# King County Housing Regression and Analysis

- Student name: Johnny Dryman
- Student pace: full time
- Scheduled project review date/time: 04/29/21, 2pm
- Instructor name: James Irving

# INTRODUCTION

> For the Phase 2 Project, we will be analyzing housing sales data for King County (Seattle, WA area). We will be using multivariate linear regression to explore which features of the data have the greatest influence on price.

## Business Problem

As home values continue to sky rocket in the pandemic era, many King County residents have inquired about how to increase the value of their homes. Fortunately, we have access to all homes sold in King County for roughly one year, from May 2014 - May 2015.

This data gives us access to a variety of important metrics both quantitative and qualitative.

After scrubbing the data and assuring quality, we will use multivariate linear regression to analyze our features and determine their relationship with sale price.

Finally, we will formulate our observations into useful recommendations to any resident interested in increasing their home value.

# OBTAIN

We will begin by importing our packages for data exploration and load our .csv data into a pandas dataframe.

In [616]:
```python
import pandas as pd
import seaborn as sns
sns.set_theme(color_codes=True)
import matplotlib.pyplot as plt
import numpy as np

df = pd.read_csv('data/kc_house_data.csv')

df.columns
```

Out[616]:
```
Index(['id', 'date', 'price', 'bedrooms', 'bathrooms', 'sqft_living',
       'sqft_lot', 'floors', 'waterfront', 'view', 'condition', 'grade',
       'sqft_above', 'sqft_basement', 'yr_built', 'yr_renovated', 'zipcode',
       'lat', 'long', 'sqft_living15', 'sqft_lot15'],
      dtype='object')
```

# SCRUB

## Data Preparation

We'll begin by getting a brief overview of our data and check for null values.

In [617]:
```python
1  print(df.info())
2
3  print(df.isna().sum())
```
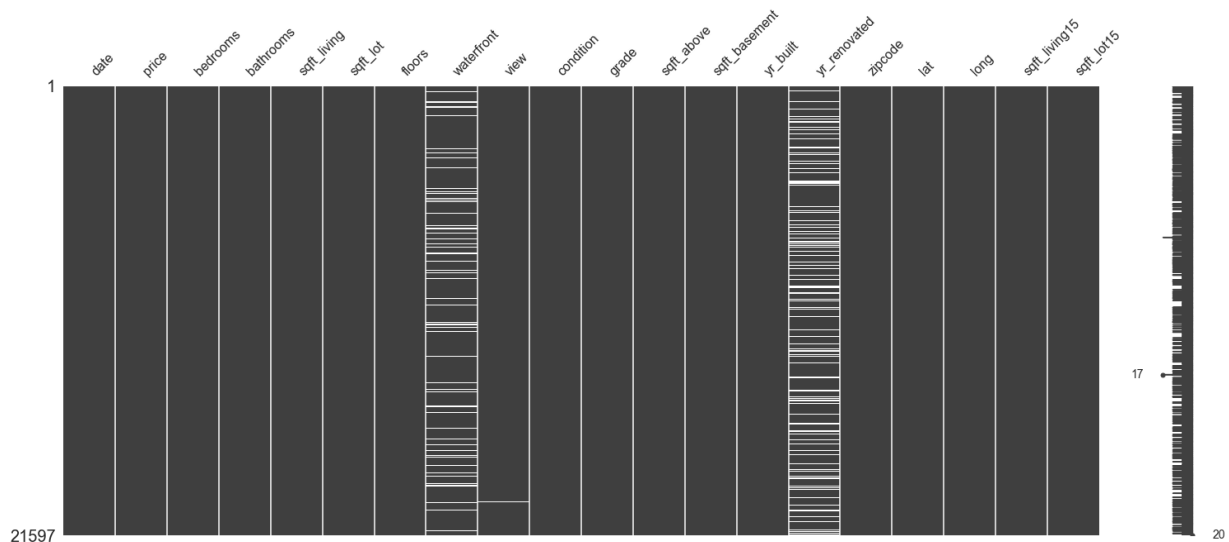
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21597 entries, 0 to 21596
Data columns (total 21 columns):
 #   Column         Non-Null Count  Dtype
---  ------         --------------  -----
 0   id             21597 non-null  int64
 1   date           21597 non-null  object
 2   price          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
None
id                  0
date                0
price               0
bedrooms            0
bathrooms           0
sqft_living         0
sqft_lot            0
floors              0
waterfront       2376
view               63
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
```

In [618]:
```python
df = df.set_index('id')
```

Using missingno package to visualize null values.

In [619]:
```python
import missingno as msno

msno.matrix(df)
```

Out[619]: <AxesSubplot:>



Creating inspect_column function, which will help us look at unique items. This will be useful for identifying any data that seems off or incorrect.

In [620]:
```python
#be cautions of naming conventions

def inspect_column(column, unique_count=10):
    column_str = str(column)
    print('Datatype: ' + str(df[column].dtypes))
    print('Total unique itms: ' + str(df[column].nunique()))
    print('Displaying first ' + str(unique_count) + ':')
    print(df[column].unique()[0:unique_count])
    return column_str

def null_count(df):
    print('---Total Entries---')
    print(df.describe())
    print('---Non-Null Values---')
    print(df.notna().describe())
```

Let's take a look at our features that have null values.

We could conceivably estimate our null values, and that might be interesting for further analysis. Mapping could be used with 'latitude' and 'longitude' and potentially calculate distance to water. However, for this project, our safest bet will be to drop the null values.

In [621]:
```
1  null_count(df['waterfront'])
```

```
---Total Entries---
count   19221.00
mean        0.01
std         0.09
min         0.00
25%         0.00
50%         0.00
75%         0.00
max         1.00
Name: waterfront, dtype: float64
---Non-Null Values---
count      21597
unique         2
top         True
freq       19221
Name: waterfront, dtype: object
```

In [622]:
```
1  inspect_column('waterfront')
2
3  df = df[df['waterfront'].notna()]
4
5  inspect_column('waterfront')
```

```
Datatype: float64
Total unique itms: 2
Displaying first 10:
[nan  0.  1.]
Datatype: float64
Total unique itms: 2
Displaying first 10:
[0. 1.]
```

Out[622]:  'waterfront'

View has very few null values, it is safe to remove them from the dataset.

In [623]:    `1` `null_count(df['view'])`

```
---Total Entries---
count    19164.00
mean         0.23
std          0.76
min          0.00
25%          0.00
50%          0.00
75%          0.00
max          4.00
Name: view, dtype: float64
---Non-Null Values---
count        19221
unique           2
top           True
freq         19164
Name: view, dtype: object
```

In [624]:    
```
1  inspect_column('view')
2
3  df = df[df['view'].notna()]
4
5  inspect_column('view')
```

```
Datatype: float64
Total unique itms: 5
Displaying first 10:
[ 0. nan  3.  4.  2.  1.]
Datatype: float64
Total unique itms: 5
Displaying first 10:
[0. 3. 4. 2. 1.]
```

Out[624]: `'view'`

Yr_renovated has ~3,500 null values. We would want to consider removing the column in this case, but yr_renovated indicates a renovation occurred with a year (e.g. 2007) and a renovation has never occurred with a zero (e.g. 0). The null values could also represent houses that have never been renovated, but we can't be sure.

In [625]: 
```
1  null_count(df['yr_renovated'])
```

```
---Total Entries---
count    15762.00
mean        82.44
std        397.21
min          0.00
25%          0.00
50%          0.00
75%          0.00
max       2015.00
Name: yr_renovated, dtype: float64
---Non-Null Values---
count        19164
unique           2
top           True
freq         15762
Name: yr_renovated, dtype: object
```

In [626]: 
```
1  inspect_column('yr_renovated', unique_count=115)
```

```
Datatype: float64
Total unique itms: 70
Displaying first 115:
[1991.    nan     0. 2002. 2010. 1992. 2013. 1994. 1978. 2005. 2003. 1984.
 1954. 2014. 2011. 1983. 1990. 1988. 1977. 1981. 1995. 2000. 1999. 1998.
 1970. 1989. 2004. 1986. 2007. 1987. 2006. 1985. 2001. 1980. 1971. 1945.
 1979. 1997. 1950. 1969. 1948. 2009. 2015. 2008. 2012. 1968. 1963. 1951.
 1962. 1953. 1993. 1955. 1996. 1982. 1956. 1940. 1976. 1946. 1975. 1964.
 1973. 1957. 1959. 1960. 1965. 1967. 1934. 1972. 1944. 1958. 1974.]
```

Out[626]: 'yr_renovated'

We will first remove rows with nan values from the dataset.

In [627]:
```python
1  df = df[df['yr_renovated'].notna()]
2
3  df.describe()
4
5
```

Out[627]:

| | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view |
|---|---|---|---|---|---|---|---|---|
| count | 15762.00 | 15762.00 | 15762.00 | 15762.00 | 15762.00 | 15762.00 | 15762.00 | 15762.00 |
| mean | 541317.18 | 3.38 | 2.12 | 2084.51 | 15280.82 | 1.50 | 0.01 | 0.23 |
| std | 372225.84 | 0.94 | 0.77 | 918.62 | 41822.88 | 0.54 | 0.09 | 0.76 |
| min | 82000.00 | 1.00 | 0.50 | 370.00 | 520.00 | 1.00 | 0.00 | 0.00 |
| 25% | 321000.00 | 3.00 | 1.75 | 1430.00 | 5048.50 | 1.00 | 0.00 | 0.00 |
| 50% | 450000.00 | 3.00 | 2.25 | 1920.00 | 7602.00 | 1.50 | 0.00 | 0.00 |
| 75% | 644875.00 | 4.00 | 2.50 | 2550.00 | 10720.00 | 2.00 | 0.00 | 0.00 |
| max | 7700000.00 | 33.00 | 8.00 | 13540.00 | 1651359.00 | 3.50 | 1.00 | 4.00 |

In [628]:
```python
1  pd.set_option('display.float_format', lambda x: '%.2f' % x)
```

In [629]:
```python
1  ren_df = df[df['yr_renovated'] != 0]
2
3  not_ren_df = df[df['yr_renovated'] == 0]
4
```

In [630]:
```python
1  ren_df['price'].describe()
```

Out[630]:
```
count        651.00
mean      760872.06
std       637150.64
min       110000.00
25%       410000.00
50%       600000.00
75%       886250.00
max      7700000.00
Name: price, dtype: float64
```

In [631]:
```python
1  not_ren_df['price'].describe()
```

Out[631]:
```
count      15111.00
mean      531858.49
std       353400.02
min        82000.00
25%       320000.00
50%       449000.00
75%       633000.00
max      6890000.00
Name: price, dtype: float64
```
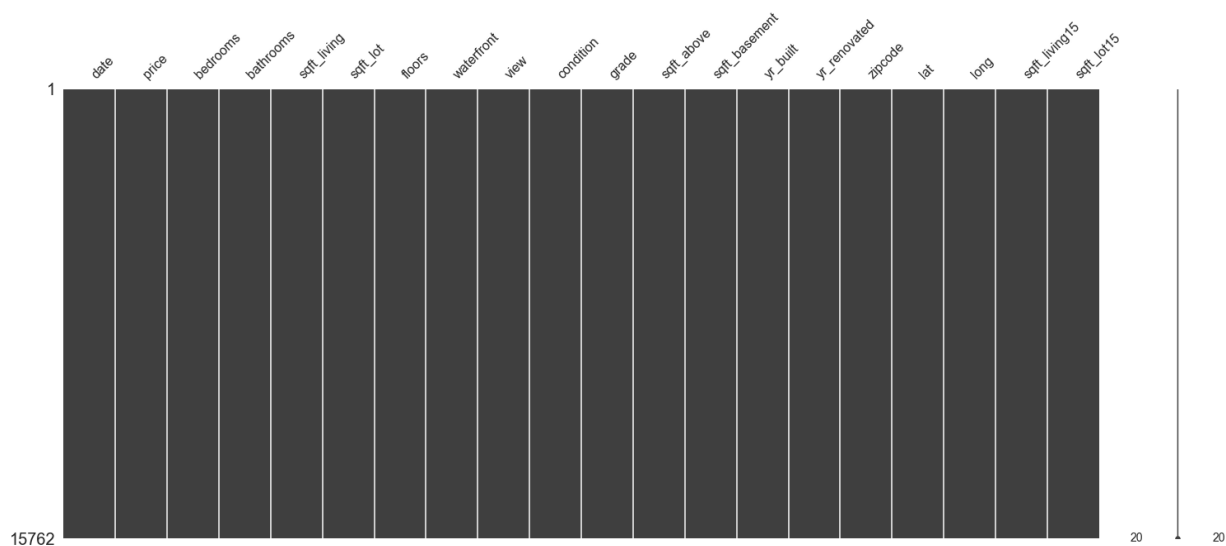
The means, standard deviations, and medians for renovated and non-renovated houses are

significant. We will revisit yr_renovated and potentially convert the column to a binary value.

Checking non-nulls again.

```
In [632]:    1  msno.matrix(df)
```

Out[632]:    <AxesSubplot:>



Now we'll take a look at each column and see if anything needs correction.

# Feature Review

In [633]:
```python
1  print(df.info())
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 15762 entries, 6414100192 to 1523300157
Data columns (total 20 columns):
 #   Column         Non-Null Count   Dtype
---  ------         --------------   -----
 0   date           15762 non-null   object
 1   price          15762 non-null   float64
 2   bedrooms       15762 non-null   int64
 3   bathrooms      15762 non-null   float64
 4   sqft_living    15762 non-null   int64
 5   sqft_lot       15762 non-null   int64
 6   floors         15762 non-null   float64
 7   waterfront     15762 non-null   float64
 8   view           15762 non-null   float64
 9   condition      15762 non-null   int64
 10  grade          15762 non-null   int64
 11  sqft_above     15762 non-null   int64
 12  sqft_basement  15762 non-null   object
 13  yr_built       15762 non-null   int64
 14  yr_renovated   15762 non-null   float64
 15  zipcode        15762 non-null   int64
 16  lat            15762 non-null   float64
 17  long           15762 non-null   float64
 18  sqft_living15  15762 non-null   int64
 19  sqft_lot15     15762 non-null   int64
dtypes: float64(8), int64(10), object(2)
memory usage: 2.5+ MB
None
```

We will define a few functions to more efficiently analyze individual features.

In [634]:
```python
 1  def inspect_column(column, unique_count=10):
 2      column_str = str(column)
 3      print('Datatype: ' + str(df[column].dtypes))
 4      print('Total unique itms: ' + str(df[column].nunique()))
 5      print('Displaying first ' + str(unique_count) + ':')
 6      print(df[column].unique()[0:unique_count])
 7      print(f"Minimum value: {df[column].min()}.  Maximum value: {df[column].m
 8      print(df[column].describe())
 9      return column_str
10
11  def regplot(column, df=df):
12      return sns.regplot(data=df, x=column, y='price')
13
14  def hist(column):
15      hist = df[column].hist()
16      return plt.show()
17
18  def displot(column):
19      return sns.displot(data=df, x=column, y='price')
```

## Date

```
In [635]:   1  df['date'] = df['date'].apply(pd.to_datetime)
            2
```

