

```
#!/usr/bin/env python3
# Save notebook as notebook is pretty large ~
var rto = 500;
console.log("%s: Increase require timeout to " + rto + " seconds");
require.resolve.js.config({waitSeconds: rto});
```

Business Problem

The Zillow Failure

As reported on a November 9, 2021 episode of the Wall Street Journal, "The Journal" podcast, Zillow attempted to transform its business in 2018 by buying up thousands of homes all over the country. By utilizing it's massive amount of search data, Zillow hoped to leverage its advantage. According to the podcast,

Ryan Knutson: So Zillow went huge and bought thousands of homes all over the country. The company would often do some light renovations and turn around and sell them quickly.

Unfortunately, for Zillow, their efforts ultimately failed (the pandemic did not help of course). Zillow has returned to its core efforts.

Will Parker: So there's a debate right now about whether what happened at Zillow is a Zillow problem or an iBuyer problem. And I think what happens with its competitors, these other companies, is really going to tell us a lot about what the answer to that question really is.

Project Goals

The problem Zillow experienced offer an opportunity for competitors in the market. In this project I will use Linear Regression techniques identify key areas where a real estate competitor could make improvements to tranches of homes in an effort to reliably secure profit. I expect a variety of factors to play a role in the analysis.

Data Understanding

The give Dataset is for King County, Washington. According to Wikipedia:

King County is located in the U.S. state of Washington. The population was 2,269,675 in the 2020 census, making it the most populous county in Washington, and the 12th-most populous in the United States. The county seat is Seattle, also the state's most populous city.

King County has a mixture of the urban Seattle area and surrounding Suburbs. Given the size and population I anticipate it will serve as an excellent case study for further modeling projects.

```
In [23]: import pandas as pd
import matplotlib.pyplot as plt

In [24]: pd.reset_option('display.float_format')

In [42]: #Importing standard packages for initial data analysis
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
data = pd.read_csv('data/kc_house_data.csv')

pd.set_option('display.float_format', lambda x: '%.5f' % x)

In [26]: data.info()

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

Column Names

The column names are are included in a markdown file that I will make into a readable dataframe.

```
In [27]: # f = open('r')
# f.readlines()
with open('data/column_name.md') as file:
    lines = file.readlines()
    lines = [line.strip() for line in lines]
    #reading a list of column names
    col_list = []
    #adding other information into an appendix list
    #including the column names and description in list
    for line in lines[1:]:
        if "==" in line:
            col_list.append(line.split("=="))
        else:
            col_list.append(line)

In [28]: for idx, col in enumerate(col_list):
    if len(col)>2:
        col_list[idx-1].append(col)

In [29]: for col in col_list:
    if len(col)>3:
        col_list.remove(col)

In [30]: pd.set_option('max_colwidth', 200)
col_names = pd.DataFrame(col_list[1:], columns=['Name', 'Description', 'Notes'])

In [31]: col_names

Out[31]:
```

	Name	Description	Notes
0	* price	Date house was sold	None
1	* price	Sale price (prediction target)	None
2	bedr ocms	* Number of bedrooms	None
3	bathr ocms	* Number of bathrooms	None
4	sqftliv ∈ g	* Square footage of living space in the home	None
5	sqftlot	* Square footage of the lot	None
6		* [x] Number of floors (levels) in house	None
7	waterfront	* Whether the house is on a waterfront	* Includes Duwamish, Elliott Bay, Puget Sound, Lake Union, Ship Canal, Lake Washington, Lake Sammamish, other lake, and rivers/creeks/waterfronts
8		* view Quality of view from the house	* Includes views of Mt. Rainier, Olympics, Cascades, Territorial, Seattle Skyline, Puget Sound, Lake Washington, Lake Sammamish, small lake / river / creek, and other
9	condition	* How good the overall condition of the house is. Related to maintenance of house.	* See the [King County Assessor Website] (https://info.kingcounty.gov/assessor/realty/Glossary.aspx?type=r) for further explanation of each condition code
10		* Overall grade of the house. Related to the construction and design of the house.	* See the [King County Assessor Website] (https://info.kingcounty.gov/assessor/realty/Glossary.aspx?type=r) for further explanation of each building grade code
11	sqftlbse	* Square footage of house apart from basement	None
12	sqftbasement	* Square footage of the basement	None
13	yr_built <	* Year when house was built	None
14	yr_renovated	* Year when house was renovated	None
15	zipcode	* zip code used by the United States Postal Service	None
16		* lat Latitude coordinate	None
17		* long Longitude coordinate	None
18	sqftliv ∈ g15	* The square footage of interior housing living space for the nearest 15 neighbors	None
19	sqftlot15	* The square footage of the land lots of the nearest 15 neighbors	None

This dataframe gives a clear explanation of each of the features that I will explore below individually.

Data Exploration

Price

```
In [32]: data['price'].describe()

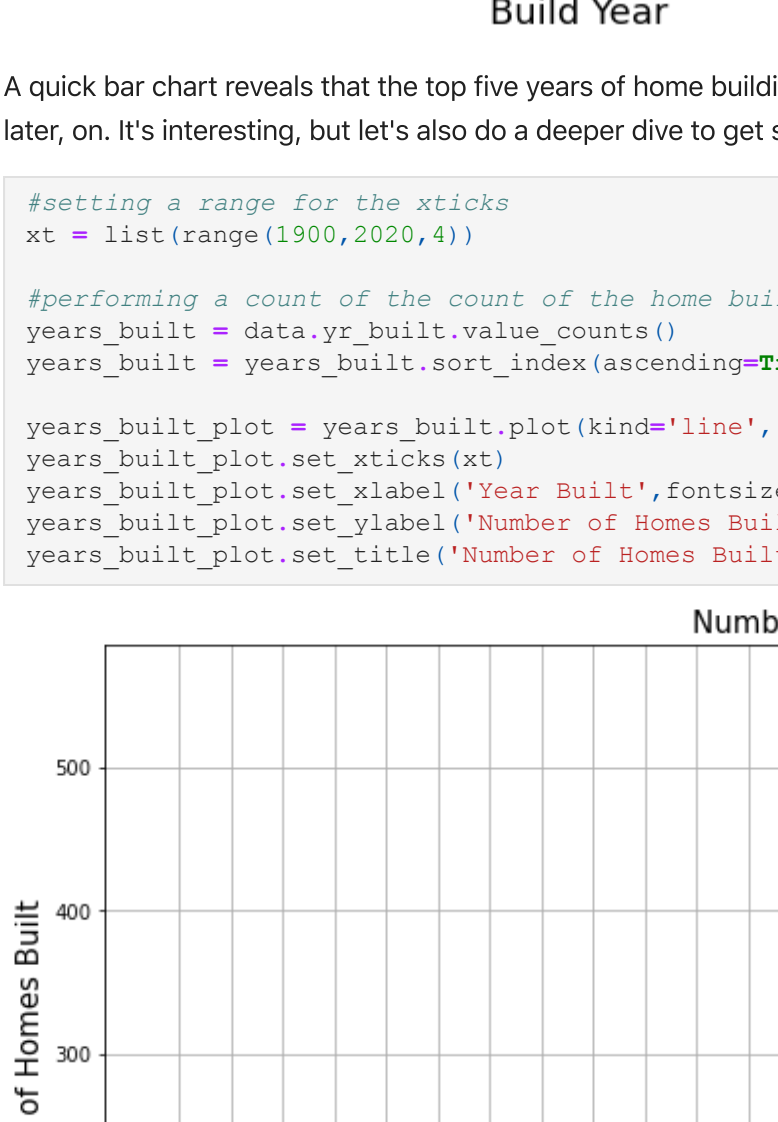
Out[32]: count    21597.000000
mean    540296.57351
std     367368.14010
min      78000.00000
25%     220000.00000
50%     450000.00000
75%     645000.00000
max     7700000.00000
Name: price, dtype: float64

In [33]: average_price = data['price'].mean()
median_price = data['price'].median()

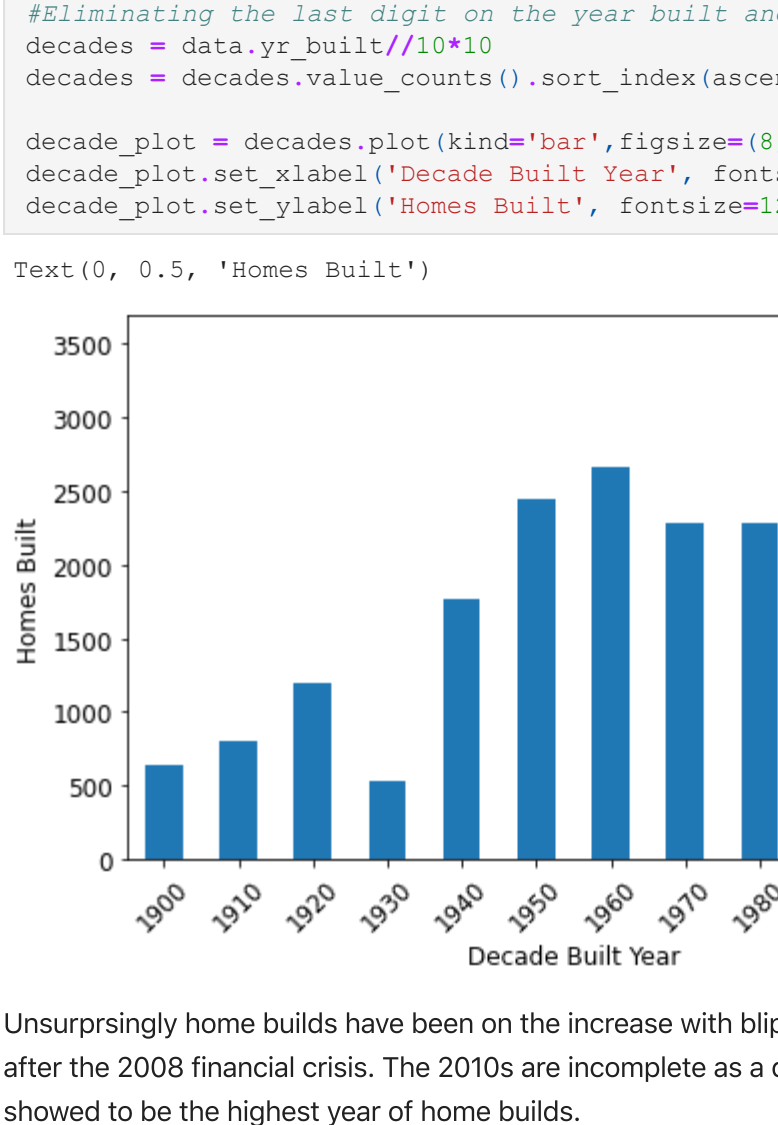
print(f'The average home price in the data set is {average_price}, the median price is {median_price}')

The average home price in the data set is 540296.5735055795, the median price is 450000.0

In [34]: sns.histplot(data=data, x='price')
plt.ticklabel_format(style='plain')
plt.xticks(rotation=75);
```



```
In [35]: sns.displot(data=data, x='price')
plt.ticklabel_format(style='plain')
plt.xticks(rotation=75);
```



The price is the target variable for the Linear Regression Analysis. It is a highly right skewed dataset in this regard.

Year Built Analysis

```
In [36]: data['yr_built'].sort_values()

Out[36]: 14069    1900
14783    1900
10973    1900
115      1900
4693     1900

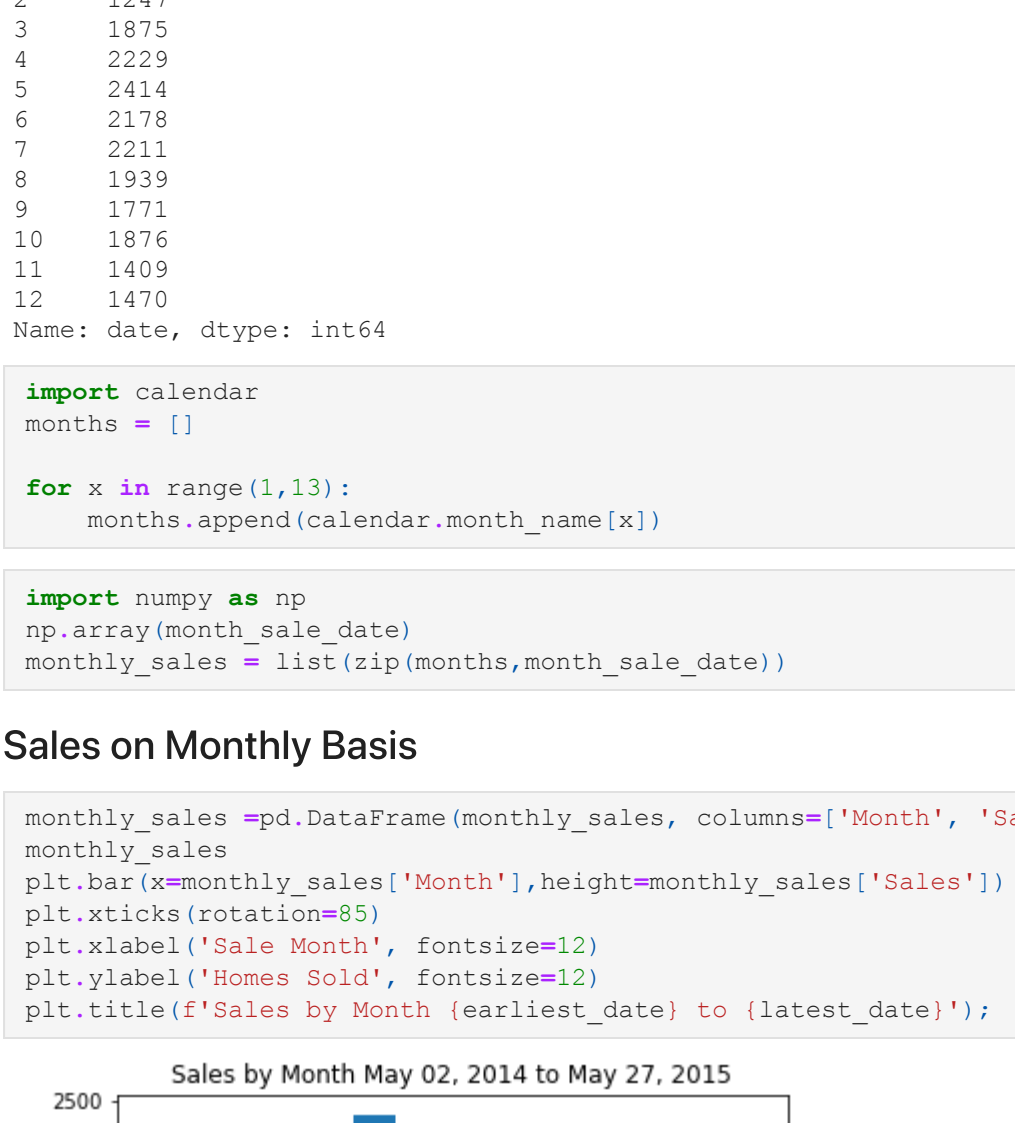
20235     201
7519     2015
14911    2015
4150     2015
19789    2015
Name: yr_built, Length: 21597, dtype: int64
```

Top Built Years

The data contains both a "date" column and "year" column. The data set starts in 1900 and extends to 2015. So we have a mixture of old homes and new homes. Let's try to get a better sense of what that split.

