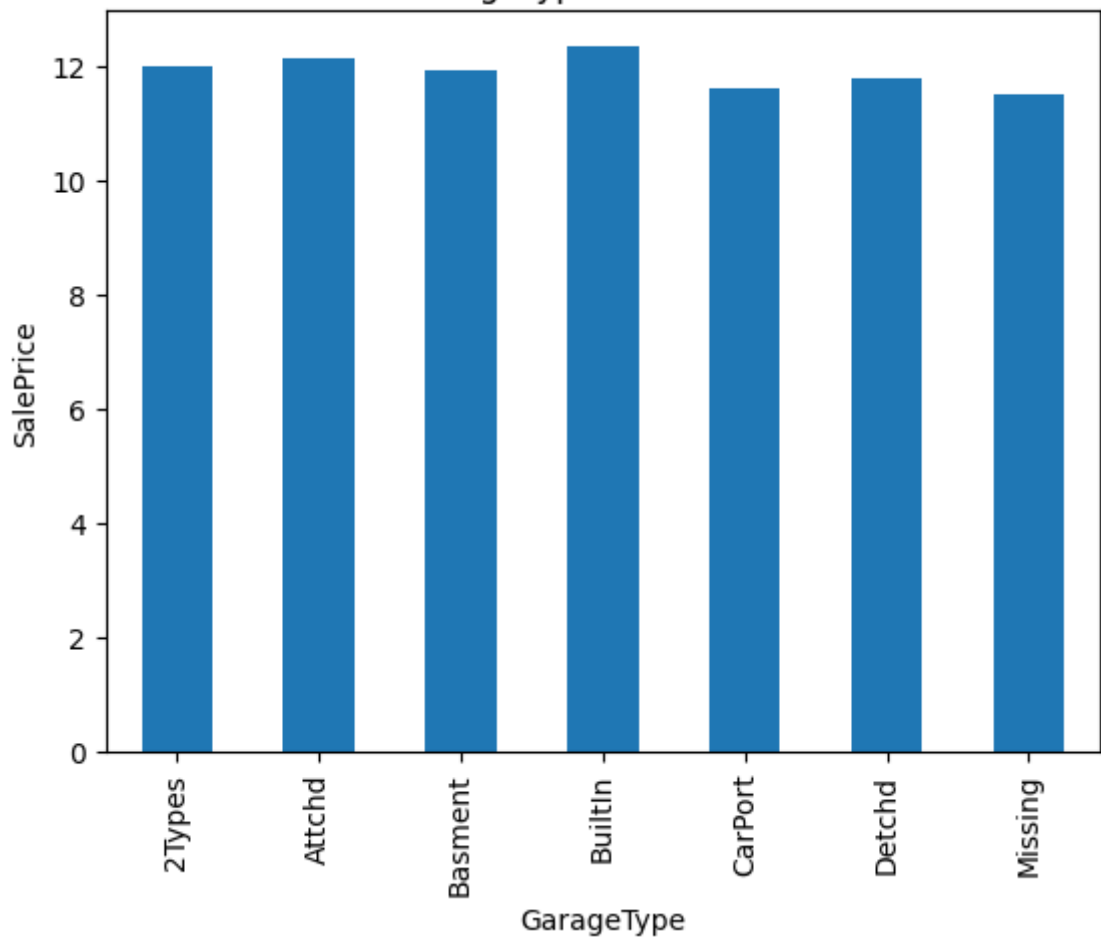
Functional vs SalePrice

