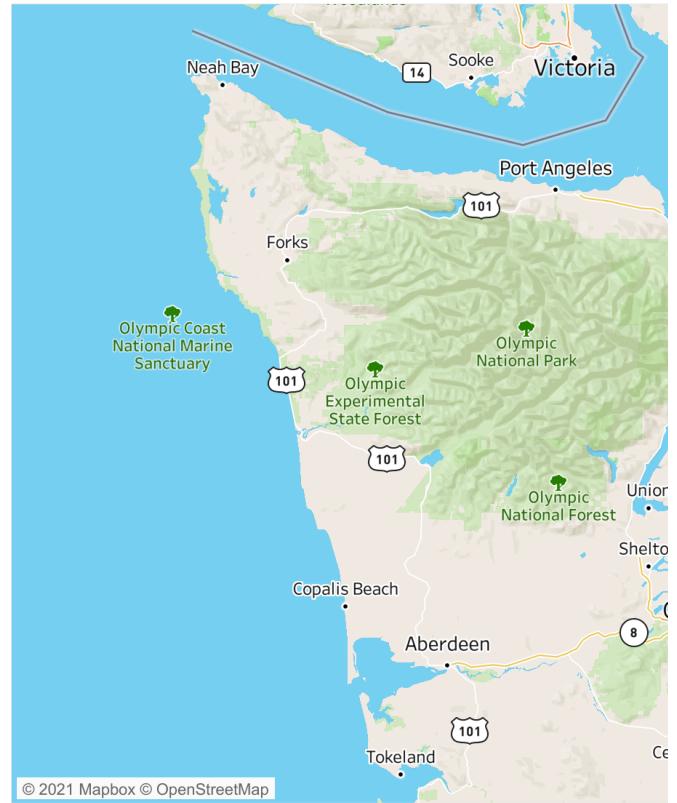


1 King County Housing Ch

King County Zip codes by Median Price



Using multiple linear regression analysis models to ir

Business problem:

King County home sales have been increasing as Se live and work. Our real estate team has been tasked ensure the price is accurate compared to the market

The model also guides clients on which features to p suggestions within the owners control.

1.1 Data

- 21,597 rows by 21 columns
- CSV Formatted

1.2 Roadmap

Contents ⚙

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' cc
 - 2.2.2 Fill in missing Values for 'yr_reno'
 - 2.2.3 Fill in missing Values for 'waterfr
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, out
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary val
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the
 - 4.3 Comparison of Square Foot living an
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearit
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Accross All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Rem
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers R
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Accross All Col
 - 7.2.2 Model 3: Z-Score All Outliers Re
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outli
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Cat
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling F
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the mos
 - 9.2.1 Digging Deep into Zip codes

- Scrub data to handle null values and duplicates
- Add additional features to better infer the price
- Check for linearity and multicollinearity to make
- Perform outlier removal methods to better meet
- Use One Hot Encoding to handle categorical vari
- Provide accompanying visualizations to support
- Circle back to how the multiple linear regression

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot
 - 4.3 Comparison of Square Foot living area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [1]:

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 pd.set_option('display.max_columns', None)
5
6 df = pd.read_csv('data/kc_house_data.csv')
7 df.head()

```

Out[1]:

	id	date	price	bedrooms	bathrooms
0	7129300520	10/13/2014	221900.0	3	2
1	6414100192	12/9/2014	538000.0	3	2
2	5631500400	2/25/2015	180000.0	2	1
3	2487200875	12/9/2014	604000.0	4	3
4	1954400510	2/18/2015	510000.0	3	2

2.1 Descriptions of columns

- **id** - unique identifier for a house
- **date** - date the house was sold
- **price** - prediction target
- **bedrooms** - number of bedrooms/house
- **bathrooms** - number of bathrooms/bedrooms
- **sqft_living** - square footage of the home
- **sqft_lot** - square footage of the lot
- **floors** - floors (levels) in house
- **waterfront** - house which has a view to a water body
- **view** - quality of view
- **condition** - how good the condition is (overall)
- **grade** - overall grade given to the housing unit, based on various factors
- **sqft_above** - square footage of house apart from basement
- **sqft_basement** - square footage of the basement
- **yr_built** - built year
- **yr_renovated** - year when house was renovated
- **zipcode** - zip code

- **lat** - Latitude coordinate
- **long** - Longitude coordinate
- **sqft_living15** - The square footage of interior home
- **sqft_lot15** - The square footage of the land lots

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot size
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Fewer Variables
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [2]:

```
1 # Evaluating if type matches column
2
3 pd.set_option('display.float_format',
4 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21597 entries, 0 to 21596
Data columns (total 21 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   id               21597 non-null   int64  
 1   date              21597 non-null   object 
 2   price             21597 non-null   float64
 3   bedrooms          21597 non-null   int64  
 4   bathrooms         21597 non-null   float64
 5   sqft_living       21597 non-null   int64  
 6   sqft_lot          21597 non-null   int64  
 7   floors             21597 non-null   float64
 8   waterfront         19221 non-null   float64
 9   view               21534 non-null   float64
 10  condition         21597 non-null   int64  
 11  grade              21597 non-null   int64  
 12  sqft_above         21597 non-null   int64  
 13  sqft_basement      21597 non-null   object 
 14  yr_built           21597 non-null   int64  
 15  yr_renovated       17755 non-null   float64
 16  zipcode            21597 non-null   int64  
 17  lat                21597 non-null   float64
 18  long               21597 non-null   float64
 19  sqft_living15      21597 non-null   int64  
 20  sqft_lot15          21597 non-null   int64  
dtypes: float64(8), int64(11), object(2)
memory usage: 3.5+ MB
```

- Date should be a datetime object
- Sqft basement should be an integer, not object

In [3]:

```
1 # Make date into datetime object
2
3 df['date'] = pd.to_datetime(df['date'])
```

In [4]:

```

1 # Observe summary statistics
2
3 df.describe()

```

Out[4]:

	id	price	bedrooms	bathrc
count	21597.00	21597.00	21597.00	21597.00
mean	4580474287.77	540296.57		3.37
std	2876735715.75	367368.14		0.93
min	1000102.00	78000.00		1.00
25%	2123049175.00	322000.00		3.00
50%	3904930410.00	450000.00		3.00
75%	7308900490.00	645000.00		4.00
max	9900000190.00	7700000.00		33.00

- ID is a random value so should not be evaluated
- Price has a large standard deviation and most likely skewed
- Waterfront is a binary variable
- Floors, view, condition, and grade are discrete variables
- Zipcode, latitude, and longitude are not continuous

2.2 Handling Null Values

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Fewer Variables
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

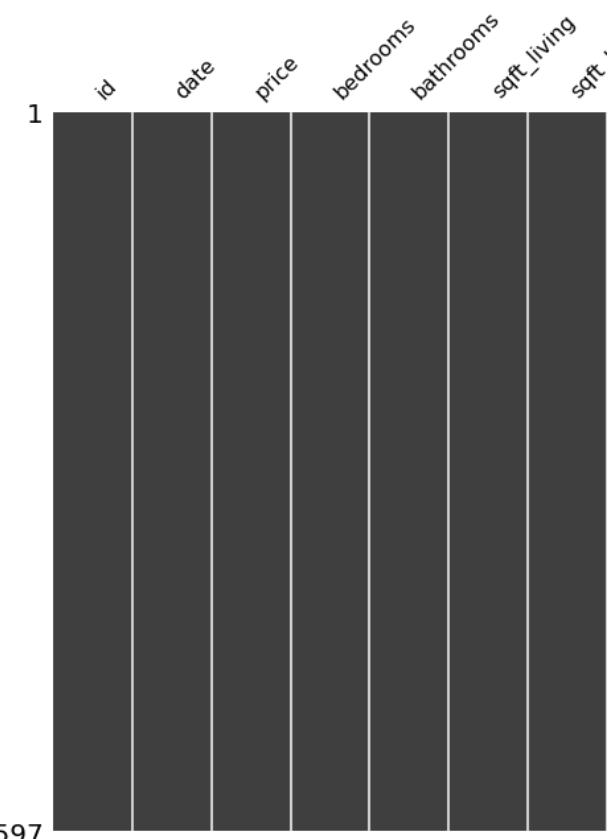
In [5]:

```

1 # Visualize which columns contain n
2
3 import missingno
4 missingno.matrix(df)

```

Out[5]: <AxesSubplot:>



Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot size
 - 4.3 Comparison of Square Foot living area vs. square foot lot area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Features
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

Waterfront, view, and yr_renovated contain null value

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Interaction
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

```
In [6]: 1 # Check how many null values are in
2
3 null = df.isna().sum()
4 null=null>1
```

```
Out[6]: waterfront    2376
view           63
yr_renovated  3842
dtype: int64
```

```
In [7]: 1 # Create formula to impute null values
2
3 def impute_cat(df, col):
4     ...
5     Impute null value with value based on mode occurring in the original column
6
7     val_prob = dict(df[col].value_counts())
8     prob = list(val_prob.values())
9     val = list(val_prob.keys())
10    np.random.choice(val, p=prob)
11    df[col].fillna(np.random.choice(val))
12
13    return df
```

2.2.1 Fill in missing Values for 'view' column

Interpreting 'view' as quality of the view from the horizon. 0 means no view, 1 means view of nature or the urban environment. A view of 0 would be a natural characteristic.

In [8]:

```

1 print('Value Counts Normalized')
2 print(df['view'].value_counts(1, dropna=False))
3 print('-----')
4 print('Value Counts Absolute')
5 print(df['view'].value_counts(dropna=False))

```

Value Counts Normalized

0.00	0.90
2.00	0.04
3.00	0.02
1.00	0.02
4.00	0.01
nan	0.00

Name: view, dtype: float64

Value Counts Absolute

0.00	19422
2.00	957
3.00	508
1.00	330
4.00	317
nan	63

Name: view, dtype: int64

```

1 I have gone ahead and made the assumption that there are 5 categories for view. I will use impute_cat to insert a value based on the distribution
2 - 90% chance of imputing a 0
3 - 2% chance of imputing a 1
4 - 4% of imputing a 2
5 - 2% of imputing a 3
6 - 1% of imputing a 4

```

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Collinearity
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

In [9]: 1 `impute_cat(df, 'view')`

Out[9]:

	id	date	price	bedrooms	ba
0	7129300520	2014-10-13	221900.00	3	
1	6414100192	2014-12-09	538000.00	3	
2	5631500400	2015-02-25	180000.00	2	
3	2487200875	2014-12-09	604000.00	4	
4	1954400510	2015-02-18	510000.00	3	
...
21592	263000018	2014-05-21	360000.00	3	
21593	6600060120	2015-02-23	400000.00	4	
21594	1523300141	2014-06-23	402101.00	2	
21595	291310100	2015-01-16	400000.00	3	
21596	1523300157	2014-10-15	325000.00	2	

21597 rows × 21 columns

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Accross All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Accross All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orde
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Fewer Categorical Variables
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

```
In [10]: 1 # Confirm that there are no more nu
          2
          3 df.isna().sum()
```

```
Out[10]: id                  0
          date                 0
          price                 0
          bedrooms                0
          bathrooms                0
          sqft_living               0
          sqft_lot                  0
          floors                   0
          waterfront              2376
          view                     0
          condition                 0
          grade                     0
          sqft_above                 0
          sqft_basement                0
          yr_built                  0
          yr_renovated              3842
          zipcode                   0
          lat                      0
          long                      0
          sqft_living15                0
          sqft_lot15                  0
          dtype: int64
```

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs Price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Selection
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

2.2.2 Fill in missing Values for 'yr_renovated'

Describes when the home was most recently renovated.

In [11]:

```

1 print('Value Counts Normalized')
2 print(df['yr_renovated'].value_coun
3 print('-----')
4 print('Value Counts Absolute')
5 print(df['yr_renovated'].value_coun

```

Value Counts Normalized

0.00	0.79
nan	0.18
2014.00	0.00
2003.00	0.00
2013.00	0.00
...	
1944.00	0.00
1948.00	0.00
1976.00	0.00
1934.00	0.00
1953.00	0.00

Name: yr_renovated, Length: 71, dtype: -----

Value Counts Absolute

0.00	17011
nan	3842
2014.00	73
2003.00	31
2013.00	31
...	
1944.00	1
1948.00	1
1976.00	1
1934.00	1
1953.00	1

Name: yr_renovated, Length: 71, dtype: -----

Most of the values in yr_renovated are either 0 or nan based on probability.

Since this value has ~79% 0 values, I will most likely

Contents ↻

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deeper Into Zip codes

In [12]: 1 df[df['yr_renovated'] > 0].describe()

Out[12]:

	id	price	bedrooms	bathrc
count	744.00	744.00	744.00	744.00
mean	4418716401.67	768901.89	3.46	1.07
std	2908265353.00	627125.79	1.00	3.00
min	3600057.00	110000.00	1.00	1.00
25%	1922984893.00	412250.00	3.00	3.00
50%	3899100167.50	607502.00	3.00	3.00
75%	7014200237.50	900000.00	4.00	4.00
max	9829200250.00	7700000.00	11.00	11.00

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot size
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

The most recent renovation took place in 2015. The distribution of renovations is skewed right.

In [13]: 1 # Going to assume that a null value
2 # This is equivalent to a 0 which is
3
4 df['yr_renovated'].fillna(0, inplace=True)

In [14]: 1 # Confirm that there are no more null values
2
3 df.isna().sum()

Out[14]:

id	0
date	0
price	0
bedrooms	0
bathrooms	0
sqft_living	0
sqft_lot	0
floors	0
waterfront	2376
view	0
condition	0
grade	0
sqft_above	0
sqft_basement	0
yr_built	0
yr_renovated	0
zipcode	0
lat	0
long	0
sqft_living15	0
sqft_lot15	0
dtype:	int64

2.2.3 Fill in missing Values for 'waterfr

Waterfront is a binary variable. 1 means the home ha

- Based on national home prices, waterfront propo a lake, and enjoy easy access to bodies of water
- I'd like to explore if homes prices at greater than
- I can then use this finding to subset the data bas

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_reno' column
 - 2.2.3 Fill in missing Values for 'waterfront' column
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary val
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling F
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

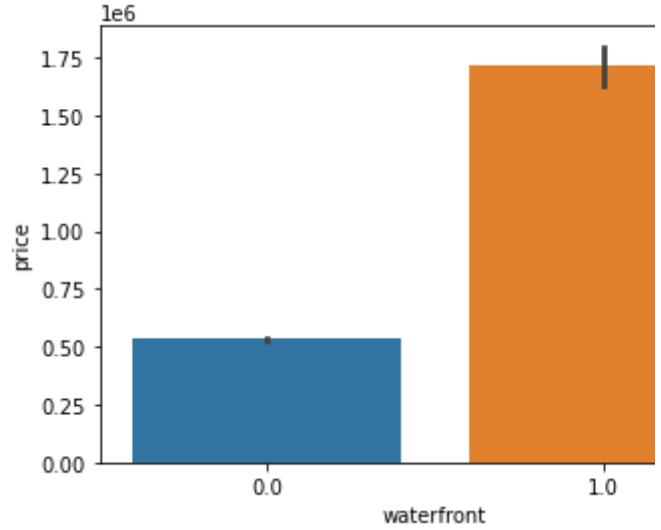
```
In [15]: 1 # Check if mean of price is greater
2
3 df.groupby('waterfront')[['price',
```

Out[15]:

	price	view
waterfront		
0.00	532641.99	0.20
1.00	1717214.73	3.76

As expected, waterfront homes have a greater mean price. This makes sense because it is easier to see the ocean, lake, or river which enhances the value of the property.

```
In [16]: 1 import seaborn as sns
2
3 sns.barplot(data=df, x='waterfront',
```



Clearly, waterfront homes are more expensive than non-waterfront homes.

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary values
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Collinearity
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

```
In [17]: 1 mplus_water = len(df[(df['price'] > 1000000)]
2 print(f'Number of houses over $1,000,000 with water front: {mplus_water}')
3 mminus_water = len(df[(df['price'] < 1000000)]
4 print(f'Number of houses under $1,000,000 with water front: {mminus_water}'
```

Number of houses over \$1,000,000 with water front: 100
Number of houses under \$1,000,000 with water front: 1000000

```
In [18]: 1 # Observe ratio of waterfront homes
2
3 pd.set_option('display.float_format', '{:.4f}'.format)
4
5 print('$1M+ with waterfront');
6 print(df.loc[df['price'] > 1000000][['waterfront']].mean())
7 print('-----')
8 # Prob of having waterfront view for $1M+
9 print('$1M- with waterfront');
10 print(df.loc[df['price'] < 1000000][['waterfront']].mean())
```

	\$1M+ with waterfront	\$1M- with waterfront
0.00000	0.92666	0.99726
1.00000	0.07334	0.00274
Name:	waterfront, dtype: float64	waterfront, dtype: float64

- 7.3% of homes priced over \$1 million have waterfront views
- 0.02% of homes priced under \$1 million have waterfront views
- As a result, I am going to subset the data by a \$1M threshold because the more expensive homes are far more likely to have a waterfront view.

```
In [19]: 1 # Subset the data into two slices based on price
2
3 df_1mplus=df.loc[df['price']>1000000]
4 df_1mmminus=df.loc[df['price']<1000000]
```

```
In [20]: 1 # Use impute_cat on homes over $1,000,000
2
3 df_1mplus =impute_cat(df_1mplus, 'waterfront')
```

/Users/ethankunin/opt/anaconda3/envs/learn_ml/lib/python3.7/site-packages/pandas/core/generic.py:517: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.

See the caveats in the documentation: <https://pandas.pydata.org/pandas-docs/stable/indexing.html#inplace-versus-copy>
return super().fillna(value)

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot size
 - 4.3 Comparison of Square Foot living area vs. square foot lot area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Collinearity
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

```
In [21]: 1 # Simply fill the missing waterfront
2 # probability of them having a waterfron
3
4 df_lmmminus['waterfront'] =df_lmmminus[
```

<ipython-input-21-3a74373b56b6>:4: Setting A value is trying to be set on a copy of a DataFrame. Try using .loc[row_indexer,col_indexer]

See the caveats in the documentation: <https://pandas.pydata.org/pandas-docs/stable/indexing.html#view-versus-copy>

```
df_lmmminus['waterfront'] =df_lmmminus[
```

```
In [22]: 1 # Join the data back together
2
3 df=pd.concat([df_lmmminus, df_limplus])
```

```
In [23]: 1 # Confirm there are no more missing values
2
3 df.isna().sum()
```

```
Out[23]: id                      0
date                     0
price                     0
bedrooms                  0
bathrooms                  0
sqft_living                 0
sqft_lot                     0
floors                     0
waterfront                   0
view                       0
condition                   0
grade                       0
sqft_above                   0
sqft_basement                  0
yr_built                     0
yr_renovated                  0
zipcode                     0
lat                         0
long                        0
sqft_living15                  0
sqft_lot15                     0
dtype: int64
```

Missing values are handled, next step is to address categorical variables.

2.3 Handling Duplicates

In [24]:

1 df[df.duplicated()]

Out[24]:

id	date	price	bedrooms	bathrooms	sqft_living
----	------	-------	----------	-----------	-------------

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Selection
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

Initially, it appears that we don't have any duplicates matching. It may be prudent to check if there are any.

In [25]:

```
1 # Check duplicates by 'id'
2 display(len(df))
3 df[df.duplicated(subset=['id']), keep=False]
```

21565

Out[25]:

	id	date	price	bedrooms
93	6021501535	2014-07-25	430000.00000	3
94	6021501535	2014-12-23	700000.00000	3
324	7520000520	2014-09-05	232000.00000	2
325	7520000520	2015-03-11	240500.00000	2
345	3969300030	2014-07-23	165000.00000	4

14294	3528000040	2014-10-01	1690000.00000	3
14295	3528000040	2015-03-26	1800000.00000	3
15999	5536100020	2015-05-12	1190000.00000	3
18976	7856400300	2014-07-02	1410000.00000	2
18977	7856400300	2015-03-22	1510000.00000	2

353 rows × 21 columns

Here, we see that when we check for duplicates by 'id', we find 21565 rows. This leads me to believe that the duplicate is reflecting the most recent change of value and accurate.

```
In [26]: 1 df=df.drop_duplicates(subset=['id'])
          2 print(len(df))
```

21388

Our dataset went from 21,565 observations to 21,388

```
In [27]: 1 # Confirm that there are no more duplicates
          2
          3 df.duplicated(subset=['id'], keep=False)
```

Out[27]: 0

3 Exploratory Data Analysis

- Explore the distribution of each variable and the outliers
- Determine if variables are discrete or continuous
- Determine if variables are categorical or numeric
- Check if their skew in the distribution or normal
- Check if there are outlier values and if they appear

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [28]:

```

1 def distr_(df, col):
2     """
3         Produces a boxplot, scatterplot
4         Produces summary statistics
5     """
6     fig, ax = plt.subplots(figsize=
7         mean=df[col].mean()
8         median=df[col].median()
9         max_=df[col].max()
10        min_=df[col].min()
11        std_=df[col].std()
12        sns.histplot(df[col], alpha=0.5,
13        sns.kdeplot(df[col], color='green')
14        ax[0].set_xlabel(col)
15        ax[0].set_title(f'{col} Distribution')
16        ax[0].axvline(mean, label=f'Mean')
17        ax[0].axvline(median, label=f'Median')
18        ax[0].legend()
19
20        sns.boxplot(data=df, x=col, ax=ax[1])
21
22        sns.scatterplot(data=df, x=df[col], y=df['price'])
23
24        fig.tight_layout();
25        print(f'{col.capitalize()} Summary')
26        print(f'Median: {median}')
27        print(f'Mean: {mean:.4f}')
28        print(f'Max: {max_}')
29        print(f'Min: {min_}')
30        print(f'StD: {std:.4f}')
31        plt.show()

```

3.1 Handling Error in Basement

When initially running the EDA check, sqft_basement

In [29]:

```

1 # Check values
2
3 df['sqft_basement'].value_counts(0)

```

Out[29]:

0.0	12701
?	452
600.0	215
500.0	205
700.0	205
...	
2130.0	1
861.0	1
3480.0	1
4130.0	1
516.0	1

Name: sqft_basement, Length: 304, dtype: int64

452 values have a question mark. I am going to assume they are missing.

data, changing it to 0 will not alter the original distribution. Turn sqft_basement into a binary variable. 0 for no basement, 1 for yes. It's difficult to interpret.

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, out
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Oral Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Collinearity
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

```
In [30]: 1 # Replace the question mark with a
2
3 df['sqft_basement'].replace(to_replace='?', value=0)
```

```
In [31]: 1 df['sqft_basement'].value_counts(1)
```

```
Out[31]: 0.0    0.61497
         600.0   0.01005
         700.0   0.00958
         500.0   0.00958
         800.0   0.00940
                  ...
         2130.0  0.00005
         861.0   0.00005
         3480.0  0.00005
         4130.0  0.00005
         516.0   0.00005
Name: sqft_basement, Length: 303, dtype: float64
```

Question mark is gone. 61% of homes do not have a basement.