```
In [636]:   1  inspect_column('date')
            2
            3  df['bedrooms'].describe()
```
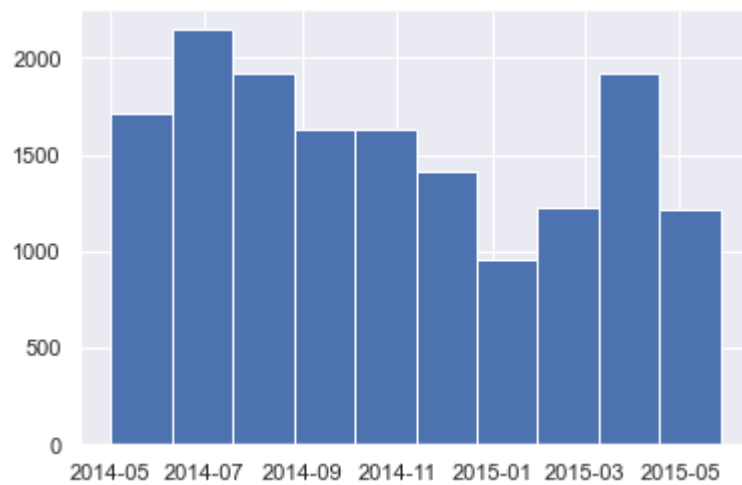
```
Datatype: datetime64[ns]
Total unique itms: 369
Displaying first 10:
['2014-12-09T00:00:00.000000000' '2015-02-18T00:00:00.000000000'
 '2014-05-12T00:00:00.000000000' '2014-06-27T00:00:00.000000000'
 '2015-04-15T00:00:00.000000000' '2015-03-12T00:00:00.000000000'
 '2014-05-27T00:00:00.000000000' '2014-10-07T00:00:00.000000000'
 '2015-01-24T00:00:00.000000000' '2014-07-31T00:00:00.000000000']
Minimum value: 2014-05-02 00:00:00.  Maximum value: 2015-05-27 00:00:00
count                       15762
unique                        369
top        2014-06-25 00:00:00
freq                          103
first      2014-05-02 00:00:00
last       2015-05-27 00:00:00
Name: date, dtype: object

<ipython-input-634-07321a3d05f6>:8: FutureWarning: Treating datetime data as ca
tegorical rather than numeric in `.describe` is deprecated and will be removed
in a future version of pandas. Specify `datetime_is_numeric=True` to silence th
is warning and adopt the future behavior now.
  print(df[column].describe())
```

```
Out[636]:  count    15762.00
           mean         3.38
           std          0.94
           min          1.00
           25%          3.00
           50%          3.00
           75%          4.00
           max         33.00
           Name: bedrooms, dtype: float64
```

In [637]:
```
1 df['date'].hist()
```

Out[637]: <AxesSubplot:>



Taking a look at dates with a histogram, we see that our sales are only from 2014-2015. It could be useful in future analysis to analyze season of sale with more years of sales data.
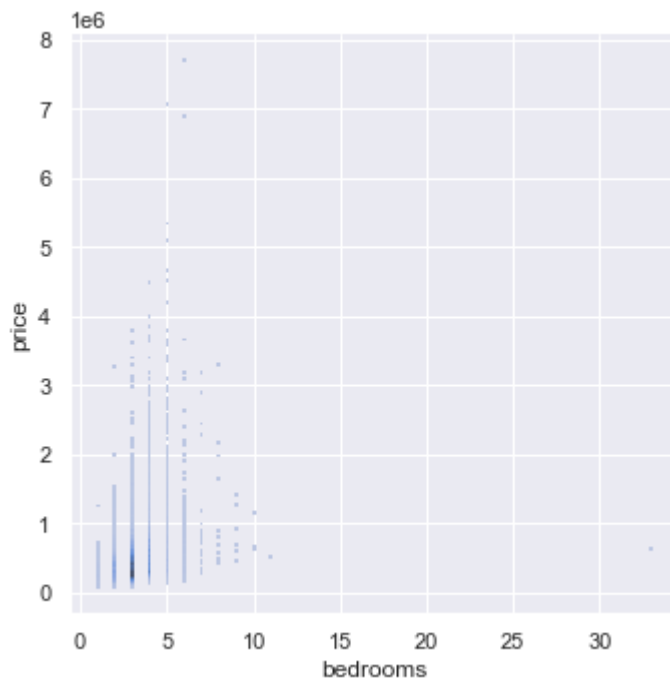
## Bedrooms

In [638]:
```
1 inspect_column('bedrooms', unique_count=20)
```

```
Datatype: int64
Total unique itms: 12
Displaying first 20:
[ 3  4  2  5  1  6  7  8  9 11 10 33]
Minimum value: 1.  Maximum value: 33
count    15762.00
mean         3.38
std          0.94
min          1.00
25%          3.00
50%          3.00
75%          4.00
max         33.00
Name: bedrooms, dtype: float64
```

Out[638]: 'bedrooms'

In [639]:
```
1 displot('bedrooms')
```

Out[639]: <seaborn.axisgrid.FacetGrid at 0x1f7a4c1c460>



In [640]:
```
1 df.loc[df['bedrooms'] == 33]
```

Out[640]:

| id | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view |
|---|---|---|---|---|---|---|---|---|---|
| 2402100895 | 2014-06-25 | 640000.00 | 33 | 1.75 | 1620 | 6000 | 1.00 | 0.00 | 0.00 |

Based on other stats, we assume the one entry with 33 bedrooms to actually be 3 bedrooms.

Correcting below.

In [641]:
```
1  df['bedrooms'] = df['bedrooms'].replace([33],3)
2
3  df.loc[df['bedrooms'] == 33]
```
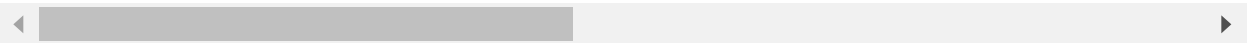
Out[641]:

| | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view | condition | gra |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **id** | | | | | | | | | | | |

In [642]:
```
1  df.loc[df['bedrooms'] == 11]
```

Out[642]:

| | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view |
|---|---|---|---|---|---|---|---|---|---|
| **id** | | | | | | | | | |
| **1773100755** | 2014-08-21 | 520000.00 | 11 | 3.00 | 3000 | 4960 | 2.00 | 0.00 | 0.00 |

The 11 bedroom house also seems unlikely based on square footage. Googling the ID '1773100755' revelas it to be a 4 bedroom house.

In [643]:
```
1  df['bedrooms'] = df['bedrooms'].replace([11],4)
2
3  df.loc[df['bedrooms'] == 11]
4
```
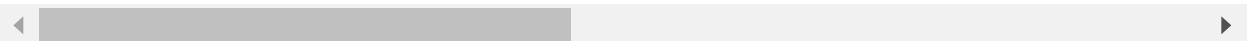
Out[643]:

| | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view | condition | gra |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **id** | | | | | | | | | | | |

In [644]:
```
1  df.loc[df['bedrooms'] == 10]
```

Out[644]:

| | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view |
|---|---|---|---|---|---|---|---|---|---|---|
| **id** | | | | | | | | | |
| **627300145** | 2014-08-14 | 1150000.00 | 10 | 5.25 | 4590 | 10920 | 1.00 | 0.00 | 2.00 |
| **5566100170** | 2014-10-29 | 650000.00 | 10 | 2.00 | 3610 | 11914 | 2.00 | 0.00 | 0.00 |
| **8812401450** | 2014-12-29 | 660000.00 | 10 | 3.00 | 2920 | 3745 | 2.00 | 0.00 | 0.00 |

Even though two of the 10 bedroom houses seem unlikely, a quick google shows that they are recorded as 9 bedroom houses on zillow. We will assume these entries were accurate at the time, and will not change.

In [645]:
```
1 displot('bedrooms')
```

Out[645]: `<seaborn.axisgrid.FacetGrid at 0x1f7a4c1c3d0>`



It looks like there is a large clump around 3 bedrooms. 3-5 bedrooms seems to be where most of the houses are concentrated.

## Bathrooms

In [646]:
```
1 inspect_column('bathrooms', unique_count=29)
```

```
Datatype: float64
Total unique itms: 27
Displaying first 29:
[2.25 3.   2.   4.5  1.   2.5  1.75 2.75 1.5  3.25 4.   3.5  0.75 5.
 4.25 3.75 1.25 5.25 4.75 0.5  5.5  6.   5.75 8.   6.75 7.5  7.75]
Minimum value: 0.5.  Maximum value: 8.0
count    15762.00
mean         2.12
std          0.77
min          0.50
25%          1.75
50%          2.25
75%          2.50
max          8.00
Name: bathrooms, dtype: float64
```

Out[646]: `'bathrooms'`

In [647]:    1  hist('bathrooms')



In [648]:    1  regplot('bathrooms')

Out[648]:  <AxesSubplot:xlabel='bathrooms', ylabel='price'>



It seems like there are some outliers and a few examples of bathrooms more than 5. In the future, it might be worthwhile to eliminate these from analysis.

## Squarefoot - Living

In [649]:
```python
1  inspect_column('sqft_living')
```

```
Datatype: int64
Total unique itms: 912
Displaying first 10:
[2570 1960 1680 5420 1715 1780 1890 1160 1370 1810]
Minimum value: 370.  Maximum value: 13540
count    15762.00
mean      2084.51
std        918.62
min        370.00
25%       1430.00
50%       1920.00
75%       2550.00
max      13540.00
Name: sqft_living, dtype: float64
```

Out[649]: 'sqft_living'

In [650]:
```python
1  regplot('sqft_living')
```

Out[650]: <AxesSubplot:xlabel='sqft_living', ylabel='price'>

In [651]:  `1  hist('sqft_living')`



There seem to be at least one unusual outlier for the price. We will want to take a look at the largest values to verify the quality of the data.

In [652]:  `1  df.sort_values(by=['sqft_living'], ascending=False).head(5)`

Out[652]:

| id | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view |
|---|---|---|---|---|---|---|---|---|---|
| 1225069038 | 2014-05-05 | 2280000.00 | 7 | 8.00 | 13540 | 307752 | 3.00 | 0.00 | 4.00 |
| 6762700020 | 2014-10-13 | 7700000.00 | 6 | 8.00 | 12050 | 27600 | 2.50 | 0.00 | 3.00 |
| 9808700762 | 2014-06-11 | 7060000.00 | 5 | 4.50 | 10040 | 37325 | 2.00 | 1.00 | 2.00 |
| 9208900037 | 2014-09-19 | 6890000.00 | 6 | 7.75 | 9890 | 31374 | 2.00 | 0.00 | 4.00 |
| 1924059029 | 2014-06-17 | 4670000.00 | 5 | 6.75 | 9640 | 13068 | 1.00 | 1.00 | 4.00 |

After reviewing the one outlier, it seems to be a compound in a rural area, and the sqft seems realistic.

## Squarefoot - Lot

In [653]:
```
1  inspect_column('sqft_lot')
2  regplot('sqft_lot')
```

Datatype: int64
Total unique itms: 7927
Displaying first 10:
[  7242    5000    8080  101930   6819    7470    6560    6000    9680    4850]
Minimum value: 520.  Maximum value: 1651359
count      15762.00
mean       15280.82
std        41822.88
min          520.00
25%         5048.50
50%         7602.00
75%        10720.00
max      1651359.00
Name: sqft_lot, dtype: float64

Out[653]:  <AxesSubplot:xlabel='sqft_lot', ylabel='price'>

In [654]:
```
1  hist('sqft_lot')
```



In [655]:
```
1  df.sort_values(by=['sqft_lot'], ascending=False).head(5)
2
3  #outlier looks like a farm, will keep
```

Out[655]:

| id | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view |
|---|---|---|---|---|---|---|---|---|---|
| 1020069017 | 2015-03-27 | 700000.00 | 4 | 1.00 | 1300 | 1651359 | 1.00 | 0.00 | 3.00 |
| 3326079016 | 2015-05-04 | 190000.00 | 2 | 1.00 | 710 | 1164794 | 1.00 | 0.00 | 0.00 |
| 2323089009 | 2015-01-19 | 855000.00 | 4 | 3.50 | 4030 | 1024068 | 2.00 | 0.00 | 0.00 |
| 722069232 | 2014-09-05 | 998000.00 | 4 | 3.25 | 3770 | 982998 | 2.00 | 0.00 | 0.00 |
| 3626079040 | 2014-07-30 | 790000.00 | 2 | 3.00 | 2560 | 982278 | 1.00 | 0.00 | 0.00 |

Given the acerage of some of these lots throughout King County, the results do not seem unreasonable.

## Floors

```
In [656]:    1  inspect_column('floors')
             2
             3  regplot('floors')
```

```
Datatype: float64
Total unique itms: 6
Displaying first 10:
[2.  1.  1.5 3.  2.5 3.5]
Minimum value: 1.0.  Maximum value: 3.5
count    15762.00
mean         1.50
std          0.54
min          1.00
25%          1.00
50%          1.50
75%          2.00
max          3.50
Name: floors, dtype: float64
```

Out[656]:  <AxesSubplot:xlabel='floors', ylabel='price'>



There seems to be some relationship between price and number of floors.

```
In [657]:   1  hist('floors')
```



3.5 floors seems within reason, nothing seems to need correction here.

## Waterfront

```
In [658]:    1  inspect_column('waterfront')
             2
             3  regplot('waterfront')
```

```
Datatype: float64
Total unique itms: 2
Displaying first 10:
[0. 1.]
Minimum value: 0.0.   Maximum value: 1.0
count    15762.00
mean         0.01
std          0.09
min          0.00
25%          0.00
50%          0.00
75%          0.00
max          1.00
Name: waterfront, dtype: float64
```

Out[658]:  <AxesSubplot:xlabel='waterfront', ylabel='price'>



There seems to be a sizeable relationship with price.

In [659]:    1  hist('waterfront')



It seems like there are very few examples of waterfront properties.

# View

In [660]:
```
1  inspect_column('view')
2
3  regplot('view')
```
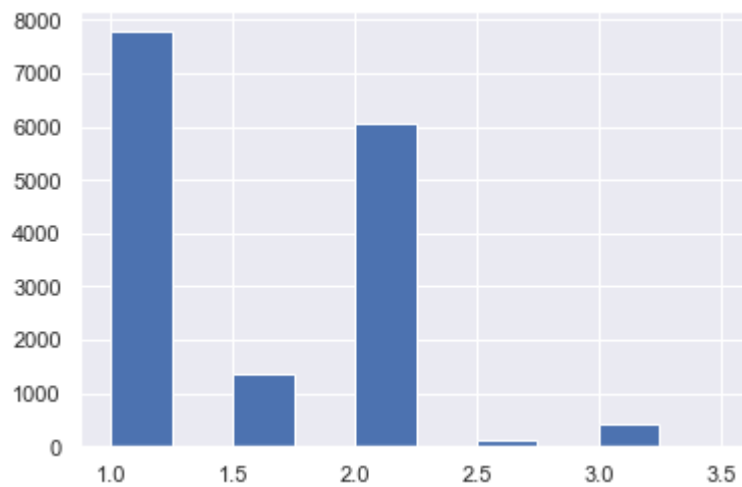
Datatype: float64
Total unique itms: 5
Displaying first 10:
[0. 3. 4. 2. 1.]
Minimum value: 0.0.  Maximum value: 4.0
count    15762.00
mean         0.23
std          0.76
min          0.00
25%          0.00
50%          0.00
75%          0.00
max          4.00
Name: view, dtype: float64

Out[660]: <AxesSubplot:xlabel='view', ylabel='price'>



There seems to be a positive relatinoship with view.


## Condition

```
In [661]:   1  inspect_column('condition')
            2
            3  regplot('condition')
```

```
Datatype: int64
Total unique itms: 5
Displaying first 10:
[3 5 4 1 2]
Minimum value: 1.  Maximum value: 5
count    15762.00
mean         3.41
std          0.65
min          1.00
25%          3.00
50%          3.00
75%          4.00
max          5.00
Name: condition, dtype: float64
```

Out[661]: <AxesSubplot:xlabel='condition', ylabel='price'>



There doesn't seem to be a very strong relatinship with price.

In [662]:    1  hist('condition')

It seems odd that there are very few examples of 1 and 2.

## Grade

In [663]:
```
1  inspect_column('grade')
2
3  regplot('grade')
```