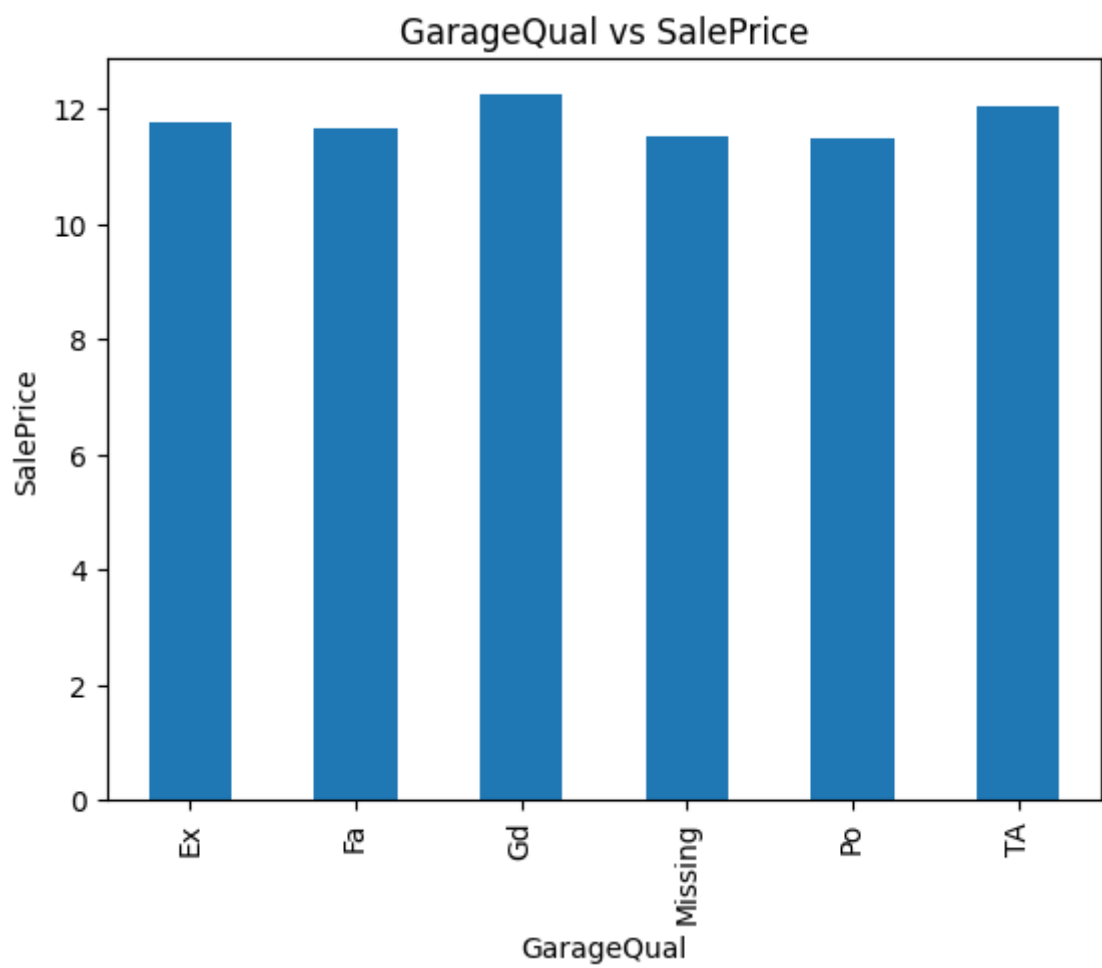
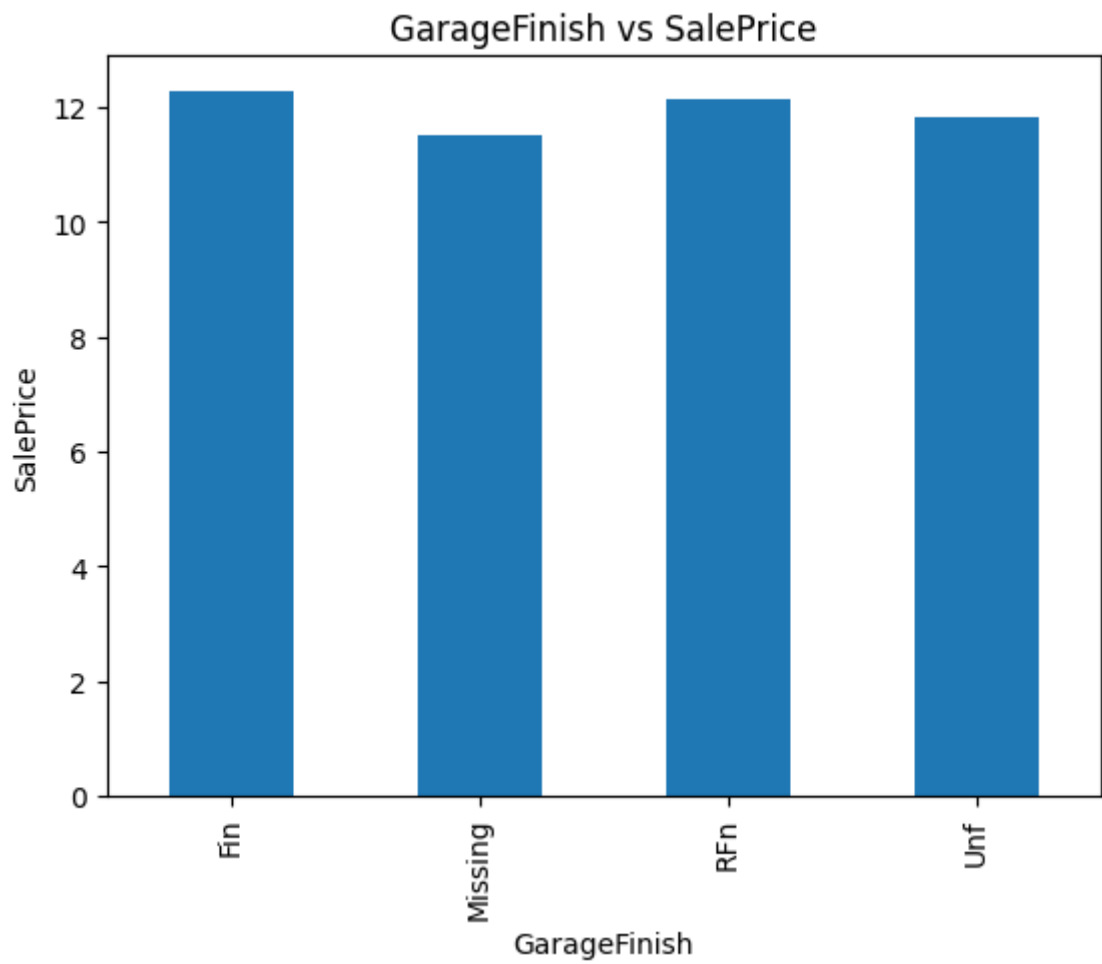