```
In [37]: #Sorting the year built to show the top years built with value counts.
top_home_built_yrs = data['yr_built'].value_counts(sort=True).head(20)

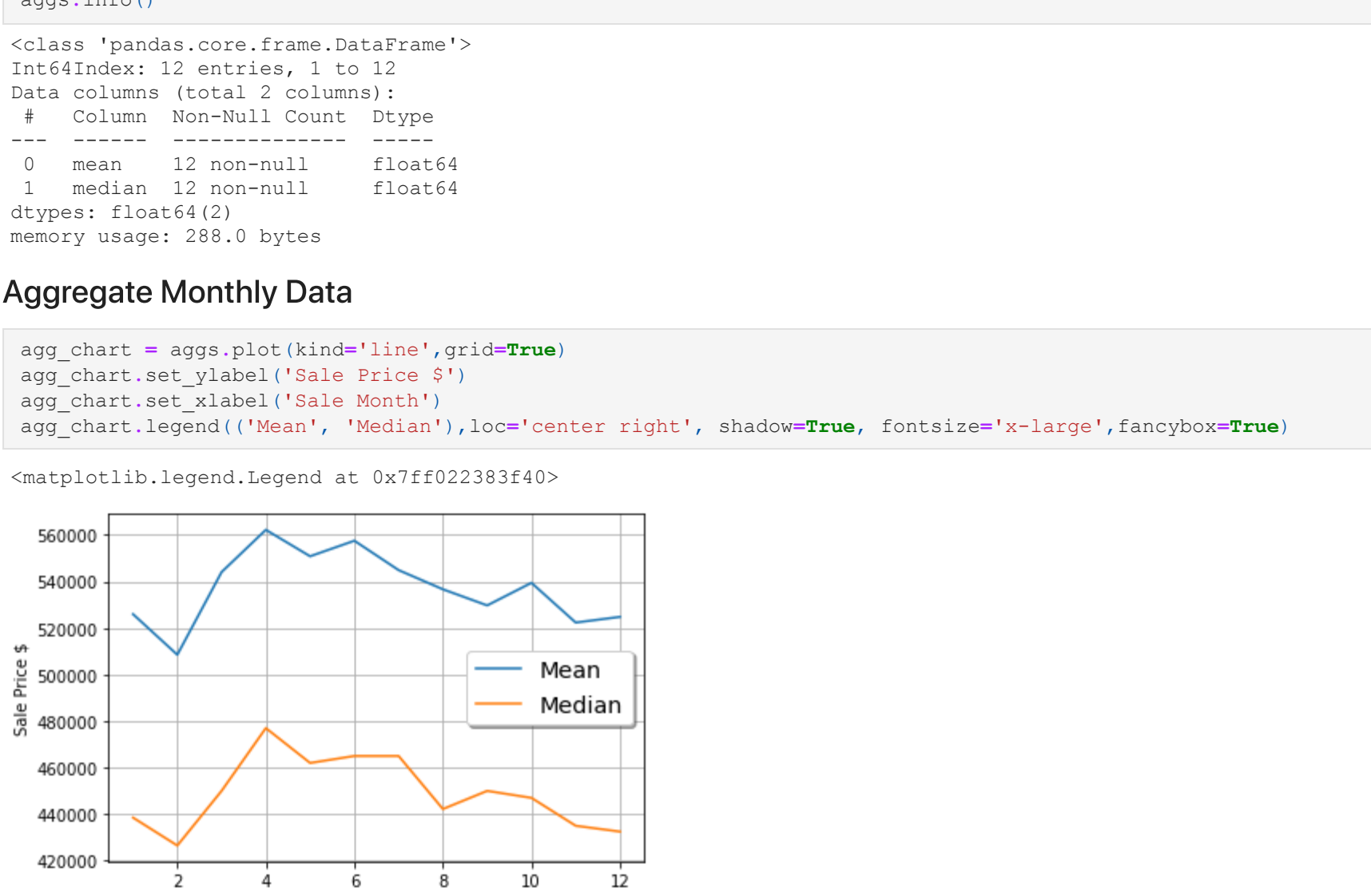
#displaying a bar chart of the years by number of homes built
top_built_yrs = top_home_built_yrs.plot(kind='bar', rot=45, figsize=(8,5), grid=True, linewidth=3)
top_built_yrs.set_xlabel('Build Year', fontsize=18)
top_built_yrs.set_ylabel('Number of Homes', fontsize=18);
```



A quick bar chart reveals that the top five years of home building occurred recently (the past 20 years). This may or may not be relevant later on, it's interesting, but let's also do a deeper dive to get some information on decade by decade home building maybe.

```
In [38]: #Sorting the last digit on the year built and multiplying by ten
decades = data.yr_built//10*10
decades = decades.value_counts().sort_index(ascending=True)

years_built_plot = years_built.plot(kind='line', figsize=(15,8), rot=45, grid=True, linewidth=3)
years_built_plot.set_xticks(xt)
years_built_plot.set_xlabel('Year Built', fontsize=15)
years_built_plot.set_ylabel('Number of Homes Built', fontsize=15)
years_built_plot.set_title('Number of Homes Built by Year', fontsize=15);
```



Unsurprisingly home builds have been on the increase with blips of slowdowns - less homes were built during the Great Depression, and after the 2008 financial crisis. The 2010s are incomplete as a dataset goes, so look to be a declining decade on the charts, but 2014 showed to be the highest year of home builds.

This isn't just a sign of when homes were built, this is a data set of the build date of homes that were sold. So Initially there is a sense that more recent homes are perhaps more desirable.

```
In [40]: data.yr_built.value_counts()[20:2015]

Out[40]: 38
```

The home built year of 38 in 2015 suggests I may consider throwing out the 2015 data, or at least accounting for the year built as incomplete in some way from the end of the data set, I may think of treating 'year built' as categorical. Considering the target of SalePrice, potentially year built will be an important coefficient in this determination.

Date Sold Information

According to the outside chart breakdown, the 'date' column indicates, the year sold information is in the 'date' column.

```
In [318]: from datetime import datetime

In [506]: #Converting date column to datetime column
data['date'] = pd.to_datetime(data['date'])

In [43]: earliest_date = datetime.strptime(data['date'].min(), '%B %d, %Y')
latest_date = datetime.strptime(data['date'].max(), '%B %d, %Y')

print(f'The earliest sale by date is {earliest_date}, the most recent sale date is {latest_date}.')

The earliest sale by date is May 02, 2014, the most recent sale date is May 27, 2015.

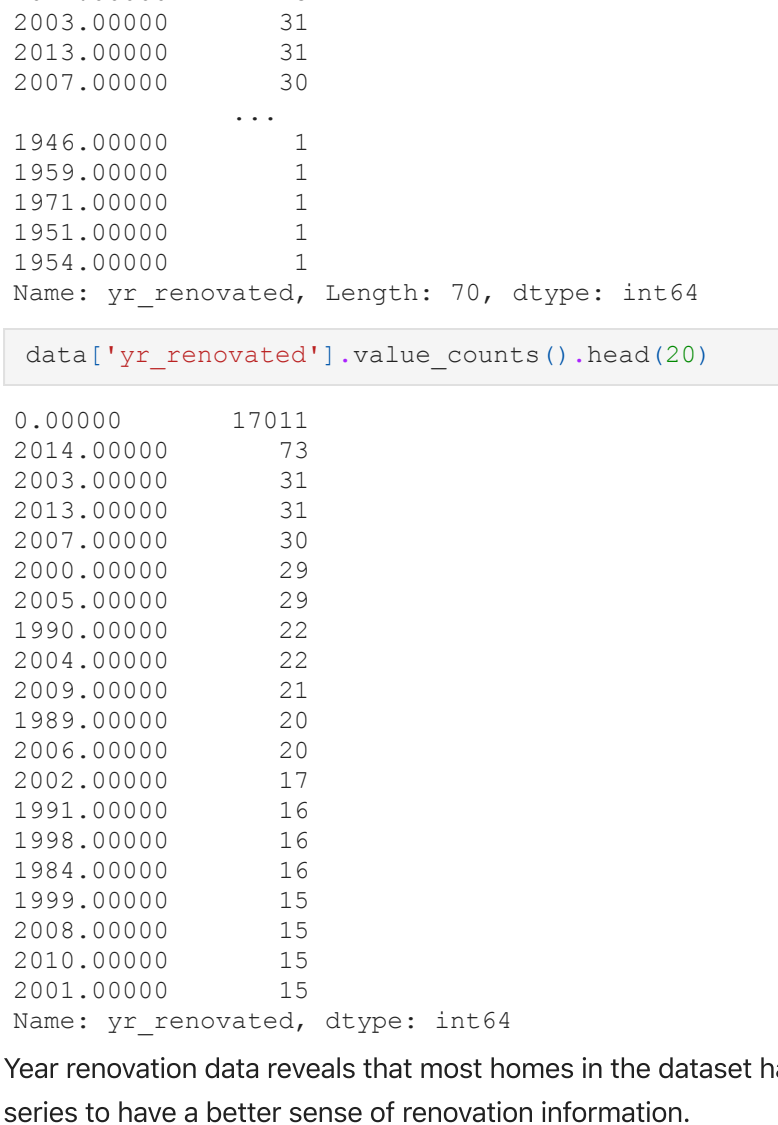
In [44]: #Top 10 correlations
corr = data.corr()[['price']].sort_values(ascending=False)[1:-1]
corr

Out[44]: sqft_living    0.70192
grade              0.66795
sqft_above         0.60337
sqft_living15     0.58324
bathrooms         0.52591
view              0.32573
bedrooms         0.30879
Name: price, dtype: float64
```

This is a very cursory look at the correlations in the dataset. I will further look at multicollinearities later in the notebook. But at first glance the best indicators look to be sqft, "grade", and number of bathrooms. Of course, this is uncleaned data, that hasn't been normalized, or investigated. I will want to separate categoricals and continuous, remove multicorrelations.

```
In [45]: corr.plot(kind='bar', rot=75)

Out[45]: <AxesSubplot>
```



```
In [46]: month_sales = data.date.dt.month.value_counts().sort_index(ascending=True)

In [47]: month_sale_date = data.date.dt.month.value_counts().sort_index()
month_sale_date

Out[47]: 1    978
2    1247
3    1875
4    1229
5    2414
6    2178
7    1211
8    1939
9    1791
10   1874
11   1409
12   1470
Name: date, dtype: int64
```

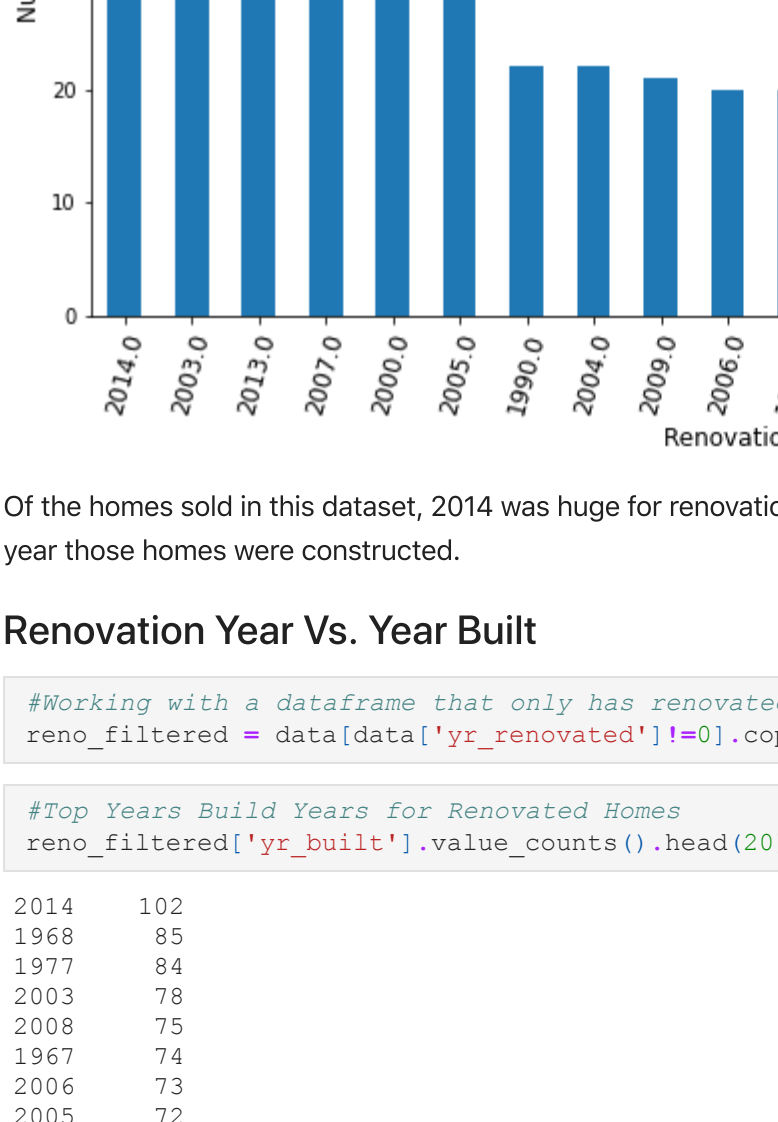
```
In [48]: import calendar
months = []

for x in range(1,13):
    months.append(calendar.month_name[x])

In [49]: import numpy as np
np.array(month_sale_date)
month_sales = list(zip(months, month_sale_date))
```

Sales on Monthly Basis

```
In [50]: monthly_sales = mpd.DataFrame(monthly_sales, columns=['Month', 'Sales'])
plt.bar(monthly_sales['Month'], height=monthly_sales['Sales'])
plt.xticks(rotation=90)
plt.xlabel('Sale Month', fontsize=12)
plt.ylabel('Homes Sold', fontsize=12)
plt.title('Sales by Month (earliest date to latest date)');
```



Homes sales in this area follow a predictable monthly pattern, the sales transactions were lowest in the January/February, tick up in May and taper off in the fall months. This may factor into the model, but so far this is merely looking at when the homes are sold with factoring into price. Let's see if we can do a simple look at sale price vs month.