```
In [32]: 1 # Use .map to make sqft_basement in binary
2
3 df['basementyes'] = (df['sqft_basement'].map(lambda x: 1 if x > 0 else 0))
4
```

```
In [33]: 1 # Confirm that original sqft_basement is still there
2
3 df['basementyes'].value_counts(1)
```

```
Out[33]: 0    0.61497
         1    0.38503
Name: basementyes, dtype: float64
```

```
In [34]: 1 # Drop sqft_basement because we are not interested in it (or it does not)
2
3
4 df=df.drop('sqft_basement', axis=1)
```

3.2 Return to checking distributions

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Oral Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Fewer Categorical Variables
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [35]:

```

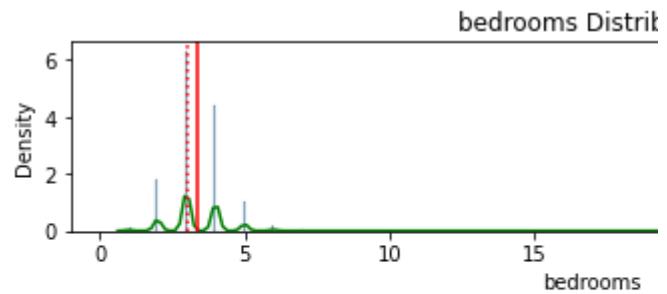
1 # Checking all variables except for
2
3 eda_check = ['price', 'bedrooms', 'sqft_lot', 'floors', 'waterfront', 'sqft_above', 'yr_built', 'year_renovated', 'lat', 'long', 'sqft_living1']
4
5 for col in eda_check:
6     print(distr_(df, col))
7
8     print('-----')
9

```

None

Bedrooms Summary

Median: 3.0
 Mean: 3.373
 Max: 33
 Min: 1
 Std: 0.9252



3.2.1 Individual EDA Analysis

1. Price

- Distribution:** Binomial, Right skewed
- Outliers:** Outliers upper IQR threshold
- Relationship with price:** NA

2. Bedrooms

- Distribution:** Bimodal
- Outliers:** Outliers upper IQR threshold. Extraneous values
- Relationship with price:** Linear until 5/6 bedrooms, then non-linear
- Discrete or Continuous:** Discrete-Possibly continuous

3. Bathrooms

- Distribution:** Bimodal
- Outliers:** Outliers upper IQR threshold.
- Relationship with price:** Linear
- Discrete or Continuous:** Discrete-Ordinal

4. Sqft_living

- Distribution:** Bimodal, skewed
- Outliers:** Outliers upper IQR threshold.
- Relationship with price:** Linear
- Discrete or Continuous:** Continuous

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or Linear Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Fewer Variables
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

5. Sqft_lot

- A. **Distribution:** Binomial, Right skewed
- B. **Outliers:** Outliers upper IQR threshold.
- C. **Relationship with price:** Non-Linear
- D. **Discrete or Continuous:** Continuous

6. Floors

- A. **Distribution:** Bimodal, Right skewed
- B. **Outliers:** None
- C. **Relationship with price:** Non-Linear-Non Continuous
- D. **Discrete or Continuous:** Discrete

7. Waterfront

- A. **Distribution:** Bernoulli
- B. **Outliers:** None
- C. **Relationship with price:** Unclear
- D. **Discrete or Continuous:** Discrete - Binary

8. View

- A. **Distribution:** Bernoulli
- B. **Outliers:** None
- C. **Relationship with price:** Unclear
- D. **Discrete or Continuous:** Discrete

9. Condition

- A. **Distribution:** Bernoulli
- B. **Outliers:** None
- C. **Relationship with price:** Unclear. Seems to have a positive correlation with price.
- D. **Discrete or Continuous:** Discrete

10. Grade

- A. **Distribution:** Bernoulli
- B. **Outliers:** None
- C. **Relationship with price:** Linear
- D. **Discrete or Continuous:** Discrete

11. Sqft_above

- A. **Distribution:** Right skewed
- B. **Outliers:** Outliers upper IQR threshold.
- C. **Relationship with price:** Linear
- D. **Discrete or Continuous:** Continuous

12. Yr_Built

- A. **Distribution:** Left skewed
- B. **Outliers:** Outliers upper IQR threshold.
- C. **Relationship with price:** None
- D. **Discrete or Continuous:** Continuous

13. Yr_Renovated

- A. **Distribution:** Bernoulli
- B. **Outliers:** None
- C. **Relationship with price:** None
- D. **Discrete or Continuous:** Continuous

14. Sqft_living15

- A. **Distribution:** Right skewed
- B. **Outliers:** Outliers upper IQR threshold.
- C. **Relationship with price:** Linear

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary values
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Collinearity
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

D. **Discrete or Continuous:** Continuous

15. Sqft_lot15

- A. **Distribution:** Right skewed
- B. **Outliers:** Outliers upper IQR threshold.
- C. **Relationship with price:** Linear
- D. **Discrete or Continuous:** Continuous

16. Basementyes

- A. **Distribution:** Binary
- B. **Outliers:** None
- C. **Relationship with price:** Unclear
- D. **Discrete or Continuous:** Discrete-Binary

3.2.2 Overall EDA Analysis

- Most of the continuous variables are right skewed
- Supported by the distribution and mean being greater than median
- High outlier values on the upper IQR threshold
- Yr_renovated has a lot of 0 values so may be imbalanced
- Bedrooms has a mistaken entry (33 bedrooms)

3.2.2.1 Handle Bedroom error

```
In [36]: 1 # Find the observation where bedrooms = 33
           2
           3 df[df['bedrooms'] > 20]
```

Out[36]:

		id	date	price	bedrooms
		15856	2402100895	2014-06-25	640000.00000 33

Appears to be a mistake because there is only 1.75 bedrooms

```
In [37]: 1 df.drop(index=15856, inplace=True)
```

```
In [38]: 1 # Confirm it has been removed
           2
           3 df[df['bedrooms'] > 20]
```

Out[38]:

	id	date	price	bedrooms	bathrooms	sqft_living

3.2.3 Turn yr_renovated into binary values

- High amount of zero values, binary encoding will be better
- Otherwise, the standard deviation will be very high

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Collinearity
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

```
In [39]: 1 df['yr_renovated'].value_counts(1)
```

```
Out[39]: 0.00000 0.96549
2014.00000 0.00341
2003.00000 0.00145
2013.00000 0.00145
2007.00000 0.00140
...
1934.00000 0.00005
1971.00000 0.00005
1954.00000 0.00005
1950.00000 0.00005
1944.00000 0.00005
Name: yr_renovated, Length: 70, dtype: float64
```

Aprox. 97% of values suggest the home has not been renovated.

```
In [40]: 1 df['renovated_yes'] = (df['yr_renovated'] > 0).astype(int)
2
3
```

```
In [41]: 1 # Confirm distribution has not changed
2
3 df['renovated_yes'].value_counts(1)
```

```
Out[41]: 0 0.96549
1 0.03451
Name: renovated_yes, dtype: float64
```

```
In [42]: 1 # Dropping yr_renovated because we have a binary representation
2
3 df.drop('yr_renovated', axis=1, inplace=True)
```

4 Feature Engineering

- Explore adding additional predictor values to the model
- Feature engineering allows for us to broaden our scope
- Be cautious of multicollinearity because features are correlated

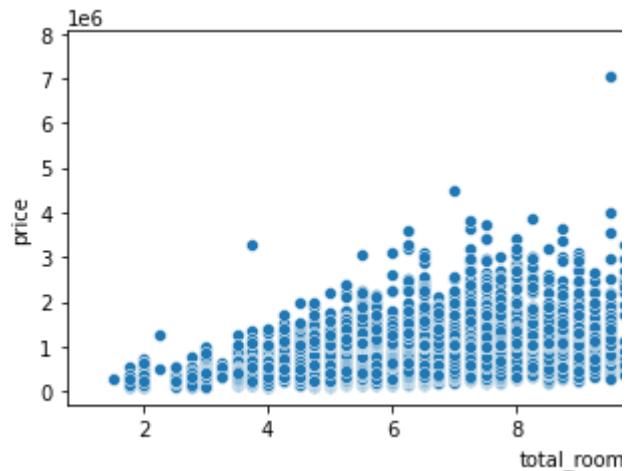
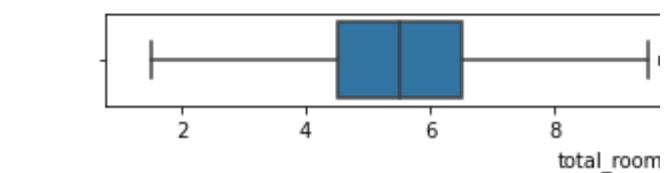
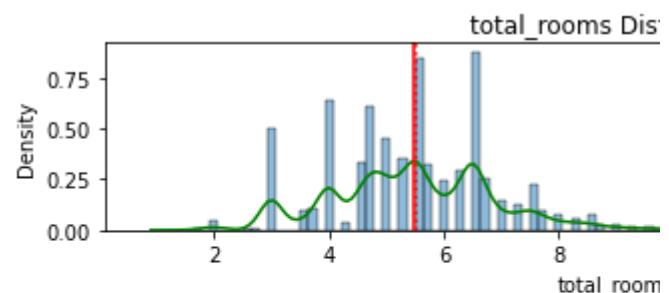
4.1 Total Rooms

- Add bedrooms and bathrooms to create new columns
- This will provide us a summary of the number of rooms

```
In [43]: 1 # Combine bedrooms and bathrooms
2
3 df['total_rooms'] = df['bedrooms']+df['bathrooms']
```

In [44]: 1 `distr_(df, 'total_rooms')`

```
Total_rooms Summary
Median: 5.5
Mean: 5.49
Max: 16.5
Min: 1.5
Std: 1.463
```



- 1 - Distribution is not normal
- 2 - Looks like there are a significant number of outliers
- 3 - Initially looks like there is a positive correlation

4.2 Backyard Size as a proportion

- Use sqft_above/sqft_lot as a proxy for backyard
- Essentially, we are capturing how big the home is relative to its lot size
- A larger value means a relatively smaller backyard

Contents ⚙️

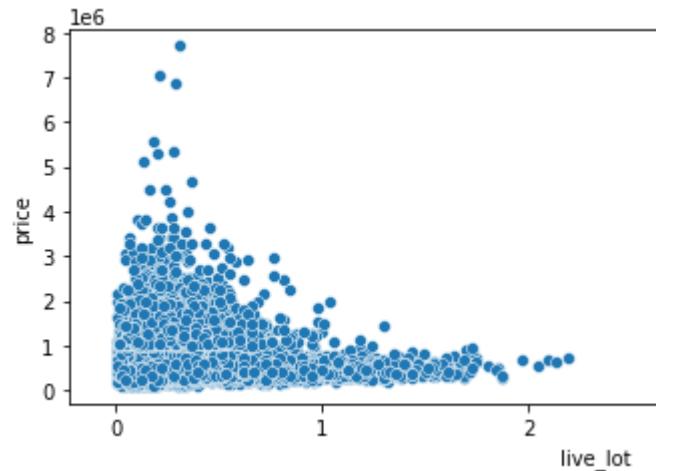
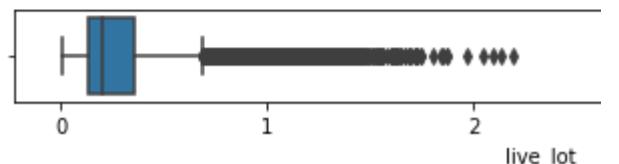
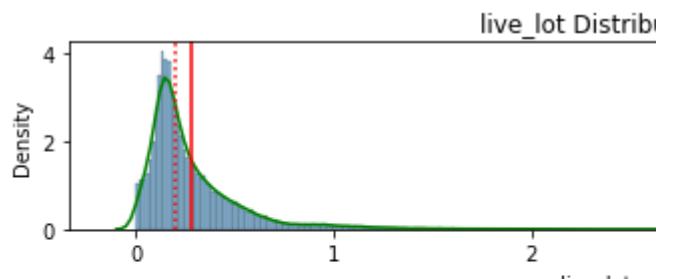
- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling F
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

```
In [45]: 1 # Divide sqft_living/sqft_lot
2 # Using sqft_above as opposed to sq
3
4 df['live_lot'] = df['sqft_above']/d
```

```
In [46]: 1 distr_(df, 'live_lot')
```

Live_lot Summary

Median: 0.20274485339987525
 Mean: 0.2808
 Max: 4.653846153846154
 Min: 0.0006095498431482305
 Std: 0.2426



- Data is right skewed
- Significant number of outliers to the right of the box plot
- Values may be greater than 1 because the living floors. In other words, it has a lot of sqft footage

4.3 Comparison of Square Footage

- Sqft_living15 represents the average living space

- Would like to compare how the living space of t

```
In [47]: 1 # SQF_living compared to neighbors
2 df['living_vs_neighbor'] = df['sqft_l
```

```
In [48]: 1 distr_(df, 'living_vs_neighbor')
```

Living_vs_neighbor Summary

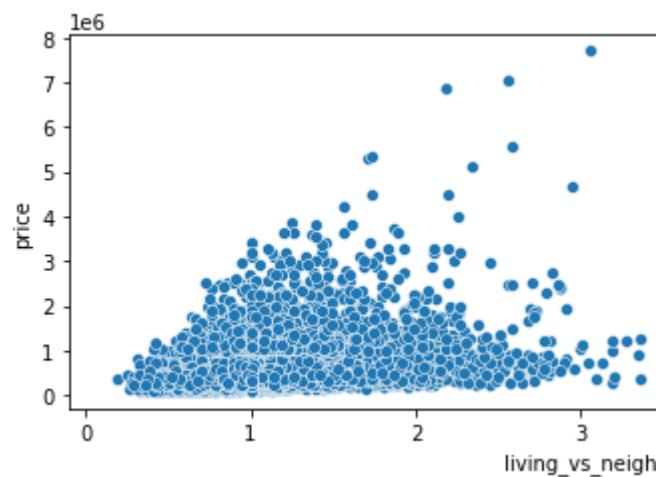
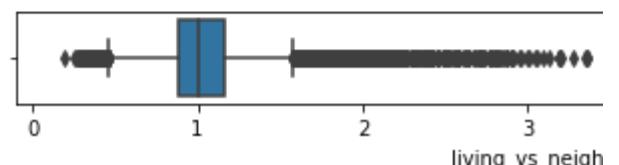
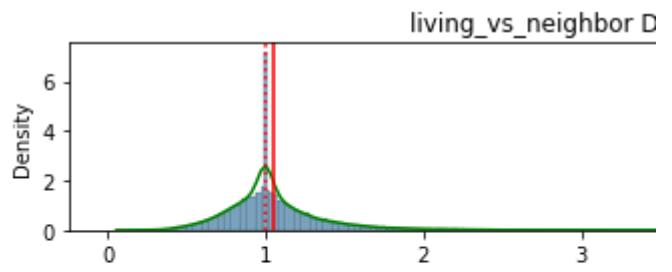
Median: 1.0

Mean: 1.054

Max: 6.0

Min: 0.1872791519434629

Std: 0.3203



- | | |
|---|-------------------------------------|
| 1 | - Data looks pretty normally distri |
| 2 | - Outliers to the right of the uppe |
| 3 | - Unclear if there is a linear rela |
| 4 | - A lot of the much larger homes in |

- Sqft_lot15 represents the average living space o
- Would like to compare how the lot size of the ob

Contents ⚙

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary val
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the house
 - 4.3 Comparison of Square Foot living area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Accross All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
- ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Accross All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
- 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling F
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

```
In [49]: 1 df['lot_vs_neighbor'] = df['sqft_lo
```

```
In [50]: 1 distr_(df, 'lot_vs_neighbor')
```

Lot_vs_neighbor Summary

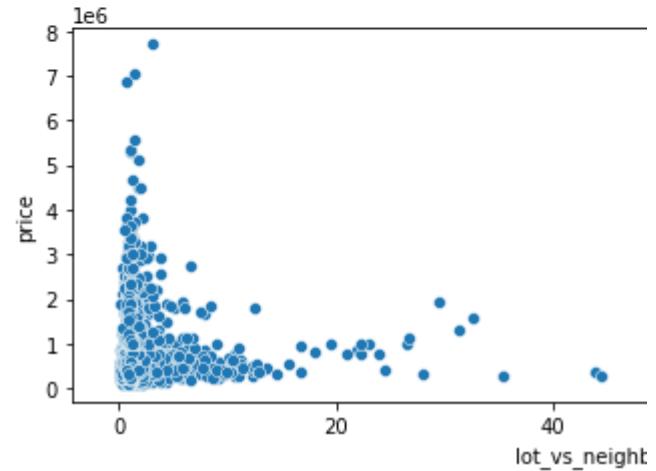
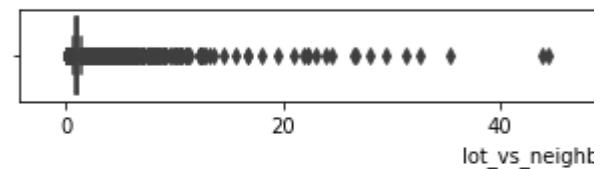
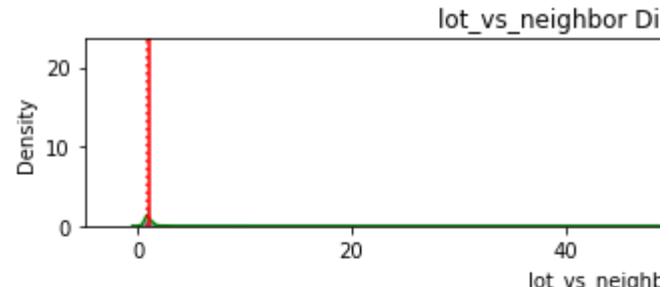
Median: 1.0

Mean: 1.134

Max: 87.52717948717948

Min: 0.054971997700810314

Std: 1.286



- Hard to tell distribution because there are clearly outliers
- There are values greater than 80 which seem unlikely
- Linearity seems unlikely because as lot_vs_neighbor increases, price increases but at a decreasing rate

5 Check Assumptions of Linear Model

- For the model to provide accurate inferences, it must satisfy certain assumptions

5.1 Check Assumption of Linearity

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot size
 - 4.3 Comparison of Square Foot living area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
- ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
- ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
- 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling F
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

- There must be a linear relationship between the
 - In our case, the predictor variable refers to the x-variable
- By linear relationship, we mean that as the x-value increases, the y-value also increases at a constant rate.
- If we do not meet assumption of linearity, our model will not be able to predict accurately.
- Must check each predictor that we are going to use in our model.

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary value
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot size
 - 4.3 Comparison of Square Foot living area vs. square foot above ground
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling F
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [51]:

```

1 def lin_check(df, cols, ncols=4, fi
2     """
3         Produces regplot for each feature
4     """
5     if ncols%4==0:
6         fig, axes = plt.subplots(nr
7             for ax, col in zip(axes.flas
8                 sns.regplot(data=df, x=
9                     ax.set_title(f'{col} vs
10                    fig.tight_layout()
11    else:
12        fig, axes = plt.subplots(nr
13            for ax, col in zip(axes.flas
14                sns.regplot(data=df, x=
15                    ax.set_title(f'{col} vs
16                    fig.tight_layout()

```

In [52]:

1 df.columns

Out[52]:

```

Index(['id', 'date', 'price', 'bedrooms',
       'sqft_lot', 'floors', 'waterfront',
       'sqft_above', 'yr_built', 'zipcode',
       'sqft_lots15', 'basementyes', 're
       'living_vs_neighbor', 'lot_vs_ne
       dtype='object')

```

In [53]:

```

# Choosing to remove latitude and longitude
# is a sufficient proxy for location
# will be easier for residents to understand
cols_to_check = ['bedrooms', 'bathrooms',
                  'sqft_lot', 'floors', 'waterfront',
                  'sqft_above', 'yr_built', 'zipcode',
                  'sqft_lots15', 'basementyes',
                  'living_vs_neighbor', 'lot_vs_ne

```

In [54]:

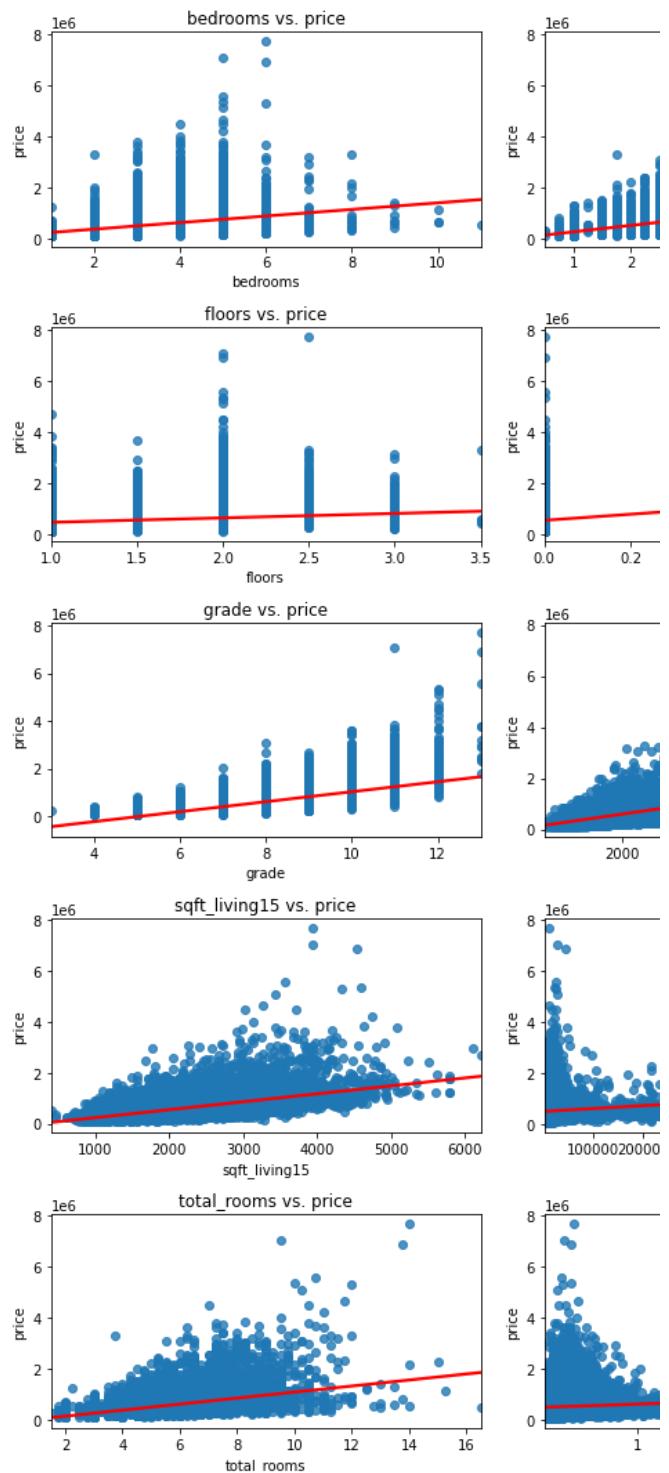
```

1 # Checking for linearity between pr
2
3 lin_check(df, cols_to_check)

```

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot size
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes



The following cells do not have a linear relationship v

- Sqft_lot
- Floors (categorical)
- View (categorical)
- Condition (categorical)
- Yr_built
- Zip code (categorical)
- Basementyes (categorical)
- Renovatedyes (categorical)
- Sqft_lot15
- Live_lot

Sqft_lot, Yr_built, sqft_lot15, and live_lot are numeric not have a linear relationship with price, will proceed price, I will One Hot Encode them

In [55]: 1 df.columns

Out[55]: Index(['id', 'date', 'price', 'bedrooms',
 'sqft_lot', 'floors', 'waterfron',
 'sqft_above', 'yr_built', 'zipco',
 'sqft_lot15', 'basementyes', 're',
 'living_vs_neighbor', 'lot_vs_ne',
 'dtype='object')]

In [56]: 1 # These were continuous variables t
 2
 3 cols_to_drop = ['sqft_lot', 'sqft_1',
 4 df_lin = df.drop(cols_to_drop, axis=0)

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot size
 - 4.3 Comparison of Square Foot living area vs square foot above ground
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [57]:

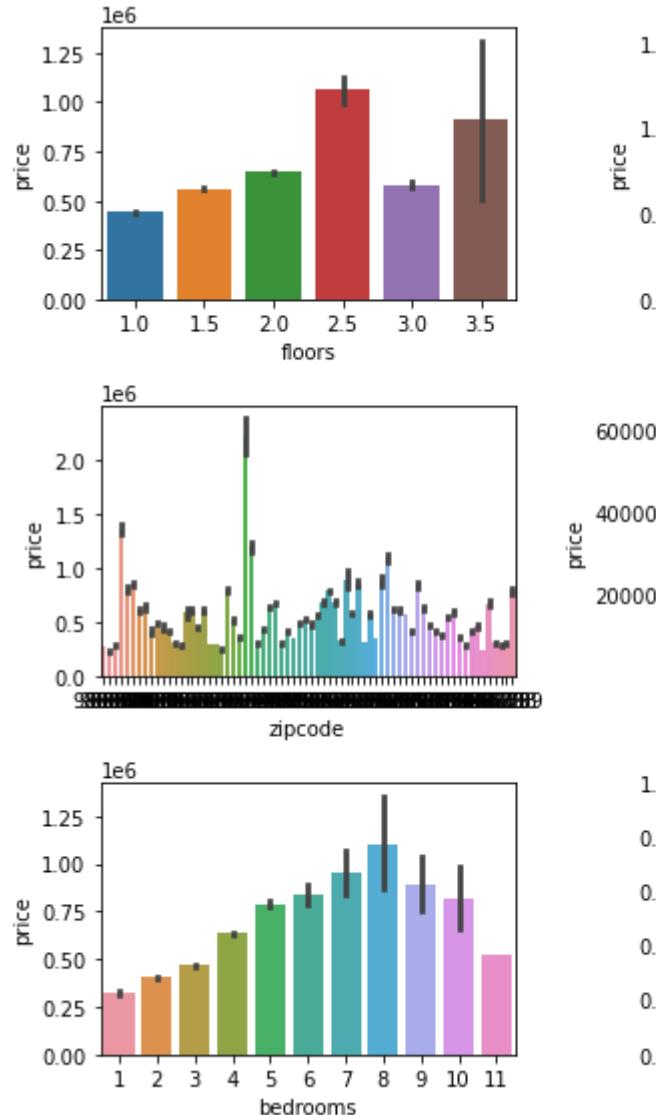
```

1 # Checking if categorical variables
2
3 cat_bars = ['floors', 'view', 'cond'
4 fig, axes = plt.subplots(nrows=3, n
5 for ax, col in zip(axes.flatten(),
6         sns.barplot(data=df, x=col, y='
7         fig.tight_layout()
8

```

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot size
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Collinearity
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes



Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
 - 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

- Floors, and zipcode do not have linear relationships
- Condition is close, but roughly has a linear relationship
- Will turn floors, zip code, and bedrooms into OH

Conclusion:

- Will be dropping: Sqft_lot, Yr_built, sqft_lots15, and condition
- Will be One Hot Encoding: Floors, zipcode, and bedrooms

5.2 Check Assumption of Multicollinearity

For a multiple linear regression model to be accurate not only should the predictors have a linear relationship, but they should also move very close together, then they are redundant for change in predictor values

RoadMap for checking assumptions of multicollinearity

1. Run initial check of correlation with price
2. Observe heatmap triangle to see which predictors move together
3. Build table to show which variables have a correlation > 0.75
 - 0.75 is the norm for determining if predictor variable is redundant

```
In [58]: 1 # Check correlation with price
           2
           3 def initial_corr_check(df, col='price'):
           4     return df.corr()[['price']].round(2)
```

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

```
In [59]: 1 initial_corr_check(df_lin.drop(columns=[
```

```
Out[59]: price           1.00000
sqft_living        0.70000
grade             0.67000
sqft_above         0.60000
sqft_living15      0.58000
bathrooms          0.53000
total_rooms         0.47000
view               0.39000
bedrooms            0.32000
lat                0.31000
living_vs_neighbor 0.30000
waterfront          0.27000
floors              0.25000
basementyes         0.18000
renovated_yes       0.12000
lot_vs_neighbor     0.04000
condition           0.03000
long                0.02000
zipcode             -0.05000
Name: price, dtype: float64
```

Initial correlation check shows that sqft_living, grade

```
In [60]: 1 # Reference: https://heartbeat.fritz.ai/using-pandas-to-create-a-correlation-heatcamp.html
2
3 def corr_triangle(df):
4     """
5     Correlation heatmap, including
6     """
7     corr2 = df.corr()
8     fig, ax = plt.subplots(figsize=(10, 10))
9     matrix = np.triu(corr2)
10    return sns.heatmap(corr2, cmap="
```

```
In [61]: 1 corr_triangle(df_lin)
```

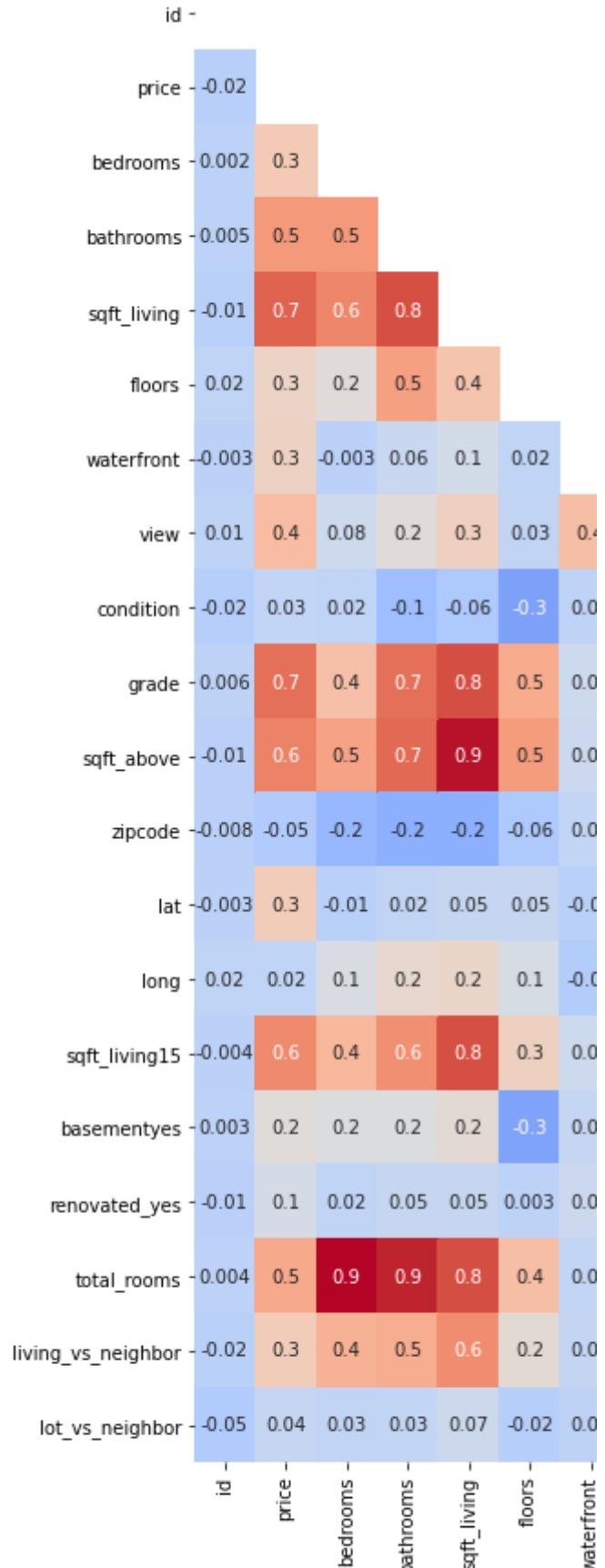
```
Out[61]: <AxesSubplot:>
```

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Features
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot
 - 4.3 Comparison of Square Foot living area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes



There are a number of features that have strong multicollinearity.

columns

- Total rooms correlates strongly with bedrooms a
- Sqft_living and sqft_living 15 are strongly corre
- Sqft_living and sqft_above and grade are strong
 - Because sqft_living and sqft_above will be t

Contents ↻

▼ 1 King County Housing Characteristics
1.1 Data
1.2 Roadmap
▼ 2 Scrub Data
2.1 Descriptions of columns
▼ 2.2 Handling Null Values
2.2.1 Fill in missing Values for 'view' cc
2.2.2 Fill in missing Values for 'yr_reno
2.2.3 Fill in missing Values for 'waterfr
2.3 Handling Duplicates
▼ 3 Exploratory Data Analysis
3.1 Handling Error in Basement encoding
▼ 3.2 Return to checking distributions, out
3.2.1 Individual EDA Analysis
▼ 3.2.2 Overall EDA Analysis
3.2.2.1 Handle Bedroom error
3.2.3 Turn yr_renovated into binary val
▼ 4 Feature Engineering
4.1 Total Rooms
4.2 Backyard Size as a proportion of the
4.3 Comparison of Square Foot living an
▼ 5 Check Assumptions of Linearity and Multi
5.1 Check Assumption of Linearity
5.2 Check Assumption of Multicollinearit
6 Model 1: Baseline Model
▼ 7 Outlier Removal: IQR + Z-Score
▼ 7.1 IQR Method
7.1.1 IQR Method Accross All Columns
7.1.1.1 Model 2: IQR All Outliers Rem
▼ 7.1.2 IQR Price Outliers Removed
7.1.2.1 Model 3: IQR Price Outliers R
▼ 7.2 Z-Score Method
7.2.1 Z-Score Method Accross All Col
7.2.2 Model 3: Z-Score All Outliers Re
▼ 7.2.3 Z-Score Price Outliers Removed
7.2.3.1 Model 4: Z-Score Price Outli
7.3 Table to Compare 4 Outlier Removal
▼ 8 Handling Categorical Variables with One
8.1 Check relationship of Non-Linear Cat
▼ 8.2 One Hot Encode Categorical Non-Or
8.2.1 Model 5: OHE Iteration 1
8.2.2 Model 5: Iteration 2 - Handling F
▼ 9 Interpretation
9.1 Standardize data for interpretation
▼ 9.2 View which predictors make the mos
9.2.1 Digging Deepr Into Zip codes

In [62]:

```

1 # Reference:https://github.com/lear
2
3 def corr_finder(df):
4     """
5         Shows pairs of features that ha
6         each other
7     """
8     df_corr = df.corr().abs().stack()
9     df_corr['pairs'] = list(zip(df_
10     df_corr.set_index(['pairs'], in
11     df_corr.drop(columns=['level_1'
12
13     # # cc for correlation coeffici
14     df_corr.columns = ['cc']
15     df_corr.drop_duplicates(inplace
16
17     return df_corr[(df_corr.cc>.75)

```

In [63]:

1 corr_finder(df_lin)

Out[63]:

cc	pairs
(total_rooms, bedrooms)	0.89510
(sqft_living, sqft_above)	0.87655
(bathrooms, total_rooms)	0.85186
(total_rooms, sqft_living)	0.76392
(sqft_living, grade)	0.76243
(sqft_living15, sqft_living)	0.75630
(grade, sqft_above)	0.75610
(bathrooms, sqft_living)	0.75581

Methodology to handle collinearity:

- For each pair, drop the feature that has the lowe
 - Maintaining the feature that has a stronger r
 - Bedrooms & Total_Rooms: Drop Total_F between suggesting adding either bedr
 - Sqft_living & Sqft_above: Drop Sqft_ab represented with a binary variable. Nua
 - Total_Rooms and Bathrooms: Already e

- Total_Rooms and Sqft_living: Already eliminated
- Sqft_living15 and Sqft_living: Already eliminated
- Sqft_above and grade: Electing to keep
- Sqft_living and bathrooms: Already dropped

Contents ⚙️

- 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- 2 Scrub Data
 - 2.1 Descriptions of columns
 - 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary values
- 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living and square foot above grade
- 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- 7 Outlier Removal: IQR + Z-Score
 - 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
- 7.3 Table to Compare 4 Outlier Removal
- 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - 8.2 One Hot Encode Categorical Non-Orde
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling F
- 9 Interpretation
 - 9.1 Standardize data for interpretation
 - 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

```
In [64]: 1 cols_to_drop = ['total_rooms', 'sqft_above']
2 df_linco = df_linco.drop(cols_to_drop)
```

```
In [65]: 1 # Confirm no multicollinearity issue
2 # Elected to keep because they represent different information
3
4 corr_finder(df_linco)
```

Out[65]:

	cc	pairs
(grade, sqft_living)	0.76243	

Our model is now closer to meeting all the necessary assumptions.

Next steps will be to handle outliers and then categorize variables.

6 Model 1: Baseline Model

- Now that assumptions of linearity and no multicollinearity are met
- Mainly observing R squared, Adjusted R squared, and P-values
- Will check if certain variables are statistically insignificant outliers

What is R^2, Adjusted R^2, QQ Plot, and Homoskedasticity?

- **R^2:** Indicated how much variance is explained by our 'goodness of fit' test. The higher the value, the better the fit of the dependent model.

- **Adjusted R^2:** Similar to R squared, but it takes into account the number of predictors, so it has a downward bias as the number of predictors increases.

- **Homoskedasticity:** For our model to be valid, the residuals (actual value - predicted value) should not have a recognizable pattern or trend.

- **QQ Plot:** Helps us measure homoskedasticity. If the data points follow a straight line, the variance becomes non-uniform.

Ref: <https://statisticsbyjim.com/regression/interpreting-regression-coefficients/>

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs Price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Oral Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Collinearity
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [66]:

```

1 import statsmodels.api as sm
2 import statsmodels.stats.api as sms
3 import statsmodels.formula.api as s
4
5 from scipy import stats
6 from sklearn.preprocessing import StandardScaler as S
7 from statsmodels.formula.api import ols
8
9 from sklearn.datasets import make_regression as make_r
10 from sklearn.linear_model import LinearRegression as Li
11 import sklearn.metrics as metrics

```

In [67]:

```

1 def model_summary(df, X_targets, y,
2                   ...):
3     """
4     Produces OLS Linear Regression
5     plot is displayed below the summary
6     ...
7     outcome = y
8     x_cols = X_targets
9     predictors = '+' .join(x_cols)
10    formula = outcome + '~' + predictors
11    model = ols(formula=formula, data=df)
12    resid1 = model.resid
13    display(model.summary())
14    if qq==True:
15        sm.graphics.qqplot(resid1,
16                            ... )
17    return model

```

In [68]:

```

1 from sklearn.datasets import make_regression as make_r
2 from sklearn.linear_model import LinearRegression as Li
3 import sklearn.metrics as metrics
4
5
6 def sked_show(df, X_cols, lr=None,
7               ...):
8     """
9     Produces scatter plot showing model fit
10    ...
11    if lr is None:
12        lr = LinearRegression()
13        lr.fit(df[X_cols], df[val])
14
15        y_hat = lr.predict(df[X_cols])
16    else:
17        y_hat = lr.predict(df)
18
19    resid = (df[val] - y_hat)
20    fig, ax= plt.subplots(figsize=(10, 6))
21    ax.scatter(x=y_hat,y=resid, alpha=0.5)
22    ax.axhline(0, color='red')
23    ax.set_xlabel('Price')
24    ax.set_ylabel('Residual')
25    return fig,ax

```

In [69]:

```

1 # Begin by using all columns as pre
2
3 x_targs = ['bedrooms', 'floors', 'w
4     'view', 'condition', 'grade'
5     'basementyes', 'renovated_ye
6     'lot_vs_neighbor', 'sqft_liv

```

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs sqft_living15
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Accross All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Accross All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orde
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [70]:

```

1 model_base = model_summary(df_linco
2 sked_show(df, x_targs, model_base)

```

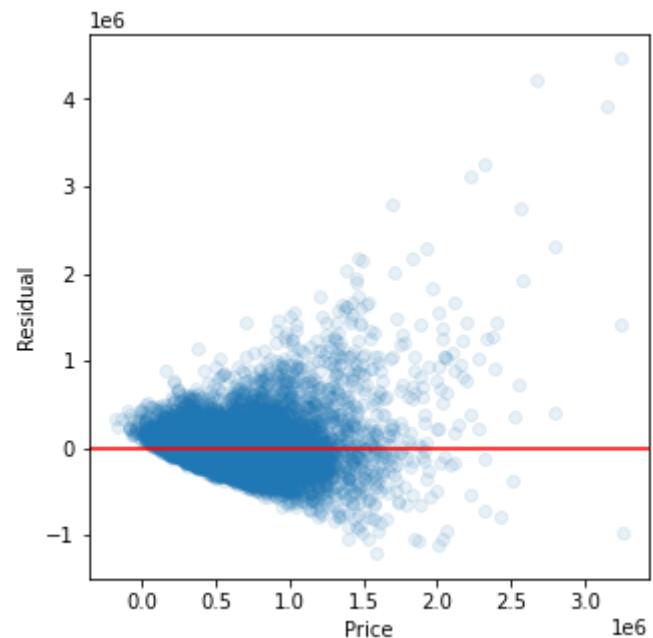
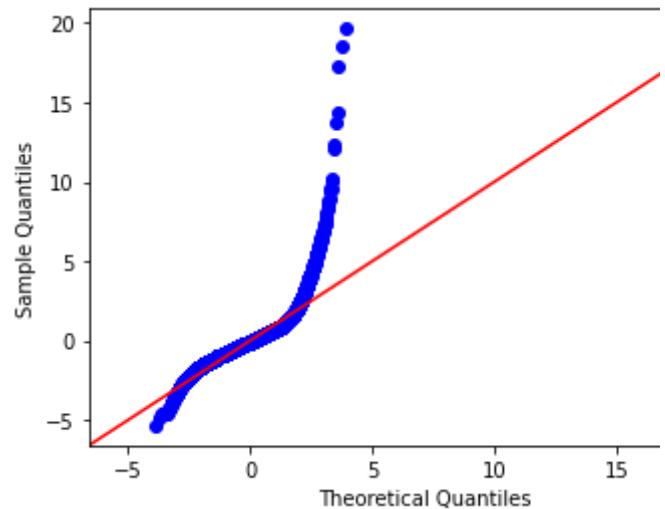
OLS Regression Results

Dep. Variable:	price	R-squared		
Model:	OLS	Adj. R-squared		
Method:	Least Squares	F-statistic		
Date:	Wed, 05 May 2021	Prob (F-statistic)		
Time:	14:10:57	Log-Likelihood		
No. Observations:	21387	AIC		
Df Residuals:	21374	BIC		
Df Model:	12			
Covariance Type:	nonrobust			
	coef	std err	t	P
Intercept	-5.254e+07	3.07e+06	-17.129	0.0
bedrooms	-3.239e+04	2204.113	-14.694	0.0
floors	-3302.0105	3555.724	-0.929	0.3
waterfront	6.124e+05	2.05e+04	29.829	0.0
view	4.962e+04	2361.401	21.012	0.0
condition	6.098e+04	2514.696	24.248	0.0
grade	9.624e+04	2284.261	42.133	0.0
zipcode	528.9805	31.267	16.918	0.0
basementyes	2.791e+04	3651.532	7.643	0.0
renovated_yes	1.524e+05	8613.485	17.697	0.0
living_vs_neighbor	-1.008e+05	6466.963	-15.591	0.0
lot_vs_neighbor	716.1896	1216.573	0.589	0.0
sqft_living	213.4715	3.621	58.950	0.0
Omnibus:	15558.225	Durbin-Watson:	1	
Prob(Omnibus):	0.000	Jarque-Bera (JB):	926735	
Skew:	2.918	Prob(JB):		
Kurtosis:	34.716	Cond. No.	1.94e	

Notes:

- [1] Standard Errors assume that the covariance matrix is correctly specified.
- [2] The condition number is large, 1.94e+08. This might indicate strong multicollinearity or other numerical problems.

Out[70]: <Figure size 360x360 with 1 Axes>,
<AxesSubplot:xlabel='Price', ylabel='R



Conclusions

- $R^2: 0.618$
- Adjusted $R^2: 0.618$
- QQ Plot: Deviates upwards at the 2nd/3rd quantiles
- Homoskedasticity: Becomes cone shaped around the mean
- Non-Statistically Significant Predictors: Floors, lotsize
- Not dropping non-statistically significant values

7 Outlier Removal: IQR + Z-Score

- Due to high outliers shown in the QQ Plot, next step is to remove them.
- Will try using Z-Score and IQR method

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the total area
 - 4.3 Comparison of Square Foot living area vs. Total Area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
- 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Fewer Outliers
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

- Will evaluate data loss to determine which method is better

7.1 IQR Method

- The IQR method uses quantiles to determine if a value is an outlier
- The process:
 1. Determine IQR which is calculated as Quantile 3 - Quantile 1
 2. Upper Thresh: Quantile 3 + 1.5 * IQR
 3. Lower Thresh: Quantile 1 - 1.5 * IQR
 4. Any values outside of the upper and lower thresholds are outliers
- This method is more strict for evaluating outliers

7.1.1 IQR Method Across All Columns

- In this approach, we are going to classify outliers across all columns
- For example, if an observation has 10 bedrooms, it would be an outlier
- This method is more strict for determining outliers
- CAUTION: we may have significant data loss

In [71]:

```

1 def find_outliers_IQR(data):
2     """Detects outliers using the IQR method.
3     Returns a boolean Series where True indicates an outlier.
4     """
5     res = data.describe()
6     q1 = res['25%']
7     q3 = res['75%']
8     thresh = 1.5*(q3-q1)
9     idx_outliers = (data < (q1-thresh) | data > (q3+thresh))
10    return idx_outliers

```

In [72]:

```
1 df_linco.columns
```

Out[72]:

```
Index(['id', 'date', 'price', 'bedrooms',
       'waterfront', 'view', 'condition',
       'basementyes', 'renovated_yes',
       'lot_vs_neighbor'],
      dtype='object')
```

In [73]:

```

1 # Cols we are going to check with the IQR method
2 # Do not perform IQR check on binary columns
3
4 iqr_check = ['price', 'bedrooms',
5               'sqft_living',
6               'living_vs_neighbor',
7               'waterfront']

```

Contents

▼ 1 King County Housing Characteristics	
1.1 Data	
1.2 Roadmap	
▼ 2 Scrub Data	
2.1 Descriptions of columns	
▼ 2.2 Handling Null Values	
2.2.1 Fill in missing Values for 'view' column	
2.2.2 Fill in missing Values for 'yr_renovated'	
2.2.3 Fill in missing Values for 'waterfront'	
2.3 Handling Duplicates	
▼ 3 Exploratory Data Analysis	
3.1 Handling Error in Basement encoding	
▼ 3.2 Return to checking distributions, outliers	
3.2.1 Individual EDA Analysis	
▼ 3.2.2 Overall EDA Analysis	
3.2.2.1 Handle Bedroom error	
3.2.2.3 Turn yr_renovated into binary variable	
▼ 4 Feature Engineering	
4.1 Total Rooms	
4.2 Backyard Size as a proportion of the lot size	
4.3 Comparison of Square Foot living area vs price	
▼ 5 Check Assumptions of Linearity and Multicollinearity	
5.1 Check Assumption of Linearity	
5.2 Check Assumption of Multicollinearity	
6 Model 1: Baseline Model	
▼ 7 Outlier Removal: IQR + Z-Score	
▼ 7.1 IQR Method	
▼ 7.1.1 IQR Method Across All Columns	
7.1.1.1 Model 2: IQR All Outliers Removed	
▼ 7.1.2 IQR Price Outliers Removed	
7.1.2.1 Model 3: IQR Price Outliers Removed	
▼ 7.2 Z-Score Method	
7.2.1 Z-Score Method Across All Columns	
7.2.2 Model 3: Z-Score All Outliers Removed	
▼ 7.2.3 Z-Score Price Outliers Removed	
7.2.3.1 Model 4: Z-Score Price Outliers Removed	
7.3 Table to Compare 4 Outlier Removal	
▼ 8 Handling Categorical Variables with One Hot Encoding	
8.1 Check relationship of Non-Linear Categorical Variables	
▼ 8.2 One Hot Encode Categorical Non-Or	
8.2.1 Model 5: OHE Iteration 1	
8.2.2 Model 5: Iteration 2 - Handling Feature	
▼ 9 Interpretation	
9.1 Standardize data for interpretation	
▼ 9.2 View which predictors make the most sense	
9.2.1 Digging Deeper Into Zip codes	

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot size
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Oral Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Fewer Variables
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

```
In [74]: 1 iqr_outliers = pd.DataFrame()
2 for col in iqr_check:
3     iqr_outliers[col]=find_outliers(df, col)
4 iqr_outliers['total'] = iqr_outlier
5 df_iqr = df_linco[~iqr_outliers['to
```

```
In [75]: 1 # Confirm that our data looks corre
2
3 df_iqr.head()
```

Out[75]:

	id	date	price	bedrooms	sqft
0	7129300520	2014-10-13	221900.00000	3	
1	6414100192	2014-12-09	538000.00000	3	
3	2487200875	2014-12-09	604000.00000	4	
4	1954400510	2015-02-18	510000.00000	3	
6	1321400060	2014-06-27	257500.00000	3	

```
In [76]: 1 df_iqr.describe()
```

Out[76]:

	id	price	bedroom
count	15475.00000	15475.00000	15475.00000
mean	4734745636.78559	467032.92775	3.3054
std	2878922190.79975	200357.29212	0.7680
min	1200019.00000	81000.00000	2.0000
25%	2297400055.00000	310000.00000	3.0000
50%	4046710050.00000	429000.00000	3.0000
75%	7504400730.00000	586750.00000	4.0000
max	9900000190.00000	1120000.00000	5.0000

With the new DataFrame, the range of values has been reduced.

- Price
 - Min: \$81,000
 - Max: \$1,120,000
- Bedrooms
 - Min: 2
 - Max: 5
- Floors
 - Min: 1
 - Max: 3.5

- Sqft_Living
 - Min: 560
 - Max: 4230

This constrains our dataset to only be able to provide 15475 observations. The QQ Plot was trailing off around \$1.25 million.

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Oral Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Interaction
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [77]:

```
1 print(f'Num observations before dropping with IQR: {len(df)}')
2 print(f'Num observations after dropping with IQR: {len(df_iqr)}')
3 print(f'Num observations removed: {len(df) - len(df_iqr)}')
4 print(f'Num observations removed as percent of total: {(len(df) - len(df_iqr)) / len(df) * 100}%)'
```

Num observations before dropping with IQR: 15475
 Num observations after dropping with IQR: 14883
 Num observations removed: 5,912
 Num observations removed as percent of total: 38.88%

With this type of outlier removal we have significant constraints. The model performs with these constraints.

7.1.1.1 Model 2: IQR All Outliers Removed

- Check to see if assumption of homoskedasticity is violated
- Check if significant p-values has changed

In [78]:

```
1 # Ensure our DataFrame only includes numerical columns
2 df_iqr = df.select_dtypes(include=[np.number])
3 df_iqr.describe()
```

Out[78]:

	id	price	bedroom
count	15475.00000	15475.00000	15475.00000
mean	4734745636.78559	467032.92775	3.3054
std	2878922190.79975	200357.29212	0.7680
min	1200019.00000	81000.00000	2.0000
25%	2297400055.00000	310000.00000	3.0000
50%	4046710050.00000	429000.00000	3.0000
75%	7504400730.00000	586750.00000	4.0000
max	9900000190.00000	1120000.00000	5.0000

In [79]:

```
1 x_targs = df_iqr.columns
2 x_targs = list(x_targs)
3 x_targs = [x for x in x_targs if x != 'id' and x != 'price' and x != 'bedroom']
```

In [80]:

```

1 model_iqra = model_summary(df_iqr,
2 sked_show(df_iqr, x_targs, model_iq

```

OLS Regression Results

Dep. Variable:	price	R-squared		
Model:	OLS	Adj. R-squared		
Method:	Least Squares	F-statistic		
Date:	Wed, 05 May 2021	Prob (F-statistic)		
Time:	14:10:57	Log-Likelihood		
No. Observations:	15475	AIC		
Df Residuals:	15460	BIC		
Df Model:	14			
Covariance Type:	nonrobust			
	coef	std err	t	P
Intercept	-2.777e+07	1.84e+06	-15.108	0.0
bedrooms	-8472.8199	1525.406	-5.554	0.0
sqft_living	134.7622	2.678	50.328	0.0
floors	1.37e+04	2184.062	6.273	0.0
waterfront	1.016e+05	3.01e+04	3.373	0.0
view	3.655e+04	1686.343	21.674	0.0
condition	4.756e+04	1535.022	30.983	0.0
grade	5.888e+04	1520.542	38.722	0.0
zipcode	-119.1025	21.897	-5.439	0.0
lat	5.857e+05	6849.178	85.512	0.0
long	-9.245e+04	8475.166	-10.909	0.0
basementyes	1.76e+04	2287.099	7.694	0.0
renovated_yes	9.912e+04	6055.472	16.368	0.0
living_vs_neighborhood	-1.105e+05	5856.715	-18.869	0.0
lot_vs_neighborhood	3044.9759	8182.017	0.372	0.0
Omnibus:	1495.077	Durbin-Watson:	1.86	
Prob(Omnibus):	0.000	Jarque-Bera (JB):	2626.11	
Skew:	0.679	Prob(JB):	0.0	
Kurtosis:	4.494	Cond. No.	1.97e+00	

Notes:

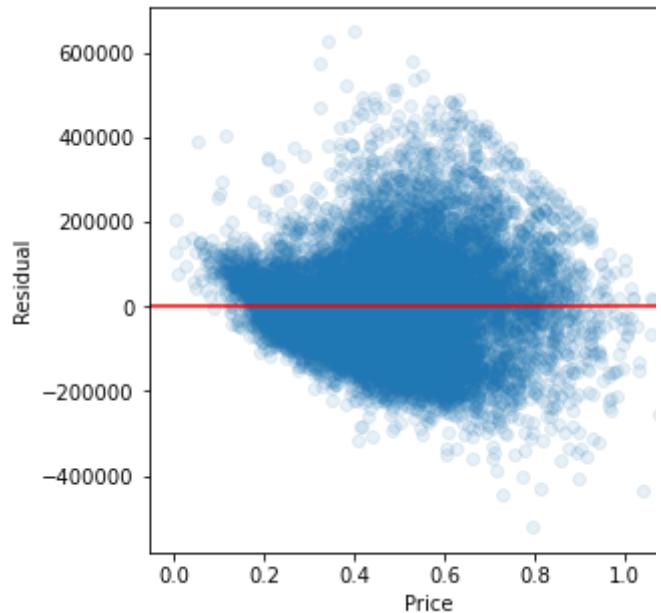
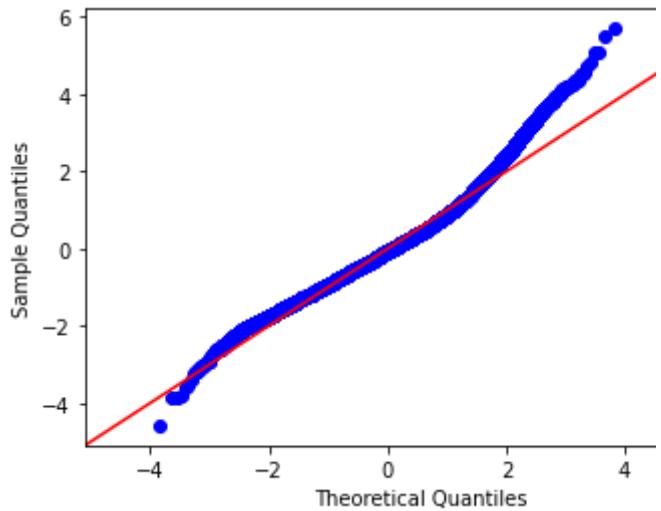
Contents ⚙

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.2 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the house size
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - ▼ 7.1.1 IQR Method Accross All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Accross All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orde
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

[1] Standard Errors assume that the covariance matrix is full rank.

[2] The condition number is large, 1.97e+08. This might indicate strong multicollinearity or other numerical problems.

Out[80]: (`<Figure size 360x360 with 1 Axes>`,
`<AxesSubplot:xlabel='Price', ylabel='R'`)



Contents ⚙

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs. total area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Fewer Variables
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

Conclusions

- R^2 : 0.676
- Adjusted R^2 : 0.676
- QQ Plot: Meets assumption but treads up slightly. The assumption scale would be -4 to 4.

- Homoskedasticity: No cone shape, passes baseline model
- Non-Statistically Significant Predictors: lot_vs_ne
- Not going to use this model because there is too

Contents ⚙️

▼ 1 King County Housing Characteristics
1.1 Data
1.2 Roadmap
▼ 2 Scrub Data
2.1 Descriptions of columns
▼ 2.2 Handling Null Values
2.2.1 Fill in missing Values for 'view' column
2.2.2 Fill in missing Values for 'yr_renovated'
2.2.3 Fill in missing Values for 'waterfront'
2.3 Handling Duplicates
▼ 3 Exploratory Data Analysis
3.1 Handling Error in Basement encoding
▼ 3.2 Return to checking distributions, outliers
3.2.1 Individual EDA Analysis
▼ 3.2.2 Overall EDA Analysis
3.2.2.1 Handle Bedroom error
3.2.3 Turn yr_renovated into binary variable
▼ 4 Feature Engineering
4.1 Total Rooms
4.2 Backyard Size as a proportion of the lot
4.3 Comparison of Square Foot living area
▼ 5 Check Assumptions of Linearity and Multicollinearity
5.1 Check Assumption of Linearity
5.2 Check Assumption of Multicollinearity
6 Model 1: Baseline Model
▼ 7 Outlier Removal: IQR + Z-Score
▼ 7.1 IQR Method
7.1.1 IQR Method Across All Columns
7.1.1.1 Model 2: IQR All Outliers Removed
▼ 7.1.2 IQR Price Outliers Removed
7.1.2.1 Model 3: IQR Price Outliers Removed
▼ 7.2 Z-Score Method
7.2.1 Z-Score Method Across All Columns
7.2.2 Model 3: Z-Score All Outliers Removed
▼ 7.2.3 Z-Score Price Outliers Removed
7.2.3.1 Model 4: Z-Score Price Outliers Removed
7.3 Table to Compare 4 Outlier Removal Methods
▼ 8 Handling Categorical Variables with One Hot Encoding
8.1 Check relationship of Non-Linear Categorical Variables
▼ 8.2 One Hot Encode Categorical Non-Or
8.2.1 Model 5: OHE Iteration 1
8.2.2 Model 5: Iteration 2 - Handling Feature Interaction
▼ 9 Interpretation
9.1 Standardize data for interpretation
▼ 9.2 View which predictors make the model perform well
9.2.1 Digging Deep into Zip codes

7.1.2 IQR Price Outliers Removed

- Rather than considering an observation an outlier
- This method reduces data loss because we are

In [81]:

```
1 # Finding the upper bound of the price
2 # Since prices cannot be negative,
3
4 res=df_linco['price'].describe()
5 thresh = res['75%'] - res['25%']
6 u_bound=res['75%']+1.5*thresh
7 u_bound
```

Out[81]: 1125564.75

Maximum price for this model will be \$1,125,564

In [82]:

```
1 # Subsetting the data to only include
2 # observations above the upper bound
3 df_iqrp = df_linco[df_linco['price'] >= u_bound]
```

In [83]:

```
1 # Ensure that observations have been
2 # removed
3 print(len(df_iqrp))
4 print(df_iqrp['price'].max())
```

20235
1120000.0

Confirmed that we have less data than the original DataFrame

In [84]:

```
1 print(f'Num observations before dropping with IQR: {len(df_linco)}')
2 print(f'Num observations after dropping with IQR: {len(df_iqrp)}')
3 print(f'Num observations removed: {len(df_linco) - len(df_iqrp)}')
4 print(f'Num observations removed as percent of total: {((len(df_linco) - len(df_iqrp)) / len(df_linco)) * 100}')
```

Num observations before dropping with IQR: 20235
Num observations after dropping with IQR: 1120000.0
Num observations removed: 1,152
Num observations removed as percent of total: 5.70%

With this type of outlier removal we have much lower R-squared and adjusted R-squared values. The model performs with these constraints.

Compared to removing values based on a single feature means that our model can provide inferences for a wider range of values.

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Interaction
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [85]:

```
1 max_ip = df_iqrp['price'].max()
2 min_ip = df_iqrp['price'].min()
3
4 print(f'Min Price: ${min_ip:,}')
5 print(f'Max Price: ${max_ip:,}')
```

```
Min Price: $78,000.0
Max Price: $1,120,000.0
```

Our model provides inferential capabilities for homes.

7.1.2.1 Model 3: IQR Price Outliers Removed

- Considering an observation an outlier if price is outside the IQR
- This method reduces data loss compared to previous methods

In [86]:

```

1 model_iqrp = model_summary(df_iqrp,
2 sked_show(df_iqrp, x_targs, model_i

```

OLS Regression Results

Dep. Variable:	price	R-squared		
Model:	OLS	Adj. R-squared		
Method:	Least Squares	F-statistic		
Date:	Wed, 05 May 2021	Prob (F-statistic)		
Time:	14:10:57	Log-Likelihood		
No. Observations:	20235	AIC		
Df Residuals:	20220	BIC		
Df Model:	14			
Covariance Type:	nonrobust			
	coef	std err	t	P
Intercept	-2.723e+07	1.67e+06	-16.304	0.0
bedrooms	-7627.4038	1221.868	-6.242	0.0
sqft_living	121.3012	2.212	54.847	0.0
floors	1.412e+04	1936.055	7.293	0.0
waterfront	1.211e+05	1.76e+04	6.880	0.0
view	3.667e+04	1428.352	25.676	0.0
condition	4.31e+04	1374.848	31.352	0.0
grade	6.184e+04	1288.038	48.009	0.0
zipcode	-134.2122	20.068	-6.688	0.0
lat	5.895e+05	6321.978	93.251	0.0
long	-9.848e+04	7632.839	-12.902	0.0
basementyes	1.008e+04	2030.602	4.962	0.0
renovated_yes	8.853e+04	4918.509	18.000	0.0
living_vs_neighborhood	-6.305e+04	3695.609	-17.062	0.0
lot_vs_neighborhood	4833.7807	678.485	7.124	0.0
Omnibus:	1847.605	Durbin-Watson:	1.83	
Prob(Omnibus):	0.000	Jarque-Bera (JB):	3189.74	
Skew:	0.655	Prob(JB):	0.0	
Kurtosis:	4.438	Cond. No.	1.96e+00	

Notes:

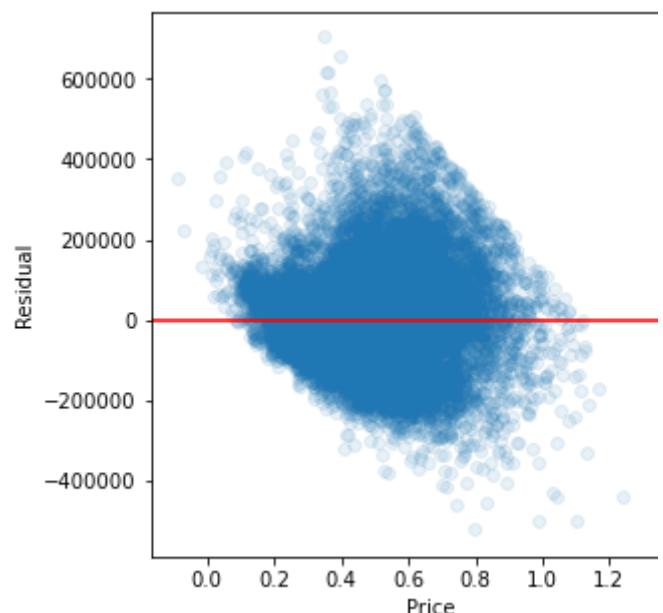
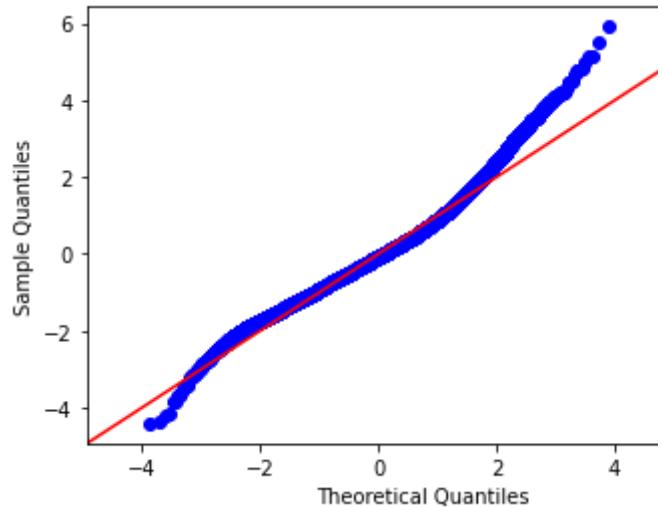
Contents ⚙

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the house
 - 4.3 Comparison of Square Foot living area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - ▼ 7.1.1 IQR Method Accross All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Accross All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orde
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

[1] Standard Errors assume that the covariance matrix is full rank.

[2] The condition number is large, 1.96e+08. This might indicate strong multicollinearity or other numerical problems.

Out[86]: (`<Figure size 360x360 with 1 Axes>`,
`<AxesSubplot:xlabel='Price', ylabel='R`



Conclusions

- $R^2: 0.667$
- Adjusted $R^2: 0.667$
- QQ Plot: Meets assumption but treads up slightly
- Homoskedasticity: No cone shape, passes assumption
- Non-Statistically Significant Predictors: None
- This model is highly preferable to removing all categorical variables for wider range of inferences

Contents ⚙

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Features
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

7.2 Z-Score Method

- Z-score outlier removal is another method for cleaning data
- To begin we standardize our values as z scores.
- This allows us to compare different features without scaling
- We will not standardize Boolean values
- From there, determine if a value is an outlier based on the 3-sigma rule
- This rule of thumb comes from the empirical rule
- Essentially, we are classifying a value as an outlier if it is more than 3 standard deviations away from the mean
- This method is less strict at classifying outliers than IQR

7.2.1 Z-Score Method Across All Columns

- In this approach, we are going to classify outliers across all columns
- For example, if an observation has bedrooms > 10, it is an outlier
- This method is more strict for determining outliers than IQR
- CAUTION: we may have significant data loss

```
In [87]: 1 # Create scaler object
          2
          3 scaler = StandardScaler()
          4 scaler
```

Out[87]: StandardScaler()

```
In [88]: 1 # Create new DF to prepare for fit
          2
          3 df_z = df_linco.copy()
```

```
In [89]: 1 df_z.columns
```

```
Out[89]: Index(['id', 'date', 'price', 'bedrooms',
       'waterfront', 'view', 'condition',
       'basementyes', 'renovated_yes',
       'lot_vs_neighborhood'],
      dtype='object')
```

```
In [90]: 1 # Not scaling binary variables such
          2 # Binary Variables are already encoded
          3 # 0 or 1
          4 # Note: Scaling does not actually change the values
          5
          6 cols_to_scale = ['price', 'bedrooms',
                            'view', 'condition', 'grade',
                            'living_vs_neighborhood',
                            'lot_vs_neighborhood']
```

Contents

▼ 1 King County Housing Characteristics
1.1 Data
1.2 Roadmap
▼ 2 Scrub Data
2.1 Descriptions of columns
▼ 2.2 Handling Null Values
2.2.1 Fill in missing Values for 'view' column
2.2.2 Fill in missing Values for 'yr_renovated'
2.2.3 Fill in missing Values for 'waterfront'
2.3 Handling Duplicates
▼ 3 Exploratory Data Analysis
3.1 Handling Error in Basement encoding
▼ 3.2 Return to checking distributions, outliers
3.2.1 Individual EDA Analysis
▼ 3.2.2 Overall EDA Analysis
3.2.2.1 Handle Bedroom error
3.2.3 Turn yr_renovated into binary variable
▼ 4 Feature Engineering
4.1 Total Rooms
4.2 Backyard Size as a proportion of the lot area
4.3 Comparison of Square Foot living area vs price
▼ 5 Check Assumptions of Linearity and Multicollinearity
5.1 Check Assumption of Linearity
5.2 Check Assumption of Multicollinearity
6 Model 1: Baseline Model
▼ 7 Outlier Removal: IQR + Z-Score
▼ 7.1 IQR Method
7.1.1 IQR Method Accross All Columns
7.1.1.1 Model 2: IQR All Outliers Removed
7.1.2 IQR Price Outliers Removed
7.1.2.1 Model 3: IQR Price Outliers Removed
▼ 7.2 Z-Score Method
7.2.1 Z-Score Method Accross All Columns
7.2.2 Model 3: Z-Score All Outliers Removed
7.2.3 Z-Score Price Outliers Removed
7.2.3.1 Model 4: Z-Score Price Outliers Removed
7.3 Table to Compare 4 Outlier Removal Methods
▼ 8 Handling Categorical Variables with One Hot Encoding
8.1 Check relationship of Non-Linear Categorical Variables
▼ 8.2 One Hot Encode Categorical Non-Or
8.2.1 Model 5: OHE Iteration 1
8.2.2 Model 5: Iteration 2 - Handling Feature Selection
▼ 9 Interpretation
9.1 Standardize data for interpretation
▼ 9.2 View which predictors make the most difference
9.2.1 Digging Deep into Zip codes

```
In [91]: 1 # Fit and transform original values
2
3 df_z[cols_to_scale] = scaler.fit_tr
4 df_z.describe()
```

Out[91]:

	id	price	bedrooms
count	21387.00000	21387.00000	21387.00000
mean	4581721940.59443	0.00000	0.00000
std	2876772841.46664	1.00002	1.00002
min	1000102.00000	-1.26066	-2.62718
25%	2124049194.50000	-0.58940	-0.41185
50%	3904930240.00000	-0.24815	-0.41185
75%	7309100170.00000	0.28260	0.69582
max	9900000190.00000	19.48493	8.44950

As we can see, the columns that we have scaled hav

```
In [92]: 1 # Create new DataFrame where we are
2 # less than 3 and greater than -3
3
4 outliers_z = pd.DataFrame()
5
6 for col in cols_to_scale:
7     outliers_z[col] = df_z[col].abs
8
9 outliers_z['total'] = outliers_z.an
10 df_za = df_z[~outliers_z['total']].
```

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

In [93]:

```
1 print(len(df_za))
2 df_za.describe()
```

19735

Out[93]:

	id	price	bedrooms
count	19735.00000	19735.00000	19735.00000
mean	4608056409.49688	-0.14098	-0.05404
std	2876186313.51690	0.64731	0.94099
min	1000102.00000	-1.24841	-2.62718
25%	2140950145.00000	-0.61559	-0.41185
50%	3918400013.00000	-0.27945	-0.41185
75%	7334550735.00000	0.17373	0.69582
max	9900000190.00000	2.99080	2.91116

- We have created a new DataFrame that only includes the columns we care about.
- Our mean and SD are no longer (0,1) because we removed the categorical variables.
- In total, we now have 19,735 observations.

In [94]:

```
1 # Formula that returns z-score back
2
3 def z_to_value(z, mu=df['price'].mean(),
4                 std=df['price'].std()):
5     """
6         Converts z-score to original value
7     """
8     x = sd*z+mu
9     return round(x,2)
```

In [95]:

```
1 print(f'Max price: ${z_to_value(2.9,
2 print(f'Min price: ${z_to_value(-1.5)}
```

Max price: \$1,639,733.25
Min price: \$82,490.54

The model will be able to infer prices of homes between \$82,490.54 and \$1,639,733.25. The Z-score is less strict in terms of classifying outliers.

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, out
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs. square foot of the house
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Accross All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Accross All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [96]:

```

1 print(f'Num observations before drop')
2 print(f'Num observations after drop')
3 print(f'Num observations removed: {')
4 print(f'Num observations removed as '

```

Num observations before dropping with I
 Num observations after dropping with IQ
 Num observations removed: 1,652
 Num observations removed as percent of

With this type of outlier removal our data loss is apro
 are met

This method will remove more data then when we or

7.2.2 Model 3: Z-Score All Outliers Rei

Contents ↻

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' cc
 - 2.2.2 Fill in missing Values for 'yr_reno
 - 2.2.3 Fill in missing Values for 'waterfr
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, out
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary val
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the
 - 4.3 Comparison of Square Foot living an
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearit
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Accross All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Rem
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers R
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Accross All Col
 - 7.2.2 Model 3: Z-Score All Outliers Rei
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outli
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Cat
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling F
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the mos
 - 9.2.1 Digging Deepr Into Zip codes

In [97]:

```

1 model_za = model_summary(df_za, x_t
2 sked_show(df_za, x_targs, model_za)

```

OLS Regression Results

Dep. Variable:	price	R-squared		
Model:	OLS	Adj. R-squared		
Method:	Least Squares	F-statistic		
Date:	Wed, 05 May 2021	Prob (F-statistic)		
Time:	14:10:58	Log-Likelihood		
No. Observations:	19735	AIC		
Df Residuals:	19720	BIC		
Df Model:	14			
Covariance Type:	nonrobust			
	coef	std err	t	P> t
Intercept	-81.9507	5.385	-15.219	0.000
bedrooms	-0.0249	0.004	-6.775	0.000
sqft_living	0.3705	0.007	55.250	0.000
floors	0.0162	0.003	4.762	0.000
waterfront	0.6130	0.155	3.962	0.000
view	0.1000	0.005	19.771	0.000
condition	0.0888	0.003	30.751	0.000
grade	0.2319	0.005	46.912	0.000
zipcode	-0.0006	6.47e-05	-9.042	0.000
lat	1.6675	0.020	81.557	0.000
long	-0.4903	0.025	-19.848	0.000
basementyes	0.0290	0.007	4.421	0.000
renovated_yes	0.2978	0.016	18.476	0.000
living_vs_neighborhood	-0.0819	0.004	-18.459	0.000
lot_vs_neighborhood	0.0457	0.009	4.909	0.000
Omnibus:	5210.929	Durbin-Watson:	1.41	
Prob(Omnibus):	0.000	Jarque-Bera (JB):	21272.01	
Skew:	1.256	Prob(JB):	0.0	
Kurtosis:	7.423	Cond. No.	1.96e+11	

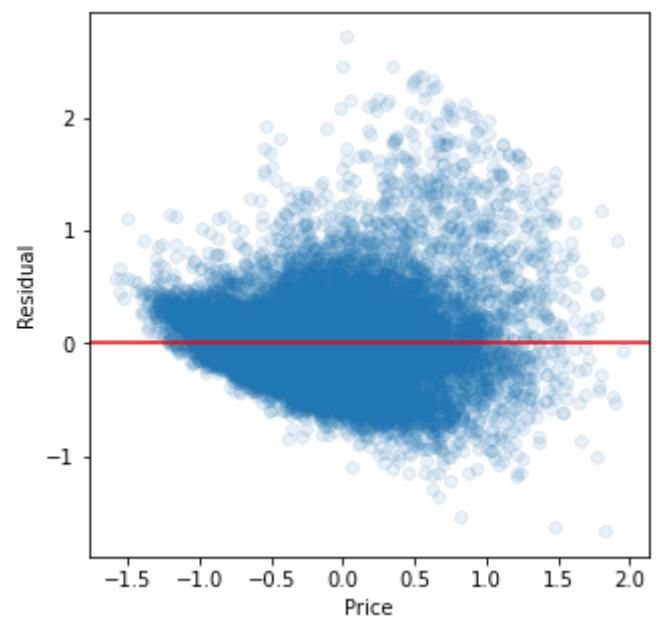
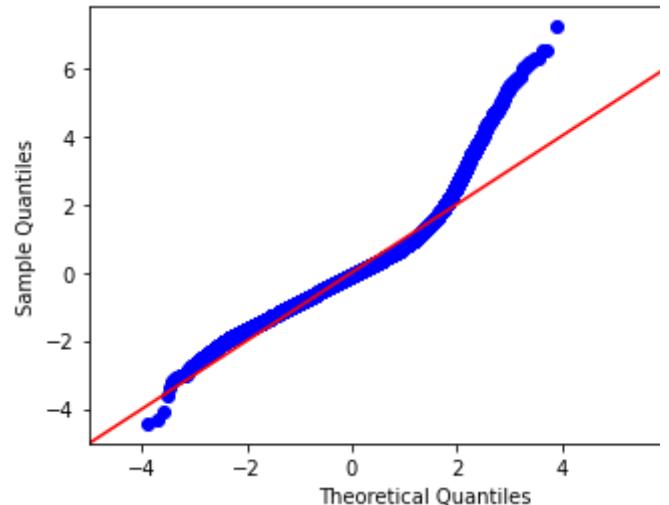
Notes:

Contents ⚙

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.2 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs. square foot of the lot
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - ▼ 7.1.1 IQR Method Accross All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Accross All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
- 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Ordeered Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Collinearity
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

- [1] Standard Errors assume that the covariance matrix is full rank.
- [2] The condition number is large, 1.96×10^8 . This might indicate strong multicollinearity or other numerical problems.

Out[97]: (<Figure size 360x360 with 1 Axes>,
<AxesSubplot:xlabel='Price', ylabel='R>



Conclusions

- $R^2: 0.660$
- Adjusted $R^2: 0.660$
- QQ Plot: Does not do a good job at meeting assumptions of normality and linearity.
- Homoskedasticity: Cone shaped, especially as price increases.
- Non-Statistically Significant Predictors: None
- Model does not meet all 4 assumptions, however.

7.2.3 Z-Score Price Outliers Removed

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs. total area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
- 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Features
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

- Rather than considering an observation an outlier
- This method reduces data loss because we are
- Will be using same scaled data from previous m

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary values
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [98]:

```

1 # Check to see if we have values greater than 5 standard deviations
2 # Only need to worry about values greater than 5 stds
3
4 df_z.describe()

```

Out[98]:

	id	price	bedrooms
count	21387.00000	21387.00000	21387.00000
mean	4581721940.59443	0.00000	0.00000
std	2876772841.46664	1.00002	1.00002
min	1000102.00000	-1.26066	-2.62718
25%	2124049194.50000	-0.58940	-0.41185
50%	3904930240.00000	-0.24815	-0.41185
75%	7309100170.00000	0.28260	0.69582
max	9900000190.00000	19.48493	8.44950

In [99]:

```

1 # Create new DataFrame only containing rows where price < 3
2
3 df_zp = df_z[df_z['price'] < 3]

```

In [100]:

```

1 # Confirm that data has been removed
2 # Confirm that max price is less than 3
3
4 print(len(df_zp))
5 df_zp['price'].max()

```

19735

Out[100]: 2.990795957515217

In [101]:

```

1 print(f'Max price: ${z_to_value(2.990795957515217)}')
2 print(f'Min price: ${z_to_value(-1.26066)}')

```

Max price: \$1,639,733.25

Min price: \$77,989.74

The model will be able to infer prices of homes between -\$1.26066 and \$1,639,733.25

In [102]:

```

1 print(f'Num observations before dro
2 print(f'Num observations after drop
3 print(f'Num observations removed: {
4 print(f'Num observations removed as

```

Num observations before dropping with I
 Num observations after dropping with IQ
 Num observations removed: 404
 Num observations removed as percent of

With this type of outlier removal our data loss is apro
 are met

This method removes the least amount of data which
 Second, Z-Score outlier removal is more strict than I

7.2.3.1 Model 4: Z-Score Price Outliers Removed

- Considering an observation an outlier if price is >
- This method reduces data loss compared to pre

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [103]:

```

1 model_zp = model_summary(df_zp, x_t
2 sked_show(df_zp, x_targs, model_zp)

```

OLS Regression Results

Dep. Variable:	price	R-squared		
Model:	OLS	Adj. R-squared		
Method:	Least Squares	F-statistic		
Date:	Wed, 05 May 2021	Prob (F-statistic)		
Time:	14:10:58	Log-Likelihood		
No. Observations:	20983	AIC		
Df Residuals:	20968	BIC		
Df Model:	14			
Covariance Type:	nonrobust			
	coef	std err	t	P> t
Intercept	-79.8958	5.424	-14.730	0.000
bedrooms	-0.0267	0.004	-7.542	0.000
sqft_living	0.3624	0.006	57.164	0.000
floors	0.0146	0.003	4.305	0.000
waterfront	0.5487	0.047	11.691	0.000
view	0.0991	0.003	30.002	0.000
condition	0.0860	0.003	29.774	0.000
grade	0.2411	0.005	50.040	0.000
zipcode	-0.0006	6.52e-05	-9.711	0.000
lat	1.7002	0.021	81.776	0.000
long	-0.4999	0.025	-20.061	0.000
basementyes	0.0215	0.007	3.271	0.001
renovated_yes	0.2953	0.016	19.038	0.000
living_vs_neighborhood	-0.0637	0.004	-17.037	0.000
lot_vs_neighborhood	0.0156	0.003	5.641	0.000
Omnibus:	4868.625	Durbin-Watson:	1.41	
Prob(Omnibus):	0.000	Jarque-Bera (JB):	17767.81	
Skew:	1.134	Prob(JB):	0.0	
Kurtosis:	6.896	Cond. No.	1.95e+1	

Notes:

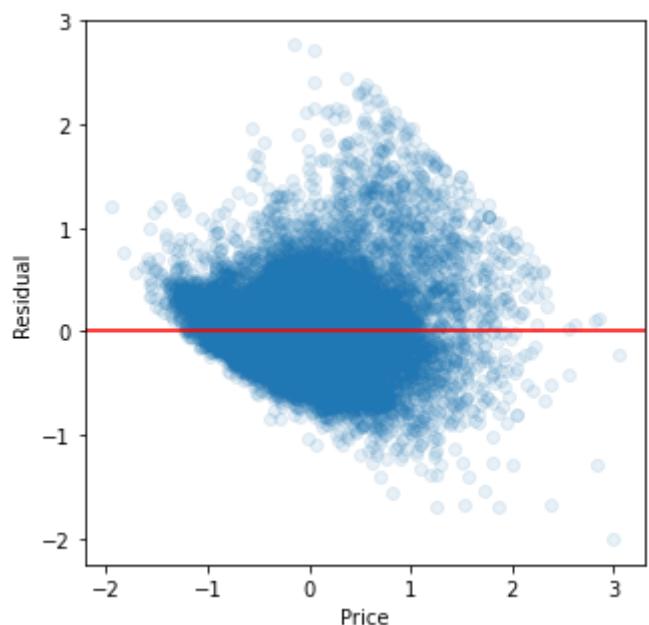
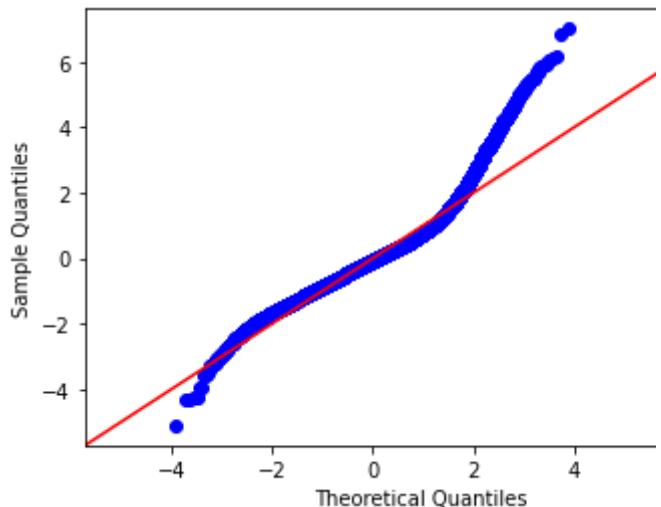
Contents ⚙

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.2 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs. square foot of the lot
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Accross All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Accross All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orde
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Collinearity
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

[1] Standard Errors assume that the covariance matrix is full rank.

[2] The condition number is large, 1.95e+08. This might indicate strong multicollinearity or other numerical problems.

Out[103]: (`<Figure size 360x360 with 1 Axes>`,
`<AxesSubplot:xlabel='Price', ylabel='R`



Conclusions

- $R^2: 0.686$
- Adjusted $R^2: 0.686$
- QQ Plot: Does not do a good job at meeting assumptions.
- Homoskedasticity: Cone shaped, especially as price increases.
- Non-Statistically Significant Predictors: Bathroom.
- Model does not meet all 4 assumptions, however.

Unfortunately, this model is not sufficient because it does not meet all four assumptions.

Contents ⚙

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs. total area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
- ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
- 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Features
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

7.3 Table to Compare 4 Outlier

- Want a simple way to evaluate all 4 models and

Contents ⚙

- 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- 2 Scrub Data
 - 2.1 Descriptions of columns
 - 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs Price
- 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- 7 Outlier Removal: IQR + Z-Score
 - 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
- 7.3 Table to Compare 4 Outlier Removal
- 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - 8.2 One Hot Encode Categorical Non-Ordinal Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Interaction
- 9 Interpretation
 - 9.1 Standardize data for interpretation
 - 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deeper Into Zip codes

In [104]:

```

1 # Create DataFrame that compares all four outlier removal types
2
3 d = {
4     'Outlier Type': ['IQR-All', 'IQR-Price', 'Z-All', 'Z-Price'],
5     'Data Loss %': [27.6, 5.4, 7.7, 1.9],
6     'R^2': [0.676, 0.677, 0.660, 0.686],
7     'Homoskedacity': ['Pass', 'Pass', 'Fail', 'Pass'],
8     'QQ Plot': ['Pass', 'Pass', 'Fail', 'Pass'],
9     'Min Price': [81000, 78000.0, 82000, 1120000],
10    'Max Price': [1120000, 1120000, 1120000, 1120000]
11 }
12 table_o = pd.DataFrame(d)
13 table_o.set_index('Outlier Type')

```

Out[104]:

	Data Loss %	R^2	Homoskedacity
Outlier Type			
IQR-All	27.60000	0.67600	Pass
IQR-Price	5.40000	0.67700	Pass
Z-All	7.70000	0.66000	Fail
Z-Price	1.90000	0.68600	Fail

As we can see, each outlier removal type creates very different assumptions about the assumption of homoskedasticity. For that reason, I am going to use the Z-Price method because the data loss is the lowest.

In conclusion, will move forward modeling with **IQR-Price**.

8 Handling Categorical Variables

- Dummy variables (one hot encoded variables) mitigate multicollinearity. In other words, one of the dummy variables must be statistically significant.
- Our next step is to One Hot Encode the ordinal variables.
- These variables, when evaluated from an ordinal perspective, are not statistically significant.
- However, we will evaluate their P-Values to determine if they are statistically significant.
 - Majority of OHE variables must be statistically significant.
 - If not, can potentially feature engineer them.

8.1 Check relationship of Non-Linear Categorical Variables

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
- 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orde
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Interaction
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [105]:

```

1 def ordinal_check(df, col, val='pri
2 """
3     Produces stripplot and barplot
4     the feature and price
5 """
6 fig, axes = plt.subplots(ncols=
7     sns.stripplot(data=df, x=col, y=
8     sns.barplot(data=df, x=col, y=val,
9
10    fig.suptitle(f'Z-{col.upper()}')
11    plt.show()
12    print('-----')
13    print(df[col].value_counts(1))

```

In [106]:

```

1 # Based on our findings in our line
2 # Don't need to check binary variables
3
4 cat_bars = ['floors', 'view', 'condition']
5 for col in cat_bars:
6     ordinal_check(df_iqrp, col)

```

Name: zipcode, Length: 70, dtype: float

1e6

Based on the results we are going to One Hot Encode the following variables:

- Floors: Increases up to 2.5 and then decreases
- Zipcode: Completely random
- Bedrooms: Increases up to 6 and then decrease
- Condition: Increases at 3 and then decreases at 4

This is because they do not appear ordinal. In other words, they are not ordered. Check if they are statistically significant.

8.2 One Hot Encode Categorical Variables

```
In [107]: 1 from sklearn.preprocessing import OneHotEncoder
2 encoder = OneHotEncoder(sparse=False)
3 encoder
```

```
Out[107]: OneHotEncoder(drop='first', sparse=False)
```

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Fewer Variables
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

```
In [108]: 1 # Separate the columns we are going to encode
2
3 cat_cols=['floors', 'zipcode', 'bedrooms']
4
```

```
In [109]: 1 # Fit and transform categorical columns
2 # Turn matrix into DataFrame
3
4 encoder.fit(df_linco[cat_cols])
5
6 ohe_vars = encoder.transform(df_iqr)
7 encoder.get_feature_names(cat_cols)
8 cat_vars = pd.DataFrame(ohe_vars,columns=ohe_vars.feature_names_)
```

```
In [110]: 1 # Confirm variables are OHE
2
3 cat_vars
```

```
Out[110]:
```

	floors_1.5	floors_2.0	floors_2.5	floors_3.0	floors_3.5
0	0.00000	0.00000	0.00000	0.00000	0.00000
1	0.00000	1.00000	0.00000	0.00000	0.00000
2	0.00000	0.00000	0.00000	0.00000	0.00000
3	0.00000	0.00000	0.00000	0.00000	0.00000
4	0.00000	0.00000	0.00000	0.00000	0.00000
...
20230	0.00000	1.00000	0.00000	0.00000	0.00000
20231	0.00000	1.00000	0.00000	0.00000	0.00000
20232	0.00000	1.00000	0.00000	0.00000	0.00000
20233	0.00000	1.00000	0.00000	0.00000	0.00000
20234	0.00000	1.00000	0.00000	0.00000	0.00000

20235 rows × 88 columns

In [111]:

```

1 # OLS Formula does not accept '.'s
2
3 name_dict = {}
4 for col in cat_vars.columns:
5     name_dict[col]=col.replace('.','')
6 name_dict

```

Out[111]:

```

{'floors_1.5': 'floors_1_5',
 'floors_2.0': 'floors_2_0',
 'floors_2.5': 'floors_2_5',
 'floors_3.0': 'floors_3_0',
 'floors_3.5': 'floors_3_5',
 'zipcode_98002': 'zipcode_98002',
 'zipcode_98003': 'zipcode_98003',
 'zipcode_98004': 'zipcode_98004',
 'zipcode_98005': 'zipcode_98005',
 'zipcode_98006': 'zipcode_98006',
 'zipcode_98007': 'zipcode_98007',
 'zipcode_98008': 'zipcode_98008',
 'zipcode_98010': 'zipcode_98010',
 'zipcode_98011': 'zipcode_98011',
 'zipcode_98014': 'zipcode_98014',
 'zipcode_98019': 'zipcode_98019',
 'zipcode_98022': 'zipcode_98022',
 'zipcode_98023': 'zipcode_98023',
 'zipcode_98024': 'zipcode_98024',
 'zipcode_98027': 'zipcode_98027',
 'zipcode_98028': 'zipcode_98028',
 'zipcode_98029': 'zipcode_98029',
 'zipcode_98030': 'zipcode_98030',
 'zipcode_98031': 'zipcode_98031',
 'zipcode_98032': 'zipcode_98032',
 'zipcode_98033': 'zipcode_98033',
 'zipcode_98034': 'zipcode_98034',
 'zipcode_98038': 'zipcode_98038',
 'zipcode_98039': 'zipcode_98039',
 'zipcode_98040': 'zipcode_98040',
 'zipcode_98042': 'zipcode_98042',
 'zipcode_98045': 'zipcode_98045',
 'zipcode_98052': 'zipcode_98052',
 'zipcode_98053': 'zipcode_98053',
 'zipcode_98055': 'zipcode_98055',
 'zipcode_98056': 'zipcode_98056',
 'zipcode_98058': 'zipcode_98058',
 'zipcode_98059': 'zipcode_98059',
 'zipcode_98065': 'zipcode_98065',
 'zipcode_98070': 'zipcode_98070',
 'zipcode_98072': 'zipcode_98072',
 'zipcode_98074': 'zipcode_98074',
 'zipcode_98075': 'zipcode_98075',
 'zipcode_98077': 'zipcode_98077',
 'zipcode_98092': 'zipcode_98092',
 'zipcode_98102': 'zipcode_98102',
 'zipcode_98103': 'zipcode_98103',
 'zipcode_98105': 'zipcode_98105',
 'zipcode_98106': 'zipcode_98106',
 'zipcode_98107': 'zipcode_98107'}

```

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
- ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
- ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
- 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deeper Into Zip codes

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs. square foot total area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Features
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

```
'zipcode_98108': 'zipcode_98108',
'zipcode_98109': 'zipcode_98109',
'zipcode_98112': 'zipcode_98112',
'zipcode_98115': 'zipcode_98115',
'zipcode_98116': 'zipcode_98116',
'zipcode_98117': 'zipcode_98117',
'zipcode_98118': 'zipcode_98118',
'zipcode_98119': 'zipcode_98119',
'zipcode_98122': 'zipcode_98122',
'zipcode_98125': 'zipcode_98125',
'zipcode_98126': 'zipcode_98126',
'zipcode_98133': 'zipcode_98133',
'zipcode_98136': 'zipcode_98136',
'zipcode_98144': 'zipcode_98144',
'zipcode_98146': 'zipcode_98146',
'zipcode_98148': 'zipcode_98148',
'zipcode_98155': 'zipcode_98155',
'zipcode_98166': 'zipcode_98166',
'zipcode_98168': 'zipcode_98168',
'zipcode_98177': 'zipcode_98177',
'zipcode_98178': 'zipcode_98178',
'zipcode_98188': 'zipcode_98188',
'zipcode_98198': 'zipcode_98198',
'zipcode_98199': 'zipcode_98199',
'bedrooms_2': 'bedrooms_2',
'bedrooms_3': 'bedrooms_3',
'bedrooms_4': 'bedrooms_4',
'bedrooms_5': 'bedrooms_5',
'bedrooms_6': 'bedrooms_6',
'bedrooms_7': 'bedrooms_7',
'bedrooms_8': 'bedrooms_8',
'bedrooms_9': 'bedrooms_9',
'bedrooms_10': 'bedrooms_10',
'bedrooms_11': 'bedrooms_11',
'condition_2': 'condition_2',
'condition_3': 'condition_3',
'condition_4': 'condition_4',
'condition_5': 'condition_5'}
```

```
In [112]: 1 # Rename cat_vars DF with new names
2
3 cat_vars.rename(columns=name_dict,
4 cat_vars
```

Out[112]:

	floors_1_5	floors_2_0	floors_2_5	floors_3_0
0	0.00000	0.00000	0.00000	0.00000
1	0.00000	1.00000	0.00000	0.00000
2	0.00000	0.00000	0.00000	0.00000
3	0.00000	0.00000	0.00000	0.00000
4	0.00000	0.00000	0.00000	0.00000
...
20230	0.00000	1.00000	0.00000	0.00000
20231	0.00000	1.00000	0.00000	0.00000
20232	0.00000	1.00000	0.00000	0.00000
20233	0.00000	1.00000	0.00000	0.00000
20234	0.00000	1.00000	0.00000	0.00000

20235 rows × 88 columns

```
In [113]: 1 # Join OHE DataFrame back with orig
2 # Ensure each DataFrame has the same index
3 # Begin by resetting index
4
5 df_iqrp=df_iqrp.reset_index()
```

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Fewer Variables
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deeper Into Zip codes

```
In [114]: 1 # Check number of rows and ensure t
2 # Should be 20235
3
4 df_iqrp
```

Out[114]:

	index	id	date	price	bedrooms
0	0	7129300520	2014-10-13	221900.00000	
1	1	6414100192	2014-12-09	538000.00000	
2	2	5631500400	2015-02-25	180000.00000	
3	3	2487200875	2014-12-09	604000.00000	
4	4	1954400510	2015-02-18	510000.00000	
...
20230	21453	1245002281	2014-05-12	1050000.00000	
20231	21461	7010700308	2014-11-12	1010000.00000	
20232	21532	8835770330	2014-08-19	1060000.00000	
20233	21577	8672200110	2015-03-17	1090000.00000	
20234	21590	7936000429	2015-03-26	1010000.00000	

20235 rows × 18 columns

```
In [115]: 1 # Check to ensure that concat was done correctly
2
3 df_ohe = pd.concat([df_iqrp, cat_variables], axis=1)
```

Contents ↗

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot size
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Accross All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Accross All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orde
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Fewer Categorical Variables
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

```
In [116]: 1 # Ensure we do not have any missing
2
3 df_ohe.isna().sum()
```

```
Out[116]: index      0
          id        0
          date      0
          price     0
          bedrooms   0
          ...
          bedrooms_11 0
          condition_2 0
          condition_3 0
          condition_4 0
          condition_5 0
Length: 106, dtype: int64
```

We now have a new DataFrame called df_ohe that in meet assumptions of no multicollinearity, and all inde from DataFrame so they are not double counted

OHE variables are interpreted as such: With respect coefficient. For example, if floors_1_5 had a coefficie

```
In [117]: 1 # Drop non-OHE variables so they ar
2
3 cols_to_drop = ['floors', 'zipcode']
4 df_ohe.drop(cols_to_drop, axis=1, i
```

```
In [118]: 1 # Ensure we do not have extra colum
2
3 df_ohe.describe()
```

Out[118]:

	index	id	price
count	20235.00000	20235.00000	20235.0000
mean	10778.21596	4605299135.81389	477281.0373
std	6231.88071	2877559932.98793	206564.7901
min	0.00000	1000102.00000	78000.0000
25%	5381.50000	2136600137.50000	316000.0000
50%	10758.00000	3905081500.00000	439000.0000
75%	16173.00000	7338220140.00000	600000.0000
max	21596.00000	9900000190.00000	1120000.0000

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs. total area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Accross All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Accross All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Selection
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

```
In [119]: 1 # Confirm non-OHE variables have been
2
3 df_ohe.head()
```

Out[119]:

	index	id	date	price	sqft_livir
0	0	7129300520	2014-10-13	221900.00000	1180
1	1	6414100192	2014-12-09	538000.00000	2570
2	2	5631500400	2015-02-25	180000.00000	710
3	3	2487200875	2014-12-09	604000.00000	1960
4	4	1954400510	2015-02-18	510000.00000	1680

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Accross All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Accross All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orde
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Fewer Categorical Variables
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs sqft_living
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Fewer Variables
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most difference
 - 9.2.1 Digging Deep into Zip codes

```
In [120]: 1 # Create list of columns excluding
2
3 x_targs = df_ohe.columns
4 x_targs = list(x_targs)
5 x_targs = [x for x in x_targs if x
6 x_targs
```

```
Out[120]: ['sqft_living',
'waterfront',
'vew',
'grade',
'lat',
'long',
'basementyes',
'renovated_yes',
'living_vs_neighborhood',
'lot_vs_neighborhood',
'floors_1_5',
'floors_2_0',
'floors_2_5',
'floors_3_0',
'floors_3_5',
'zipcode_98002',
'zipcode_98003',
'zipcode_98004',
'zipcode_98005',
'zipcode_98006',
'zipcode_98007',
'zipcode_98008',
'zipcode_98010',
'zipcode_98011',
'zipcode_98014',
'zipcode_98019',
'zipcode_98022',
'zipcode_98023',
'zipcode_98024',
'zipcode_98027',
'zipcode_98028',
'zipcode_98029',
'zipcode_98030',
'zipcode_98031',
'zipcode_98032',
'zipcode_98033',
'zipcode_98034',
'zipcode_98038',
'zipcode_98039',
'zipcode_98040',
'zipcode_98042',
'zipcode_98045',
'zipcode_98052',
'zipcode_98053',
'zipcode_98055',
'zipcode_98056',
'zipcode_98058',
'zipcode_98059',
'zipcode_98065',
'zipcode_98070',
```

```
'zipcode_98072',
'zipcode_98074',
'zipcode_98075',
'zipcode_98077',
'zipcode_98092',
'zipcode_98102',
'zipcode_98103',
'zipcode_98105',
'zipcode_98106',
'zipcode_98107',
'zipcode_98108',
'zipcode_98109',
'zipcode_98112',
'zipcode_98115',
'zipcode_98116',
'zipcode_98117',
'zipcode_98118',
'zipcode_98119',
'zipcode_98122',
'zipcode_98125',
'zipcode_98126',
'zipcode_98133',
'zipcode_98136',
'zipcode_98144',
'zipcode_98146',
'zipcode_98148',
'zipcode_98155',
'zipcode_98166',
'zipcode_98168',
'zipcode_98177',
'zipcode_98178',
'zipcode_98188',
'zipcode_98198',
'zipcode_98199',
'bedrooms_2',
'bedrooms_3',
'bedrooms_4',
'bedrooms_5',
'bedrooms_6',
'bedrooms_7',
'bedrooms_8',
'bedrooms_9',
'bedrooms_10',
'bedrooms_11',
'condition_2',
'condition_3',
'condition_4',
'condition_5']
```

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs. square foot total area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Selection
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

8.2.1 Model 5: OHE Iteration 1

In [121]:	1	model_ohe = model_summary(df_ohe, x	
	2	sked_show(df_ohe, x_targs, model_oh	
		sqft_living	135.6679
		waterfront	1.644e+05
		view	3.255e+04
		grade	4.433e+04
		lat	1.483e+05
		long	-5.204e+04
		basementyes	-2.049e+04
		renovated_yes	4.836e+04
		living_vs_neighbor	-4.941e+04
		lot_vs_neighbor	5681.0665
		floors_1_5	1.535e+04
		floors_2_0	-7407.2962

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary values
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot
 - 4.3 Comparison of Square Foot living area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Fractional Floors
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deeper Into Zip codes

Conclusions

- R^2: 0.828
- Adjusted R^2: 0.827
- QQ Plot: Does a good job of meeting assumptions
- Homoskedasticity: Meets assumptions
- Non-Statistically Significant Predictors:
 - Maintaining bedroom because a majority of houses have one
 - Floors are majority statistically significant. Very few houses have more than 2 floors.
 - Majority of zipcodes are statistically significant.

Next step will be to turn floors into integer values

8.2.2 Model 5: Iteration 2 - Handling Fractional Floors

- In this iteration we are going to turn half floors into integer values. This may not classify as an additional floor.

In [122]:	1	df_iqrp['floors'].value_counts(1)	
Out[122]:	1.00000	0.50971	
	2.00000	0.36768	
	1.50000	0.08891	
	3.00000	0.02797	
	2.50000	0.00544	
	3.50000	0.00030	
		Name: floors, dtype: float64	

- ~51% of homes are 1 story
- ~37% are 2 stories
- ~3% are 3 stories

Going to add the half story to their respective integer

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Features
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

```
In [123]: 1 # Function to adjust floors to their integer value
2
3 df_iqrp['floors'] = df_iqrp['floors'].apply(lambda x: int(x) + int(x) % 1 == 0.5)
```

```
In [124]: 1 # Ensure transformation worked
2
3 df_iqrp['floors'].value_counts(1)
```

```
Out[124]: 1 0.59862
2 0.37312
3 0.02827
Name: floors, dtype: float64
```

As we can see, the 0.5 values were added to the respective integer values.

Now, we can try OHE with floors and the number of bedrooms.

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs Price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Fewer Variables
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [125]:

```

1 from sklearn.preprocessing import OneHotEncoder
2 encoder = OneHotEncoder(sparse=False)
3 encoder
4
5 # Separate the columns we are going to encode
6 cat_cols = []
7
8
9 # Fit and transform categorical columns
10 # Turn matrix into DataFrame
11
12 def onehotencoder(df, cat_cols):
13     encoder.fit(df[cat_cols])
14
15     ohe_vars = encoder.transform(df[cat_cols])
16     encoder.get_feature_names(cat_cols)
17     cat_vars = pd.DataFrame(ohe_vars, columns=encoder.get_feature_names(cat_cols))
18
19     # OLS Formula does not accept 'None' as a column name
20     name_dict = {}
21     for col in cat_vars.columns:
22         name_dict[col] = col.replace('None', 'None_')
23
24     # Rename cat_vars DF with new names
25     cat_vars.rename(columns=name_dict, inplace=True)
26
27     # Join OHE DataFrame back with original DataFrame
28     # Ensure each DataFrame has the same index
29     # Begin by resetting index
30
31     df = df.reset_index()
32
33     df_ohe = pd.concat([df, cat_vars], axis=1)
34
35     # Drop original column names
36     cols_to_drop = cat_cols
37     df_ohe.drop(cols_to_drop, axis=1, inplace=True)
38
39     return df_ohe
40

```

In [126]:

```

1 # Use function to re-OHE with new feature names
2
3
4 cat_cols = ['floors', 'zipcode', 'bedrooms']
5 df_ohefloors = onehotencoder(df_iqr, cat_cols)

```

```
In [127]: 1 # Should see floors_2 and floors_3
2
3 df_ohefloors.head()
```

Out[127]:

	ated_yes	living_vs_neighbor	lot_vs_neighbor	floors_2
	0	0.88060	1.00000	0.00000
	1	1.52071	0.94803	1.00000
	0	0.28309	1.24039	0.00000
	0	1.44118	1.00000	0.00000
	0	0.93333	1.07690	0.00000

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs. square foot lot area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Features
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

```
In [128]: 1 # Confirm that we do not have any null values
2
3 df_ohefloors.isna().sum().any()
```

Out[128]: False

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs sqft_basement
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Collinearity
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

In [129]:

```

1 # Create list of columns excluding
2 # Remove additional indices created
3 # Removing lat and long because zip
4
5 x_targs = df_ohefloors.columns
6 x_targs = list(x_targs)
7 x_targs = [x for x in x_targs if x
8 x_targs

```

Out[129]:

```

['sqft_living',
 'waterfront',
 'view',
 'grade',
 'basementyes',
 'renovated_yes',
 'living_vs_neighborhood',
 'lot_vs_neighborhood',
 'floors_2',
 'floors_3',
 'zipcode_98002',
 'zipcode_98003',
 'zipcode_98004',
 'zipcode_98005',
 'zipcode_98006',
 'zipcode_98007',
 'zipcode_98008',
 'zipcode_98010',
 'zipcode_98011',
 'zipcode_98014',
 'zipcode_98019',
 'zipcode_98022',
 'zipcode_98023',
 'zipcode_98024',
 'zipcode_98027',
 'zipcode_98028',
 'zipcode_98029',
 'zipcode_98030',
 'zipcode_98031',
 'zipcode_98032',
 'zipcode_98033',
 'zipcode_98034',
 'zipcode_98038',
 'zipcode_98039',
 'zipcode_98040',
 'zipcode_98042',
 'zipcode_98045',
 'zipcode_98052',
 'zipcode_98053',
 'zipcode_98055',
 'zipcode_98056',
 'zipcode_98058',
 'zipcode_98059',
 'zipcode_98065',
 'zipcode_98070',
 'zipcode_98072',
 'zipcode_98074',
 'zipcode_98075',

```