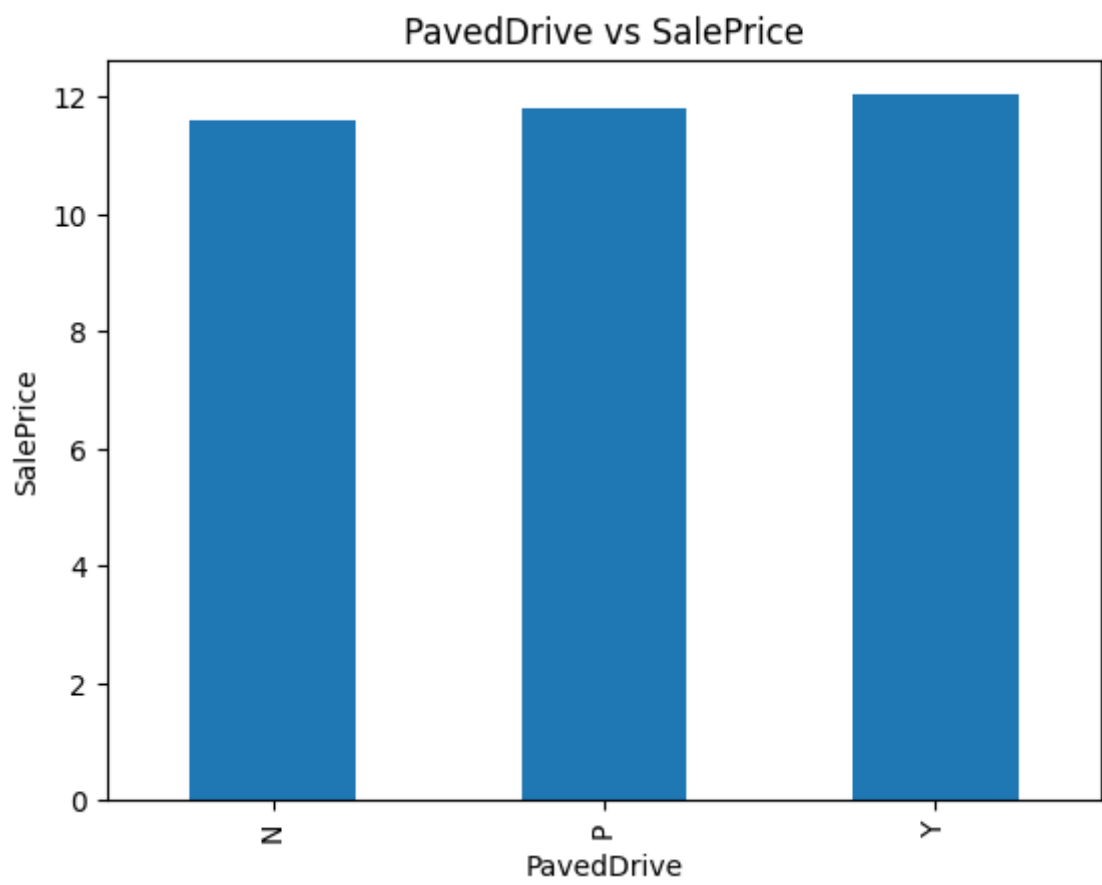
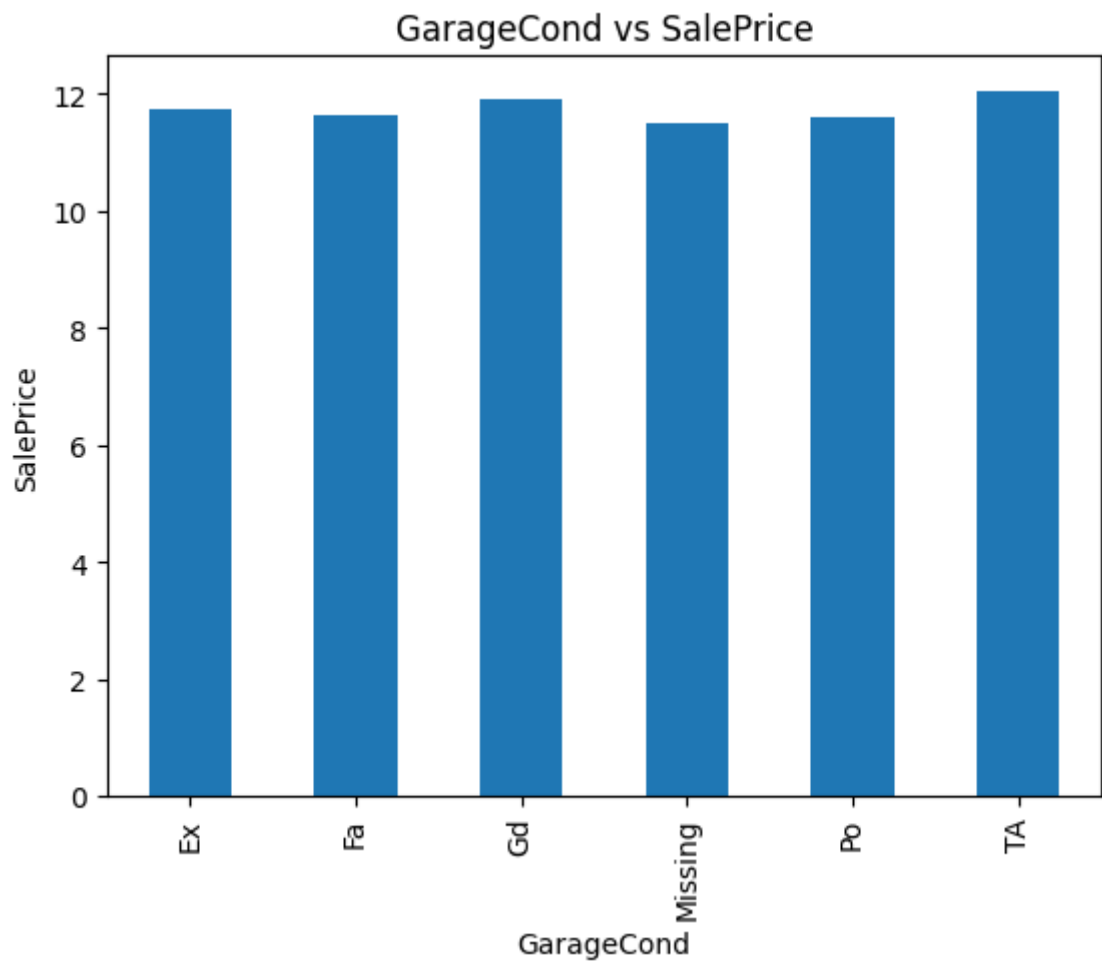
FireplaceQu vs SalePrice