```
In [51]: import numpy as np
data['month'] = data['date'].map(lambda x: x.month)
agg = data.groupby(data.month)['price'].agg([np.mean, np.median])

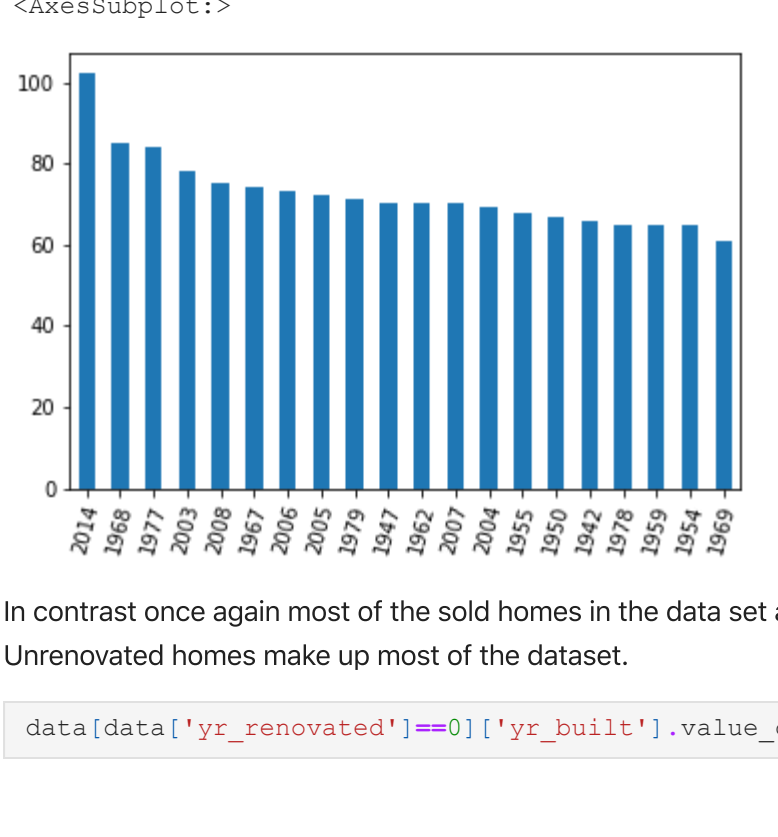
In [52]: aggs.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 12 entries, 1 to 12
Data columns (total 2 columns):
 #   Column      Non-Null Count  Dtype  ---
 0   mean       12 non-null      float64
 1   median     12 non-null      float64
dtypes: float64(2)
memory usage: 288.0 bytes
```

Aggregate Monthly Data

```
In [53]: agg_chart = aggs.plot(kind='line', grid=True)
agg_chart.set_ylabel('Sale Price $')
agg_chart.set_xlabel('Sale Month')
agg_chart.legend(['Mean', 'Median'], loc='center right', shadow=True, fontsize='x-large', fancybox=True)

Out[53]: <matplotlib.legend.Legend at 0x7ff02338f40>
```



Mean and Median sale price follow a similar but not exact pattern as well. Sale Price does seem to peak in April, and gradually lowers until December. Not only are fewer homes sold in the colder months, but the average price and the median price declines as well. Perhaps sellers lower the price as the year goes on, or the pricier homes are removed/bought up?

```
In [54]: monthly_sales

Out[54]:
```

Month	Sales
0	January 978
1	February 1247
2	March 1875
3	April 2229
4	May 2414
5	June 2178
6	July 1211
7	August 1939
8	September 1771
9	October 1876
10	November 1409
11	December 1470

```
In [55]: sns.regplot(x=monthly_sales.Sales, y=aggs['median'])
plt.xlabel('Monthly Sales')
plt.ylabel('Median Monthly Home Price')
plt.title('Monthly Sales Count Vs. Median Sale Price');
```


The above chart indicates a some linear relationship between Monthly Sales and the Median Sale Price. Perhaps this will be helpful later on for the regression, but, as sales volume is higher, sale price indeed seems to be higher. Let's see if we can look at this on not just a monthly basis.

```
In [56]: sns.regplot(data=data, x='month', y='price');
```


This scatterplot displays some relationship of month to sale price. The same sort of relationship with some outliers is show. This will be further explored later in the notebook - the potential relationship between various features in the dataset and the sale price on a monthly price.

Renovation Year

```
In [57]: data['yr_renovated'].describe()

Out[57]: count    17755.000000
mean      83.63678
std       399.94641
min        0.000000
25%        0.000000
50%        0.000000
75%        0.000000
max       2015.000000
Name: yr_renovated, dtype: float64
```

```
In [58]: data['yr_renovated'].value_counts()

Out[58]: 0.000000    17011
1.000000      73
2003.00000    31
2013.00000    31
2005.00000    30
1346.00000     1
1954.00000     1
1991.00000     1
1971.00000     1
1984.00000     1
1954.00000     1
Name: yr_renovated, Length: 70, dtype: int64
```

```
In [59]: data['yr_renovated'].value_counts().head(20)

Out[59]: 0.000000    17011
2014.00000    73
2003.00000    31
2013.00000    31
2007.00000    30
2008.00000    29
2005.00000    29
1990.00000    22
2004.00000    22
2009.00000    20
1987.00000    20
2006.00000    20
2002.00000    17
1994.00000    16
1998.00000    16
1984.00000    16
1999.00000    15
2008.00000    15
2010.00000    15
2001.00000    15
Name: yr_renovated, dtype: int64
```

Year renovation data reveals that most homes in the dataset haven't been renovated, for this initial dive, let's remove them from the series to have a better sense of renovation information.

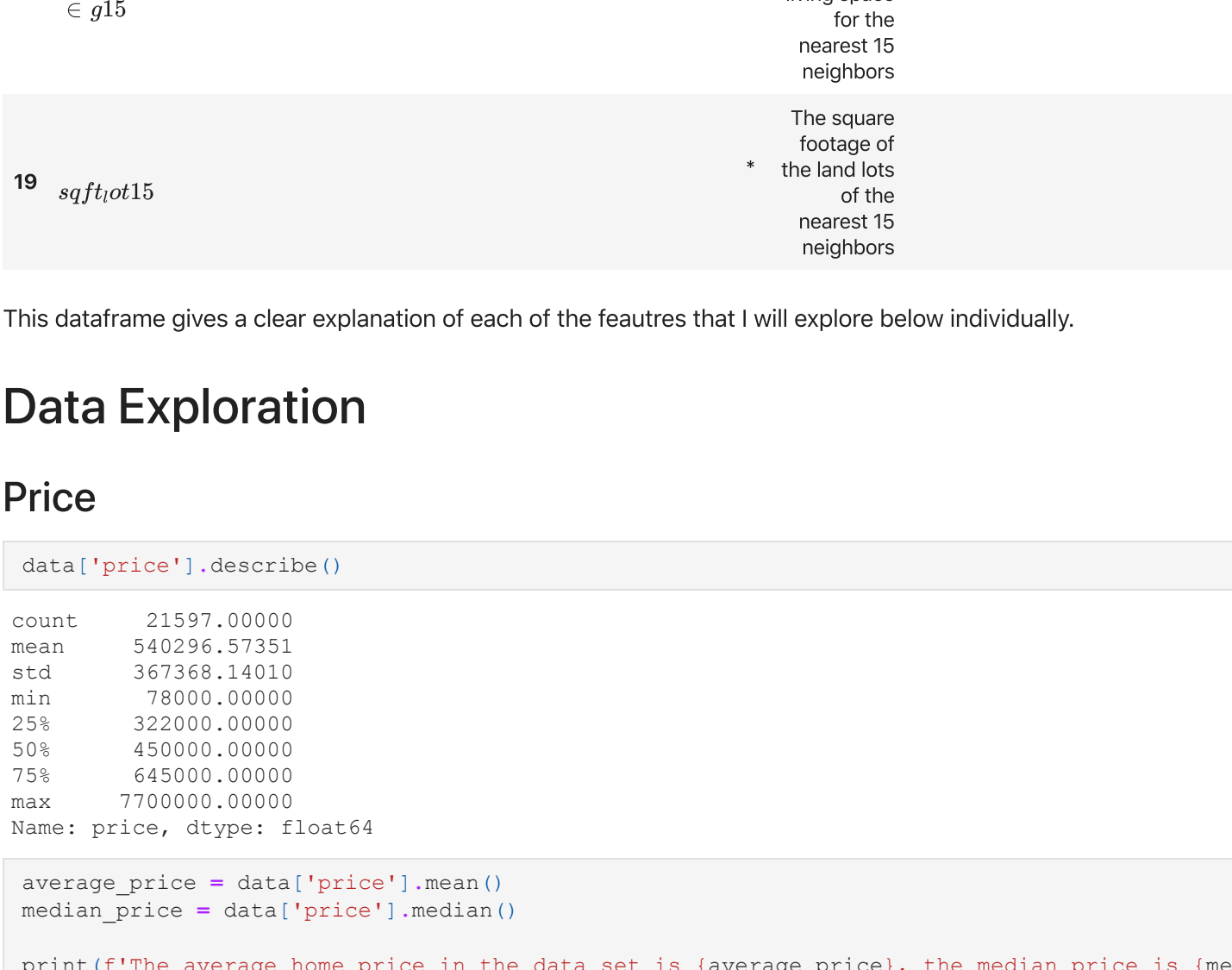
```
In [60]: #Removing where the "year renovated" is 0
renovation_filtered = data.yr_renovated['yr_renovated'!=0].copy()

In [61]: renovation_filtered = pd.to_datetime(renovation_filtered, format='%Y').dt.year

In [62]: top_20_renos = renovation_filtered.value_counts().head(20)
```

```
In [63]: top_20_renos.plot(kind='bar', rot=75, figsize=(12,8))
plt.xticks(fontsize=12);
plt.xlabel('Renovation Year', fontsize=12)
plt.ylabel('Number of Renovations', fontsize=12)
plt.title('Top 20 Renovations Per Year', fontsize=12)
```

```
Out[63]: Text(0.5, 1.0, 'Top 20 Renovations Per Year')
```



Of the homes sold in this dataset, 2014 was huge for renovations. This may not be surprising, what might be interesting to look at is what year houses were constructed.

Renovation Year Vs. Year Built

```
In [64]: #Working with a dataframe that only has renovated homes
reno_filtered = data[data['yr_renovated']!=0].copy()

In [65]: #Top years Build Years for Renovated Homes
reno_filtered['yr_built'].value_counts().head(20)
```

```
Out[65]: 2014    102
1985     85
1977     84
2003     78
2008     75
1967     74
2006     73
2005     72
1979     71
2004     69
1962     70
2007     70
2009     69
1955     68
1950     67
2005     65
1978     65
1944     65
1954     65
1969     61
Name: yr_built, dtype: int64
```

```
In [66]: #Average year for renovated home
average_year = str(reno_filtered['yr_built'].mean())[0:4]

#Median year for renovated home
median_year = str(reno_filtered['yr_built'].median())[0:4]

print(f'The average build year for renovated homes is {average_year}, the median year is {median_year}')

The average build year for renovated homes is 1965, the median year is 1967
```

Perhaps unsurprisingly renovated homes date back to around World War II. Let's look at the average time between build and renovation.

```
In [67]: int(reno_filtered['yr_renovated'].mean())-reno_filtered['yr_built'].mean()