```
Datatype: int64
Total unique itms: 11
Displaying first 10:
[ 7  8 11  9  6  5 10 12  4  3]
Minimum value: 3.  Maximum value: 13
count    15762.00
mean         7.66
std          1.17
min          3.00
25%          7.00
50%          7.00
75%          8.00
max         13.00
Name: grade, dtype: float64
```

Out[663]:  <AxesSubplot:xlabel='grade', ylabel='price'>



There is a strong relationship with grade.

In [664]:　1　hist('grade')



There is a clump in the middle, few examples of 13. Also, very few examples of 1-4.

In [665]:
```
1  df[df['grade'] == 13]
```

Out[665]:

| id | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view |
|---|---|---|---|---|---|---|---|---|---|
| 9831200500 | 2015-03-04 | 2480000.00 | 5 | 3.75 | 6810 | 7500 | 2.50 | 0.00 | 0.00 |
| 7237501190 | 2014-10-10 | 1780000.00 | 4 | 3.25 | 4890 | 13402 | 2.00 | 0.00 | 0.00 |
| 1725059316 | 2014-11-20 | 2390000.00 | 4 | 4.00 | 6330 | 13296 | 2.00 | 0.00 | 2.00 |
| 853200010 | 2014-07-01 | 3800000.00 | 5 | 5.50 | 7050 | 42840 | 1.00 | 0.00 | 2.00 |
| 6762700020 | 2014-10-13 | 7700000.00 | 6 | 8.00 | 12050 | 27600 | 2.50 | 0.00 | 3.00 |
| 1068000375 | 2014-09-23 | 3200000.00 | 6 | 5.00 | 7100 | 18200 | 2.50 | 0.00 | 0.00 |
| 9208900037 | 2014-09-19 | 6890000.00 | 6 | 7.75 | 9890 | 31374 | 2.00 | 0.00 | 4.00 |
| 3303850390 | 2014-12-12 | 2980000.00 | 5 | 5.50 | 7400 | 18898 | 2.00 | 0.00 | 3.00 |
| 2426039123 | 2015-01-30 | 2420000.00 | 5 | 4.75 | 7880 | 24250 | 2.00 | 0.00 | 2.00 |
| 4139900180 | 2015-04-20 | 2340000.00 | 4 | 2.50 | 4500 | 35200 | 1.00 | 0.00 | 0.00 |
| 2303900100 | 2014-09-11 | 3800000.00 | 3 | 4.25 | 5510 | 35000 | 2.00 | 0.00 | 4.00 |

The grade 13 has fairly high prices, which aligns with expectations.

## Squarefoot Above

In [666]:
```python
1  inspect_column('sqft_above')
2
3  regplot('sqft_above')
```

```
Datatype: int64
Total unique itms: 835
Displaying first 10:
[2170 1050 1680 3890 1715 1890  860 1370 1810 1980]
Minimum value: 370.  Maximum value: 9410
count    15762.00
mean      1792.78
std        828.40
min        370.00
25%       1200.00
50%       1570.00
75%       2220.00
max       9410.00
Name: sqft_above, dtype: float64
```

Out[666]: <AxesSubplot:xlabel='sqft_above', ylabel='price'>

In [667]:

```
1 hist('sqft_above')
```



It looks like there is a strong relationship between price and sqft above.

## Squarefoot Basement

In [668]:

```
1 inspect_column('sqft_basement')
2
```

```
Datatype: object
Total unique itms: 283
Displaying first 10:
['400.0' '910.0' '0.0' '1530.0' '?' '730.0' '300.0' '970.0' '760.0'
 '720.0']
Minimum value: 0.0.  Maximum value: ?
count     15762
unique      283
top         0.0
freq       9362
Name: sqft_basement, dtype: object
```

Out[668]: 'sqft_basement'

It seems there are some errors with question marks. Let's take a look.

In [669]:
```
1 df[df['sqft_basement'] == '?']
```

Out[669]:

| id | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront |
|---|---|---|---|---|---|---|---|---|
| 1321400060 | 2014-06-27 | 257500.00 | 3 | 2.25 | 1715 | 6819 | 2.00 | 0.00 |
| 16000397 | 2014-12-05 | 189000.00 | 2 | 1.00 | 1200 | 9850 | 1.00 | 0.00 |
| 7203220400 | 2014-07-07 | 861990.00 | 5 | 2.75 | 3595 | 5639 | 2.00 | 0.00 |
| 1531000030 | 2015-03-23 | 720000.00 | 4 | 2.50 | 3450 | 39683 | 2.00 | 0.00 |
| 2525310310 | 2014-09-16 | 272500.00 | 3 | 1.75 | 1540 | 12600 | 1.00 | 0.00 |
| 1909600046 | 2014-07-03 | 445838.00 | 3 | 2.50 | 2250 | 5692 | 2.00 | 0.00 |

It might be best to go ahead and make a "True" and "False" boolean column for 'has_basement.' We will also change all '?' values to zero (0) and conver the values into floats.

In [670]:
```
1 df['sqft_basement'] = df['sqft_living'] - df['sqft_above']
2
3 df.loc[df['sqft_basement'] == '?']
```

Out[670]:

| id | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view | condition | gra |
|---|---|---|---|---|---|---|---|---|---|---|---|

In [671]:
```
1 df['sqft_basement'] = df['sqft_basement'].astype(float)
```

In [672]:
```python
1  has_basement = np.where(df['sqft_basement'] > 0, 1, 0)
2
3  df.insert (12, 'has_basement', has_basement)
4
5  df.head(5)
```

Out[672]:

| id | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view |
|---|---|---|---|---|---|---|---|---|---|
| 6414100192 | 2014-12-09 | 538000.00 | 3 | 2.25 | 2570 | 7242 | 2.00 | 0.00 | 0.00 |
| 2487200875 | 2014-12-09 | 604000.00 | 4 | 3.00 | 1960 | 5000 | 1.00 | 0.00 | 0.00 |
| 1954400510 | 2015-02-18 | 510000.00 | 3 | 2.00 | 1680 | 8080 | 1.00 | 0.00 | 0.00 |
| 7237550310 | 2014-05-12 | 1230000.00 | 4 | 4.50 | 5420 | 101930 | 1.00 | 0.00 | 0.00 |
| 1321400060 | 2014-06-27 | 257500.00 | 3 | 2.25 | 1715 | 6819 | 2.00 | 0.00 | 0.00 |

5 rows × 21 columns

In [673]:
```python
1  regplot(has_basement)
```

Out[673]: <AxesSubplot:ylabel='price'>



There isn't the strongest relationship with price, but there is a somewhat noticeable slope.

## Year Built

In [674]:
```
1  inspect_column('yr_built')
2
3  regplot('yr_built')
```

```
Datatype: int64
Total unique itms: 116
Displaying first 10:
[1951 1965 1987 2001 1995 1960 2003 1942 1977 1900]
Minimum value: 1900.  Maximum value: 2015
count    15762.00
mean      1971.11
std         29.34
min       1900.00
25%       1952.00
50%       1975.00
75%       1997.00
max       2015.00
Name: yr_built, dtype: float64
```

Out[674]:  <AxesSubplot:xlabel='yr_built', ylabel='price'>



In [675]:
```
1  hist('yr_built')
```

We will convert this to age to more easily interpret this feature in our model.

In [676]:
```
1  df['age'] = abs(df['yr_built'] - 2015)
2
3  df.head()
```

Out[676]:

| id | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view |
|---|---|---|---|---|---|---|---|---|---|
| 6414100192 | 2014-12-09 | 538000.00 | 3 | 2.25 | 2570 | 7242 | 2.00 | 0.00 | 0.00 |
| 2487200875 | 2014-12-09 | 604000.00 | 4 | 3.00 | 1960 | 5000 | 1.00 | 0.00 | 0.00 |
| 1954400510 | 2015-02-18 | 510000.00 | 3 | 2.00 | 1680 | 8080 | 1.00 | 0.00 | 0.00 |
| 7237550310 | 2014-05-12 | 1230000.00 | 4 | 4.50 | 5420 | 101930 | 1.00 | 0.00 | 0.00 |
| 1321400060 | 2014-06-27 | 257500.00 | 3 | 2.25 | 1715 | 6819 | 2.00 | 0.00 | 0.00 |

5 rows × 22 columns

```
In [677]:    1  inspect_column('age')
             2
             3  hist('age')
```

```
Datatype: int64
Total unique itms: 116
Displaying first 10:
[ 64  50  28  14  20  55  12  73  38 115]
Minimum value: 0.   Maximum value: 115
count    15762.00
mean        43.89
std         29.34
min          0.00
25%         18.00
50%         40.00
75%         63.00
max        115.00
Name: age, dtype: float64
```



```
In [678]:    1  del df['yr_built']
```

Ultimately, this will have the same effect in our model, but age is a little easier to interpret.

## Year Renovated

In [679]:
```
1  inspect_column('yr_renovated')
2
3  regplot('yr_renovated')
```

```
Datatype: float64
Total unique itms: 70
Displaying first 10:
[1991.    0. 2002. 2010. 1992. 2013. 1994. 1978. 2005. 2003.]
Minimum value: 0.0.  Maximum value: 2015.0
count   15762.00
mean       82.44
std       397.21
min         0.00
25%         0.00
50%         0.00
75%         0.00
max      2015.00
Name: yr_renovated, dtype: float64
```

Out[679]:  <AxesSubplot:xlabel='yr_renovated', ylabel='price'>



For the purposes of this model, we will simplify yr_renovated to a binary column denoting whether or not a house has been renovated. It might be useful in the future to analyze yr_renovated as it is presented.

```
In [680]:   1  renovated = np.where(df['yr_renovated'] > 0, 1, 0)
            2
            3  df['renovated'] = renovated
            4
            5  del df['yr_renovated']
            6
            7  df.head(5)
```

Out[680]:

| id | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view |
|---|---|---|---|---|---|---|---|---|---|
| 6414100192 | 2014-12-09 | 538000.00 | 3 | 2.25 | 2570 | 7242 | 2.00 | 0.00 | 0.00 |
| 2487200875 | 2014-12-09 | 604000.00 | 4 | 3.00 | 1960 | 5000 | 1.00 | 0.00 | 0.00 |
| 1954400510 | 2015-02-18 | 510000.00 | 3 | 2.00 | 1680 | 8080 | 1.00 | 0.00 | 0.00 |
| 7237550310 | 2014-05-12 | 1230000.00 | 4 | 4.50 | 5420 | 101930 | 1.00 | 0.00 | 0.00 |
| 1321400060 | 2014-06-27 | 257500.00 | 3 | 2.25 | 1715 | 6819 | 2.00 | 0.00 | 0.00 |

5 rows × 21 columns

```
In [681]:   1  df.head()
```

Out[681]:

| id | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view |
|---|---|---|---|---|---|---|---|---|---|
| 6414100192 | 2014-12-09 | 538000.00 | 3 | 2.25 | 2570 | 7242 | 2.00 | 0.00 | 0.00 |
| 2487200875 | 2014-12-09 | 604000.00 | 4 | 3.00 | 1960 | 5000 | 1.00 | 0.00 | 0.00 |
| 1954400510 | 2015-02-18 | 510000.00 | 3 | 2.00 | 1680 | 8080 | 1.00 | 0.00 | 0.00 |
| 7237550310 | 2014-05-12 | 1230000.00 | 4 | 4.50 | 5420 | 101930 | 1.00 | 0.00 | 0.00 |
| 1321400060 | 2014-06-27 | 257500.00 | 3 | 2.25 | 1715 | 6819 | 2.00 | 0.00 | 0.00 |

5 rows × 21 columns

In [682]:
```
1 regplot('renovated')
```

Out[682]: `<AxesSubplot:xlabel='renovated', ylabel='price'>`



There seems to be a relationship between renovated and price.

## Zipcode

In [683]:
```
1 inspect_column('zipcode')
```

```
Datatype: int64
Total unique itms: 70
Displaying first 10:
[98125 98136 98074 98053 98003 98146 98038 98115 98107 98126]
Minimum value: 98001.   Maximum value: 98199
count    15762.00
mean     98077.56
std         53.41
min      98001.00
25%      98033.00
50%      98065.00
75%      98117.00
max      98199.00
Name: zipcode, dtype: float64
```

Out[683]: `'zipcode'`

Zipcode should be integer for now, since there will be no decimals. It might be worth considering conversion to string as well further in the project.

```
In [684]:   1  df['zipcode'] = df['zipcode'].astype(int)
```

```
In [685]:   1  inspect_column('zipcode')
            2
            3  regplot('zipcode')
```

```
Datatype: int32
Total unique itms: 70
Displaying first 10:
[98125 98136 98074 98053 98003 98146 98038 98115 98107 98126]
Minimum value: 98001.  Maximum value: 98199
count    15762.00
mean     98077.56
std         53.41
min      98001.00
25%      98033.00
50%      98065.00
75%      98117.00
max      98199.00
Name: zipcode, dtype: float64
```

Out[685]:  <AxesSubplot:xlabel='zipcode', ylabel='price'>

In [686]:
```
1  hist('zipcode')
```



The regplot and histogram are not particularly useful here since each zip code is actually an independent variable.

## 'sqft_living15'

The square footage of interior housing living space for the nearest 15 neighbors

In [687]:
```
1  inspect_column('sqft_living15')
2
3  regplot('sqft_living15')
```

Datatype: int64
Total unique itms: 694
Displaying first 10:
[1690 1360 1800 4760 2238 1780 2390 1330 1370 2140]
Minimum value: 399.  Maximum value: 6210
count    15762.00
mean      1990.22
std        684.14
min        399.00
25%       1490.00
50%       1846.00
75%       2370.00
max       6210.00
Name: sqft_living15, dtype: float64

Out[687]: <AxesSubplot:xlabel='sqft_living15', ylabel='price'>

In [688]:

```
1 hist('sqft_living15')
```



There is a relationship between price and sqft_living15, but this factor detracts from the uniqueness of the home itself (in our opinion). We could consider reviewing this in future analyses.

## 'sqft_lot15'

The square footage of the land lots of the nearest 15 neighbors

In [689]:
```
1  inspect_column('sqft_lot15')
2
3  regplot('sqft_lot15')
```

```
Datatype: int64
Total unique itms: 7126
Displaying first 10:
[  7639    5000    7503 101930    6819    8113    7570    6000   10208    4850]
Minimum value: 659.   Maximum value: 871200
count     15762.00
mean      12900.42
std       27977.23
min         659.00
25%        5100.00
50%        7620.00
75%       10107.50
max      871200.00
Name: sqft_lot15, dtype: float64
```

Out[689]:  <AxesSubplot:xlabel='sqft_lot15', ylabel='price'>



In [690]:
```
1  hist('sqft_lot15')
```

Similar to sqft_living15, we think this would be better saved for future analysis.

# Removing Outliers

In [691]:
```python
1  fig, ax = plt.subplots(figsize=(12, 4))
2  ax = sns.boxplot(df['price'])
```

```
C:\Users\johnn\anaconda3\envs\learn-env\lib\site-packages\seaborn\_decorators.p
y:36: FutureWarning: Pass the following variable as a keyword arg: x. From vers
ion 0.12, the only valid positional argument will be `data`, and passing other
arguments without an explicit keyword will result in an error or misinterpretat
ion.
  warnings.warn(
```



There seem to be quite a few outliers. We should consider removing some with the IQR method.

## Target Variable (price)

In [692]:
```python
1  def iqr_df(column):
2
3      column = column
4
5      describe = df.describe()[column]
6
7      q1 = describe['25%']
8      q3 = describe['75%']
9
10     iqr = q3 - q1
11
12     outlier_index = (df[column] > (q3 + 1.5 * iqr)) | (df[column] < (q1 - 1.
13
14     return df[~outlier_index]
```

In [693]:
```
1  iqr_price = iqr_df('price')
2
3  print(df.shape)
4  print(iqr_price.shape)
```

```
(15762, 21)
(14931, 21)
```

In [694]:
```
1  regplot('sqft_living', df=df)
```

Out[694]: `<AxesSubplot:xlabel='sqft_living', ylabel='price'>`



In [695]:
```
1  regplot('sqft_living', df=iqr_price)
```

Out[695]: `<AxesSubplot:xlabel='sqft_living', ylabel='price'>`



This shaves quite a few examples from our dataset, but it will be helpful in normalizing our dataset.

In [696]:
```
1  df = iqr_price
```

## Squarefoot Living

```
In [697]:    1  iqr_sqft_living = iqr_df('sqft_living')
             2
             3  print(df.shape)
             4  print(iqr_sqft_living.shape)
```

(14931, 21)
(14715, 21)

```
In [698]:    1  regplot('sqft_living', df=df)
```

Out[698]:  <AxesSubplot:xlabel='sqft_living', ylabel='price'>



```
In [699]:    1  regplot('sqft_living', df=iqr_sqft_living)
```

Out[699]:  <AxesSubplot:xlabel='sqft_living', ylabel='price'>



Only removes ~200 examples and gives us a more normal dataset.