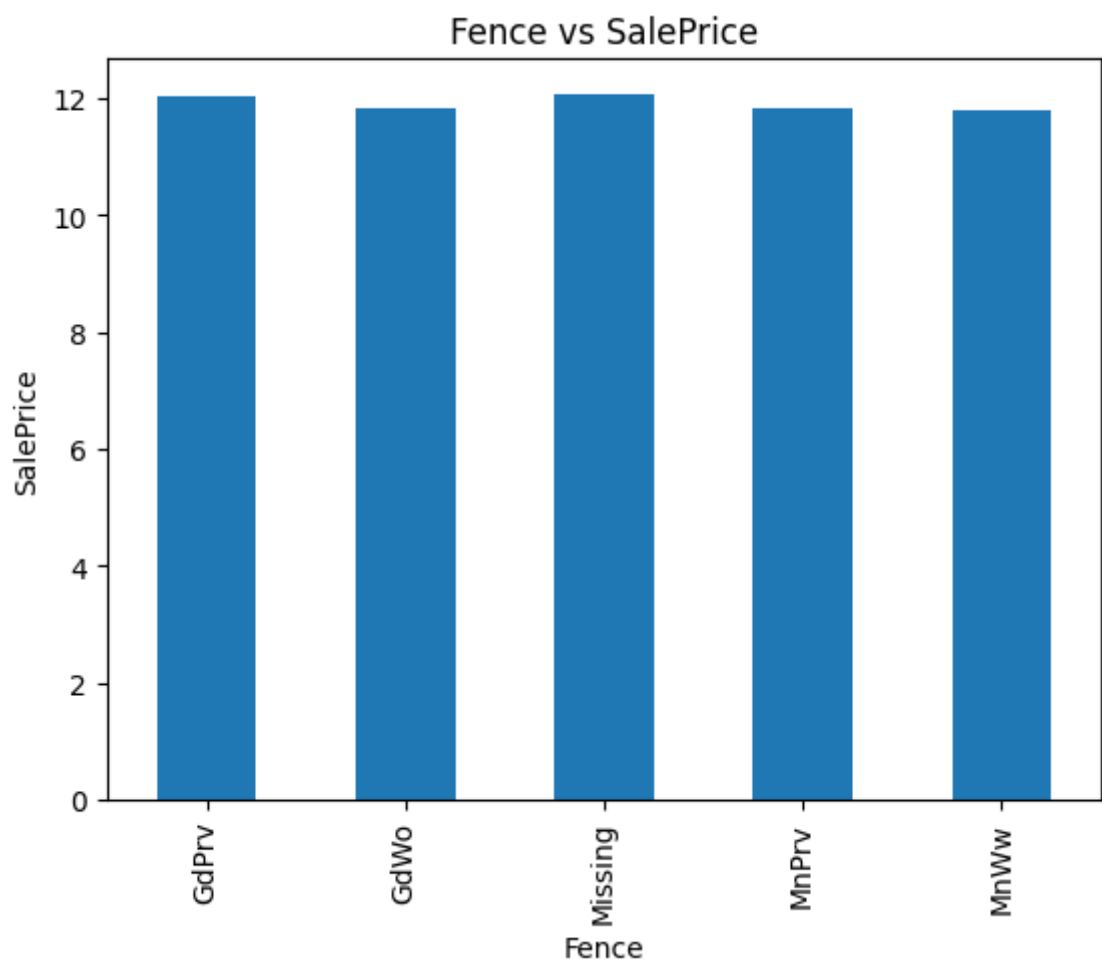
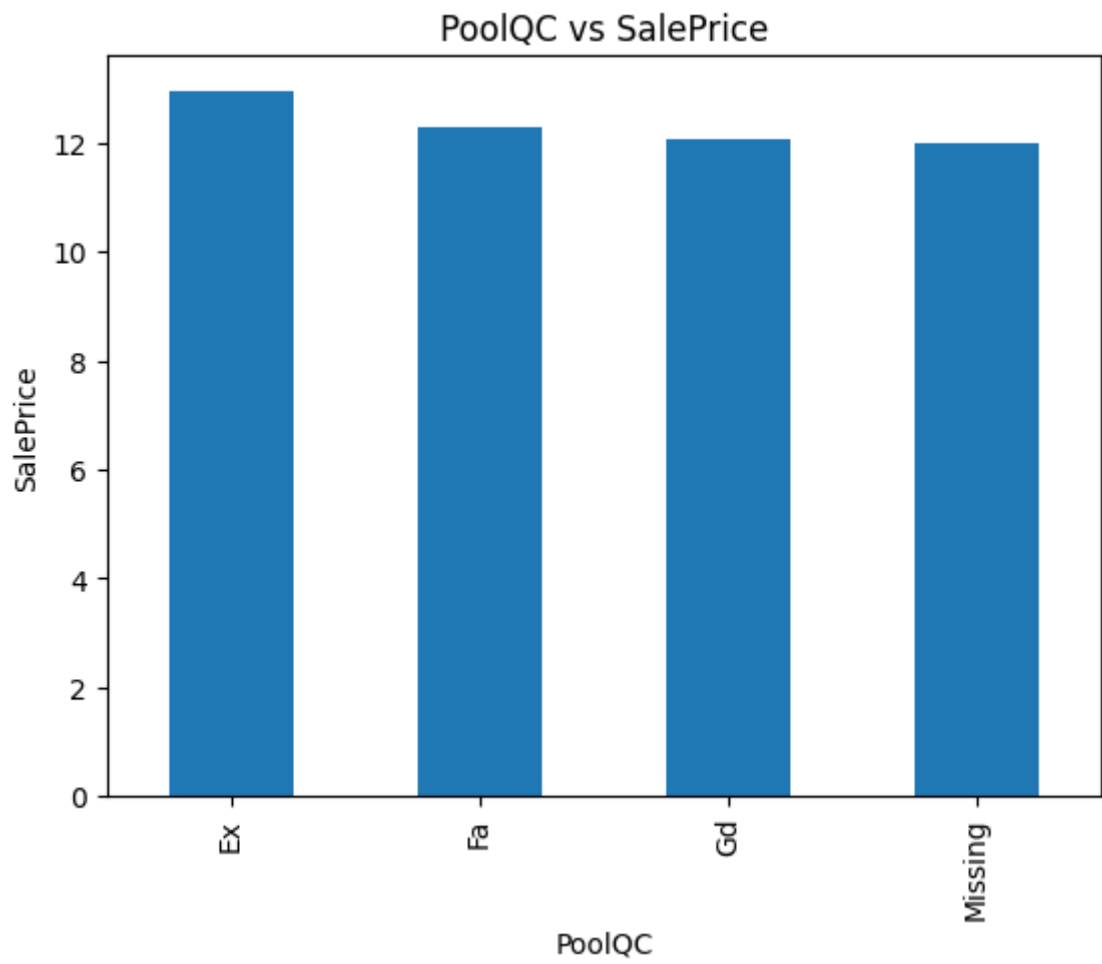