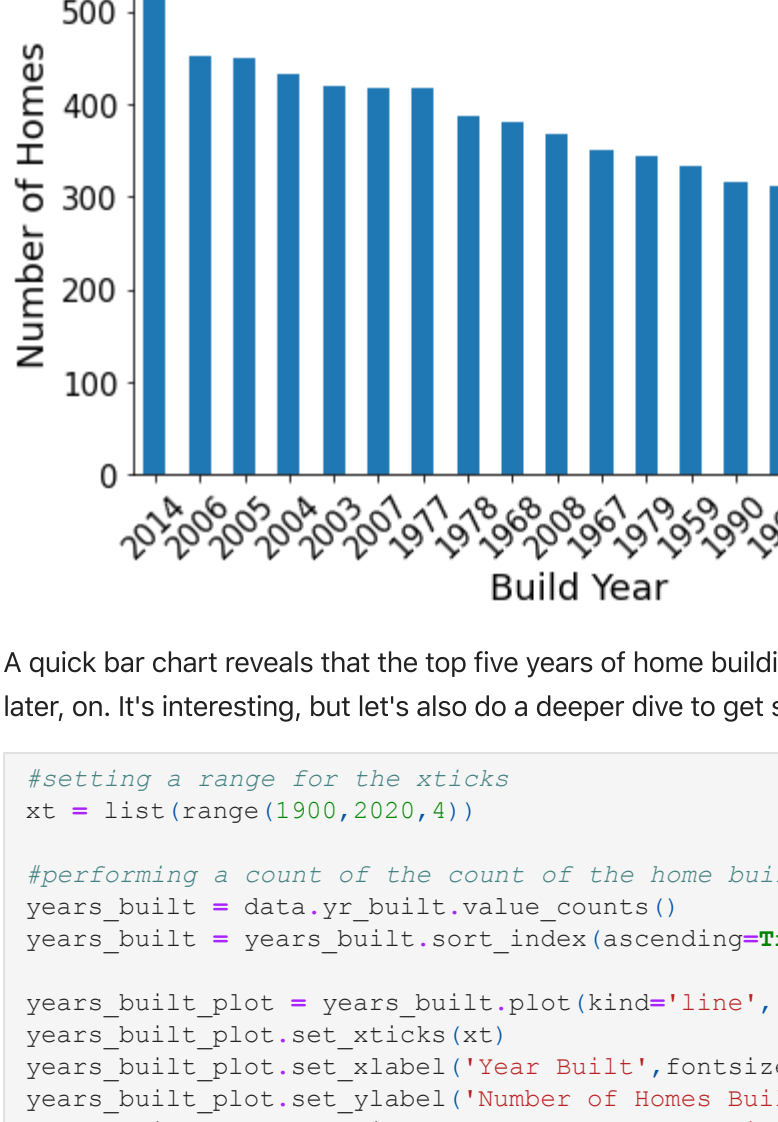
Out[67]: 30
```

On average it was 56 years between construction and renovation for the renovated properties.

The following chart shows us the top build years of the renovated homes that were sold. No real trends seem to be gleaned here, 1950 tops the chart, but then homes built in 1924 are third on the list.

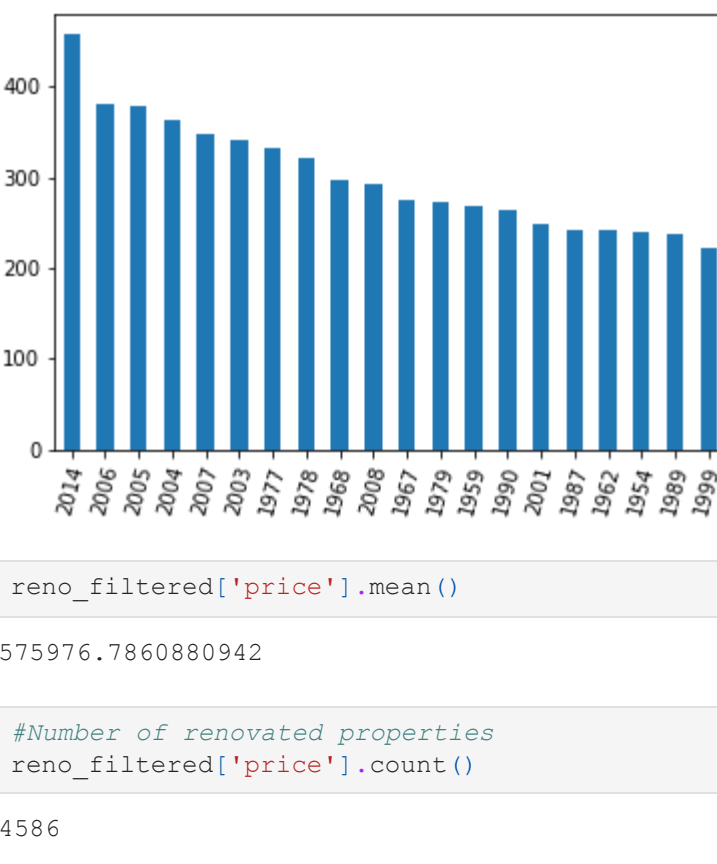
```
In [68]: reno_filtered['yr_built'].value_counts().head(20).plot(kind='bar', rot=75)

Out[68]: <AxesSubplot>
```



In contrast once again most of the sold homes in the data set are from 2014, so recent homes seem to be the top homes on the list. Unrenovated homes make up most of the dataset.

```
In [69]: data[data['yr_renovated']==0]['yr_built'].value_counts().head(20).plot(kind='bar', rot=75);
```

```
In [70]: reno_filtered['price'].mean()
Out[70]: 575976.786080942
```

```
In [71]: #Number of renovated properties
reno_filtered['price'].count()
Out[71]: 4586
```

The average price for renovated home was 768,901. I think I need to be careful reading into this number, because there is a question of why these homes have been selected for renovation.

```
In [72]: #Number of unrenovated properties
data[data['yr_renovated']==0]['price'].count()
Out[72]: 17041
```

The renovated vs. unrenovated properties contrast sharply. Renovated properties are much older, constructed around World War Two. The top properties unrenovated are recent builds - maybe part of large developments, or perhaps new homes are that much more desirable.

```
In [73]: reno_filtered[reno_filtered['yr_built']<=1980]
Out[73]:
```

	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view	...	sqft_above	sqft
1	641410192	2014-12-09	538000.00000	3	2	25000	2570	7242	2.00000	0.00000	0.00000	...	2170
2	5631500400	2015-02-25	180000.00000	2	1	100000	770	10000	1.00000	0.00000	0.00000	...	770
12	114101516	2014-06-28	310000.00000	3	1	100000	1430	19901	1.50000	0.00000	0.00000	...	1430
26	1794500383	2014-06-26	937000.00000	3	1	175000	2450	2691	2.00000	0.00000	0.00000	...	1750
28	5101402488	2014-06-24	438000.00000	3	1	175000	1520	6380	1.00000	0.00000	0.00000	...	790
...
20428	430660360	2015-02-25	500012.00000	4	2	250000	2400	9612	1.00000	0.00000	0.00000	...	1230
20431	3319500628	2015-02-12	356999.00000	3	1	150000	1010	1546	2.00000	nan	0.00000	...	1010
20946	1279000210	2015-03-11	110000.00000	2	1	100000	828	4524	1.00000	0.00000	0.00000	...	828
21027	9276200220	2014-07-17	375000.00000	1	1	100000	720	3166	1.00000	0.00000	0.00000	...	720
21224	7174800094	2015-04-20	525000.00000	1	1	150000	1030	5923	1.00000	nan	0.00000	...	1030

2990 rows x 22 columns

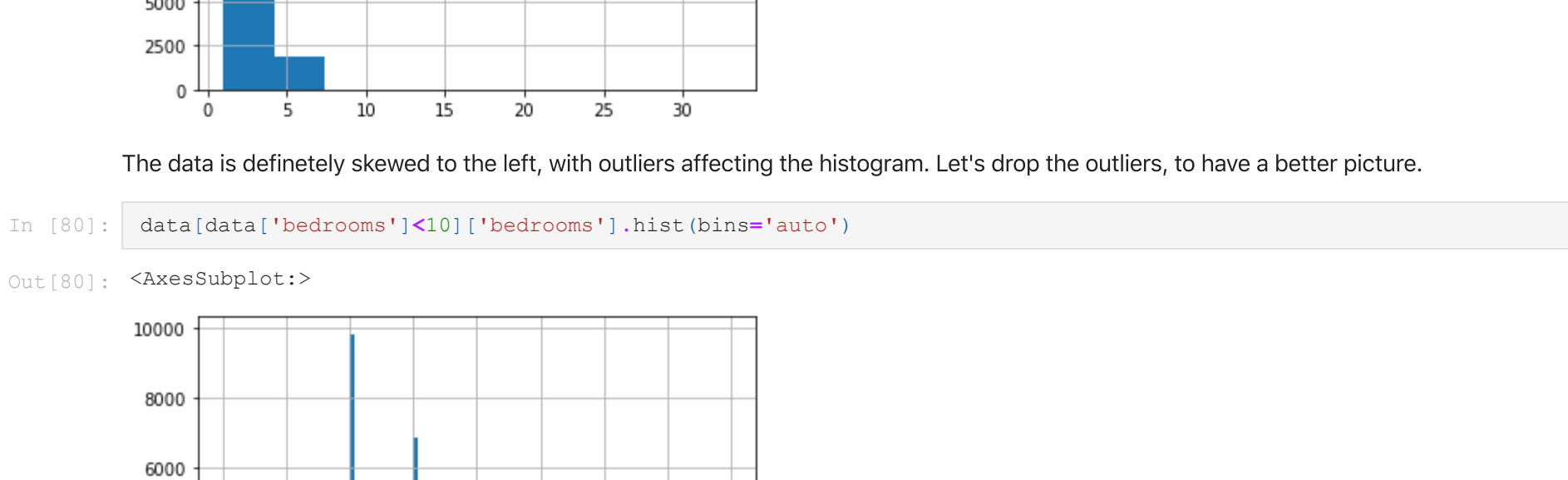
The vast majority of renovated homes were built before 1980. Just for fun right now let's take a look at homes built before 1980 renovated vs. unrenovated.

```
In [74]: reno_filtered[reno_filtered['yr_built']<1980]['price'].mean()
Out[74]: 563035.372242959
```

```
In [75]: #creating a dataframe of unrenovated properties
unrenovated_filtered = data[data['yr_renovated']==0].copy()
In [76]: #filtering unrenovated properties to find the average price before 1980
unrenovated_filtered[unrenovated_filtered['yr_built']<1980]['price'].mean()
Out[76]: 486308.9730359377
```

The plot below shows build year alone doesn't have a strong correlation to an increase in price.

```
In [77]: sns.regplot(data=data, x='yr_built', y='price', line_kws={"color": "red"}, scatter_kws={"alpha":0.3})
plt.ticklabel_format(style='plain')
Out[77]:
```

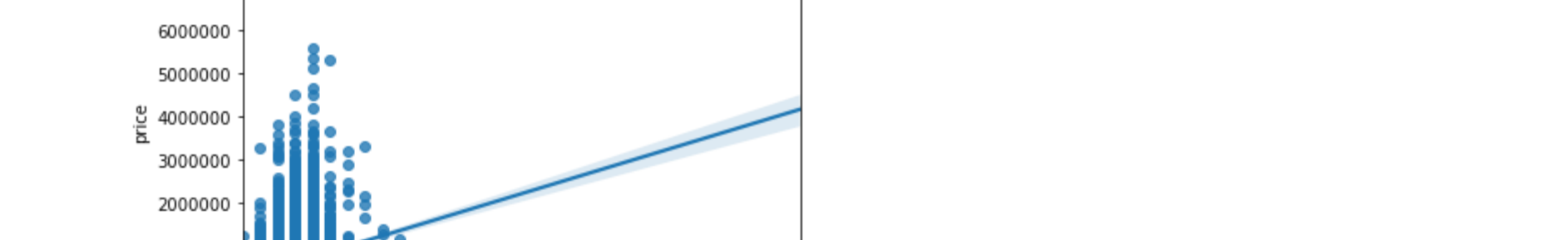


Bedrooms

```
In [78]: #Number of bedrooms per apartment sold
data['bedrooms'].value_counts()
Out[78]:
```

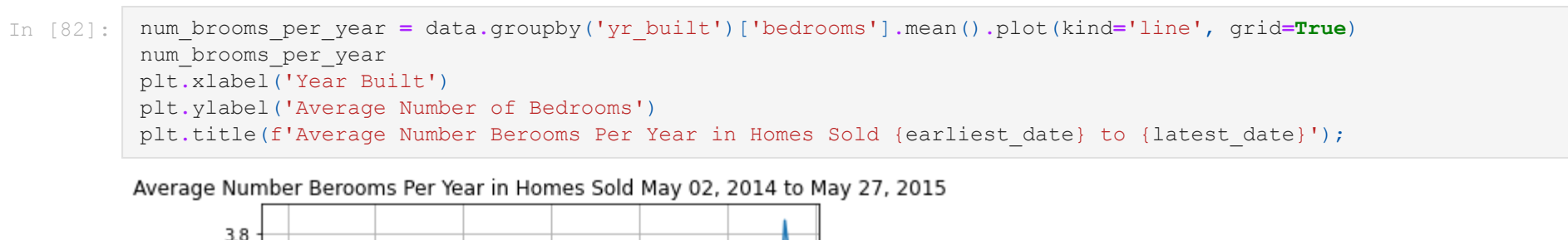
```
3    9824
4    8885
2    2760
5    1601
6    272
1     196
7     36
8     13
9      6
10     3
11     1
33     1
Name: bedrooms, dtype: int64
```

```
In [79]: data['bedrooms'].hist()
Out[79]: <AxesSubplot>
```



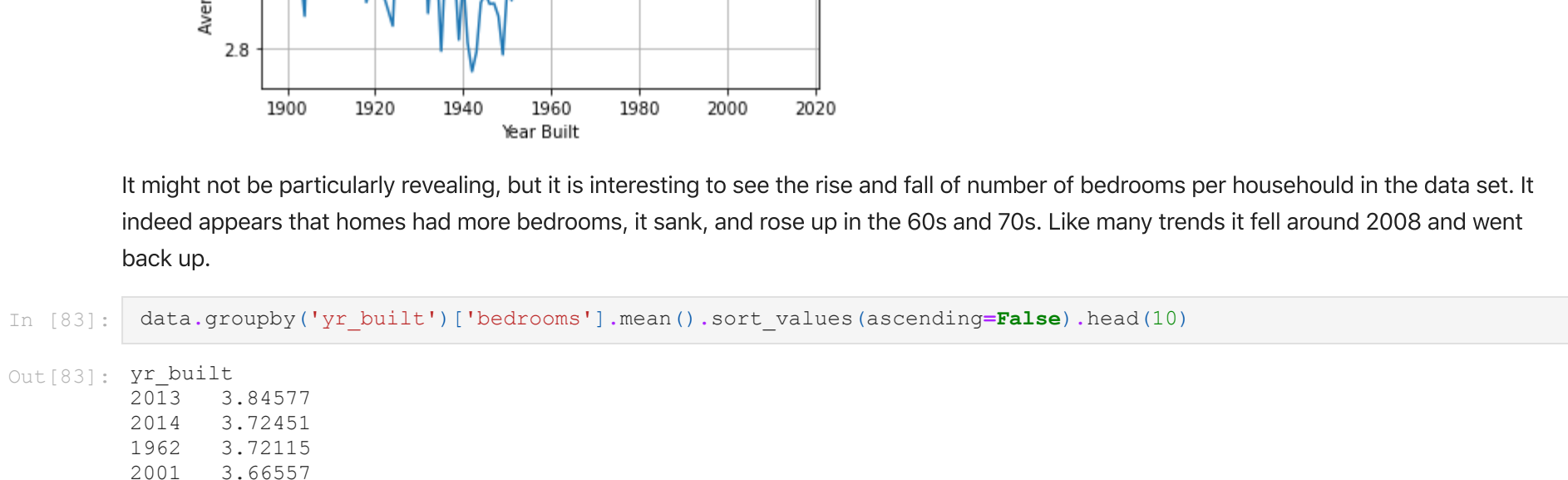
The data is definitely skewed to the left, with outliers affecting the histogram. Let's drop the outliers, to have a better picture.