```
In [700]:   1  df = iqr_sqft_living
```

## Squarefoot Lot

```
In [701]:   1  iqr_sqft_lot = iqr_df('sqft_lot')
            2
            3  print(df.shape)
            4  print(iqr_sqft_lot.shape)
```

```
(14715, 21)
(13128, 21)
```

```
In [702]:   1  regplot('sqft_lot', df=df)
```

Out[702]: <AxesSubplot:xlabel='sqft_lot', ylabel='price'>

```
In [703]:    1  regplot('sqft_lot', df=iqr_sqft_lot)
```

Out[703]:  <AxesSubplot:xlabel='sqft_lot', ylabel='price'>



This will remove a lot of the farms and might have a small impact on how our model interprets sqft_lot, but it will make our model more useful for average home owners.

```
In [704]:    1  df = iqr_sqft_lot
```

# EXPLORE

Now that we're comfortable that we have quality data, it's time to determine which columns we'll want to analyze for our primary analysis.

First we'll review which columns we have to work with:

## Feature Selection

Based on prior anaysis and scrubbing, we'll categorize our columns into three sections:

Continuous variables:

- price
- sqft_living
- sqft_lot
- sqft_above
- sqft_basement
- yr_built
- sqft_living15
- sqft_lot15

Categorical variables - while some of these may appear continuous, their values represent integers and fractions that are more categorical even if they are for a specific count.

- bedrooms
- bathrooms
- floors
- condition
- grade
- waterfront
- renovated
- zipcode
- has_basement

Remove from model:

- date - date of sale could be interesting to analyze if we had a longer time horizon. Home prices could sell for more less based on season, and this could be interesting for further analysis
- lat - will not have a linear relationship
- long - will not have a linear relationship
- view - while we have a range of values, the column description reads "Has been viewed" which should be binary. Seems like there could be an error, further review could make this column eligible for future analysis

```
In [705]:    1  df.drop(['date', 'lat', 'long', 'view'],axis=1,inplace=True)
```

```
C:\Users\johnn\anaconda3\envs\learn-env\lib\site-packages\pandas\core\frame.py:
4163: 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/sta
ble/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pyd
ata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-c
opy)
  return super().drop(
```

```
In [706]:    1  df.head()
```

Out[706]:

| id | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | condition | g |
|---|---|---|---|---|---|---|---|---|---|
| 6414100192 | 538000.00 | 3 | 2.25 | 2570 | 7242 | 2.00 | 0.00 | 3 | |
| 2487200875 | 604000.00 | 4 | 3.00 | 1960 | 5000 | 1.00 | 0.00 | 5 | |
| 1954400510 | 510000.00 | 3 | 2.00 | 1680 | 8080 | 1.00 | 0.00 | 3 | |
| 1321400060 | 257500.00 | 3 | 2.25 | 1715 | 6819 | 2.00 | 0.00 | 3 | |
| 2414600126 | 229500.00 | 3 | 1.00 | 1780 | 7470 | 1.00 | 0.00 | 3 | |

# Multicollinearity

We will create a heat map to identify multicollinearity.

```
In [707]:
1  def heatmap(df_name, figsize=(15,10), cmap='Reds'):
2      corr = df_name.drop('price',axis=1).corr()
3      mask = np.zeros_like(corr)
4      mask[np.triu_indices_from(mask)] = True
5      fig, ax = plt.subplots(figsize=figsize)
6      sns.heatmap(corr, annot=True, cmap=cmap, mask=mask)
7      return fig, ax
8
9  heatmap(df)
```

Out[707]: (<Figure size 1080x720 with 2 Axes>, <AxesSubplot:>)



We will drop the following:

sqft_above + sqft_basement - these are duplicative of sqft_living.

sqft_lot15 and sqft_living15 - these could be more interesting for broader analysis of areas. Since there is high multicolinearity, we can save these for when we look at zip, lat, and long.

```
In [708]:   1  del df['sqft_above']
            2  del df['sqft_basement']
            3  del df['sqft_lot15']
            4  del df['sqft_living15']
```

```
In [709]:   1  heatmap(df)
```

Out[709]:  (<Figure size 1080x720 with 2 Axes>, <AxesSubplot:>)

Sqft_living, bathrooms, and grade appear to have potential for multicollinearity. This issue should be remedied by encoding grade and bathrooms as categorical variables, which we will do next in the modeling stage.

# MODEL

## Data Modeling

Describe and justify the process for analyzing or modeling the data.

---

Questions to consider:

- How did you analyze or model the data?
- How did you iterate on your initial approach to make it better?
- Why are these choices appropriate given the data and the business problem?

In [710]:
```
1 df.head()
```

Out[710]:

| id | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | condition | g |
|---|---|---|---|---|---|---|---|---|---|
| 6414100192 | 538000.00 | 3 | 2.25 | 2570 | 7242 | 2.00 | 0.00 | 3 | |
| 2487200875 | 604000.00 | 4 | 3.00 | 1960 | 5000 | 1.00 | 0.00 | 5 | |
| 1954400510 | 510000.00 | 3 | 2.00 | 1680 | 8080 | 1.00 | 0.00 | 3 | |
| 1321400060 | 257500.00 | 3 | 2.25 | 1715 | 6819 | 2.00 | 0.00 | 3 | |
| 2414600126 | 229500.00 | 3 | 1.00 | 1780 | 7470 | 1.00 | 0.00 | 3 | |

Installing stats and modeling packages:

Note: uncomment first line to install -U fsds.

In [711]:
```
1  # !pip install -U fsds
2  from scipy import stats
3  from fsds.imports import *
4
5  import statsmodels.api as sm
6  import statsmodels.stats.api as sms
7  import statsmodels.formula.api as smf
8  import scipy.stats as stats
9
10 import scipy.stats as stats
11 import statsmodels.api as sms
```

# Initial Model

We'll go ahead and define our categorical variables so that we can implement the code into our model function:

In [712]:
```
1  categoricals = ['bedrooms',
2                  'bathrooms',
3                  'floors',
4                  'waterfront',
5                  'condition',
6                  'grade',
7                  'has_basement',
8                  'zipcode',
9                  'renovated']
10
```

Function to draw a QQ plot and a homoscedasticity check.

```
In [713]:    1  def check_model(model):
             2
             3      resids = model.resid
             4
             5      fig,ax = plt.subplots(ncols=2,figsize=(12,5))
             6      sms.qqplot(resids, stats.distributions.norm, fit=True, line='45',ax=ax[0
             7      xs = np.linspace(0,1,len(resids))
             8
             9      y_hat = model.predict(df)
            10      y = df['price']
            11      resid = y - y_hat
            12      plot = plt.scatter(x=y_hat, y=resid)
            13      plt.axhline(0)
            14
            15      ax[1].scatter(x=y_hat,y=resid)
            16
            17      return fig,ax
            18
            19  # check_model(model1)
```

Function to run the model and output summary statistics and graphs.

```
In [714]:    1  def make_model(df_name, categoricals=categoricals):
             2
             3      features = ' + '.join(df.drop('price',axis=1).columns)
             4      for variable in categoricals:
             5          features = features.replace(variable, ("C(" + variable + ")"))
             6
             7      f  = "price~"+features
             8
             9      model = smf.ols(f, df_name).fit()
            10      display(model.summary())
            11
            12      fig,ax = check_model(model)
            13      plt.show()
            14
            15      return model
            16
            17  model1 = make_model(df)
```

OLS Regression Results

| | | | |
|---|---|---|---|
| **Dep. Variable:** | price | **R-squared:** | 0.831 |
| **Model:** | OLS | **Adj. R-squared:** | 0.829 |
| **Method:** | Least Squares | **F-statistic:** | 528.1 |
| **Date:** | Sat, 01 May 2021 | **Prob (F-statistic):** | 0.00 |
| **Time:** | 11:42:03 | **Log-Likelihood:** | -1.6728e+05 |
| **No. Observations:** | 13128 | **AIC:** | 3.348e+05 |
| **Df Residuals:** | 13006 | **BIC:** | 3.357e+05 |
| **Df Model:** | 121 | | |
| **Covariance Type:** | nonrobust | | |

| | coef | std err | t | P>\|t\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| **Intercept** | -9.876e+04 | 1.01e+05 | -0.982 | 0.326 | -2.96e+05 | 9.83e+04 |

Our first model has a fairly strong R-squared at 0.831. The QQ plot indicates that there might be some outliers that we could remove to further refine our model. The homoscedasticity graph also shows some outliers, but the graph has a noticeable cone shape indicating we are mostly on track with our current refinement of the overall dataset.

## Reviewing P-values

Next, we'll want to look at the features that have a P-value greater than 0.05. Removing these features will help us isolate the most statistically significant variables of our model.

In [715]:
```python
1  model1.pvalues
2
3  pvals = model1.pvalues
4
5  pvals[pvals > 0.05]
6  # pvals[pvals > 0.05].index
```

Out[715]:
```
Intercept                 0.33
C(bedrooms)[T.5]          0.69
C(bedrooms)[T.6]          0.35
C(bedrooms)[T.8]          0.15
C(bedrooms)[T.9]          0.47
C(bedrooms)[T.10]         0.15
C(bathrooms)[T.0.75]      0.14
C(bathrooms)[T.1.0]       0.26
C(bathrooms)[T.1.25]      0.90
C(bathrooms)[T.1.5]       0.26
C(bathrooms)[T.1.75]      0.20
C(bathrooms)[T.2.0]       0.25
C(bathrooms)[T.2.25]      0.15
C(bathrooms)[T.2.5]       0.12
C(bathrooms)[T.2.75]      0.08
C(bathrooms)[T.3.0]       0.11
C(bathrooms)[T.3.25]      0.06
C(bathrooms)[T.4.5]       0.38
C(bathrooms)[T.4.75]      0.63
C(bathrooms)[T.5.0]       0.39
C(bathrooms)[T.5.25]      0.34
C(bathrooms)[T.5.75]      0.94
C(floors)[T.2.0]          0.07
C(floors)[T.3.5]          0.10
C(grade)[T.4]             0.22
C(grade)[T.5]             0.11
C(grade)[T.6]             0.13
C(grade)[T.7]             0.24
C(grade)[T.8]             0.55
C(grade)[T.9]             0.70
C(grade)[T.10]            0.40
C(zipcode)[T.98002]       0.25
C(zipcode)[T.98003]       0.39
C(zipcode)[T.98022]       0.61
C(zipcode)[T.98030]       0.55
C(zipcode)[T.98031]       0.05
C(zipcode)[T.98032]       0.93
C(zipcode)[T.98042]       0.10
dtype: float64
```

It seems that certain bedroom numbers don't have a significant effect. Bathrooms have very little effect. 1.5 and 3.5 floors might not have an effect, likely due to low representaion in dataset. Some conditions seems important, grade seems negligible, and a 12 of the 69 zip codes are not significant.

We will try running our model again and convert the following to numerical values:

- bedrooms

- bathrooms
- grade

# Second model

We will convert bedrooms, bathrooms, and grade to numerical values.