```
'zipcode_98077',
'zipcode_98092',
'zipcode_98102',
'zipcode_98103',
'zipcode_98105',
'zipcode_98106',
'zipcode_98107',
'zipcode_98108',
'zipcode_98109',
'zipcode_98112',
'zipcode_98115',
'zipcode_98116',
'zipcode_98117',
'zipcode_98118',
'zipcode_98119',
'zipcode_98122',
'zipcode_98125',
'zipcode_98126',
'zipcode_98133',
'zipcode_98136',
'zipcode_98144',
'zipcode_98146',
'zipcode_98148',
'zipcode_98155',
'zipcode_98166',
'zipcode_98168',
'zipcode_98177',
'zipcode_98178',
'zipcode_98188',
'zipcode_98198',
'zipcode_98199',
'bedrooms_2',
'bedrooms_3',
'bedrooms_4',
'bedrooms_5',
'bedrooms_6',
'bedrooms_7',
'bedrooms_8',
'bedrooms_9',
'bedrooms_10',
'bedrooms_11',
'condition_2',
'condition_3',
'condition_4',
'condition_5']
```

Targets look correct with new OHE variables and no

Next step is to run Model 5 Iteration 2 and check if fl

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view'
 - 2.2.2 Fill in missing Values for 'yr_reno'
 - 2.2.3 Fill in missing Values for 'waterfr
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, out
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary val
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the
 - 4.3 Comparison of Square Foot living an
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearit
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Accross All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers R
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Accross All Col
 - 7.2.2 Model 3: Z-Score All Outliers Re
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outli
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Cat
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling F
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the mos
 - 9.2.1 Digging Deep into Zip codes

In [130]:

1	model_ohefloors = model_summary(df_ohefloors)
2	sked_show(df_ohefloors, x_targs, model_ohefloors)
	floors_3 -6.452e+04 4129.581 -15.623 0.1
	zipcode_98002 1.167e+04 7660.497 1.523 0.1
	zipcode_98003 -4413.3594 6900.304 -0.640 0.1
	zipcode_98004 5.169e+05 8431.130 61.312 0.1
	zipcode_98005 3.309e+05 8446.732 39.169 0.1
	zipcode_98006 2.714e+05 6362.838 42.655 0.1
	zipcode_98007 2.625e+05 8732.050 30.061 0.1
	zipcode_98008 2.44e+05 6999.033 34.861 0.1
	zipcode_98010 9.473e+04 9781.390 9.685 0.1
	zipcode_98011 1.474e+05 7690.627 19.160 0.1
	zipcode_98014 1.249e+05 9099.142 13.730 0.1
	zipcode_98019 1.06e+05 7742.104 13.693 0.1

As we can see, floors are statistically significant now

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the house
 - 4.3 Comparison of Square Foot living area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or Linear Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Fewer Categorical Variables
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deeper Into Zip codes

Conclusions

- R²: 0.827
- Adjusted R²: 0.827
- QQ Plot: Does a good job of meeting assumptions
- Homoskedasticity: Meets assumptions
- Non-Statistically Significant Predictors:
 - Maintaining bedroom because a majority of houses have one
 - Majority of zipcodes are statistically significant

- **Model 5 Iteration 2 will act as our final model**
 - It does not have any statistically insignificant predictors
 - R² of 0.83 which means that the predictor explains 83% of the variance
 - Meets all assumptions:
 - Homoskedasticity
 - Predictor variables have a linear relationship with the outcome
 - No multicollinearity between predictor variables

9 Interpretation

- Share what the results of our multiple linear regression mean
- Support these points with visualizations from the notebook

9.1 Standardize data for interpretation

For our data to be on the same unit, we must use a **standardization**. Standardizing data will change the range and distribution across predictors.

Contents

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot
 - 4.3 Comparison of Square Foot living area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Fewer Variables
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most difference
 - 9.2.1 Digging Deeper Into Zip codes

```
In [131]: 1 df_zf = df_ohefloors.copy()
```

```
In [132]: 1 # Predictor values, need to evaluate
2 # standard scaler
3
4 x_scale = df_zf.columns
5 x_scale
```

```
Out[132]: Index(['level_0', 'index', 'id', 'date',
       'view', 'grade', 'lat', 'long',
       'living_vs_neighbor', 'lot_vs_neighbo
       'zipcode_98002', 'zipcode_98003',
       'zipcode_98006', 'zipcode_98007',
       'zipcode_98011', 'zipcode_98014',
       'zipcode_98023', 'zipcode_98024',
       'zipcode_98029', 'zipcode_98030',
       'zipcode_98033', 'zipcode_98034',
       'zipcode_98040', 'zipcode_98042',
       'zipcode_98053', 'zipcode_98055',
       'zipcode_98059', 'zipcode_98065',
       'zipcode_98074', 'zipcode_98075',
       'zipcode_98102', 'zipcode_98103',
       'zipcode_98107', 'zipcode_98108',
       'zipcode_98115', 'zipcode_98116',
       'zipcode_98119', 'zipcode_98122',
       'zipcode_98133', 'zipcode_98136',
       'zipcode_98148', 'zipcode_98155',
       'zipcode_98177', 'zipcode_98178',
       'zipcode_98199', 'bedrooms_2',
       'bedrooms_6', 'bedrooms_7',
       'bedrooms_11', 'condition_2',
       'condition_5'],
      dtype='object')
```

We don't scale our target variable (price) for interpretation.

```
In [133]: 1 # Not going to scale OHE variables
2 # Not going to scale binary variables
3
4 x_scale = [x for x in x_targs if x != 'price']
5
```

```
In [134]: 1 x_scale
```

```
Out[134]: ['sqft_living', 'view', 'grade', 'living_vs_neighbor', 'lot_vs_neighbo
       'zipcode_98002', 'zipcode_98003', 'zipcode_98006', 'zipcode_98007', 'zipcode_98011', 'zipcode_98014', 'zipcode_98023', 'zipcode_98024', 'zipcode_98029', 'zipcode_98030', 'zipcode_98033', 'zipcode_98034', 'zipcode_98040', 'zipcode_98042', 'zipcode_98053', 'zipcode_98055', 'zipcode_98059', 'zipcode_98065', 'zipcode_98074', 'zipcode_98075', 'zipcode_98102', 'zipcode_98103', 'zipcode_98107', 'zipcode_98108', 'zipcode_98115', 'zipcode_98116', 'zipcode_98119', 'zipcode_98122', 'zipcode_98133', 'zipcode_98136', 'zipcode_98148', 'zipcode_98155', 'zipcode_98177', 'zipcode_98178', 'zipcode_98199', 'bedrooms_2', 'bedrooms_6', 'bedrooms_7', 'bedrooms_11', 'condition_2', 'condition_5']
```

We are ready to scale all of our numeric data.

```
In [135]: 1 df_zf=df_zf.drop(['level_0', 'index'])
```

Dropping lat, and long because they are difficult to interpret

```
In [136]: 1 # Fit and transform original values
2
3 df_zf[x_scale] = scaler.fit_transform(df_zf)
4 df_zf.describe()
```

Out[136]:

	id	price	sqft_living
count	20235.00000	20235.00000	20235.00000
mean	4605299135.81389	477281.03736	0.00000
std	2877559932.98793	206564.79018	1.00000
min	1000102.00000	78000.00000	-2.0760
25%	2136600137.50000	316000.00000	-0.7451
50%	3905081500.00000	439000.00000	-0.1507
75%	7338220140.00000	600000.00000	0.5986
max	9900000190.00000	1120000.00000	7.1108

As we can observe, the column that we have used to standardize has been converted into Z-scores which allows us to compare the effect of a 1 bedroom increase and a 1 sqft_living increase now.

9.2 View which predictors make the most difference

- When evaluating coefficients, we are looking at the following:
 - As we increase this predictor variable, does the price increase?
 - How large is the coefficient. As we increase the predictor by one standard deviation, how much does the price increase?
 - Coefficient values are interpreted as 'with respect to'

Contents

- 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- 2 Scrub Data
 - 2.1 Descriptions of columns
 - 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living and price
- 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- 7 Outlier Removal: IQR + Z-Score
 - 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
- 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
- 7.3 Table to Compare 4 Outlier Removal
- 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling F
- 9 Interpretation
 - 9.1 Standardize data for interpretation
 - 9.2 View which predictors make the most difference
 - 9.2.1 Digging Deep into Zip codes

In [137]:	1	x_targs
-----------	---	---------

Out[137]:	['sqft_living', 'waterfront', 'view', 'grade', 'basementyes', 'renovated_yes', 'living_vs_neighborhood', 'lot_vs_neighborhood', 'floors_2', 'floors_3', 'zipcode_98002', 'zipcode_98003', 'zipcode_98004', 'zipcode_98005', 'zipcode_98006', 'zipcode_98007', 'zipcode_98008', 'zipcode_98010', 'zipcode_98011', 'zipcode_98014', 'zipcode_98019', 'zipcode_98022', 'zipcode_98023', 'zipcode_98024', 'zipcode_98027', 'zipcode_98028', 'zipcode_98029', 'zipcode_98030', 'zipcode_98031', 'zipcode_98032', 'zipcode_98033', 'zipcode_98034', 'zipcode_98038', 'zipcode_98039', 'zipcode_98040', 'zipcode_98042', 'zipcode_98045', 'zipcode_98052', 'zipcode_98053', 'zipcode_98055', 'zipcode_98056', 'zipcode_98058', 'zipcode_98059', 'zipcode_98065', 'zipcode_98070', 'zipcode_98072', 'zipcode_98074', 'zipcode_98075', 'zipcode_98077', 'zipcode_98092', 'zipcode_98102', 'zipcode_98103', 'zipcode_98105', 'zipcode_98106', 'zipcode_98107',]
-----------	------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Collinearity
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

```
'zipcode_98108',
'zipcode_98109',
'zipcode_98112',
'zipcode_98115',
'zipcode_98116',
'zipcode_98117',
'zipcode_98118',
'zipcode_98119',
'zipcode_98122',
'zipcode_98125',
'zipcode_98126',
'zipcode_98133',
'zipcode_98136',
'zipcode_98144',
'zipcode_98146',
'zipcode_98148',
'zipcode_98155',
'zipcode_98166',
'zipcode_98168',
'zipcode_98177',
'zipcode_98178',
'zipcode_98188',
'zipcode_98198',
'zipcode_98199',
'bedrooms_2',
'bedrooms_3',
'bedrooms_4',
'bedrooms_5',
'bedrooms_6',
'bedrooms_7',
'bedrooms_8',
'bedrooms_9',
'bedrooms_10',
'bedrooms_11',
'condition_2',
'condition_3',
'condition_4',
'condition_5']
```

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs. square foot total area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Accross All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Accross All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orde
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Selection
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

```
In [138]: 1 # Run model to get parameter (coeff
2
3 model_final = model_summary(df_zf,
4 sked_show(df_zf, x_targs, model_fin
```

OLS Regression Results

Dep. Variable:	price	R-squared	
Model:	OLS	Adj. R-squared	
Method:	Least Squares	F-statistic	
Date:	Wed, 05 May 2021	Prob (F-statistic)	
Time:	14:11:05	Log-Likelihood	
No. Observations:	20235	AIC	
Df Residuals:	20141	BIC	
Df Model:	93		
Covariance Type:	nonrobust		
coef	std err	t	P
Intercept	2.121e+05	1.882e-04	1.122e-04

Confirmed that assumptions and values look the same

Scaling data does not make a difference

```
In [139]: 1 # Use statsmodels .params to extract
2 # Include coefficient and absolute
3
4 df_coeff=pd.DataFrame({'coeff': mod
```

Contents ↗

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [140]:

1 df_coeff

Out[140]:

		coeff	abs_coeff
	Intercept	243052.98038	243052.98038
	sqft_living	105564.12643	105564.12643
	waterfront	164587.53385	164587.53385
	view	20736.15964	20736.15964
	grade	45704.32075	45704.32075

	bedrooms_11	11786.85452	11786.85452
	condition_2	48276.34358	48276.34358
	condition_3	66106.77111	66106.77111
	condition_4	87080.87703	87080.87703
	condition_5	122763.55463	122763.55463

94 rows × 2 columns

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the total area
 - 4.3 Comparison of Square Foot living area vs Price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Accross All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Accross All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orde
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Fewer Categories
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

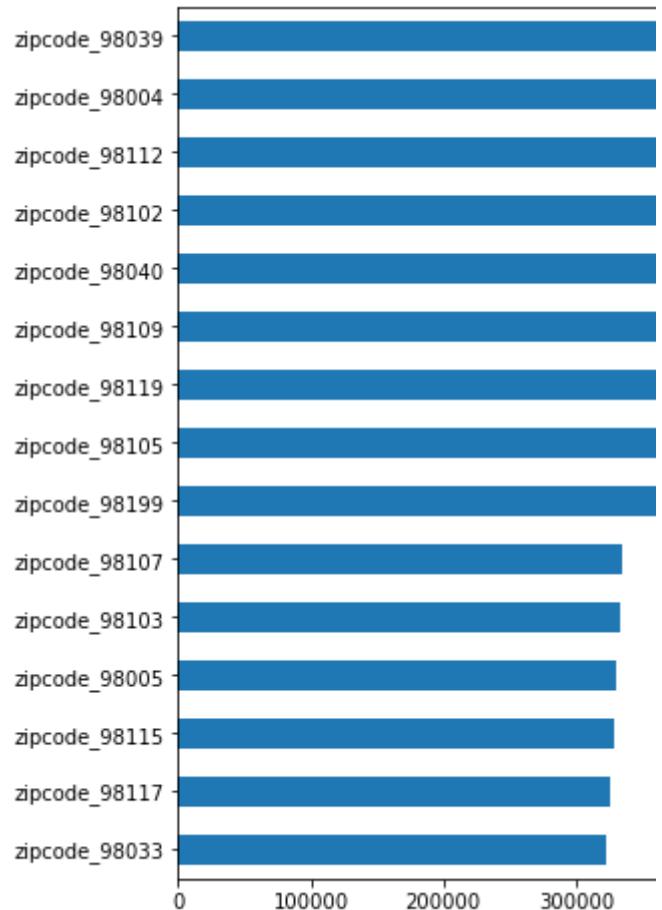
In [141]:

```

1 # Create horizontal bar plot to view
2 # absolute perspective
3
4 df_coeff.drop('Intercept', axis=0).

```

Out[141]: <AxesSubplot:>



```

1 As we can see, the 15 most impactful zipcodes can explain ~$600,000
2
3 98039, 98004, and 98112 are the top
4 Check that later
5 Almost all zip codes are statistically

```

9.2.1 Digging Deep Into Zip codes

- Evaluating how the top zipcodes compare to the rest

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs. total area
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Features
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep Into Zip codes

In [142]:

```
#1 Gather data on the top 10 zipcode
2
top_10_zips = df_iqrp.groupby('zipcode')
top_10_zips
```

Out[142]:

	zipcode	price
0	98039	901250.00000
1	98040	850000.00000
2	98004	825000.00000
3	98005	740000.00000
4	98075	725393.00000
5	98112	714250.00000
6	98109	700000.00000
7	98006	691100.00000
8	98119	667500.00000
9	98102	667475.00000

Using median in the groupby statement to handle outliers

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Interaction
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

Contents ⚙️

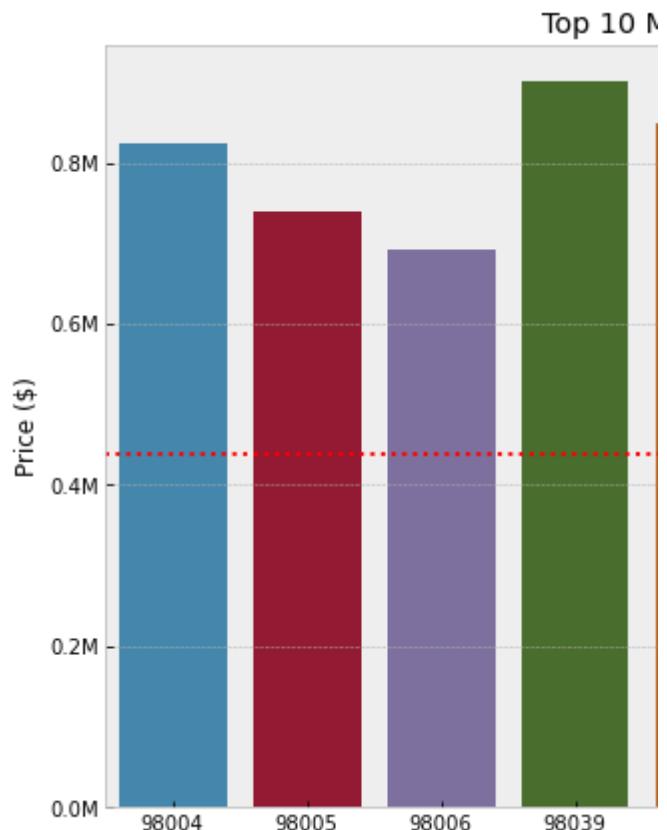
- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs Price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Selection
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [143]:

```

1 # Comparing top 10 zip codes by median home price
2
3 # Function to display values in millions
4 from matplotlib.ticker import FuncFormatter
5
6 def millions(x, pos):
7     return '%1.1fM' % (x * 1e-6)
8
9 # Kings County median home price
10 median = int(df_iqrp['price'].median())
11
12 # Construct visualization
13 with plt.style.context('bmh'):
14     fig, ax = plt.subplots(figsize=(10, 6))
15     sns.barplot(data=top_10_zips, x='zip_code', y='price')
16     ax.set_xlabel('Zip Code')
17     ax.set_ylabel('Price ($)')
18     ax.set_title('Top 10 Most Expensive Zip Codes in King County')
19     ax.axhline(median, color='r', linestyle='dashed')
20     ax.legend()
21
22     formatter = FuncFormatter(millions)
23     ax.yaxis.set_major_formatter(formatter)

```



The median home price in King County is \$439,000. The most expensive zip codes tend to be situated around the northern part of the county.

8.2 Evaluating the predictors

- Evaluate how the predictor values outside of zip

Contents

▼ 1	King County Housing Characteristics
1.1	Data
1.2	Roadmap
▼ 2	Scrub Data
2.1	Descriptions of columns
▼ 2.2	Handling Null Values
2.2.1	Fill in missing Values for 'view' column
2.2.2	Fill in missing Values for 'yr_renovated'
2.2.3	Fill in missing Values for 'waterfront'
2.3	Handling Duplicates
▼ 3	Exploratory Data Analysis
3.1	Handling Error in Basement encoding
▼ 3.2	Return to checking distributions, outliers
3.2.1	Individual EDA Analysis
▼ 3.2.2	Overall EDA Analysis
3.2.2.1	Handle Bedroom error
3.2.2.3	Turn yr_renovated into binary variable
▼ 4	Feature Engineering
4.1	Total Rooms
4.2	Backyard Size as a proportion of the lot size
4.3	Comparison of Square Foot living area vs square foot lot area
▼ 5	Check Assumptions of Linearity and Multicollinearity
5.1	Check Assumption of Linearity
5.2	Check Assumption of Multicollinearity
6	Model 1: Baseline Model
▼ 7	Outlier Removal: IQR + Z-Score
▼ 7.1	IQR Method
7.1.1	IQR Method Accross All Columns
7.1.1.1	Model 2: IQR All Outliers Removed
▼ 7.1.2	IQR Price Outliers Removed
7.1.2.1	Model 3: IQR Price Outliers Removed
▼ 7.2	Z-Score Method
7.2.1	Z-Score Method Accross All Columns
7.2.2	Model 3: Z-Score All Outliers Removed
▼ 7.2.3	Z-Score Price Outliers Removed
7.2.3.1	Model 4: Z-Score Price Outliers Removed
7.3	Table to Compare 4 Outlier Removal Methods
▼ 8	Handling Categorical Variables with One Hot Encoding
8.1	Check relationship of Non-Linear Categorical Variables
▼ 8.2	One Hot Encode Categorical Non-Or
8.2.1	Model 5: OHE Iteration 1
8.2.2	Model 5: Iteration 2 - Handling Feature Interaction
▼ 9	Interpretation
9.1	Standardize data for interpretation
▼ 9.2	View which predictors make the most difference
9.2.1	Digging Deep into Zip codes

In [144]:

```

1 # Create list of features excluding
2
3 indices = df_coeff.index
4 non_zips = []
5 for ind in indices:
6     if ind.startswith('zip'):
7         pass
8     else:
9         non_zips.append(ind)

```

In [145]:

```

1 # These are the predictors that are
2
3 non_zips

```

Out[145]:

```

['Intercept',
 'sqft_living',
 'waterfront',
 'view',
 'grade',
 'basementyes',
 'renovated_yes',
 'living_vs_neighbo',
 'lot_vs_neighbo',
 'floors_2',
 'floors_3',
 'bedrooms_2',
 'bedrooms_3',
 'bedrooms_4',
 'bedrooms_5',
 'bedrooms_6',
 'bedrooms_7',
 'bedrooms_8',
 'bedrooms_9',
 'bedrooms_10',
 'bedrooms_11',
 'condition_2',
 'condition_3',
 'condition_4',
 'condition_5']

```

In [146]:

```

1 # Select the non-zip rows
2
3 df_coeff2 = df_coeff.loc[non_zips]

```

In [147]:

```

1 # Remove intercept because it is no longer needed
2
3 df_coeff2.drop('Intercept', axis=0,

```

In [148]:

```

1 # Reset index so that we can sort by value
2
3 df_coeff2=df_coeff2.reset_index()

```

In [149]:

```

1 # Sort in descending order so the values are sorted by coefficient
2
3 df_coeff2=df_coeff2.sort_values(by=-1)

```

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Accross All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Accross All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orde
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [150]:

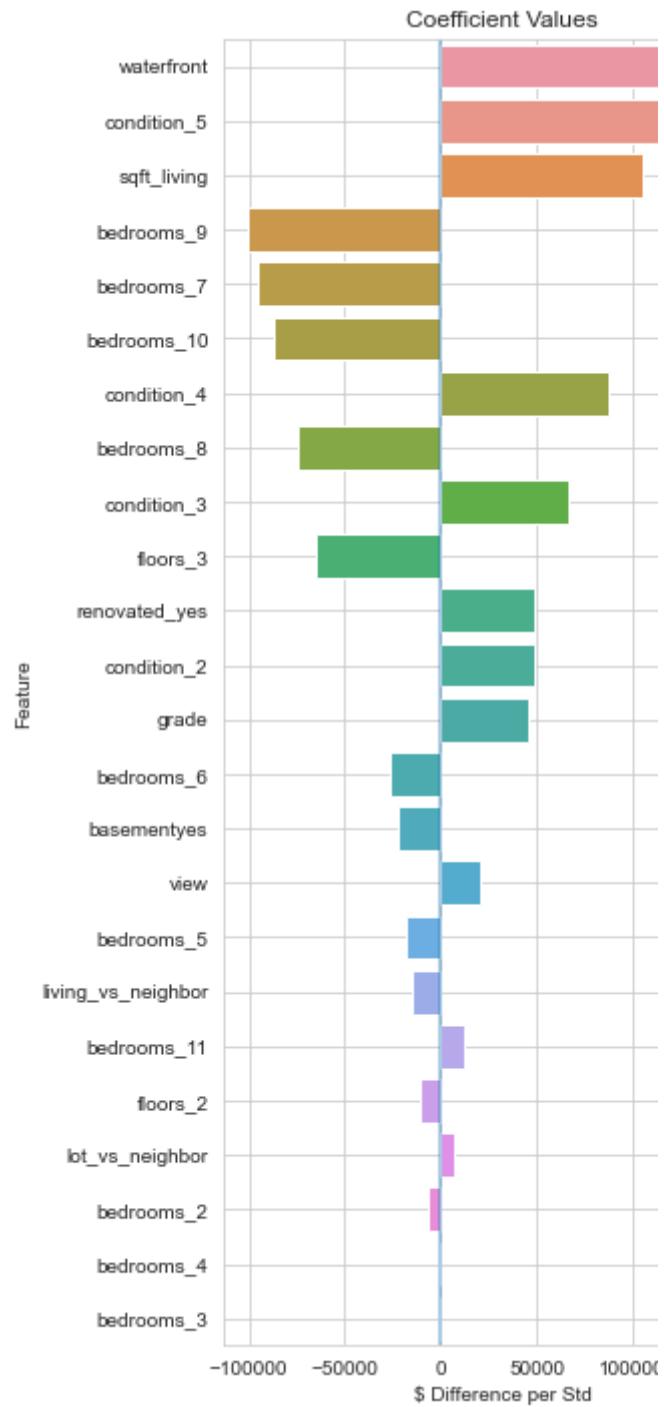
```

1 sns.set_style('whitegrid')
2 fig, ax = plt.subplots(figsize=(5,1
3
4
5 sns.barplot(data=df_coeff2, y='inde
6 ax.set_title('Coefficient Values')
7 ax.set_xlabel('$ Difference per Std')
8 ax.set_ylabel('Feature')
9
10
11 ax.axvline(0, alpha=0.5)
12 ax.grid(True)

```

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs. square foot of the entire house
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Collinearity
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most difference
 - 9.2.1 Digging Deep into Zip codes



Based on the graph above, waterfront (binary variable) adds ~\$150,000 more. Next comes condition, all condition. Condition represents the overall state of the house. It has structural integrity, be generally well kept, and it price. A 1 Standard Deviation increase in sqft_living results in a standard deviation above the mean sqft_living will contribute ~\$150,000 more.

Bedrooms 9, 7, and 10 all have negative impacts on price. We will need to OHE them.

All floors with respect to 1 floor have a negative relationship with price.

Grade has a positive impact on price, which is how to good condition.

Homes with basements have a lower price on average.

View has a positive relationship with price.

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

Contents ⚙️

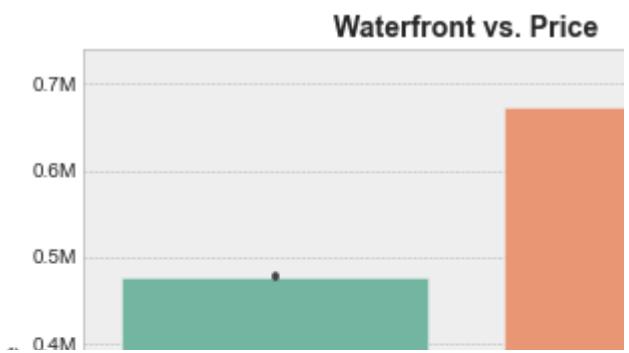
- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs Price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

In [158]:

```

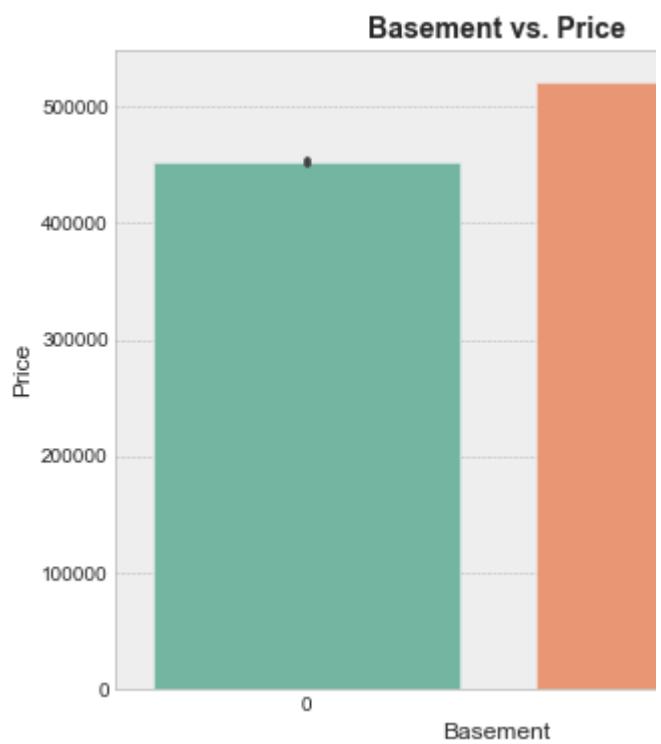
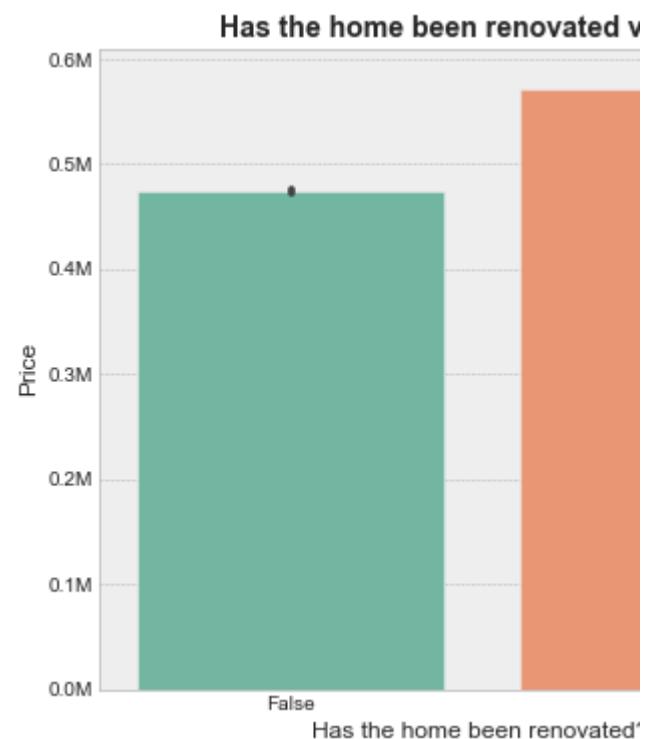
1 with plt.style.context('bmh'):
2
3     fig, axs=plt.subplots(nrows=3,
4
5         sns.barplot(data=df_ohefloors,
6         sns.barplot(data=df_ohefloors,
7         sns.barplot(data=df_ohefloors,
8         sns.barplot(data=df_ohefloors,
9         sns.barplot(data=df_ohefloors,
10        sns.barplot(data=df_iqrp, x='co'
11
12        tf_labs = [ 'False', 'True' ]
13
14        axs[0,0].set_title('Waterfront vs. Price')
15        axs[0,0].set_xlabel('Waterfront')
16        axs[0,0].set_ylabel('Price')
17        axs[0,0].set_xticklabels(tf_labs)
18
19        axs[0,1].set_title('View vs. Price')
20        axs[0,1].set_xlabel('View')
21        axs[0,1].set_ylabel('Price')
22
23        axs[1,0].set_title('Has the home vs. Price')
24        axs[1,0].set_xlabel('Has the home')
25        axs[1,0].set_ylabel('Price')
26        axs[1,0].set_xticklabels(tf_labs)
27
28        axs[1,1].set_title('Grade vs. Price')
29        axs[1,1].set_xlabel('Grade')
30        axs[1,1].set_ylabel('Price')
31
32        axs[2,0].set_title('Basement vs. Price')
33        axs[2,0].set_xlabel('Basement')
34        axs[2,0].set_ylabel('Price')
35
36        axs[2,1].set_title('Condition vs. Price')
37        axs[2,1].set_xlabel('Condition')
38        axs[2,1].set_ylabel('Price')
39
40
41        formatter = FuncFormatter(milli)
42        axs[0,0].yaxis.set_major_formatter(formatter)
43        axs[0,1].yaxis.set_major_formatter(formatter)
44        axs[1,0].yaxis.set_major_formatter(formatter)
45        axs[1,1].yaxis.set_major_formatter(formatter)
46        axs[2,0].yaxis.set_major_formatter(formatter)
47        axs[2,1].yaxis.set_major_formatter(formatter)

```



Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs Price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Features
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes



- Waterfront homes are more expensive than non-
- View increases in a semi linear fashion
- Homes that have been renovated are more expensive
- Condition has a linear relationship with price
- Homes with basements are marginally more expensive
- Condition does not have a linear relationship with price

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs Price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes

10 Models for Presentation

In [152]:

```

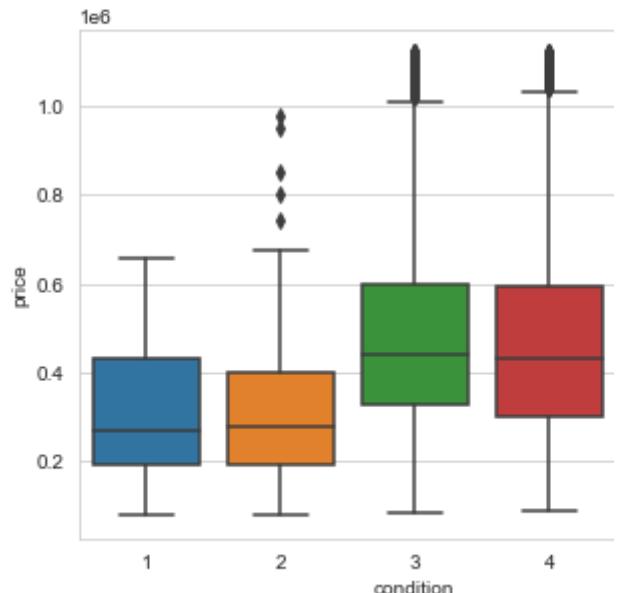
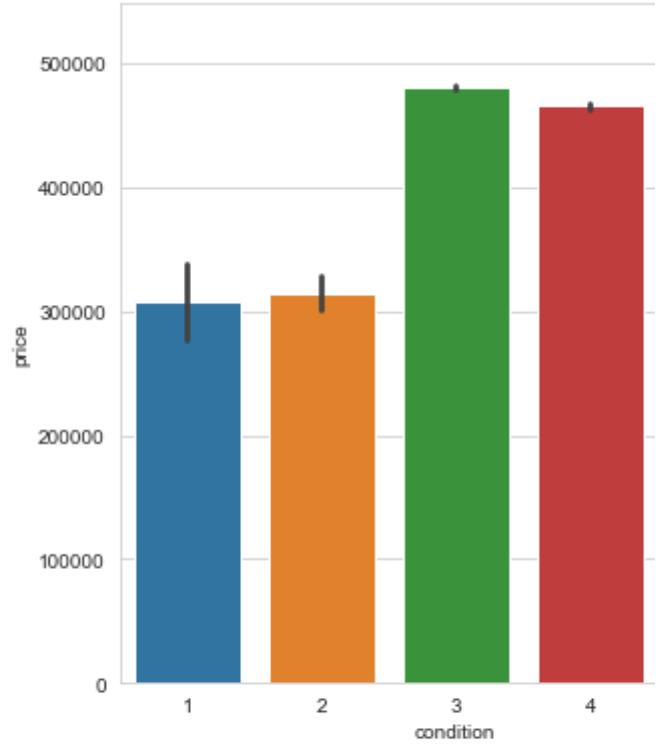
1 fig, axs=plt.subplots(nrows=2, ncol=1)
2
3 sns.barplot(data=df_iqrp, x='condition', y='price')
4 sns.boxplot(data=df_iqrp, x='condition', y='price')

```

Out[152]: <AxesSubplot:xlabel='condition', ylabel='price'

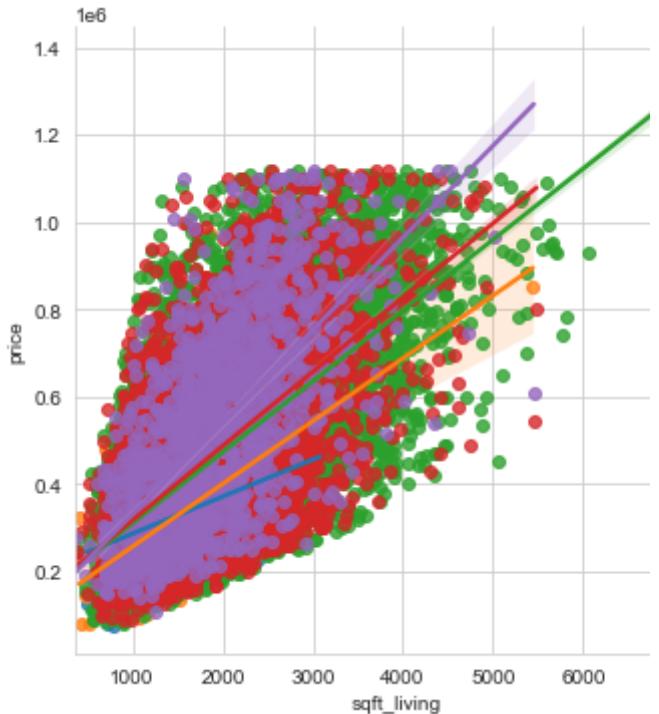
Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes



In [153]:

1 g = sns.lmplot(data=df_iqrp, x="sqf



Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot size
 - 4.3 Comparison of Square Foot living area vs Price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
dered Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Features
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

In [154]:

```

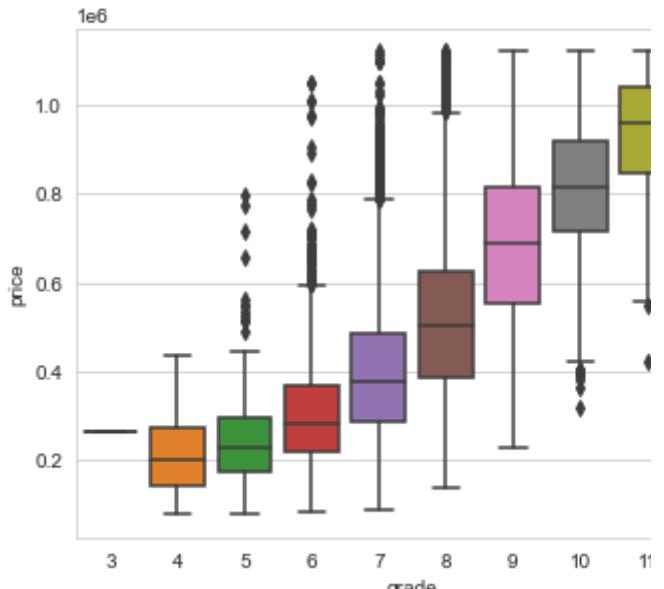
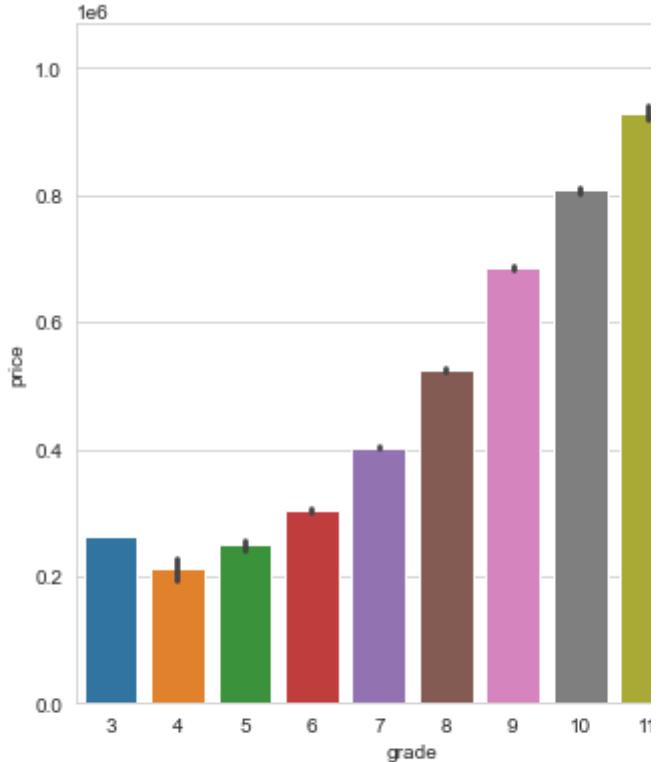
1 fig, axs=plt.subplots(nrows=2, ncol=1)
2
3 sns.barplot(data=df_iqrp, x='grade', y='price')
4 sns.boxplot(data=df_iqrp, x='grade', y='price')

```

Out[154]: <AxesSubplot:xlabel='grade', ylabel='pr

Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot size
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.2 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Or
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Feature Interaction
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most impact
 - 9.2.1 Digging Deep into Zip codes



Contents ⚙️

- ▼ 1 King County Housing Characteristics
 - 1.1 Data
 - 1.2 Roadmap
- ▼ 2 Scrub Data
 - 2.1 Descriptions of columns
 - ▼ 2.2 Handling Null Values
 - 2.2.1 Fill in missing Values for 'view' column
 - 2.2.2 Fill in missing Values for 'yr_renovated'
 - 2.2.3 Fill in missing Values for 'waterfront'
 - 2.3 Handling Duplicates
- ▼ 3 Exploratory Data Analysis
 - 3.1 Handling Error in Basement encoding
 - ▼ 3.2 Return to checking distributions, outliers
 - 3.2.1 Individual EDA Analysis
 - ▼ 3.2.2 Overall EDA Analysis
 - 3.2.2.1 Handle Bedroom error
 - 3.2.2.3 Turn yr_renovated into binary variable
- ▼ 4 Feature Engineering
 - 4.1 Total Rooms
 - 4.2 Backyard Size as a proportion of the lot area
 - 4.3 Comparison of Square Foot living area vs price
- ▼ 5 Check Assumptions of Linearity and Multicollinearity
 - 5.1 Check Assumption of Linearity
 - 5.2 Check Assumption of Multicollinearity
- 6 Model 1: Baseline Model
- ▼ 7 Outlier Removal: IQR + Z-Score
 - ▼ 7.1 IQR Method
 - 7.1.1 IQR Method Across All Columns
 - 7.1.1.1 Model 2: IQR All Outliers Removed
 - ▼ 7.1.2 IQR Price Outliers Removed
 - 7.1.2.1 Model 3: IQR Price Outliers Removed
 - ▼ 7.2 Z-Score Method
 - 7.2.1 Z-Score Method Across All Columns
 - 7.2.1.1 Model 3: Z-Score All Outliers Removed
 - ▼ 7.2.3 Z-Score Price Outliers Removed
 - 7.2.3.1 Model 4: Z-Score Price Outliers Removed
 - 7.3 Table to Compare 4 Outlier Removal Methods
- ▼ 8 Handling Categorical Variables with One Hot Encoding
 - 8.1 Check relationship of Non-Linear Categorical Variables
 - ▼ 8.2 One Hot Encode Categorical Non-Orchard Variables
 - 8.2.1 Model 5: OHE Iteration 1
 - 8.2.2 Model 5: Iteration 2 - Handling Features
- ▼ 9 Interpretation
 - 9.1 Standardize data for interpretation
 - ▼ 9.2 View which predictors make the most sense
 - 9.2.1 Digging Deep into Zip codes

In []:

1