```
In [80]: data[data['bedrooms']<10]['bedrooms'].hist(bins='auto')
Out[80]: <AxesSubplot>
```



The number of rooms actually has a somewhat normal distribution in this histogram, I believe we are dealing with a fairly normally distributed set of bedrooms in each home.

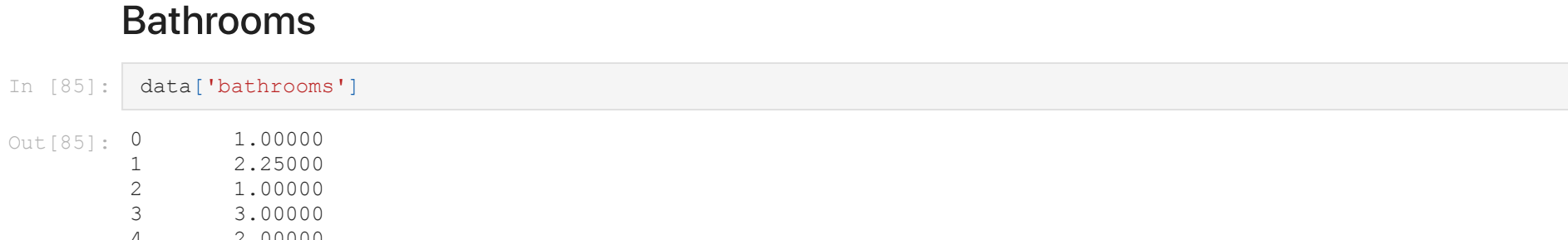
```
In [81]: sns.regplot(data=data, x='bedrooms', y='price')
plt.ticklabel_format(style='plain')
Out[81]:
```



This does seem to show somewhat of linear relationship, as bedrooms go up, price does look to go up a little bit.

```
In [82]: num_brooms_per_year = data.groupby('yr_built')['bedrooms'].mean().plot(kind='line', grid=True)
num_brooms_per_year
plt.xlabel('Year Built')
plt.ylabel('Average Number of Bedrooms')
plt.title('Average Number Bedrooms Per Year in Homes Sold (earliest_date) to (latest_date)');
Out[82]:
```

Average Number Bedrooms Per Year in Homes Sold May 02, 2014 to May 27, 2015



It might not be particularly revealing, but it is interesting to see the rise and fall of number of bedrooms per household in the data set. It indeed appears that homes had more bedrooms, it sank, and rose up in the 60s and 70s. Like many trends it fell around 2008 and went back up.

```
In [83]: data.groupby('yr_built')['bedrooms'].mean().sort_values(ascending=False).head(10)
Out[83]:
```

```
yr_built
2013    3.84577
2014    3.72451
1962    3.72115
2001    3.66557
1963    3.66310
1966    3.65863
2015    3.65789
1963    3.65098
1964    3.63953
1913    3.29793
Name: bedrooms, dtype: float64
```

```
In [84]: data.groupby('yr_built')['bedrooms'].mean()[2009]
Out[84]: 3.139130434782609
```

Bathrooms

```
In [85]: data['bathrooms'].hist()
Out[85]:
```

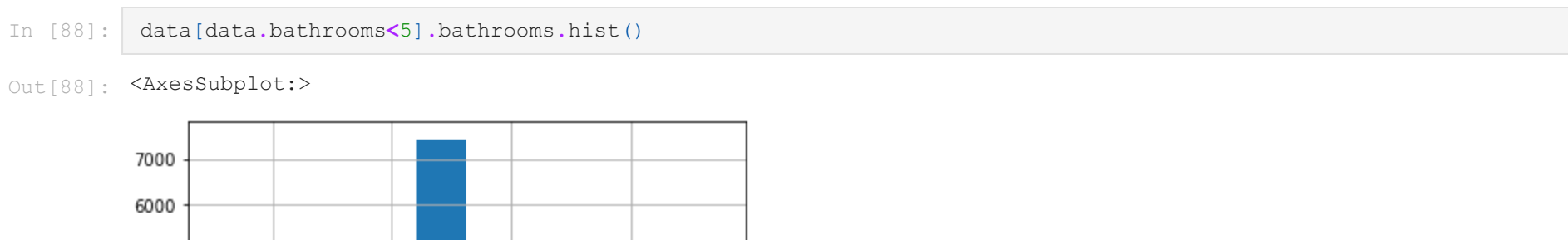
```
0      1.00000
1      2.25000
2      1.00000
3      3.00000
4      2.00000
...
21592  2.50000
21593  2.50000
21594  0.75000
21595  2.50000
21596  0.75000
Name: bathrooms, Length: 21597, dtype: float64
```

```
In [86]: data['bathrooms'].value_counts()
Out[86]:
```

```
1.00000    3851
1.75000    3048
2.25000    2047
2.00000    1930
1.50000    1445
2.75000    1185
3.00000     753
3.50000     731
3.25000     589
3.75000     155
4.00000     136
4.50000     100
4.25000     79
4.75000     71
4.75000     23
5.00000     21
5.25000     13
5.50000     10
1.25000      9
6.00000      6
5.75000      4
6.50000      4
8.00000      2
6.25000      2
6.75000      2
6.50000      2
7.50000      1
7.75000      1
Name: bathrooms, dtype: int64
```

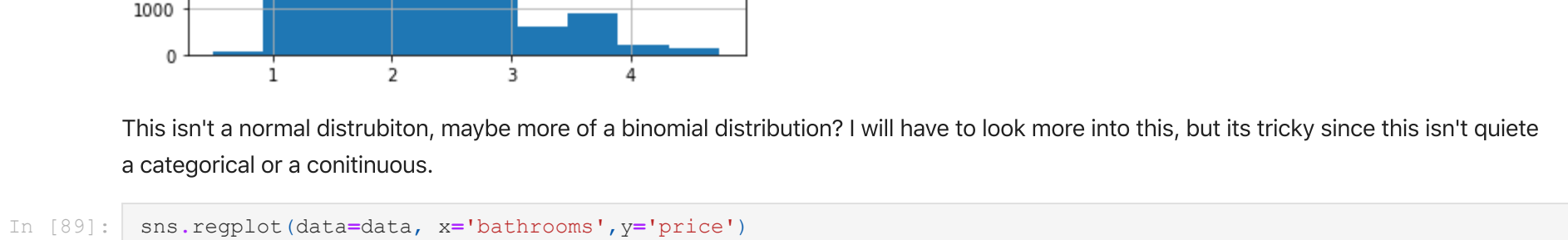
Value counts on the bathrooms reveals that the top bathroom amount is 2.5 bathrooms. Half and quarter bathrooms are a feature in some living situations, so this does seem to be an accurate dataset of some kind. Let's look at a histogram.

```
In [87]: data.bathrooms.hist()
Out[87]: <AxesSubplot>
```



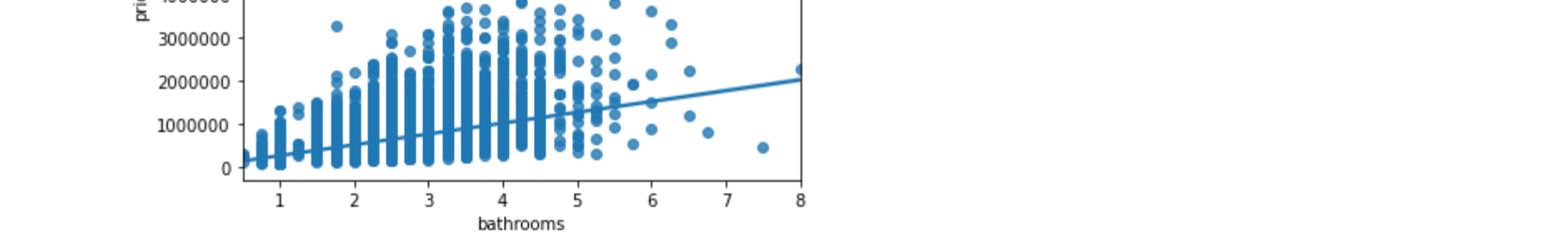
Once again this isn't perfectly normal, but let's also throw out some of the outliers for now.

```
In [88]: data[data.bathrooms<5].bathrooms.hist()
Out[88]: <AxesSubplot>
```



This isn't a normal distribution, maybe more of a binomial distribution? I will have to look more into this, but its tricky since this isn't quite a categorical or a continuous.

```
In [89]: sns.regplot(data=data, x='bathrooms', y='price')
plt.ticklabel_format(style='plain')
Out[89]:
```



This is consistent with the initial correlations we saw previously, there is some linearity in regards to number of bathrooms and the price of a home. This will be a feature to keep in mind, interestingly it has a higher initial correlation than the number of bedrooms. Bathrooms are certainly a necessity, but also an amenity. Maybe we can look at home number of home affects homes with similar number of bedrooms? Is that an interaction?

I'm not really sure what to do on that right now, but for now I am going to look at how bathrooms affect homes of a set number of bedrooms.

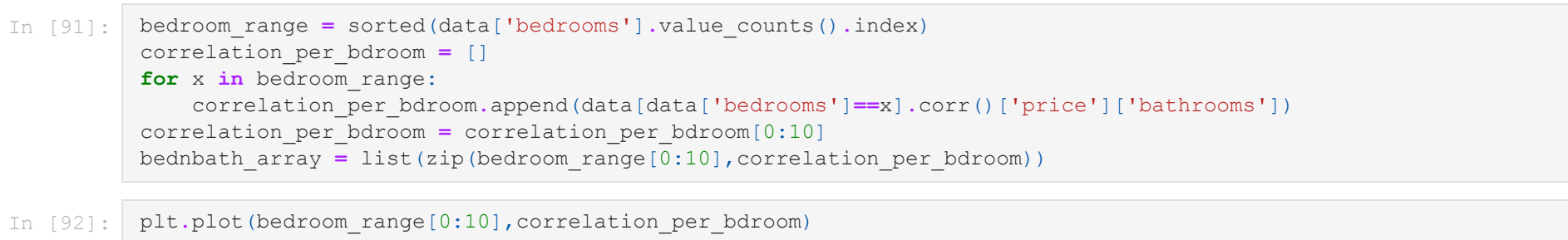
```
In [90]: sns.regplot(data=data, x='bedrooms', y='price')
Out[90]: <AxesSubplot>
```



```
In [91]: bedroom_range = sorted(data['bedrooms'].value_counts().index)
for x in bedroom_range:
    correlation_per_bedroom = data[data['bedrooms']==x].corr()['price']['bathrooms']
    correlation_per_bedroom = correlation_per_bedroom[0:10]
bednbath_array = list(zip(bedroom_range[0:10], correlation_per_bedroom))
Out[91]:
```

```
plt.plot(bedroom_range[0:10], correlation_per_bedroom)
plt.xlabel('Number of Bedrooms')
plt.ylabel('Bathroom Correlation to Sale Price')
plt.title('Correlation of Bathrooms to Sale Price as Bedrooms Increase')
Out[91]:
```

Text(0.5, 1.0, 'Correlation of Bathrooms to Sale Price as Bedrooms Increase')



This is some window into how bathrooms correlation changes as the number of bedrooms change. What I am looking for is to see if the value add of bathroom changes as each bedroom number increases. What I do see is it is a better correlation until about 6 bedrooms, but of course we do not have many homes on those higher end in the data set.

I think there is more exploring to be done in the bathroom area and will return to it later.

sqft_living

The Square Foot Living represents a major feature as it is the actual living space of the homes that were sold. It also has the highest correlation on our simple correlation list.

```
In [93]: data.corr()['price']['sqft_living']
Out[93]: 0.7019173021377597
```

```
In [94]: data['sqft_living'].describe()
Out[94]:
```

```
count    21597.00000
mean      2080.32185
std       918.10613
min       370.00000
25%      1430.00000
50%      1920.00000
75%      2550.00000
max      13540.00000
Name: sqft_living, dtype: float64
```

```
In [95]: data['sqft_living'].value_counts()
Out[95]:
```

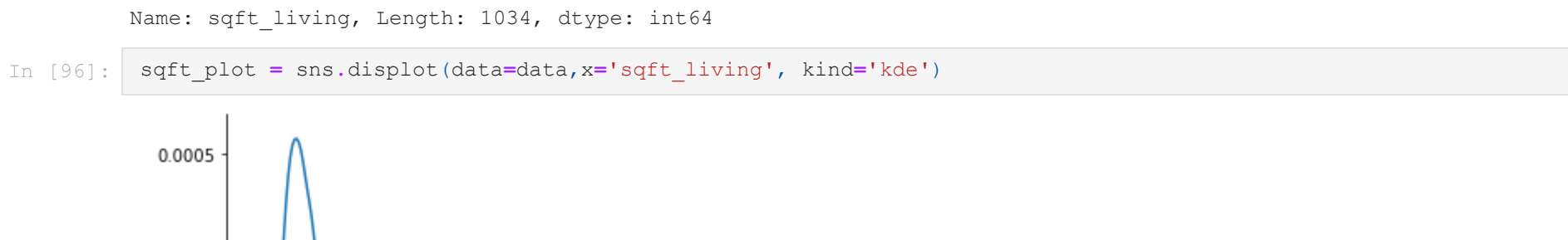
```
1300    138
1400     135
1460     133
1660     129
1010     125
...
4970      1
2905      1
2793      1
4810      1
1975      1
Name: sqft_living, Length: 1034, dtype: int64
```