```
In [716]:   1  categoricals = [
            2  #               'bedrooms',
            3  #               'bathrooms',
            4                  'floors',
            5                  'waterfront',
            6                  'condition',
            7  #               'grade',
            8                  'has_basement',
            9                  'zipcode',
           10                  'renovated'
           11                  ]
```

In [717]:
```
1  model2 = make_model(df, categoricals)
```

OLS Regression Results

| | | | |
|---|---|---|---|
| Dep. Variable: | price | R-squared: | 0.825 |
| Model: | OLS | Adj. R-squared: | 0.824 |
| Method: | Least Squares | F-statistic: | 705.5 |
| Date: | Sat, 01 May 2021 | Prob (F-statistic): | 0.00 |
| Time: | 11:42:04 | Log-Likelihood: | -1.6752e+05 |
| No. Observations: | 13128 | AIC: | 3.352e+05 |
| Df Residuals: | 13040 | BIC: | 3.359e+05 |
| Df Model: | 87 | | |
| Covariance Type: | nonrobust | | |

| | coef | std err | t | P>\|t\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| Intercept | -4.66e+05 | 2.58e+04 | -18.065 | 0.000 | -5.17e+05 | -4.15e+05 |
| C(floors)[T.1.5] | 6715.7196 | 3061.388 | 2.194 | 0.028 | 714.953 | 1.27e+04 |
| C(floors)[T.2.0] | -1360.5062 | 2604.330 | -0.522 | 0.601 | -6465.373 | 3744.360 |
| C(floors)[T.2.5] | -1.692e+04 | 1.08e+04 | -1.570 | 0.116 | -3.81e+04 | 4204.357 |
| C(floors)[T.3.0] | -4.032e+04 | 5601.636 | -7.199 | 0.000 | -5.13e+04 | -2.93e+04 |
| C(floors)[T.3.5] | -6.423e+04 | 3.81e+04 | -1.686 | 0.092 | -1.39e+05 | 1.04e+04 |
| C(waterfront)[T.1.0] | 3.258e+05 | 1.88e+04 | 17.305 | 0.000 | 2.89e+05 | 3.63e+05 |
| C(condition)[T.2] | 5.806e+04 | 2.51e+04 | 2.311 | 0.021 | 8820.975 | 1.07e+05 |
| C(condition)[T.3] | 8.728e+04 | 2.37e+04 | 3.688 | 0.000 | 4.09e+04 | 1.34e+05 |
| C(condition)[T.4] | 1.056e+05 | 2.37e+04 | 4.462 | 0.000 | 5.92e+04 | 1.52e+05 |
| C(condition)[T.5] | 1.327e+05 | 2.38e+04 | 5.587 | 0.000 | 8.62e+04 | 1.79e+05 |
| C(has_basement)[T.1] | -1.854e+04 | 1963.480 | -9.445 | 0.000 | -2.24e+04 | -1.47e+04 |
| C(zipcode)[T.98002] | 1.259e+04 | 9076.948 | 1.387 | 0.165 | -5198.997 | 3.04e+04 |
| C(zipcode)[T.98003] | 2474.7576 | 8286.881 | 0.299 | 0.765 | -1.38e+04 | 1.87e+04 |
| C(zipcode)[T.98004] | 5.212e+05 | 1.04e+04 | 50.232 | 0.000 | 5.01e+05 | 5.42e+05 |
| C(zipcode)[T.98005] | 3.361e+05 | 1.09e+04 | 30.887 | 0.000 | 3.15e+05 | 3.57e+05 |
| C(zipcode)[T.98006] | 2.808e+05 | 7919.132 | 35.456 | 0.000 | 2.65e+05 | 2.96e+05 |
| C(zipcode)[T.98007] | 2.486e+05 | 1.05e+04 | 23.653 | 0.000 | 2.28e+05 | 2.69e+05 |
| C(zipcode)[T.98008] | 2.446e+05 | 8494.599 | 28.796 | 0.000 | 2.28e+05 | 2.61e+05 |
| C(zipcode)[T.98010] | 1.075e+05 | 1.68e+04 | 6.408 | 0.000 | 7.46e+04 | 1.4e+05 |
| C(zipcode)[T.98011] | 1.442e+05 | 9554.627 | 15.097 | 0.000 | 1.26e+05 | 1.63e+05 |
| C(zipcode)[T.98014] | 1.003e+05 | 1.61e+04 | 6.240 | 0.000 | 6.88e+04 | 1.32e+05 |
| C(zipcode)[T.98019] | 8.918e+04 | 1.07e+04 | 8.345 | 0.000 | 6.82e+04 | 1.1e+05 |
| C(zipcode)[T.98022] | 8315.2112 | 1e+04 | 0.829 | 0.407 | -1.14e+04 | 2.8e+04 |

| | | | | | |
|---|---|---|---|---|---|
| C(zipcode)[T.98023] | -1.689e+04 | 7388.109 | -2.286 | 0.022 | -3.14e+04 | -2410.859 |
| C(zipcode)[T.98024] | 1.324e+05 | 1.98e+04 | 6.690 | 0.000 | 9.36e+04 | 1.71e+05 |
| C(zipcode)[T.98027] | 2.369e+05 | 8682.293 | 27.281 | 0.000 | 2.2e+05 | 2.54e+05 |
| C(zipcode)[T.98028] | 1.352e+05 | 8480.907 | 15.945 | 0.000 | 1.19e+05 | 1.52e+05 |
| C(zipcode)[T.98029] | 2.287e+05 | 8167.594 | 27.999 | 0.000 | 2.13e+05 | 2.45e+05 |
| C(zipcode)[T.98030] | 3966.8977 | 8613.197 | 0.461 | 0.645 | -1.29e+04 | 2.09e+04 |
| C(zipcode)[T.98031] | 1.401e+04 | 8477.743 | 1.653 | 0.098 | -2606.716 | 3.06e+04 |
| C(zipcode)[T.98032] | 725.3311 | 1.06e+04 | 0.069 | 0.945 | -2e+04 | 2.14e+04 |
| C(zipcode)[T.98033] | 3.231e+05 | 7853.812 | 41.141 | 0.000 | 3.08e+05 | 3.39e+05 |
| C(zipcode)[T.98034] | 1.959e+05 | 7279.182 | 26.913 | 0.000 | 1.82e+05 | 2.1e+05 |
| C(zipcode)[T.98038] | 4.285e+04 | 7317.478 | 5.856 | 0.000 | 2.85e+04 | 5.72e+04 |
| C(zipcode)[T.98039] | 5.87e+05 | 4.92e+04 | 11.940 | 0.000 | 4.91e+05 | 6.83e+05 |
| C(zipcode)[T.98040] | 4.281e+05 | 1.01e+04 | 42.399 | 0.000 | 4.08e+05 | 4.48e+05 |
| C(zipcode)[T.98042] | 1.159e+04 | 7444.168 | 1.557 | 0.120 | -3002.747 | 2.62e+04 |
| C(zipcode)[T.98045] | 1.045e+05 | 1.01e+04 | 10.295 | 0.000 | 8.46e+04 | 1.24e+05 |
| C(zipcode)[T.98052] | 2.588e+05 | 7316.720 | 35.375 | 0.000 | 2.44e+05 | 2.73e+05 |
| C(zipcode)[T.98053] | 2.608e+05 | 8737.921 | 29.849 | 0.000 | 2.44e+05 | 2.78e+05 |
| C(zipcode)[T.98055] | 4.52e+04 | 8471.855 | 5.335 | 0.000 | 2.86e+04 | 6.18e+04 |
| C(zipcode)[T.98056] | 1.032e+05 | 7714.064 | 13.374 | 0.000 | 8.8e+04 | 1.18e+05 |
| C(zipcode)[T.98058] | 3.321e+04 | 7697.188 | 4.314 | 0.000 | 1.81e+04 | 4.83e+04 |
| C(zipcode)[T.98059] | 1.03e+05 | 7771.563 | 13.248 | 0.000 | 8.77e+04 | 1.18e+05 |
| C(zipcode)[T.98065] | 1.507e+05 | 8461.132 | 17.808 | 0.000 | 1.34e+05 | 1.67e+05 |
| C(zipcode)[T.98070] | 5.365e+04 | 2.03e+04 | 2.638 | 0.008 | 1.38e+04 | 9.35e+04 |
| C(zipcode)[T.98072] | 1.537e+05 | 1.01e+04 | 15.211 | 0.000 | 1.34e+05 | 1.74e+05 |
| C(zipcode)[T.98074] | 2.204e+05 | 8015.268 | 27.503 | 0.000 | 2.05e+05 | 2.36e+05 |
| C(zipcode)[T.98075] | 2.39e+05 | 8733.668 | 27.367 | 0.000 | 2.22e+05 | 2.56e+05 |
| C(zipcode)[T.98077] | 1.812e+05 | 1.65e+04 | 10.954 | 0.000 | 1.49e+05 | 2.14e+05 |
| C(zipcode)[T.98092] | -1.837e+04 | 8456.354 | -2.172 | 0.030 | -3.49e+04 | -1791.967 |
| C(zipcode)[T.98102] | 4.065e+05 | 1.33e+04 | 30.557 | 0.000 | 3.8e+05 | 4.33e+05 |
| C(zipcode)[T.98103] | 3.213e+05 | 7664.977 | 41.922 | 0.000 | 3.06e+05 | 3.36e+05 |
| C(zipcode)[T.98105] | 3.797e+05 | 9745.042 | 38.959 | 0.000 | 3.61e+05 | 3.99e+05 |
| C(zipcode)[T.98106] | 1.247e+05 | 8161.089 | 15.275 | 0.000 | 1.09e+05 | 1.41e+05 |
| C(zipcode)[T.98107] | 3.283e+05 | 8719.777 | 37.644 | 0.000 | 3.11e+05 | 3.45e+05 |
| C(zipcode)[T.98108] | 1.172e+05 | 9478.141 | 12.366 | 0.000 | 9.86e+04 | 1.36e+05 |
| C(zipcode)[T.98109] | 4.1e+05 | 1.26e+04 | 32.587 | 0.000 | 3.85e+05 | 4.35e+05 |
| C(zipcode)[T.98112] | 4.336e+05 | 9932.311 | 43.658 | 0.000 | 4.14e+05 | 4.53e+05 |
| C(zipcode)[T.98115] | 3.23e+05 | 7532.540 | 42.883 | 0.000 | 3.08e+05 | 3.38e+05 |

| | | | | | | |
|---|---|---|---|---|---|---|
| C(zipcode)[T.98116] | 3.013e+05 | 8295.181 | 36.317 | 0.000 | 2.85e+05 | 3.18e+05 |
| C(zipcode)[T.98117] | 3.139e+05 | 7549.873 | 41.577 | 0.000 | 2.99e+05 | 3.29e+05 |
| C(zipcode)[T.98118] | 1.754e+05 | 7538.869 | 23.266 | 0.000 | 1.61e+05 | 1.9e+05 |
| C(zipcode)[T.98119] | 4.156e+05 | 1.04e+04 | 39.846 | 0.000 | 3.95e+05 | 4.36e+05 |
| C(zipcode)[T.98122] | 2.998e+05 | 8813.270 | 34.015 | 0.000 | 2.83e+05 | 3.17e+05 |
| C(zipcode)[T.98125] | 1.984e+05 | 7749.216 | 25.602 | 0.000 | 1.83e+05 | 2.14e+05 |
| C(zipcode)[T.98126] | 2.037e+05 | 8097.744 | 25.158 | 0.000 | 1.88e+05 | 2.2e+05 |
| C(zipcode)[T.98133] | 1.486e+05 | 7475.323 | 19.884 | 0.000 | 1.34e+05 | 1.63e+05 |
| C(zipcode)[T.98136] | 2.757e+05 | 8701.258 | 31.680 | 0.000 | 2.59e+05 | 2.93e+05 |
| C(zipcode)[T.98144] | 2.481e+05 | 8298.718 | 29.892 | 0.000 | 2.32e+05 | 2.64e+05 |
| C(zipcode)[T.98146] | 1.164e+05 | 8366.998 | 13.915 | 0.000 | 1e+05 | 1.33e+05 |
| C(zipcode)[T.98148] | 5.139e+04 | 1.43e+04 | 3.594 | 0.000 | 2.34e+04 | 7.94e+04 |
| C(zipcode)[T.98155] | 1.387e+05 | 7643.726 | 18.147 | 0.000 | 1.24e+05 | 1.54e+05 |
| C(zipcode)[T.98166] | 1.023e+05 | 9096.673 | 11.245 | 0.000 | 8.45e+04 | 1.2e+05 |
| C(zipcode)[T.98168] | 5.321e+04 | 8810.583 | 6.039 | 0.000 | 3.59e+04 | 7.05e+04 |
| C(zipcode)[T.98177] | 2.232e+05 | 9115.407 | 24.490 | 0.000 | 2.05e+05 | 2.41e+05 |
| C(zipcode)[T.98178] | 6.778e+04 | 8610.518 | 7.872 | 0.000 | 5.09e+04 | 8.47e+04 |
| C(zipcode)[T.98188] | 3.039e+04 | 1.07e+04 | 2.845 | 0.004 | 9455.739 | 5.13e+04 |
| C(zipcode)[T.98198] | 4.082e+04 | 8550.046 | 4.774 | 0.000 | 2.41e+04 | 5.76e+04 |
| C(zipcode)[T.98199] | 3.581e+05 | 8770.515 | 40.824 | 0.000 | 3.41e+05 | 3.75e+05 |
| C(renovated)[T.1] | 3.018e+04 | 4280.528 | 7.050 | 0.000 | 2.18e+04 | 3.86e+04 |
| bedrooms | -8951.1767 | 1153.274 | -7.762 | 0.000 | -1.12e+04 | -6690.591 |
| bathrooms | 1.25e+04 | 1896.048 | 6.591 | 0.000 | 8780.852 | 1.62e+04 |
| sqft_living | 123.7263 | 2.142 | 57.767 | 0.000 | 119.528 | 127.925 |
| sqft_lot | 2.1906 | 0.313 | 6.992 | 0.000 | 1.576 | 2.805 |
| grade | 5.287e+04 | 1285.228 | 41.133 | 0.000 | 5.03e+04 | 5.54e+04 |
| age | 513.1744 | 48.990 | 10.475 | 0.000 | 417.147 | 609.202 |

| | | | |
|---|---|---|---|
| Omnibus: | 1251.902 | Durbin-Watson: | 1.990 |
| Prob(Omnibus): | 0.000 | Jarque-Bera (JB): | 4136.115 |
| Skew: | 0.478 | Prob(JB): | 0.00 |
| Kurtosis: | 5.578 | Cond. No. | 5.92e+05 |

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 5.92e+05. This might indicate that there are strong multicollinearity or other numerical problems.

In [718]:
```python
model2.pvalues

pvals = model2.pvalues

pvals[pvals > 0.05]
# pvals[pvals > 0.05].index
```

Out[718]:
```
C(floors)[T.2.0]      0.60
C(floors)[T.2.5]      0.12
C(floors)[T.3.5]      0.09
C(zipcode)[T.98002]   0.17
C(zipcode)[T.98003]   0.77
C(zipcode)[T.98022]   0.41
C(zipcode)[T.98030]   0.65
C(zipcode)[T.98031]   0.10
C(zipcode)[T.98032]   0.95
C(zipcode)[T.98042]   0.12
dtype: float64
```

It looks like bedrooms, bathrooms, and grade are now significant. Let's take a closer look at the p-value ranking (removing zipcodes from ranks to more easily interpret results).

In [719]:
```python
pvals = model2.pvalues
pvals_list = pvals.sort_values(ascending=True)

pvals_df = pvals_list.to_frame()

pd.options.display.max_rows = 999
pvals_df = pvals_df.reset_index()
pvals_df = pvals_df.rename(columns={'index': 'Variable', 0: 'P_Value'})

pvals_df[~pvals_df['Variable'].str.contains("zipcode")]
```

Out[719]:

| | Variable | P_Value |
|---|---|---|
| 9 | sqft_living | 0.00 |
| 11 | grade | 0.00 |
| 36 | Intercept | 0.00 |
| 38 | C(waterfront)[T.1.0] | 0.00 |
| 50 | age | 0.00 |
| 52 | C(has_basement)[T.1] | 0.00 |
| 55 | bedrooms | 0.00 |
| 56 | C(floors)[T.3.0] | 0.00 |
| 57 | C(renovated)[T.1] | 0.00 |
| 58 | sqft_lot | 0.00 |
| 60 | bathrooms | 0.00 |
| 65 | C(condition)[T.5] | 0.00 |
| 68 | C(condition)[T.4] | 0.00 |
| 70 | C(condition)[T.3] | 0.00 |
| 74 | C(condition)[T.2] | 0.02 |
| 76 | C(floors)[T.1.5] | 0.03 |
| 78 | C(floors)[T.3.5] | 0.09 |
| 80 | C(floors)[T.2.5] | 0.12 |
| 84 | C(floors)[T.2.0] | 0.60 |

## Interpretation and future analysis

When we first ran this model, we mistakenly developed our analyses and insights before running a true second model with bedrooms, bathrooms, and grade as continuous variables. Upon review of our second true model, we would have opted to continue with this approach and potentially accept this model for generating insights.

We now realize that grade, bedrooms, and bathrooms are all significant features to our model, and these should have been left in for analysis. Importantly, this also increased our adjusted R-squared significantly from 0.798 to 0.824.

For a third model, we might also consider treating floors as a continuous variable, since the half floors don't seem to provide any useful insight. Each floor other than 1.5 floors seems to decrease value, meaning it might provide a decent linear coefficient.

Let's also take a look at our coefficients.