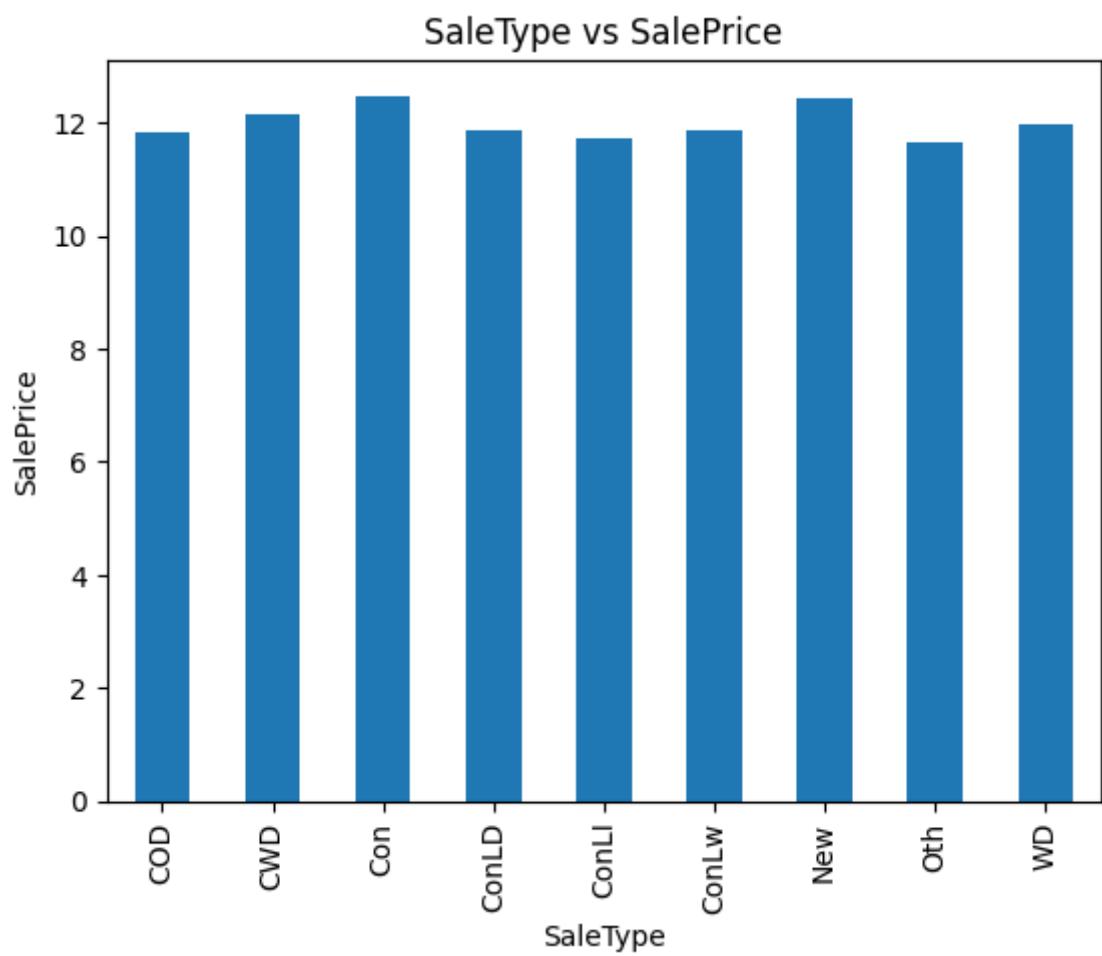