```
In [96]: sqft_plot = sns.distplot(data=data, x=sqft_living, kind='kde')
Out[96]:
```



We have a somewhat normal distribution of homes, with outliers affecting the information.

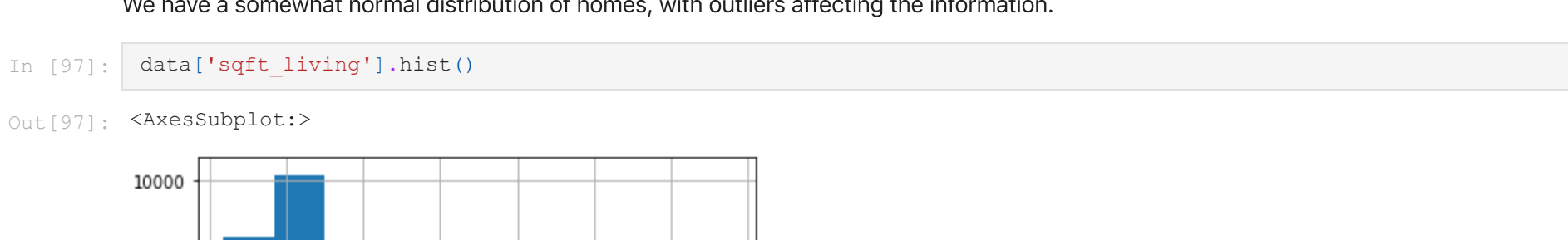
```
In [97]: data['sqft_living'].hist()
Out[97]: <AxesSubplot>
```



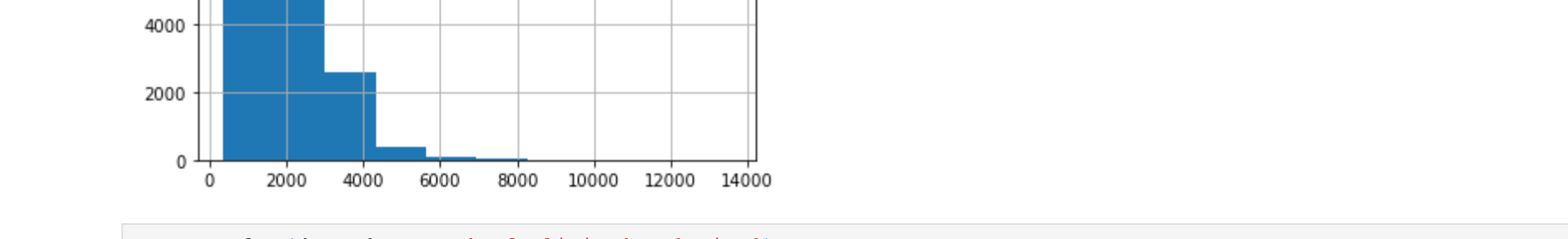
```
In [98]: sns.regplot(data=data, x='sqft_living', y='price')
Out[98]: <AxesSubplot>
```



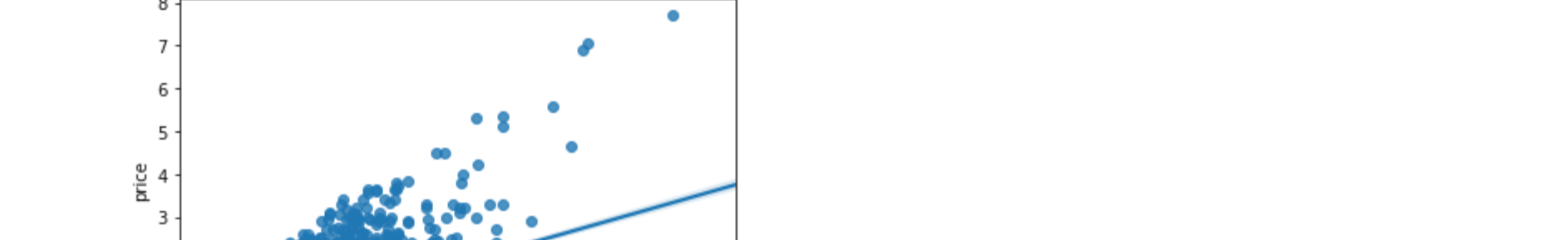
```
In [99]: lower = data.sqft_living.quantile(.05)
upper = data.sqft_living.quantile(.95)
sqft_filtered = data[(data['sqft_living']>lower) & (data['sqft_living']<upper)]
sns.distplot(data=sqft_filtered, x='sqft_living', kind='kde')
Out[99]:
```



```
In [100]: sqft_filtered['sqft_living'].hist()
Out[100]:
```



```
In [101]: sns.regplot(data=sqft_filtered, x='sqft_living', y='price', line_kws={"color": "red"}, scatter_kws={"alpha":0.3})
Out[101]: <AxesSubplot>
```



```
In [102]: sqft_filtered.corr()['price']['sqft_living']
Out[102]: 0.5785532312360123
```

Dropping the outliers of the dataset showed less of correlation with price. Perhaps square footage is a bigger drive of price at the tail of the dataset.

Square Foot Lot

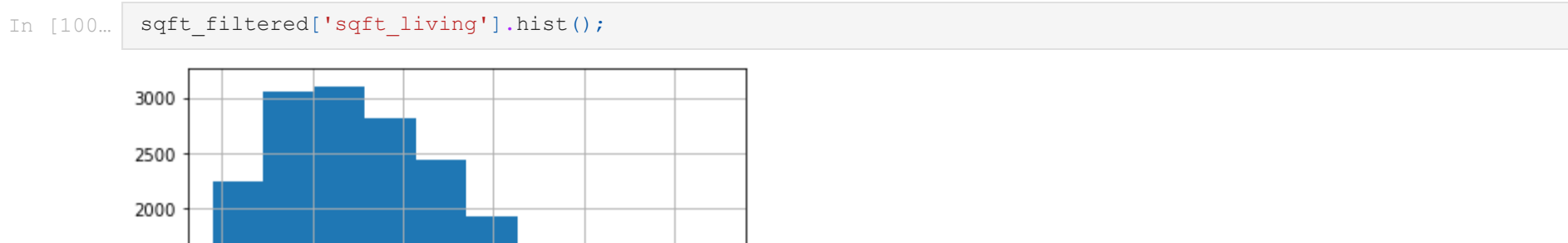
```
In [103]: data['sqft_lot'].hist()
Out[103]:
```

```
0      5650
1      7242
2     10000
3      5000
4      6080
...
21592    1131
21593    3813
21594    1350
21595    2389
21596    1076
Name: sqft_lot, Length: 21597, dtype: int64
```

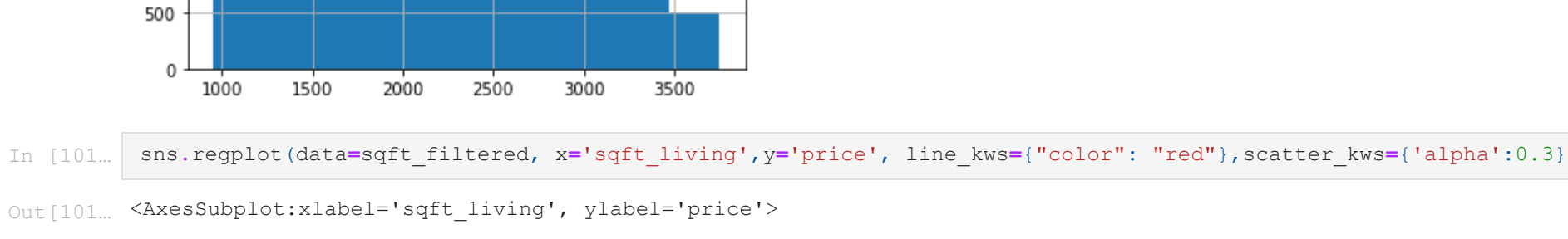
```
In [104]: data['sqft_lot'].describe()
Out[104]:
```

```
count      21597.00000
mean      15099.40876
std       41412.63688
min        520.00000
25%      5040.00000
50%      7618.00000
75%      10688.00000
max      1651359.00000
Name: sqft_lot, dtype: float64
```

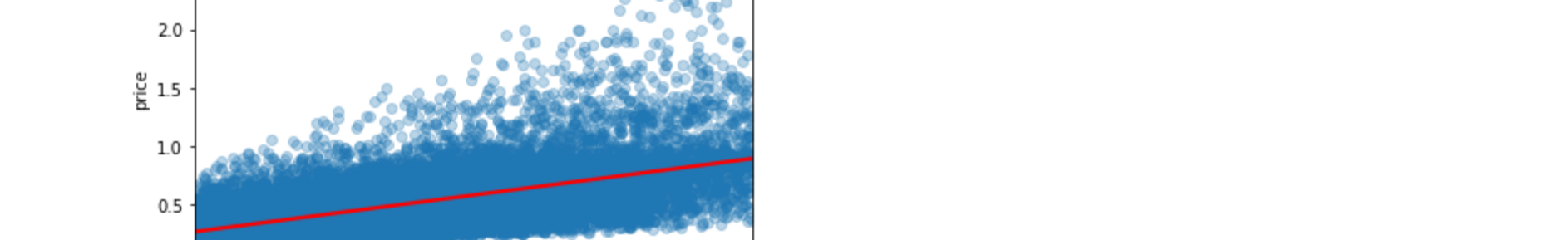
```
In [105]: data['sqft_lot'].hist()
Out[105]:
```



```
In [106]: sns.boxplot(data=data, x='sqft_lot')
plt.ticklabel_format(style='plain', axis='x')
plt.xticks(rotation=75);
Out[106]:
```

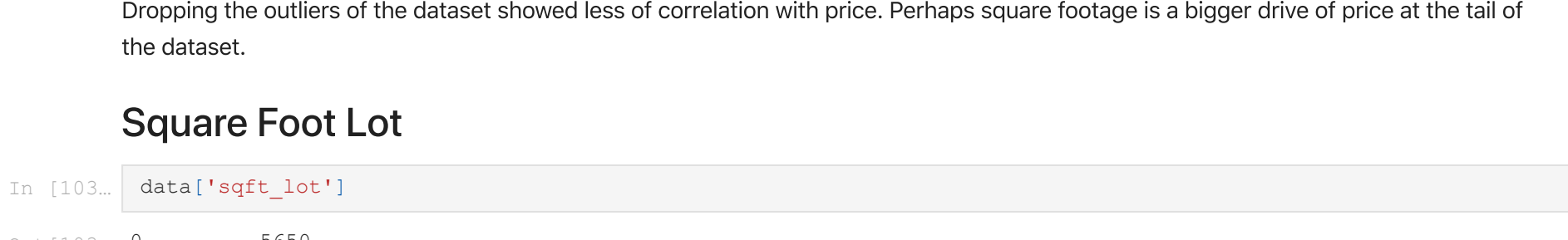


```
In [107]: sns.distplot(x=data['sqft_lot'])
Out[107]:
```



This data set has an extremely long tail and a high peak, most of the values are in the low end of lot square footage. This means most homes are on smaller lots, but we have a long trail off of larger and larger homes.

```
In [108]: sns.regplot(data=data, x='sqft_lot', y='price')
plt.ticklabel_format(style='plain', axis='x')
plt.xticks(rotation=75);
Out[108]:
```



I think this scatterplot is interesting, because we have such variability in the low end of lot square footage. Some other factor is affecting the price in those instances, controlling for square footage, location, bathrooms, bedrooms, there is another factor at play.

Floors

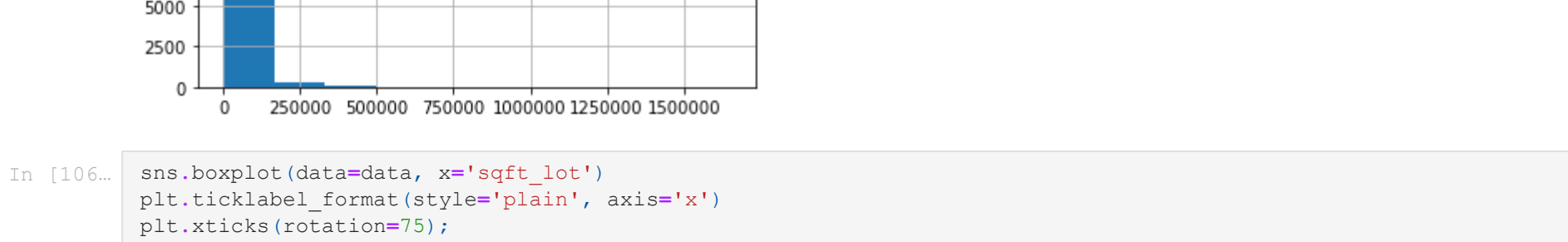
```
In [109]: data['floors'].hist()
Out[109]:
```

```
0      1.00000
1      2.00000
2      1.00000
3      1.00000
4      1.00000
...
21592    3.00000
21593    2.00000
21594    2.00000
21595    2.00000
21596    2.00000
Name: floors, Length: 21597, dtype: float64
```

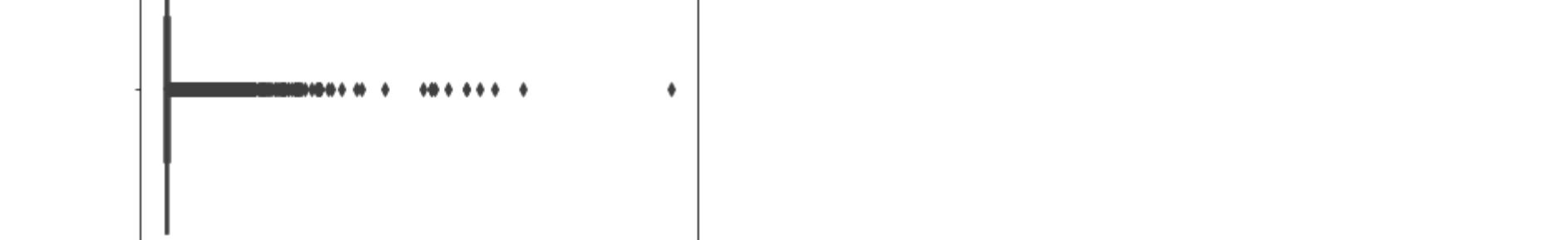
```
In [110]: data['floors'].describe()
Out[110]:
```

```
count      21597.00000
mean       1.49410
std        0.53668
min         1.00000
25%         1.00000
50%         1.00000
75%         2.00000
max         3.50000
Name: floors, dtype: float64
```

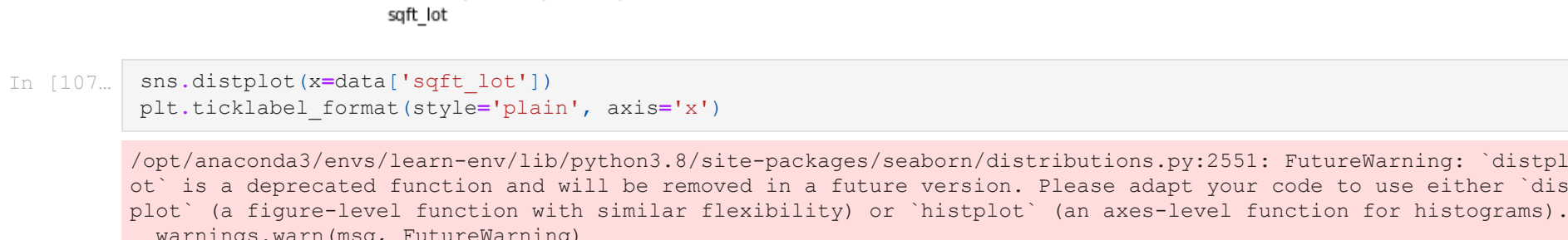
```
In [111]: data['floors'].hist()
Out[111]: <AxesSubplot>
```