```
In [720]:
1  coeffs = model2.params
2  coeffs_list = coeffs.sort_values(ascending=False).round(2)
3
4  coeff_df = coeffs_list.to_frame()
5
6  pd.options.display.max_rows = 999
7  coeff_df = coeff_df.reset_index()
8  coeff_df = coeff_df.rename(columns={'index': 'Variable', 0: 'Dollar Impact'}
9
10 coeff_df['Dollar Impact'] = coeff_df['Dollar Impact'].apply(lambda x: "{:,}"
11
12 coeff_df[~coeff_df['Variable'].str.contains("zipcode")]
```

Out[720]:

| | Variable | Dollar Impact |
|---|---|---|
| 11 | C(waterfront)[T.1.0] | 325,785.65 |
| 41 | C(condition)[T.5] | 132,720.04 |
| 47 | C(condition)[T.4] | 105,603.05 |
| 54 | C(condition)[T.3] | 87,276.68 |
| 56 | C(condition)[T.2] | 58,062.04 |
| 59 | grade | 52,865.59 |
| 66 | C(renovated)[T.1] | 30,177.93 |
| 69 | bathrooms | 12,497.38 |
| 72 | C(floors)[T.1.5] | 6,715.72 |
| 76 | age | 513.17 |
| 77 | sqft_living | 123.73 |
| 78 | sqft_lot | 2.19 |
| 79 | C(floors)[T.2.0] | -1,360.51 |
| 80 | bedrooms | -8,951.18 |
| 82 | C(floors)[T.2.5] | -16,923.14 |
| 84 | C(has_basement)[T.1] | -18,544.2 |
| 85 | C(floors)[T.3.0] | -40,323.46 |
| 86 | C(floors)[T.3.5] | -64,229.11 |
| 87 | Intercept | -465,955.84 |

The coefficients tell a bit of a different story. Grade is definitely significant. However, based on our

histogram from grade, we wouldn't expect each step in grade to be equivalent to $52k. Perhaps it should be rerun and reinterpreted as categorical. Bathrooms gives us something useful to work with, and bedrooms is a bit surprising, removing $9k in value for each additional bedroom.

Our square footage living change from $159 to $124. which would have significant alternative values four our recommendation to add additional sqft_living.

## Initial model with bedrooms, bathrooms, and grade removed

Below is the rest of the project, which was completed using the initial model and by removing insignificant p-values.

In [721]:
```python
1  df = df.drop(['bedrooms', 'bathrooms', 'grade'], axis=1)
```

In [722]:
```python
1  df.head()
```

Out[722]:

|  | price | sqft_living | sqft_lot | floors | waterfront | condition | has_basement | zipcode |
|---|---|---|---|---|---|---|---|---|
| id |  |  |  |  |  |  |  |  |
| 6414100192 | 538000.00 | 2570 | 7242 | 2.00 | 0.00 | 3 | 1 | 98125 |
| 2487200875 | 604000.00 | 1960 | 5000 | 1.00 | 0.00 | 5 | 1 | 98136 |
| 1954400510 | 510000.00 | 1680 | 8080 | 1.00 | 0.00 | 3 | 0 | 98074 |
| 1321400060 | 257500.00 | 1715 | 6819 | 2.00 | 0.00 | 3 | 0 | 98003 |
| 2414600126 | 229500.00 | 1780 | 7470 | 1.00 | 0.00 | 3 | 1 | 98146 |

In [723]:
```python
1  categoricals = [
2  #                 'bedrooms',
3  #                 'bathrooms',
4                  'floors',
5                  'waterfront',
6                  'condition',
7  #                 'grade',
8                  'has_basement',
9                  'zipcode',
10                 'renovated'
11                 ]
```

In [724]:

```
1 model1 = make_model(df, categoricals)
```

OLS Regression Results

| | | | |
|---|---|---|---|
| Dep. Variable: | price | R-squared: | 0.799 |
| Model: | OLS | Adj. R-squared: | 0.798 |
| Method: | Least Squares | F-statistic: | 617.7 |
| Date: | Sat, 01 May 2021 | Prob (F-statistic): | 0.00 |
| Time: | 11:42:05 | Log-Likelihood: | -1.6841e+05 |
| No. Observations: | 13128 | AIC: | 3.370e+05 |
| Df Residuals: | 13043 | BIC: | 3.376e+05 |
| Df Model: | 84 | | |
| Covariance Type: | nonrobust | | |

| | coef | std err | t | P>\|t\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| Intercept | -1.882e+05 | 2.64e+04 | -7.127 | 0.000 | -2.4e+05 | -1.36e+05 |
| C(floors)[T.1.5] | 9218.9133 | 3242.995 | 2.843 | 0.004 | 2862.170 | 1.56e+04 |
| C(floors)[T.2.0] | 1.145e+04 | 2699.623 | 4.241 | 0.000 | 6157.285 | 1.67e+04 |
| C(floors)[T.2.5] | 1225.8619 | 1.15e+04 | 0.107 | 0.915 | -2.13e+04 | 2.38e+04 |
| C(floors)[T.3.0] | -2.247e+04 | 5929.341 | -3.790 | 0.000 | -3.41e+04 | -1.08e+04 |
| C(floors)[T.3.5] | -4.043e+04 | 4.07e+04 | -0.992 | 0.321 | -1.2e+05 | 3.94e+04 |
| C(waterfront)[T.1.0] | 3.378e+05 | 2.02e+04 | 16.763 | 0.000 | 2.98e+05 | 3.77e+05 |
| C(condition)[T.2] | 8.736e+04 | 2.69e+04 | 3.250 | 0.001 | 3.47e+04 | 1.4e+05 |
| C(condition)[T.3] | 1.288e+05 | 2.53e+04 | 5.089 | 0.000 | 7.92e+04 | 1.78e+05 |
| C(condition)[T.4] | 1.458e+05 | 2.53e+04 | 5.762 | 0.000 | 9.62e+04 | 1.95e+05 |
| C(condition)[T.5] | 1.741e+05 | 2.54e+04 | 6.854 | 0.000 | 1.24e+05 | 2.24e+05 |
| C(has_basement)[T.1] | -2.363e+04 | 2044.381 | -11.561 | 0.000 | -2.76e+04 | -1.96e+04 |
| C(zipcode)[T.98002] | 3273.3380 | 9714.062 | 0.337 | 0.736 | -1.58e+04 | 2.23e+04 |
| C(zipcode)[T.98003] | 1.968e+04 | 8861.380 | 2.221 | 0.026 | 2315.287 | 3.71e+04 |
| C(zipcode)[T.98004] | 5.533e+05 | 1.11e+04 | 49.936 | 0.000 | 5.32e+05 | 5.75e+05 |
| C(zipcode)[T.98005] | 3.733e+05 | 1.16e+04 | 32.140 | 0.000 | 3.5e+05 | 3.96e+05 |
| C(zipcode)[T.98006] | 3.195e+05 | 8420.462 | 37.937 | 0.000 | 3.03e+05 | 3.36e+05 |
| C(zipcode)[T.98007] | 2.767e+05 | 1.12e+04 | 24.653 | 0.000 | 2.55e+05 | 2.99e+05 |
| C(zipcode)[T.98008] | 2.677e+05 | 9070.409 | 29.519 | 0.000 | 2.5e+05 | 2.86e+05 |
| C(zipcode)[T.98010] | 8.794e+04 | 1.79e+04 | 4.901 | 0.000 | 5.28e+04 | 1.23e+05 |
| C(zipcode)[T.98011] | 1.592e+05 | 1.02e+04 | 15.578 | 0.000 | 1.39e+05 | 1.79e+05 |
| C(zipcode)[T.98014] | 8.193e+04 | 1.72e+04 | 4.766 | 0.000 | 4.82e+04 | 1.16e+05 |
| C(zipcode)[T.98019] | 8.668e+04 | 1.14e+04 | 7.579 | 0.000 | 6.43e+04 | 1.09e+05 |
| C(zipcode)[T.98022] | 1.061e+04 | 1.07e+04 | 0.988 | 0.323 | -1.04e+04 | 3.17e+04 |

| | coef | std err | t | P>\|t\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| C(zipcode)[T.98023] | 25.9815 | 7897.722 | 0.003 | 0.997 | -1.55e+04 | 1.55e+04 |
| C(zipcode)[T.98024] | 1.245e+05 | 2.12e+04 | 5.876 | 0.000 | 8.3e+04 | 1.66e+05 |
| C(zipcode)[T.98027] | 2.665e+05 | 9265.219 | 28.765 | 0.000 | 2.48e+05 | 2.85e+05 |
| C(zipcode)[T.98028] | 1.475e+05 | 9073.443 | 16.253 | 0.000 | 1.3e+05 | 1.65e+05 |
| C(zipcode)[T.98029] | 2.666e+05 | 8694.197 | 30.664 | 0.000 | 2.5e+05 | 2.84e+05 |
| C(zipcode)[T.98030] | 1.129e+04 | 9218.780 | 1.225 | 0.221 | -6775.191 | 2.94e+04 |
| C(zipcode)[T.98031] | 2.115e+04 | 9073.600 | 2.331 | 0.020 | 3361.926 | 3.89e+04 |
| C(zipcode)[T.98032] | 1.05e+04 | 1.13e+04 | 0.929 | 0.353 | -1.17e+04 | 3.27e+04 |
| C(zipcode)[T.98033] | 3.461e+05 | 8388.089 | 41.262 | 0.000 | 3.3e+05 | 3.63e+05 |
| C(zipcode)[T.98034] | 2.101e+05 | 7784.164 | 26.984 | 0.000 | 1.95e+05 | 2.25e+05 |
| C(zipcode)[T.98038] | 3.795e+04 | 7830.954 | 4.846 | 0.000 | 2.26e+04 | 5.33e+04 |
| C(zipcode)[T.98039] | 6.28e+05 | 5.26e+04 | 11.934 | 0.000 | 5.25e+05 | 7.31e+05 |
| C(zipcode)[T.98040] | 4.706e+05 | 1.08e+04 | 43.755 | 0.000 | 4.5e+05 | 4.92e+05 |
| C(zipcode)[T.98042] | 1.173e+04 | 7969.154 | 1.472 | 0.141 | -3890.631 | 2.74e+04 |
| C(zipcode)[T.98045] | 1.078e+05 | 1.09e+04 | 9.926 | 0.000 | 8.65e+04 | 1.29e+05 |
| C(zipcode)[T.98052] | 2.877e+05 | 7800.158 | 36.881 | 0.000 | 2.72e+05 | 3.03e+05 |
| C(zipcode)[T.98053] | 2.695e+05 | 9311.825 | 28.938 | 0.000 | 2.51e+05 | 2.88e+05 |
| C(zipcode)[T.98055] | 5.664e+04 | 9063.361 | 6.250 | 0.000 | 3.89e+04 | 7.44e+04 |
| C(zipcode)[T.98056] | 1.042e+05 | 8256.293 | 12.624 | 0.000 | 8.8e+04 | 1.2e+05 |
| C(zipcode)[T.98058] | 4.524e+04 | 8233.756 | 5.494 | 0.000 | 2.91e+04 | 6.14e+04 |
| C(zipcode)[T.98059] | 1.079e+05 | 8317.840 | 12.967 | 0.000 | 9.16e+04 | 1.24e+05 |
| C(zipcode)[T.98065] | 1.438e+05 | 9042.380 | 15.899 | 0.000 | 1.26e+05 | 1.61e+05 |
| C(zipcode)[T.98070] | 4.448e+04 | 2.18e+04 | 2.044 | 0.041 | 1818.902 | 8.71e+04 |
| C(zipcode)[T.98072] | 1.624e+05 | 1.08e+04 | 15.016 | 0.000 | 1.41e+05 | 1.84e+05 |
| C(zipcode)[T.98074] | 2.635e+05 | 8514.085 | 30.948 | 0.000 | 2.47e+05 | 2.8e+05 |
| C(zipcode)[T.98075] | 2.832e+05 | 9281.052 | 30.510 | 0.000 | 2.65e+05 | 3.01e+05 |
| C(zipcode)[T.98077] | 2.27e+05 | 1.77e+04 | 12.846 | 0.000 | 1.92e+05 | 2.62e+05 |
| C(zipcode)[T.98092] | -5509.6749 | 9046.424 | -0.609 | 0.543 | -2.32e+04 | 1.22e+04 |
| C(zipcode)[T.98102] | 4.643e+05 | 1.42e+04 | 32.770 | 0.000 | 4.37e+05 | 4.92e+05 |
| C(zipcode)[T.98103] | 3.629e+05 | 8141.622 | 44.574 | 0.000 | 3.47e+05 | 3.79e+05 |
| C(zipcode)[T.98105] | 4.273e+05 | 1.04e+04 | 41.231 | 0.000 | 4.07e+05 | 4.48e+05 |
| C(zipcode)[T.98106] | 1.339e+05 | 8733.545 | 15.326 | 0.000 | 1.17e+05 | 1.51e+05 |
| C(zipcode)[T.98107] | 3.734e+05 | 9268.057 | 40.288 | 0.000 | 3.55e+05 | 3.92e+05 |
| C(zipcode)[T.98108] | 1.323e+05 | 1.01e+04 | 13.047 | 0.000 | 1.12e+05 | 1.52e+05 |
| C(zipcode)[T.98109] | 4.731e+05 | 1.34e+04 | 35.366 | 0.000 | 4.47e+05 | 4.99e+05 |
| C(zipcode)[T.98112] | 4.945e+05 | 1.05e+04 | 46.980 | 0.000 | 4.74e+05 | 5.15e+05 |
| C(zipcode)[T.98115] | 3.588e+05 | 8014.771 | 44.765 | 0.000 | 3.43e+05 | 3.74e+05 |

| | | | | | | |
|---|---|---|---|---|---|---|
| C(zipcode)[T.98116] | 3.428e+05 | 8820.213 | 38.871 | 0.000 | 3.26e+05 | 3.6e+05 |
| C(zipcode)[T.98117] | 3.531e+05 | 8023.029 | 44.006 | 0.000 | 3.37e+05 | 3.69e+05 |
| C(zipcode)[T.98118] | 1.958e+05 | 8054.593 | 24.307 | 0.000 | 1.8e+05 | 2.12e+05 |
| C(zipcode)[T.98119] | 4.814e+05 | 1.1e+04 | 43.571 | 0.000 | 4.6e+05 | 5.03e+05 |
| C(zipcode)[T.98122] | 3.562e+05 | 9332.522 | 38.170 | 0.000 | 3.38e+05 | 3.75e+05 |
| C(zipcode)[T.98125] | 2.166e+05 | 8282.370 | 26.147 | 0.000 | 2e+05 | 2.33e+05 |
| C(zipcode)[T.98126] | 2.317e+05 | 8636.709 | 26.826 | 0.000 | 2.15e+05 | 2.49e+05 |
| C(zipcode)[T.98133] | 1.668e+05 | 7987.400 | 20.889 | 0.000 | 1.51e+05 | 1.83e+05 |
| C(zipcode)[T.98136] | 3.127e+05 | 9266.445 | 33.746 | 0.000 | 2.95e+05 | 3.31e+05 |
| C(zipcode)[T.98144] | 2.811e+05 | 8844.537 | 31.783 | 0.000 | 2.64e+05 | 2.98e+05 |
| C(zipcode)[T.98146] | 1.198e+05 | 8955.507 | 13.382 | 0.000 | 1.02e+05 | 1.37e+05 |
| C(zipcode)[T.98148] | 6.074e+04 | 1.53e+04 | 3.969 | 0.000 | 3.07e+04 | 9.07e+04 |
| C(zipcode)[T.98155] | 1.496e+05 | 8178.491 | 18.294 | 0.000 | 1.34e+05 | 1.66e+05 |
| C(zipcode)[T.98166] | 1.174e+05 | 9730.631 | 12.066 | 0.000 | 9.83e+04 | 1.36e+05 |
| C(zipcode)[T.98168] | 4.942e+04 | 9428.239 | 5.241 | 0.000 | 3.09e+04 | 6.79e+04 |
| C(zipcode)[T.98177] | 2.539e+05 | 9727.865 | 26.098 | 0.000 | 2.35e+05 | 2.73e+05 |
| C(zipcode)[T.98178] | 6.844e+04 | 9217.735 | 7.425 | 0.000 | 5.04e+04 | 8.65e+04 |
| C(zipcode)[T.98188] | 3.22e+04 | 1.14e+04 | 2.816 | 0.005 | 9786.049 | 5.46e+04 |
| C(zipcode)[T.98198] | 4.897e+04 | 9149.900 | 5.352 | 0.000 | 3.1e+04 | 6.69e+04 |
| C(zipcode)[T.98199] | 4.053e+05 | 9314.796 | 43.511 | 0.000 | 3.87e+05 | 4.24e+05 |
| C(renovated)[T.1] | 4.226e+04 | 4533.183 | 9.322 | 0.000 | 3.34e+04 | 5.11e+04 |
| sqft_living | 159.4767 | 1.605 | 99.364 | 0.000 | 156.331 | 162.623 |
| sqft_lot | 3.4502 | 0.333 | 10.352 | 0.000 | 2.797 | 4.103 |
| age | -168.9231 | 47.874 | -3.528 | 0.000 | -262.764 | -75.083 |

| | | | |
|---|---|---|---|
| Omnibus: | 1167.900 | Durbin-Watson: | 2.001 |
| Prob(Omnibus): | 0.000 | Jarque-Bera (JB): | 3658.229 |
| Skew: | 0.459 | Prob(JB): | 0.00 |
| Kurtosis: | 5.418 | Cond. No. | 5.91e+05 |

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 5.91e+05. This might indicate that there are
strong multicollinearity or other numerical problems.