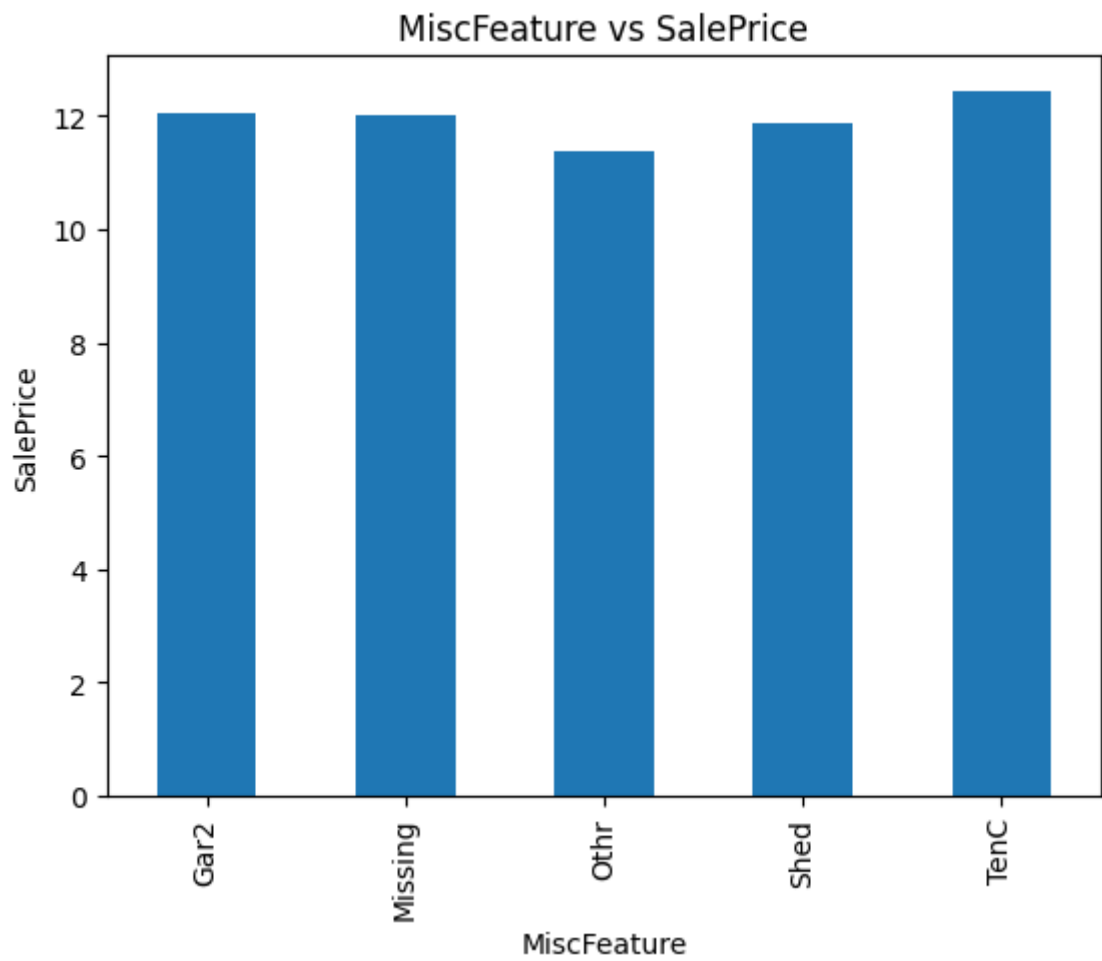
GarageType vs SalePrice