```
In [112]: sns.distplot(x=data['floors'])
Out[112]:
```



```
In [113]: sns.regplot(data=data, x='floors', y='price')
Out[113]: <AxesSubplot>
```



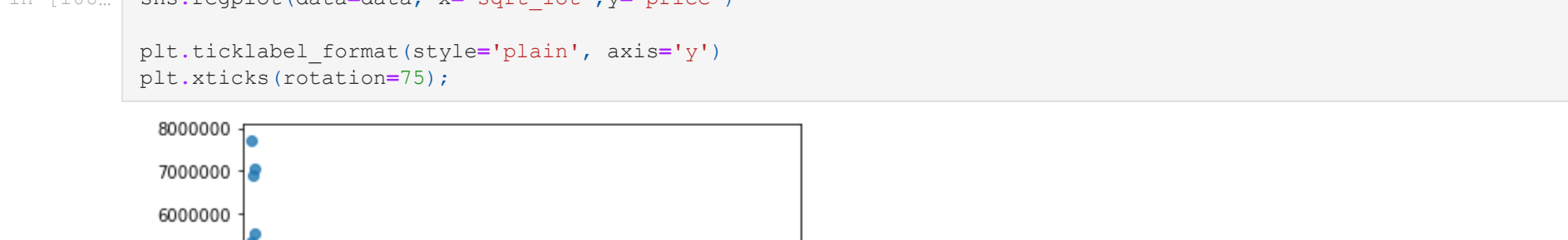
There is a slight increase in price as floors go up.

Waterfront

```
In [114]: data['waterfront'].hist()
Out[114]:
```

```
0      0.00000
1      0.00000
2      0.00000
3      0.00000
4      0.00000
...
21592    0.00000
21593    0.00000
21594    0.00000
21595    0.00000
21596    0.00000
Name: waterfront, Length: 21597, dtype: float64
```

```
In [115]: data['waterfront'].hist()
Out[115]: <AxesSubplot>
```



This column being categorical, a property either is or isn't on the waterfront, it doesn't offer as much chance for exploration and it will be thought of as an add on value when modeling later on I believe.

View

```
In [116]: data['view'].value_counts()
Out[116]:
```



```
count    21534.000000
mean      0.23386
std       0.76569
min       0.000000
25%       0.000000
50%       0.000000
75%       0.000000
max       4.000000
Name: view, dtype: float64
```

The mean isn't helpful here, the mode is obviously 0 as no view. So that's probably the best to go with at this time.

```
np.sum(data['view']==0)/len(data)
Out[437]: 0.8992915682733713
```

```
data['view']==data['view'].fillna(0)
Out[438]:
```

```
np.sum(data['view']==0)/len(data)
Out[439]: 0.90220844008397012
```

This is slightly more dramatic change to the dataset, but we still only slightly altered the distribution of the data.

Basement Square Footage

The basement square footage is currently encoded as an object, which it should not be. Let's look at why this is the case.

```
data['sqft_basement'].value_counts()
Out[440]:
```

```
0.0      12826
600.0     434
500.0     217
500.0     209
700.0     208
700.0     201
...
1275.0     1
1525.0     1
1275.0     1
1281.0     1
1005.0     1
Name: sqft_basement, Length: 304, dtype: int64
```

```
np.sum(data['sqft_basement']!=0)/len(data)
Out[441]: 0.021021438162707785
```

```
np.sum(data['sqft_basement']!=0.0)/len(data)
Out[442]: 0.5938779794601103
```

```
np.sum(data['sqft_basement']!=0.0') & (data['sqft_basement']!=0')/len(data)
Out[443]: 0.385099782377192
```

```
data['sqft_basement'] = data['sqft_basement'].replace('?', '0.0')
data['sqft_basement'].value_counts()
Out[444]:
```

```
0.0      12380
600.0     217
500.0     209
700.0     208
800.0     201
...
1275.0     1
1281.0     1
935.0     1
1816.0     1
1920.0     1
Name: sqft_basement, Length: 303, dtype: int64
```

```
data['sqft_basement'].pd.to_numeric(data['sqft_basement'])
Out[445]:
```

Creating a Square Foot Basement Computed Column

```
data['sqft_basement_computed'] = data['sqft_living']-data['sqft_above']
Out[446]:
```

```
data['sqft_basement_computed']
Out[447]:
```

```
0      0
1      0
2      0
3     910
4      0
...
21592     0
21593     0
21594     0
21595     0
21596     0
Name: sqft_basement_computed, Length: 21597, dtype: int64
```

```
filtered_values = np.where((data['sqft_basement']!=data['sqft_basement_computed']) & (data['sqft_basement']!=0))
Out[448]:
```

```
loc_filter = data.loc[filtered_values].index
Out[449]:
```

```
data[data['sqft_basement_computed']!=0 & (data['sqft_basement']!=0)]
Out[450]:
```

```
   id  date  price  bedrooms  bathrooms  sqft_living  sqft_lot  floors  waterfront  view ... sqft_above
112  2525310310  9/16/2014  272500.00000  3  175000  1540  12600  1.00000  0.00000  0.00000 ... 1160
115  3826039325  11/21/2014  740500.00000  3  350000  4380  6350  2.00000  0.00000  0.00000 ... 2280
309  3204800200  1/8/2015  665000.00000  4  275000  3320  10574  2.00000  0.00000  0.00000 ... 2720
584  7135000030  7/28/2014  1350000.00000  5  350000  4800  14984  2.00000  0.00000  2.00000 ... 3480
308  5113400431  5/8/2014  615000.00000  2  100000  1540  6872  1.00000  0.00000  0.00000 ... 820
...
21000  291310780  6/13/2014  379500.00000  3  225000  1410  1287  2.00000  0.00000  0.00000 ... 1290
21109  3438500250  6/23/2014  519500.00000  5  325000  2910  5027  2.00000  0.00000  0.00000 ... 2040
21210  3278600680  6/27/2014  235000.00000  1  150000  1170  1456  2.00000  0.00000  0.00000 ... 1070
21356  616990185  5/20/2014  490000.00000  5  350000  4460  2975  3.00000  0.00000  2.00000 ... 3280
21442  3226049565  7/11/2014  504600.00000  5  300000  2360  5000  1.00000  0.00000  0.00000 ... 1390
```

170 rows x 22 columns

```
for loc in loc_filter:
    data.at[loc, 'sqft_basement'] = data.at[loc, 'sqft_basement_computed']
Out[451]:
```

```
data['sqft_basement'].value_counts()
Out[452]:
```

```
0.00000  13110
600.00000  221
700.00000  218
500.00000  214
800.00000  206
...
2390.00000  1
602.00000  1
295.00000  1
1281.00000  1
906.00000  1
Name: sqft_basement, Length: 306, dtype: int64
```

```
data[data['sqft_basement']!=data['sqft_basement_computed']]
Out[453]:
```

```
id date price bedrooms bathrooms sqft_living sqft_lot floors waterfront view ... sqft_above sqft_basement yr_built yr_renovate
```

0 rows x 22 columns

```
data = data.drop('sqft_basement_computed', axis=1)
Out[454]:
```

```
#creating a variable to which columns are null in the dataframe
nan_values = data.isna()
#creating a list of columns that contain nulls
nan_columns = nan_values.any()
columns_with_nan = data.columns[nan_columns].tolist()
print(columns_with_nan)
['yr_renovated']
Out[455]:
```

Year Renovated

```
np.sum(data['yr_renovated'].isna())/len(data)
Out[456]: 0.17789507802009538
```

```
np.sum(data['yr_renovated']==0)/len(data)
Out[457]: 0.767653692920313
```

The year renovated value counts and a look at the dataframe suggests that the nulls are probably homes that haven't been renovated, I wouldn't want to make up years for renovation, and the value counts suggest that most homes haven't been renovated. I will set these nulls to 0'.

```
#replacing null yr_renovated column to 0'.
data['yr_renovated'].fillna(0, inplace = True)
Out[458]:
```

```
#checking that null values have been replaced
data.isna()
Out[459]:
```

```
<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 21597 non-null float64
9 view 21597 non-null int64
10 condition 21597 non-null int64
11 grade 21597 non-null int64
12 sqft_above 21597 non-null float64
13 sqft_basement 21597 non-null float64
14 yr_built 21597 non-null int64
15 yr_renovated 21597 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(9), int64(11), object(1)
memory usage: 3.5+ MB
```

Building a Baseline Model

For the baseline model I will only be including the continuous house feature variables separate from the location based-features of zip code, latitude and longitude and date features of the house including build year, renovation year and date of sale.

```
data.columns
Out[460]:
```

```
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'])
Out[461]:
```

```
house_features = ['bedrooms', 'bathrooms', 'sqft_living', 'sqft_lot', 'floors',
                  'waterfront', 'view', 'condition', 'grade', 'sqft_above', 'sqft_basement',
                  'sqft_living15', 'sqft_lot15']
Out[462]:
```

```
#selecting price to be the target of the dataset
y = data['price']
#dropping price for the feature set
X = data[house_features]
Out[463]:
```

Checking for Multicollinearity

In order to make sure the baseline model meets the assumption of independence, it's important to see which variables are multi-collinear.

```
# Create a df with the target as the first column,
# then compute the correlation matrix
heatmap_data = pd.concat([y, X], axis=1)
corr = heatmap_data.corr()
Out[464]:
```

```
# Set up figure and axes
fig, ax = plt.subplots(figsize=(20, 12))
# Plot a heatmap of the correlation matrix, with both
# numbers and colors indicating the correlations
sns.heatmap(
    # Specifies the data to be plotted
    data=corr,
    # The mask means we only show half the values,
    # instead of showing duplicates. It's optional.
    mask=np.triu(np.ones_like(corr, dtype=bool)),
    # Specifies that we want to use the existing axes
    ax=ax,
    # Specifies that we should labels, not just colors
    annot=True,
    cmap = 'Blues',
    # Customizes colorbar appearance
    cbar_kws={"label": "Correlation", "orientation": "horizontal", "pad": -.2, "extend": "both"}
)
# Customize the plot appearance
ax.set_title("Heatmap of Correlation Between Attributes (Including Target)");
Out[465]:
```



The heat matrix shows there are several multi-collinear variables over a .75 threshold. I can decide which features to initially drop.

```
df=X.corr().iabs().stack().reset_index(i.sort_values(0, ascending=False))
df['pairs'] = list(zip(df.level_0, df.level_1))
df.set_index(['pairs'], inplace = True)
df.drop_duplicates(inplace=True)
df.columns = ['cc']
df.drop_duplicates(inplace=True)
df
Out[466]:
```

```
cc
pairs
(bedrooms, bedrooms) 1.00000
(sqft_above, sqft_living) 0.87645
(grade, sqft_living) 0.76278
(sqft_living15, sqft_living) 0.75640
(grade, sqft_above) 0.75607
...
(floors, sqft_lot15) 0.01072
(sqft_lot, floors) 0.00481
(sqft_lot15, condition) 0.00983
(bedrooms, waterfront) 0.00213
Out[467]:
```

79 rows x 1 columns

```
abs(df.corr()) > 0.75
df=X.corr().iabs().stack().reset_index(i.sort_values(0, ascending=False))
df[(df[0]>.75) & (df[0] < 1)]
Out[468]:
```

```
cc
level_0 level_1 0
319 sqft_above sqft_living 0.87645
106 grade sqft_living 0.76278
34 sqft_living grade 0.76278
145 sqft_living15 sqft_living 0.75640
37 sqft_living sqft_living15 0.75607
113 grade sqft_above 0.75607
125 sqft_above grade 0.75607
27 sqft_living bathrooms 0.75576
15 bathrooms sqft_living 0.75576
Out[469]:
```

Examining this correlation set, it seems time to remove some of the features based on their correlations. Square Feet Living is a very understandable concept, so I will be dropping "Sqft_above". "Sqft_living15" also shows a high level of correlation to square_feet living.

Removing Multicollinearity

```
to_drop = ['sqft_above', 'sqft_living15']
house_features = [element for element in house_features if element not in to_drop]
Out[469]:
```

```
X = data[house_features]
Out[470]:
```

```
X.info()
Out[471]:
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21597 entries, 0 to 21596
Data columns (total 18 columns):
# Column Non-Null Count Dtype
---
0 id 21597 non-null int64
1 bedrooms 21597 non-null float64
2 sqft_living 21597 non-null int64
3 sqft_lot 21597 non-null int64
4 floors 21597 non-null float64
5 waterfront 21597 non-null float64
6 view 21597 non-null int64
7 condition 21597 non-null int64
8 sqft_basement 21597 non-null float64
9 sqft_lot15 21597 non-null int64
dtypes: float64(5), int64(6)
memory usage: 1.8 MB
```

Train Test Split

Now I can perform a a train/test split of the data in order to ensure the model is not overfitting.

```
#import sklearn train/test split
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y)
Out[472]:
```

```
print(f'X_train is a DataFrame with {X_train.shape[0]} rows and {X_train.shape[1]} columns')
print(f'y_train is a Series with {y_train.shape[0]} values')
# We always should have the same number of rows in X as values in y
assert X_train.shape[0] == y_train.shape[0]
X_train is a DataFrame with 16197 rows and 11 columns
y_train is a Series with 16197 values
Out[473]:
```

Scaling Data

In order to gage the relative impact of the house features I will scaling the data using a standard scaler.

```
#Importing the Standard Scaler
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
scaled_data = scaler.fit_transform(X_train)
Out[474]:
```

```
scaled_data
array([[0.69682284, -0.14882144, 0.1725273, ..., -0.55667713,
       -1.8697709, -1.2393118],
       [-1.50971471, -1.44337612, -1.53811406, ..., -1.40801312,
       -0.64925384, -0.26017592],
       [-0.40644533, 0.4984559, -0.55287125, ..., -0.55667713,
       0.43890535, -0.27483768],
       ...,
       [-0.40644533, -0.14882144, -0.55287125, ..., 0.29465886,
       -0.64925384, -0.34416363],
       [-0.69682284, 0.4984559, 0.43237156, ..., 1.99733084,
       -0.64925384, -0.13664851],
       [-0.40644533, -0.26017592, 0.54064, ..., 1.14599485,
       0.40644533, -1.18509771],
       ...])
```

The scaled data produces a series of arrays that will not be as interpretable without the column labels.