```
In [725]:    1  model1.pvalues
             2
             3  pvals = model1.pvalues
             4
             5  pvals[pvals > 0.05]
             6  # pvals[pvals > 0.05].index
```

```
Out[725]:  C(floors)[T.2.5]     0.92
           C(floors)[T.3.5]     0.32
           C(zipcode)[T.98002]  0.74
           C(zipcode)[T.98022]  0.32
           C(zipcode)[T.98023]  1.00
           C(zipcode)[T.98030]  0.22
           C(zipcode)[T.98032]  0.35
           C(zipcode)[T.98042]  0.14
           C(zipcode)[T.98092]  0.54
           dtype: float64
```

```
In [726]:    1  pd.set_option('display.float_format', lambda x: '%.10f' % x)
             2
             3  type(model1.pvalues)
```

```
Out[726]:  pandas.core.series.Series
```

In [727]:
```python
1  dfp = model1.pvalues.to_frame()
2
3  dfp.sort_values(by=[0])
```

Out[727]:

|  | 0 |
|---|---|
| C(zipcode)[T.98103] | 0.0000000000 |
| sqft_living | 0.0000000000 |
| C(zipcode)[T.98105] | 0.0000000000 |
| C(zipcode)[T.98199] | 0.0000000000 |
| C(zipcode)[T.98117] | 0.0000000000 |
| C(zipcode)[T.98107] | 0.0000000000 |
| C(zipcode)[T.98116] | 0.0000000000 |
| C(zipcode)[T.98033] | 0.0000000000 |
| C(zipcode)[T.98004] | 0.0000000000 |
| C(zipcode)[T.98115] | 0.0000000000 |
| C(zipcode)[T.98040] | 0.0000000000 |
| C(zipcode)[T.98119] | 0.0000000000 |
| C(zipcode)[T.98112] | 0.0000000000 |
| C(zipcode)[T.98122] | 0.0000000000 |
| C(zipcode)[T.98006] | 0.0000000000 |
| C(zipcode)[T.98052] | 0.0000000000 |
| C(zipcode)[T.98109] | 0.0000000000 |
| C(zipcode)[T.98136] | 0.0000000000 |
| C(zipcode)[T.98102] | 0.0000000000 |
| C(zipcode)[T.98005] | 0.0000000000 |
| C(zipcode)[T.98144] | 0.0000000000 |
| C(zipcode)[T.98074] | 0.0000000000 |
| C(zipcode)[T.98029] | 0.0000000000 |
| C(zipcode)[T.98075] | 0.0000000000 |
| C(zipcode)[T.98008] | 0.0000000000 |
| C(zipcode)[T.98053] | 0.0000000000 |
| C(zipcode)[T.98027] | 0.0000000000 |
| C(zipcode)[T.98034] | 0.0000000000 |
| C(zipcode)[T.98126] | 0.0000000000 |
| C(zipcode)[T.98125] | 0.0000000000 |
| C(zipcode)[T.98177] | 0.0000000000 |
| C(zipcode)[T.98007] | 0.0000000000 |
| C(zipcode)[T.98118] | 0.0000000000 |

|  | 0 |
| --- | --- |
| C(zipcode)[T.98133] | 0.0000000000 |
| C(zipcode)[T.98155] | 0.0000000000 |
| C(waterfront)[T.1.0] | 0.0000000000 |
| C(zipcode)[T.98028] | 0.0000000000 |
| C(zipcode)[T.98065] | 0.0000000000 |
| C(zipcode)[T.98011] | 0.0000000000 |
| C(zipcode)[T.98106] | 0.0000000000 |
| C(zipcode)[T.98072] | 0.0000000000 |
| C(zipcode)[T.98146] | 0.0000000000 |
| C(zipcode)[T.98108] | 0.0000000000 |
| C(zipcode)[T.98059] | 0.0000000000 |
| C(zipcode)[T.98077] | 0.0000000000 |
| C(zipcode)[T.98056] | 0.0000000000 |
| C(zipcode)[T.98166] | 0.0000000000 |
| C(zipcode)[T.98039] | 0.0000000000 |
| C(has_basement)[T.1] | 0.0000000000 |
| sqft_lot | 0.0000000000 |
| C(zipcode)[T.98045] | 0.0000000000 |
| C(renovated)[T.1] | 0.0000000000 |
| C(zipcode)[T.98019] | 0.0000000000 |
| C(zipcode)[T.98178] | 0.0000000000 |
| Intercept | 0.0000000000 |
| C(condition)[T.5] | 0.0000000000 |
| C(zipcode)[T.98055] | 0.0000000004 |
| C(zipcode)[T.98024] | 0.0000000043 |
| C(condition)[T.4] | 0.0000000085 |
| C(zipcode)[T.98058] | 0.0000000400 |
| C(zipcode)[T.98198] | 0.0000000883 |
| C(zipcode)[T.98168] | 0.0000001618 |
| C(condition)[T.3] | 0.0000003648 |
| C(zipcode)[T.98010] | 0.0000009651 |
| C(zipcode)[T.98038] | 0.0000012756 |
| C(zipcode)[T.98014] | 0.0000018998 |
| C(floors)[T.2.0] | 0.0000224120 |
| C(zipcode)[T.98148] | 0.0000725704 |
| C(floors)[T.3.0] | 0.0001514949 |

|  | 0 |
| --- | --- |
| age | 0.0004193930 |
| C(condition)[T.2] | 0.0011577263 |
| C(floors)[T.1.5] | 0.0044800397 |
| C(zipcode)[T.98188] | 0.0048690388 |
| C(zipcode)[T.98031] | 0.0197862920 |
| C(zipcode)[T.98003] | 0.0263393414 |
| C(zipcode)[T.98070] | 0.0410018990 |
| C(zipcode)[T.98042] | 0.1410626590 |
| C(zipcode)[T.98030] | 0.2205173493 |
| C(floors)[T.3.5] | 0.3210375106 |
| C(zipcode)[T.98022] | 0.3230609814 |
| C(zipcode)[T.98032] | 0.3529145341 |
| C(zipcode)[T.98092] | 0.5425055560 |
| C(zipcode)[T.98002] | 0.7361456692 |
| C(floors)[T.2.5] | 0.9151773929 |
| C(zipcode)[T.98023] | 0.9973752208 |

```
In [728]:   1  pvals = model1.pvalues
            2  pvals_list = pvals.sort_values(ascending=True)
            3
            4  pvals_df = pvals_list.to_frame()
            5
            6  pd.options.display.max_rows = 999
            7  pvals_df = pvals_df.reset_index()
            8  pvals_df = pvals_df.rename(columns={'index': 'Variable', 0: 'P_Value'})
            9
           10  pvals_df[~pvals_df['Variable'].str.contains("zipcode")]
```

Out[728]:

|     | Variable | P_Value |
|-----|----------|---------|
| 1   | sqft_living | 0.0000000000 |
| 35  | C(waterfront)[T.1.0] | 0.0000000000 |
| 48  | C(has_basement)[T.1] | 0.0000000000 |
| 49  | sqft_lot | 0.0000000000 |
| 51  | C(renovated)[T.1] | 0.0000000000 |
| 54  | Intercept | 0.0000000000 |
| 55  | C(condition)[T.5] | 0.0000000000 |
| 58  | C(condition)[T.4] | 0.0000000085 |
| 62  | C(condition)[T.3] | 0.0000003648 |
| 66  | C(floors)[T.2.0] | 0.0000224120 |
| 68  | C(floors)[T.3.0] | 0.0001514949 |
| 69  | age | 0.0004193930 |
| 70  | C(condition)[T.2] | 0.0011577263 |
| 71  | C(floors)[T.1.5] | 0.0044800397 |
| 78  | C(floors)[T.3.5] | 0.3210375106 |
| 83  | C(floors)[T.2.5] | 0.9151773929 |

After running our p-value check again, some zip codes are still insignificant, but not enough to remove zip codes from the model. The 2.5 and 3.5 floors are insignificant, but that is likely due to half-floors having little representation in our dataset.

In [729]:
```python
coeffs = model1.params
coeffs_list = coeffs.sort_values(ascending=False).round(2)

coeff_df = coeffs_list.to_frame()

pd.options.display.max_rows = 999
coeff_df = coeff_df.reset_index()
coeff_df = coeff_df.rename(columns={'index': 'Variable', 0: 'Dollar Impact'}

coeff_df['Dollar Impact'] = coeff_df['Dollar Impact'].apply(lambda x: "{:,}"

coeff_df[~coeff_df['Variable'].str.contains("zipcode")]
```

Out[729]:

| | Variable | Dollar Impact |
|---|---|---|
| 17 | C(waterfront)[T.1.0] | 337,780.33 |
| 35 | C(condition)[T.5] | 174,135.47 |
| 41 | C(condition)[T.4] | 145,846.97 |
| 45 | C(condition)[T.3] | 128,815.16 |
| 53 | C(condition)[T.2] | 87,360.03 |
| 63 | C(renovated)[T.1] | 42,260.43 |
| 69 | C(floors)[T.2.0] | 11,448.94 |
| 73 | C(floors)[T.1.5] | 9,218.91 |
| 75 | C(floors)[T.2.5] | 1,225.86 |
| 76 | sqft_living | 159.48 |
| 78 | sqft_lot | 3.45 |
| 79 | age | -168.92 |
| 81 | C(floors)[T.3.0] | -22,470.55 |
| 82 | C(has_basement)[T.1] | -23,634.14 |
| 83 | C(floors)[T.3.5] | -40,433.65 |
| 84 | Intercept | -188,190.8 |

# INTERPRET

Before looking at zipcode, let's take a look at our feature coefficients, which represent price impact.

Waterfront is the most impactful, adding $338k to price.

Condition lines up with our expectations. The greater the condition, the more valuable the home. Improving the condition from 1 to 5 would add an estimated $174,135 to a home owner's value.

Renovated homes seem to fetch a larger price of approximately $42,260, which aligns with expectations.

Floors is a bit counterintuitive. While 2 floors seems to increase the value by $11.5k, a third floor decreases value by $22.5k, 3.5 floors decreases by $40.5k. Considering the cost of adding an additional floor would likely be much more expensive than these coefficients, this might indicate that expanding the square footage of a home within floors that already exist might be a more sensible investment.

Sqft_living gives us an estimated value of $159 for every additional square foot of space.

On the surface, sqft_lot looks like it has a relatively lower impact on price. However, it is still relevant when comparing properties with significant differences in size. One acre is 43,560 square feet. Our model predicts that with a $3.45 impact to price for every square foot, an additional acre would add $150,282 to the value of two otherwise identical properties.

Age doesn't seem to have a great impact. Despite having a P-value greater than 0.05, a house will lose $168 in value every year. Even in the case of our oldest houses, age can only have a maximum price impact of $19,425.

Perhaps counterintuitively, the presence of a basement decreases the value of a home by $23,634. This might require further examination.

In [730]:
```python
1  print('Most valuable zip codes:')
2  print(coeff_df[coeff_df['Variable'].str.contains("zipcode")].head(5))
3  print('Least valuable zip codes:')
4  print(coeff_df[coeff_df['Variable'].str.contains("zipcode")].tail(5))
```

```
Most valuable zip codes:
              Variable Dollar Impact
0  C(zipcode)[T.98039]    628,000.93
1  C(zipcode)[T.98004]    553,256.56
2  C(zipcode)[T.98112]     494,473.0
3  C(zipcode)[T.98119]    481,396.09
4  C(zipcode)[T.98109]    473,148.05
Least valuable zip codes:
              Variable Dollar Impact
71  C(zipcode)[T.98022]    10,614.08
72  C(zipcode)[T.98032]     10,500.9
74  C(zipcode)[T.98002]     3,273.34
77  C(zipcode)[T.98023]        25.98
80  C(zipcode)[T.98092]     -5,509.67
```

Depending on the location, zip codes can have the most dramatic impact on price. The most valuable zip codes are those closest to the metropolitan city center (Seattle, Bellevue, and Mercer Island). The impact on price in the top 5 zip codes is an estimated $473-628k.

Other than the least valuable zip code, our model functions in a way that doesn't subtract estimated value from homes. The bottom 5 zip codes are located in Kent, near the southern end of King County. While not the furthest from the city center, they are significantly further than our most valuable zip codes.

# Model Evaluation

Our model has a semi-strong fit with an adjusted R-squared of 0.798. This means it has a predictive power of roughly 79.8%.

Additional steps could be taken to improve predictive power. Standardization and logistic normalization would theoretically improve R-squared and allow us to make more accurate predictions. We did not incorporate these processes into our model because we were more interested in the practical recommendations it could provide, and inferences are difficult to interpret after normalization.

There are likely improvements that we could make to hone in on accuracy. For our residential clients interest in improving the value of their homes, a 79.8% confidence level seems strong enough to make at least some base line recommendations.

# CONCLUSIONS & RECOMMENDATIONS

Our model generated some interesting insights about what drives price in the King County housing market. Here are our major takeaways about the most influential factors in determining a house's price:

## Insights

- Location is the most prized quality of a property. Certain zip codes are highly sought after. The top 5 most valuable zip codes will influence property value by an average of $473k-$628k. These zip codes are generally closer to the metropolitan area. Homes located further from the city to the south are less valuable.
- Similar to location, waterfront properties are also much more more valuable and add an average $337k to property value.
- One might assume that additional bedrooms and bathrooms are more valuable. However, according to our model, what actually drives value is total living area square footage. Understanding this, we can intuitively assume that with additional square footage comes additional bedrooms and bathrooms (on average), but our model does not see the bed/bath count as significant.
- The home condition also has a significant impact on price. Before analysis, we assumed that King County's 'Grade' system might behave similarly, but our model determine that the grade system was not a driver of price.

## Recommendations to Home Owners

Many of the insights generated by analyzing our model did not lead to practical recommendations for home owners. It isn't exactly practical or possible in most cases to uproot a home and move it to a new area or by the water. But we did notice two key ways that an owner can improve their value:

- Adding square footage through home construction is the most practical recommendation we can offer to improve value. Each additional square foot of living space adds an estimated $159.48 in home value. Adding a second floor gives a small bonus and adding a basement

gives a small penalty. However, when factoring in the added square footage of projects like these, the penalties will most likely be absorbed by the added value.

- Renovating also gives a noticeable bump to price, especially if that renovation improves the condition. Home owners should maintain the condition of their home, or it will decrease in value.

## Further Analysis and Modeling

The goal of this project was to develop a very general understanding of the most influential factors in property value. Given more time for data review, we might be able to implement the 'view' feature if we can get a better understanding of what it represents. Sqft_living15, sqft_lot15, and Year Renovated might be interesting to explore. Lat and long can be used to heatmap our dataset to visualize home values on a map of King County.

We could implement standardization and normalization to improve our model's predictive quality. We would also like to implement a train / test split for similar purposes.

It might be helpful to build dynamic splitting of our data. For example, how specifically could the owner of a 2 story, 4 bedroom house in Bellevue improve their home value? Would the coefficients of our features change if we ran our model using only houses that matched that criteria? Dynamic splitting could be useful for generating tailored recommendations to clients who might be willing to pay a premium for such services.

# VISUALS

## Zipcode Graph

```
In [731]:   1  coeffs = model1.params
            2  coeffs_list = coeffs.sort_values(ascending=False).round(2)
            3
            4  coeff_df = coeffs_list.to_frame()
            5
            6  # type(model1.params)
            7  pd.options.display.max_rows = 999
            8  coeff_df = coeff_df.reset_index()
            9  coeff_df = coeff_df.rename(columns={'index': 'Variable', 0: 'Dollar Impact'}
           10
           11  # coeff_df[coeff_df['Variable'].str.contains("zipcode")].head()
```

In [732]:
```python
1  zip_df = coeff_df[coeff_df['Variable'].str.contains("zipcode")]
2  zip_df.head()
```

Out[732]:

| | Variable | Dollar Impact |
|---|---|---|
| 0 | C(zipcode)[T.98039] | 628000.9300000001 |
| 1 | C(zipcode)[T.98004] | 553256.5600000001 |
| 2 | C(zipcode)[T.98112] | 494473.0000000000 |
| 3 | C(zipcode)[T.98119] | 481396.0900000000 |
| 4 | C(zipcode)[T.98109] | 473148.0500000000 |

In [733]:
```python
 1  zipcodes = []
 2
 3  for row in zip_df['Variable']:
 4      old = row
 5      old = old.replace("C(zipcode)[T.", "")
 6      old = old.replace("]", "")
 7      zipcodes.append(old)
 8
 9  zip_df['Zip Code'] = zipcodes
10
11  zip_df.head()
```