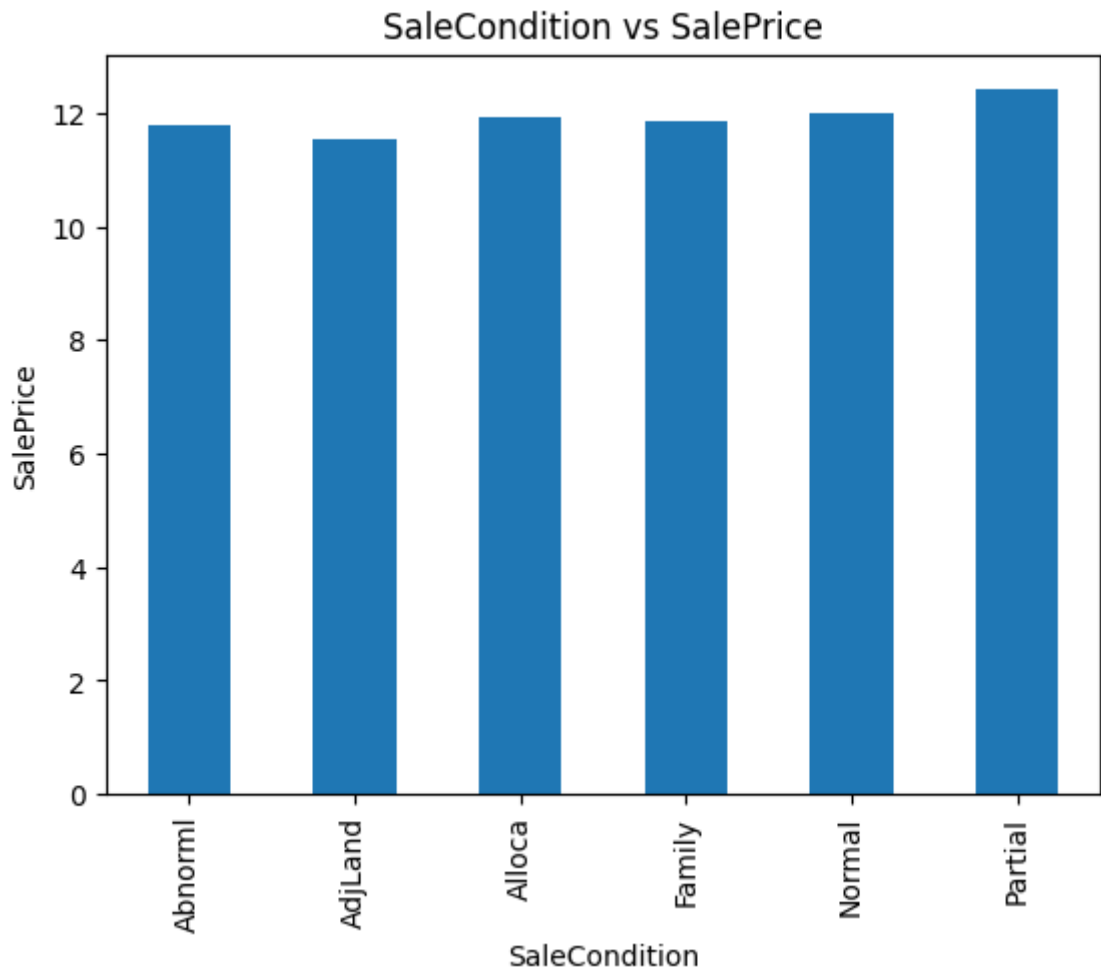












RARE CATEGORICAL VARIABLES

We will replace categorical variables that are present less than 1% of the observations by 'Rare_var'.

```
In [443... for feature in categorical_features:

    # group the dataset by each categorical feature and calculate the count of each
    # divide the count by the length of the dataset to get the relative frequency
    temp = dataset.groupby(feature)['SalePrice'].count() / len(dataset)

    # create a new DataFrame (temp_df) by filtering the categories that have a relative frequency > 0.01
    # The .index attribute retrieves the index values of the filtered categories.
    temp_df = temp[temp > 0.01].index

    dataset[feature] = np.where(dataset[feature].isin(temp_df), dataset[feature], 'Rare_var')

In [444... dataset.head(100)
```

Out[444]:

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContc
0	1	60	RL	4.174387	9.041922	Pave	Missing	Reg	
1	2	20	RL	4.382027	9.169518	Pave	Missing	Reg	
2	3	60	RL	4.219508	9.328123	Pave	Missing	IR1	
3	4	70	RL	4.094345	9.164296	Pave	Missing	IR1	
4	5	60	RL	4.430817	9.565214	Pave	Missing	IR1	
...	
95	96	60	RL	4.234107	9.186560	Pave	Missing	IR2	
96	97	20	RL	4.356709	9.236398	Pave	Missing	IR1	
97	98	20	RL	4.290459	9.298443	Pave	Missing	Reg	1
98	99	30	RL	4.442651	9.270965	Pave	Missing	Reg	
99	100	20	RL	4.343805	9.139918	Pave	Missing	IR1	

100 rows × 84 columns

```
In [445... for feature in categorical_features:
    labels_ordered = dataset.groupby([feature])['SalePrice'].mean().sort_values(
    labels_ordered = {k:i for i, k in enumerate(labels_ordered, 0)}
    dataset[feature] = dataset[feature].map(labels_ordered)
```

```
In [446... dataset.head(10)
```

Out[446]:

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour
0	1	60	3	4.174387	9.041922	1	2	0	1
1	2	20	3	4.382027	9.169518	1	2	0	1
2	3	60	3	4.219508	9.328123	1	2	1	1
3	4	70	3	4.094345	9.164296	1	2	1	1
4	5	60	3	4.430817	9.565214	1	2	1	1
5	6	50	3	4.442651	9.554993	1	2	1	1
6	7	20	3	4.317488	9.218705	1	2	0	1
7	8	60	3	4.234107	9.247829	1	2	1	1
8	9	50	1	3.931826	8.719317	1	2	0	1
9	10	190	3	3.912023	8.911934	1	2	0	1

FEATURE SCALING

(will not be performed on 'Id' column because it can be dropped later and 'SalePrice' column because it is the

dependennt variable.)

```
In [447... feature_scale = []
for feature in dataset.columns:
    if feature not in ['Id', 'SalePerice']:
        feature_scale.append(feature)

print(len(feature_scale))
dataset[feature_scale].head()
```