```
<ipython-input-733-1481b7f7e852>:9: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/sta
ble/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pyd
ata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-c
opy)
  zip_df['Zip Code'] = zipcodes
```

Out[733]:

| | Variable | Dollar Impact | Zip Code |
|---|---|---|---|
| 0 | C(zipcode)[T.98039] | 628000.9300000001 | 98039 |
| 1 | C(zipcode)[T.98004] | 553256.5600000001 | 98004 |
| 2 | C(zipcode)[T.98112] | 494473.0000000000 | 98112 |
| 3 | C(zipcode)[T.98119] | 481396.0900000000 | 98119 |
| 4 | C(zipcode)[T.98109] | 473148.0500000000 | 98109 |

In [734]:

```python
1  zip_df_top5 = zip_df.head()
2
3  zip_df_bottom5 = zip_df.tail()
4
5  zip_df['Zip Code'] = zip_df['Zip Code'].astype(int)
6  zip_df['Dollar Impact'] = zip_df['Dollar Impact'].astype(float)
7
8  print(zip_df['Zip Code'].describe())
```

```
count       69.0000000000
mean     98078.4057971015
std         56.2707004631
min      98002.0000000000
25%      98030.0000000000
50%      98070.0000000000
75%      98118.0000000000
max      98199.0000000000
Name: Zip Code, dtype: float64

<ipython-input-734-a822e190af78>:5: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/sta
ble/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pyd
ata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-c
opy)
  zip_df['Zip Code'] = zip_df['Zip Code'].astype(int)
<ipython-input-734-a822e190af78>:6: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/sta
ble/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pyd
ata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-c
opy)
  zip_df['Dollar Impact'] = zip_df['Dollar Impact'].astype(float)
```
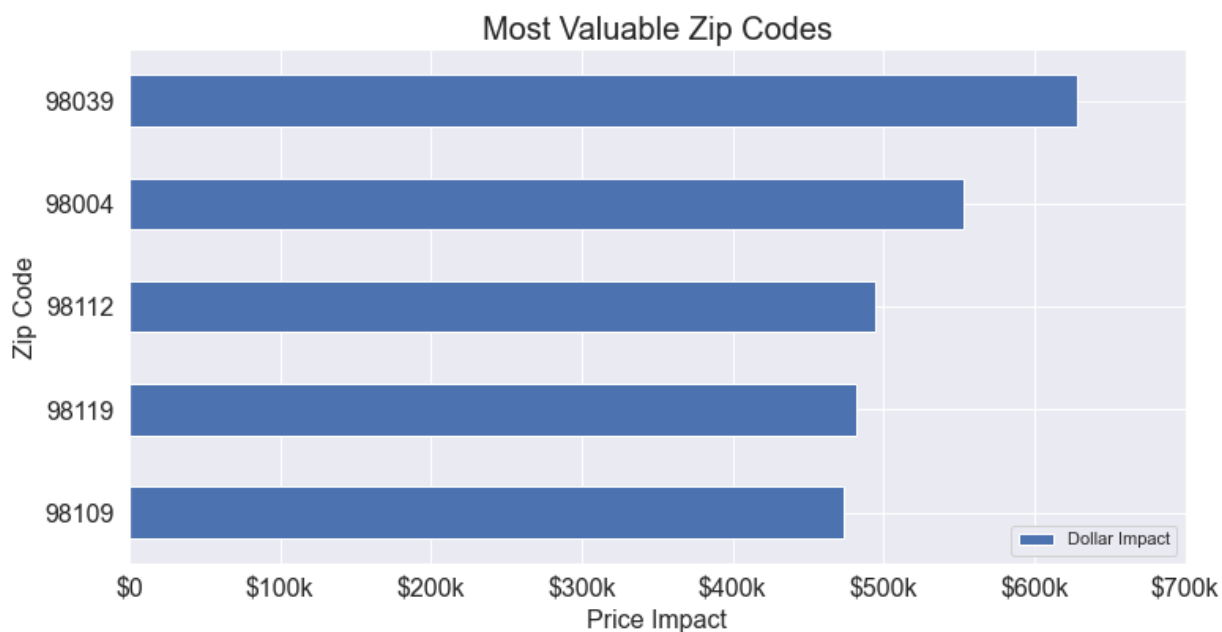
In [735]:

```python
import pandas as pd
import matplotlib.pyplot as plt

data = zip_df_top5
data.plot(x="Zip Code", y="Dollar Impact", kind="barh",figsize=(12, 6))

# plt.legend(([["Film Budget ($)", "Box Office Revenue ($)"]))
ax = plt.gca()
ax.set_xticks([0, 100000, 200000, 300000, 400000, 500000, 600000, 700000])
ax.set_xticklabels(['$0', '$100k', '$200k', '$300k', '$400k', '$500k', '$600
ax.set_yticklabels(zip_df_top5['Zip Code'], fontsize=16)
ax.set_title('Most Valuable Zip Codes', fontsize=20)
plt.ylabel('Zip Code', fontsize = 16)
plt.xlabel('Price Impact',fontsize = 16)

ax.invert_yaxis()

plt.savefig('images/top_zips.png')

plt.show()
```

In [736]:
```python
import pandas as pd
import matplotlib.pyplot as plt

data = zip_df_bottom5
data.plot(x="Zip Code", y="Dollar Impact", kind="barh",figsize=(12, 6), colo

# plt.legend(([["Film Budget ($)", "Box Office Revenue ($)"]))
ax = plt.gca()
ax.set_xticks([0, 100000, 200000, 300000, 400000, 500000, 600000, 700000])
ax.set_xticklabels(['$0', '$100k', '$200k', '$300k', '$400k', '$500k', '$600
ax.set_yticklabels(zip_df_top5['Zip Code'], fontsize=16)
ax.set_title('Least Valuable Zip Codes', fontsize=20)
plt.ylabel('Zip Code', fontsize = 16)
plt.xlabel('Price Impact',fontsize = 16)

ax.invert_yaxis()

plt.savefig('images/bottom_zips.png')

plt.show()
```

# Renovation & Condition Improvements

In [737]:
```python
coeff_df[~coeff_df['Variable'].str.contains("zipcode")]
```

Out[737]:

|  | Variable | Dollar Impact |
|---|---|---|
| 17 | C(waterfront)[T.1.0] | 337780.3300000000 |
| 35 | C(condition)[T.5] | 174135.4700000000 |
| 41 | C(condition)[T.4] | 145846.9700000000 |
| 45 | C(condition)[T.3] | 128815.1600000000 |
| 53 | C(condition)[T.2] | 87360.0300000000 |
| 63 | C(renovated)[T.1] | 42260.4300000000 |
| 69 | C(floors)[T.2.0] | 11448.9400000000 |
| 73 | C(floors)[T.1.5] | 9218.9100000000 |
| 75 | C(floors)[T.2.5] | 1225.8600000000 |
| 76 | sqft_living | 159.4800000000 |
| 78 | sqft_lot | 3.4500000000 |
| 79 | age | -168.9200000000 |
| 81 | C(floors)[T.3.0] | -22470.5500000000 |
| 82 | C(has_basement)[T.1] | -23634.1400000000 |
| 83 | C(floors)[T.3.5] | -40433.6500000000 |
| 84 | Intercept | -188190.8000000000 |

In [738]:
```python
labels = ['Reonvated Bonus', 'Condition 1 to 2', 'Condition 2 to 3', 'Condit
labels
```

Out[738]:
```
['Reonvated Bonus',
 'Condition 1 to 2',
 'Condition 2 to 3',
 'Condition 3 to 4',
 'Condition 4 to 5']
```

In [739]:
```python
renovated = 42260.43
onetwo = 87360.03
twothree = 128815.16 - 87360.03
threefour = 145846.97 - 128815.16
fourfive = 174135.47 - 145846.97

print(twothree)
print(threefour)
print(fourfive)
```
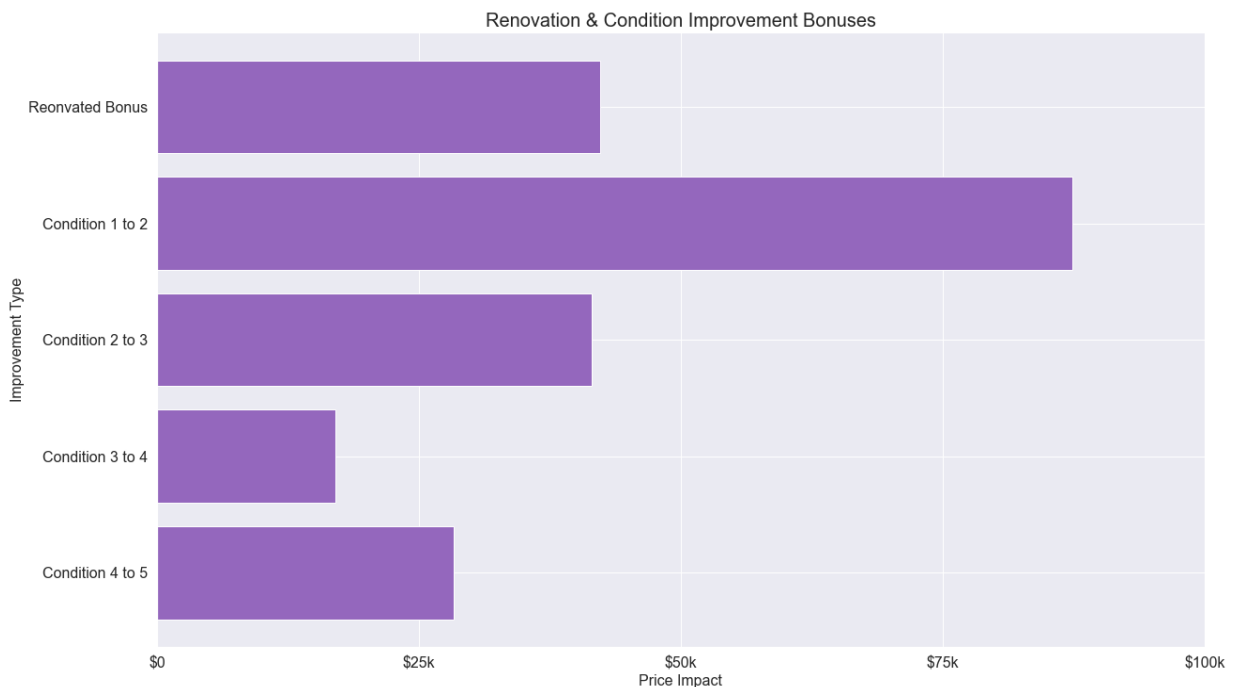
```
41455.130000000005
17031.809999999998
28288.5
```

```
In [740]:    1  values = [42260.43, 87360.03, 41455.13, 17031.80, 28288.5]
```

```
In [741]:    1  fig = plt.figure(figsize=(20, 12))
             2  ax = plt.gca()
             3
             4  ax.barh(labels,values, color='tab:purple')
             5
             6  ax.set_title('Renovation & Condition Improvement Bonuses', fontsize=20)
             7  plt.ylabel('Improvement Type', fontsize=16)
             8  plt.xlabel('Price Impact', fontsize=16)
             9
            10  ax.set_xticks([0, 25000, 50000, 75000, 100000])
            11  ax.set_xticklabels(['$0', '$25k', '$50k', '$75k', '$100k'], fontsize=16)
            12  ax.set_yticklabels(labels, fontsize=16)
            13
            14  ax.invert_yaxis()
            15
            16  plt.savefig('images/renovation.png')
            17
            18  plt.show()
```

```
<ipython-input-741-50a986fac75a>:12: UserWarning: FixedFormatter should only be
used together with FixedLocator
  ax.set_yticklabels(labels, fontsize=16)
```



## Home Addition Bonuses

In [742]:

```
1  coeff_df[~coeff_df['Variable'].str.contains("zipcode")]
```

Out[742]:

| | Variable | Dollar Impact |
|---|---|---|
| 17 | C(waterfront)[T.1.0] | 337780.3300000000 |
| 35 | C(condition)[T.5] | 174135.4700000000 |
| 41 | C(condition)[T.4] | 145846.9700000000 |
| 45 | C(condition)[T.3] | 128815.1600000000 |
| 53 | C(condition)[T.2] | 87360.0300000000 |
| 63 | C(renovated)[T.1] | 42260.4300000000 |
| 69 | C(floors)[T.2.0] | 11448.9400000000 |
| 73 | C(floors)[T.1.5] | 9218.9100000000 |
| 75 | C(floors)[T.2.5] | 1225.8600000000 |
| 76 | sqft_living | 159.4800000000 |
| 78 | sqft_lot | 3.4500000000 |
| 79 | age | -168.9200000000 |
| 81 | C(floors)[T.3.0] | -22470.5500000000 |
| 82 | C(has_basement)[T.1] | -23634.1400000000 |
| 83 | C(floors)[T.3.5] | -40433.6500000000 |
| 84 | Intercept | -188190.8000000000 |

In [743]:

```
1  labels = ['500 Sqft', '1000 Sqft', 'Second Floor', 'Finished Basement',]
2
3  labels
```

Out[743]: ['500 Sqft', '1000 Sqft', 'Second Floor', 'Finished Basement']

For this graph, we we will be making recommendations to home owners with a one story house with an unfinished basement. The 'Finished Basement' idea is a bit flawed, since we cannot tell from our data whether or not a basement is finished or unfinished. But we will assume that this add-on will toggle 'has_basement' from 0 to 1 and add the additional living space to 'sqft_living.'

In [744]:
```python
1  f1_df = df[(df['floors'] == 1.00) & (df['has_basement'] == 0)]
2
3  f1_df['sqft_living'].describe()
```

Out[744]:
```
count    3315.0000000000
mean     1284.1710407240
std       402.6186829600
min       370.0000000000
25%       990.0000000000
50%      1240.0000000000
75%      1510.0000000000
max      3430.0000000000
Name: sqft_living, dtype: float64
```

The median sqft_living for a 1 floor house that currently has no basement is 1240 sqft. For the sake of example, we will assume that a basement or second floor addon will be the same sqft as the first floor.

In [745]:
```python
1   sqft = 159.48
2
3   sqft500 = sqft * 500
4   sqft1000 = sqft * 1000
5   second_floor = (sqft * 1240) + 11448.94
6   finished_basement = (sqft * 1240) - 23634.14
7
8   print('sqft500 = ' + str(sqft500))
9   print('sqft1000 = ' + str(sqft1000))
10  print('second_floor = ' + str(second_floor))
11  print('finished_basement = ' + str(finished_basement))
12
13  1240 * sqft
14
```

```
sqft500 = 79740.0
sqft1000 = 159480.0
second_floor = 209204.13999999998
finished_basement = 174121.06
```

Out[745]: 197755.19999999998

In [746]:
```python
1  values = [79740.0, 159480.0, 209204.13999999998, 174121.06]
2
3  values
```

Out[746]: [79740.0, 159480.0, 209204.13999999998, 174121.06]

In [747]:
```python
fig = plt.figure(figsize=(20, 12))
ax = plt.gca()

ax.barh(labels,values, color='tab:green')

ax.set_title('Home Addition Bonuses', fontsize=20)
plt.ylabel('Addition Type', fontsize=16)
plt.xlabel('Price Impact', fontsize=16)

ax.set_xticks([0, 50000, 100000, 150000, 200000, 250000])
ax.set_xticklabels(['$0', '$50k', '$100k', '$150k', '$200k', '$250k'], fonts
ax.set_yticklabels(labels, fontsize=16)

ax.invert_yaxis()

plt.savefig('images/additions.png')

plt.show()
```

```
<ipython-input-747-8fa3b030f5d6>:12: UserWarning: FixedFormatter should only be
used together with FixedLocator
  ax.set_yticklabels(labels, fontsize=16)
```