83

```
Out[447]:
```

	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	Utili
0	60	3	4.174387	9.041922	1	2	0	1	
1	20	3	4.382027	9.169518	1	2	0	1	
2	60	3	4.219508	9.328123	1	2	1	1	
3	70	3	4.094345	9.164296	1	2	1	1	
4	60	3	4.430817	9.565214	1	2	1	1	

```
In [448... from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
scaler.fit(dataset[feature_scale])
scaler.transform(dataset[feature_scale])
```

```
Out[448]: array([[0.23529412, 0.75      , 0.41820812, ..., 0.      , 0.      ,
                  0.      ],
                 [0.      , 0.75      , 0.49506375, ..., 0.      , 0.      ,
                  0.      ],
                 [0.23529412, 0.75      , 0.434909  , ..., 0.      , 0.      ,
                  0.      ],
                 ...,
                 [0.29411765, 0.75      , 0.42385922, ..., 0.      , 0.      ,
                  0.      ],
                 [0.      , 0.75      , 0.434909  , ..., 0.      , 0.      ,
                  0.      ],
                 [0.      , 0.75      , 0.47117546, ..., 0.      , 0.      ,
                  0.      ]])
```

```
In [449... # transform the train and test set, and add on the Id and SalePrice variables
data = pd.concat([dataset[['Id', 'SalePrice']].reset_index(drop=True),
                  pd.DataFrame(scaler.transform(dataset[feature_scale]), column
                               axis=1)
```

```
In [450... data.head()
```

```
Out[450]:
```

	Id	SalePrice	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	Lan
0	1	12.247694	0.235294	0.75	0.418208	0.366344	1.0	1.0	0.000000	
1	2	12.109011	0.000000	0.75	0.495064	0.391317	1.0	1.0	0.000000	
2	3	12.317167	0.235294	0.75	0.434909	0.422359	1.0	1.0	0.333333	
3	4	11.849398	0.294118	0.75	0.388581	0.390295	1.0	1.0	0.333333	
4	5	12.429216	0.235294	0.75	0.513123	0.468761	1.0	1.0	0.333333	

```
In [451... data.to_csv('X_train.csv', index = False)
```

FEATURE SELECTION - ADVANCE HOUSE PRICE PREDICTION

The main aim of this project is to predict the house price based on various features.

IMPORT LIBRARIES

```
In [452... from sklearn.linear_model import Lasso
from sklearn.feature_selection import SelectFromModel
from sklearn.preprocessing import OneHotEncoder
```

READ DATA

```
In [453... dataset_fs = pd.read_csv('X_train.csv')
```

```
In [454... dataset_fs.head()
```

```
Out[454]:
```

	Id	SalePrice	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	Lan
0	1	12.247694	0.235294	0.75	0.418208	0.366344	1.0	1.0	0.000000	
1	2	12.109011	0.000000	0.75	0.495064	0.391317	1.0	1.0	0.000000	
2	3	12.317167	0.235294	0.75	0.434909	0.422359	1.0	1.0	0.333333	
3	4	11.849398	0.294118	0.75	0.388581	0.390295	1.0	1.0	0.333333	
4	5	12.429216	0.235294	0.75	0.513123	0.468761	1.0	1.0	0.333333	

```
In [455... # Capture the dependent feature
y_train = dataset_fs[['SalePrice']]
```

```
In [456... # drop dependent feature from dataset
X_train = dataset_fs.drop(['Id', 'SalePrice'], axis = 1)
```

```
In [457... ## Apply Feature Selection
# first, specify the Lasso Regression model, and
# select a suitable alpha (equivalent of penalty).
# The bigger the alpha the less features that will be selected.

# Then use the selectFromModel object from sklearn, which
# will select the features whose coefficients are non-zero

feature_sel_model = SelectFromModel(Lasso(alpha=0.005, random_state=0)) # rememb
feature_sel_model.fit(X_train, y_train)
```

```
Out[457]: ▸ SelectFromModel
          ▸ estimator: Lasso
              ▸ Lasso
```

```
In [458... feature_sel_model.get_support()
```

```
Out[458]: array([False, False, False, False, False, False, False, False, False,
                False, False, True, False, False, False, False, False, False,
                False, True, False, False, False, False, False, False, False,
                False, False, False, False, False, False, False, False, False,
                False, False, False, False, False, False, False, False, False,
                False, False, True, False, False, True, False, False, False,
                False, False, False, False, False, False, False, False, False,
                False, False, False, False, False, False, False, False, True,
                False, False, False])
```

```
In [459... # print the number of total and selected features

# make a list of the selected features
selected_feat = X_train.columns[(feature_sel_model.get_support())]

# print stats
print('total features: {}'.format((X_train.shape[1])))
print('selected features: {}'.format(len(selected_feat)))
print('features with coefficients shrank to zero: {}'.format(
    np.sum(feature_sel_model.estimator_.coef_ == 0)))

total features: 83
selected features: 5
features with coefficients shrank to zero: 78
```

```
In [460... selected_feat
```

```
Out[460]: Index(['Neighborhood', 'YearRemodAdd', 'FireplaceQu', 'GarageFinish',
                'SalePrice.1'],
                dtype='object')
```

```
In [461... X_train = X_train[selected_feat]
```

```
In [462... X_train.head()
```


Out[462]:

	Neighborhood	YearRemodAdd	FireplaceQu	GarageFinish	SalePrice.1
0	0.636364	0.098361	0.2	0.666667	0.581431
1	0.500000	0.524590	0.6	0.666667	0.536319
2	0.636364	0.114754	0.6	0.666667	0.604029
3	0.727273	0.606557	0.8	0.333333	0.451871
4	1.000000	0.147541	0.6	0.666667	0.640477

In [462...]