```
#Creating a new data frame with the scaled data. Getting the column names from the X_train dataframe.
X_train = pd.DataFrame(scaled_data, columns = X_train.columns, index=X_train.index)
Out[475]:
```

Fitting the Baseline Model

```
from sklearn.linear_model import LinearRegression
#Instantiating the Baseline Model
model = LinearRegression()
model.fit(X_train, y_train)
Out[476]:
```

Scoring the Model

```
# Run this cell without changes
from sklearn.model_selection import cross_val_score
cross_val_score = cross_val_score(model, X_train, y_train, cv=8)
cross_val_score.mean()
Out[477]:
```

```
0.5975518518659969
```

```
#Obtaining the model score for the Training data.
model.score(X_train, y_train)
Out[478]:
```

0.602614291489027

```
#Scaling the test data
X_test = scaler.transform(X_test)
Out[479]:
```

```
#Scoring the model on test data.
model.score(X_test, y_test)
Out[479]:
```

0.6197702056233946

The model's cross-validated score, training score, and test score are all around 60%, meaning 60% of the variance of the data is explained by the model - good not great.

```
prediction = model.predict(X_train)
residuals = (y_train - prediction)
Out[480]:
```

The baseline model is performing around 60% on the training data, and 59% on the test data. The model is not overfitting greatly, but will it adhere to the necessary assumptions.

Analyzing the Model in Stats Model

```
import statsmodels.api as sm
stats_model = sm.OLS(y_train, sm.add_constant(X_train)).fit()
stats_model.summary()
Out[481]:
```

OLS Regression Results					
Dep. Variable:	price	R-squared:	0.603		
Model:	OLS	Adj. R-squared:	0.602		
Method:	Least Squares	F-statistic:	2231.		
Date:	Fr, 07 Jan 2022	Log Likelihood:	0.00		
Time:	05:01:48	Prob (F-statistic):	-2.2303e+05		
No. Observations:	16197	AIC:	4.46e+05		
Df Residuals:	16185	BIC:	4.462e+05		
Df Model:	11				
Covariance Type:	nonrobust				
	coef	std err	t	P> t	[0.025 0.975]
const	5.406e+05	1816.772	297.567	0.000	5.37e+05 5.44e+05
bedrooms	-3.421e+04	2361.937	-14.485	0.000	-3.88e+04 -2.96e+04
bathrooms	-9306.9448	3150.693	-3.001	0.003	-1.56e+04 -3259.441
sqft_living	1.763e+05	3992.125	44.162	0.000	1.66e+05 1.84e+05
sqft_lot	-69.9142	2614.111	-0.027	0.979	-8786.382 5054.033
floors	-2269.9310	2457.236	-0.924	0.356	-7086.385 2546.524
waterfront	4.644e+04	1957.229	23.726	0.000	4.26e+04 5.03e+04
view	4.582e+04	2083.603	21.993	0.000	4.17e+04 4.99e+04
condition	3.408e+04	1904.274	17.898	0.000	3.03e+04 3.78e+04
grade	1.204e+05	3083.215	39.062	0.000	1.14e+05 1.26e+05
sqft_basement	1.257e+04	2400.399	5.238	0.000	7867.064 1.73e+04
sqft_lot15	-1.904e+04	2626.997	-7.248	0.000	-2.42e+04 -1.39e+04
Omnibus:	12349.218	Durbin-Watson:	2.000		
Prob(Omnibus):	0.000	Jarque-Bera (JB):	891005.788		
Skewn:	3.076	Prob(JB):	0.00		
Kurtosis:	38.811	Cond. No.	4.82		

Notes:

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

An initial examination of the Stats Model summary indicates, low p_values for many of the feature coefficients. I will consider dropping Sqft_Lot from the model as it is well above a p_value of .05 as well as the floors coefficient.

```
to_drop = ['sqft_lot']
Out[482]:
```

Examining Baseline Residuals

In order to check the assumptions of normality of residuals I will obtain the predictions with model on the test data.

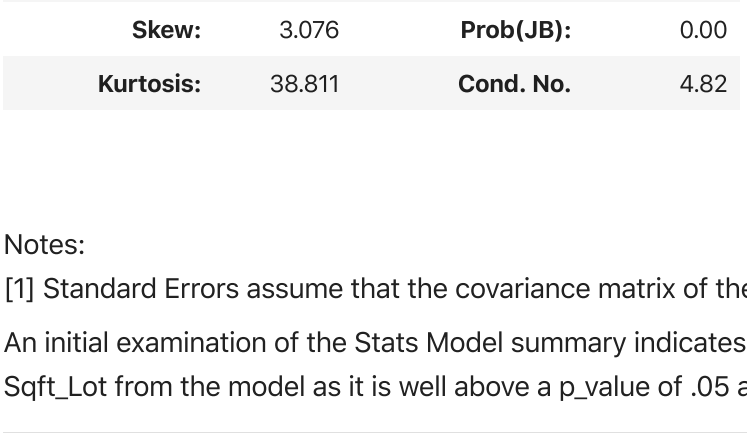
Normality of Residuals

```
from scipy import stats
sns.graphics.qqplot(stats_model.resid, dist=stats.norm, line='45', fit=True);
Out[483]:
```



Initially, according to this test it appears the distribution of the residuals is certainly not normally distributed. Plotting the residuals on a histogram.

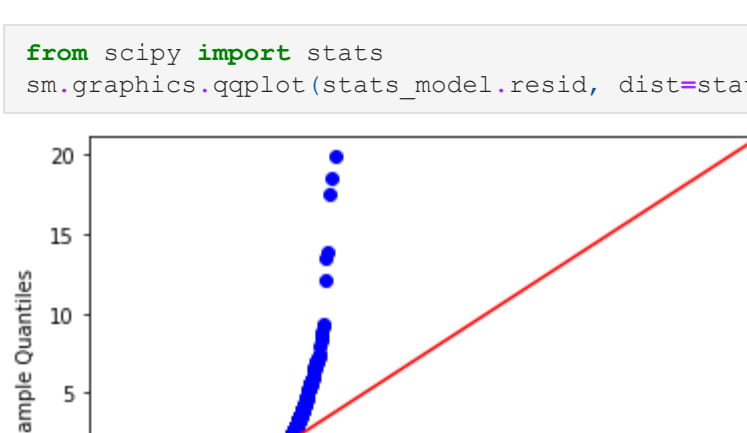
```
plt.hist(stats_model.resid, bins=100)
Out[484]:
```



Again the residuals do not show a normal distribution in this histogram, the dataset included a long tail of the target variable and of the square foot living feature.

Homoscedasticity

```
# Run this cell without changes
ax.scatter(prediction, residuals, alpha=0.5)
ax.plot(prediction, [0 for i in range(len(X_train))])
ax.set_xlabel("Predicted Value")
ax.set_ylabel("Residuals")
ax.set_title("Homoscedasticity Test");
Out[485]:
```



In addition to a problem with the normal distribution of residuals, the scatter plot of Predicted Values against Residuals, shows cone-like shape indicating that as the predicted values increase, the residuals appear to be increasing, though this seems to be more diffuse along the x-axis. This may be reflective of the skew of the data set as well.

The data is highly right-skewed as is the size of the many of the homes in the data set. For the second iteration of the model I will look at the data and perhaps exclude outliers.

```
from statsmodels.stats.diagnostic import het_breuschpagan
stats_test, stats_pvalue, stats_model_resid, stats_model_exog = statsmodels.stats.diagnostic.het_breuschpagan(
    labels = ["LM Statistic", "LM-Test p-value", "F-Statistic", "F-Test p-value"])
dic(dict(zip(labels,stats_test)))
Out[486]:
```

```
["LM Statistic": 1867.855759467683,
 "LM-Test p-value": 0.0,
 "F-Statistic": 191.79758484543374,
 "F-Test p-value": 0.0]
```

Second Model Iteration

Dropping High P Value Features

```
to_drop
Out[487]:
```

```
['sqft_lot']
Out[487]:
```

```
house_features = [element for element in house_features if element not in to_drop]
Out[488]:
```

```
house_features
Out[489]:
```

```
['bedrooms',
 'bathrooms',
 'sqft_living',
 'floors',
 'waterfront',
 'view',
 'condition',
 'grade',
 'sqft_basement',
 'sqft_living15',
 'yr_built',
 'yr_renovated',
 'zipcode',
 'lat',
 'long',
 'sqft_living15',
 'sqft_lot15',
 dtype='object']
Out[490]:
```

```
data['sqft_living_compared'] = data['sqft_living']/data['sqft_living15']*100
Out[491]:
```

```
data['sqft_lot_compared'] = data['sqft_lot']/data['sqft_lot15']*100
Out[492]:
```

```
house_features.append('sqft_living_compared')
Out[493]:
```

```
house_features.append('sqft_lot_compared')
Out[494]:
```

```
to_drop = ['sqft_living15', 'sqft_lot15']
house_features = [element for element in house_features
```



```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 20439 entries, 0 to 21596
Data columns (total 23 columns):
# Column Non-Null Count Dtype
---  ---
0 id 20439 non-null int64
1 date 20439 non-null object
2 price 20439 non-null float64
3 bedrooms 20439 non-null int64
4 bathrooms 20439 non-null float64
5 sqft_living 20439 non-null int64
6 sqft_lot 20439 non-null int64
7 floors 20439 non-null float64
8 waterfront 20439 non-null float64
9 view 20439 non-null float64
10 condition 20439 non-null int64
11 grade 20439 non-null int64
12 sqft_above 20439 non-null float64
13 sqft_basement 20439 non-null float64
14 yr_built 20439 non-null int64
15 yr_renovated 20439 non-null int64
16 zipcode 20439 non-null int64
17 lat 20439 non-null float64
18 long 20439 non-null float64
19 sqft_living15 20439 non-null int64
20 sqft_lot15 20439 non-null int64
21 sqft_living_compared 20439 non-null float64
22 sqft_lot_compared 20439 non-null float64
dtypes: float64(11), int64(11), object(1)
memory usage: 3.7+ MB
```

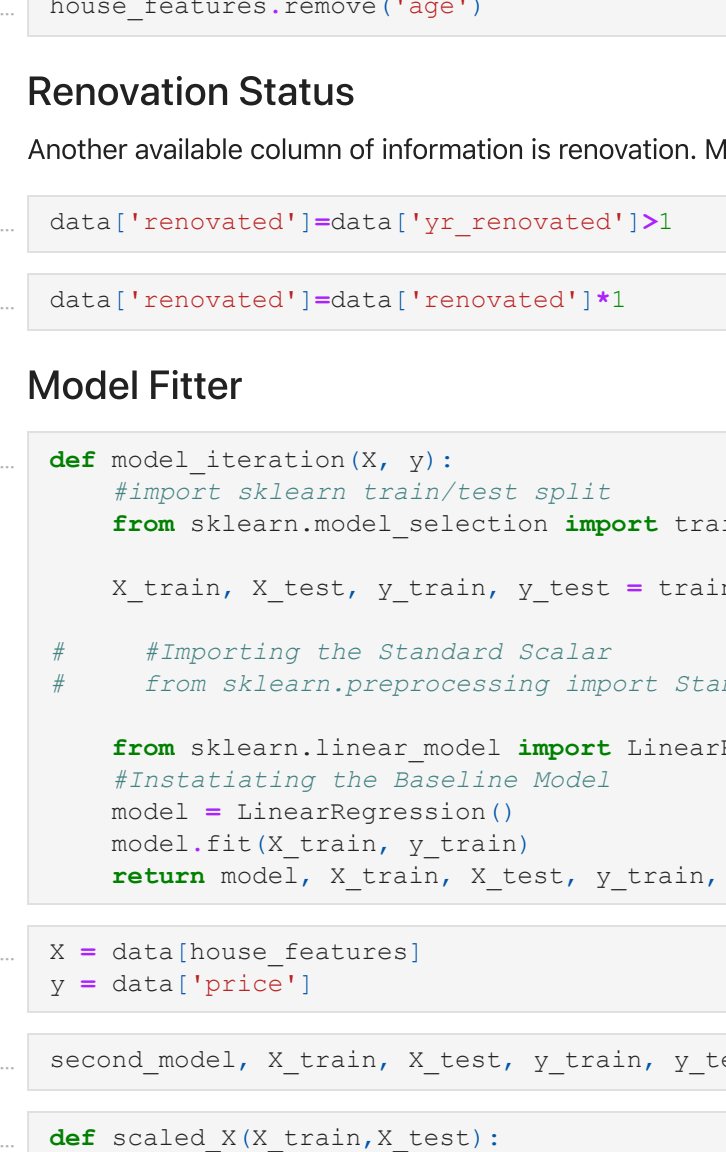
```
In [507]: data['age']=data['date'].dt.year - data['yr_built']

In [508]: house_features.append('age')
```

```
In [509]: sm.graphics.qqplot(data['age'],data['price'])

FutureWarning: Pass the following variables as keyword arguments to the only valid positional argument: the "data",
and passing other arguments without an explicit keyword will result in an error or misinterpretation.
```

```
Out[509]: <Figure with 1 Axes>
```



I don't believe age meets the standard of linearity that is needed for performing regression.

```
In [510]: house_features.remove('age')
```

Renovation Status

Another available column of information is renovation. Many of the homes were not renovated, but I can still create a renovatin feature.

```
In [511]: data['renovated']=data['yr_renovated']>1
```

```
In [512]: data['renovated']=data['renovated']*1
```

Model Fitter

```
In [513]: def model_iteration(X, y):
    #import sklearn train/test split
    from sklearn.model_selection import train_test_split
    X_train, X_test, y_train, y_test = train_test_split(X, y)
    # Importing the StandardScaler
    from sklearn.preprocessing import StandardScaler
    #Instantiating the Baseline Model
    model = LinearRegression()
    model.fit(X_train, y_train)
    return model, X_train, X_test, y_train, y_test

In [514]: X = data[house_features]
y = data['price']

second_model, X_train, X_test, y_train, y_test = model_iteration(X,y)

In [515]:

In [516]: def scaled_X(X_train,X_test):
    scaler = StandardScaler()
    scaled_data = scaler.fit_transform(X_train)
    # Creating a new data frame with the scaled data. Getting the column names from the X train dataframe.
    scaled_X_train = pd.DataFrame(scaled_data, columns = X_train.columns,index=X_train.index)
    scaled_X_test = scaler.transform(X_test)
    return scaled_X_train, scaled_X_test
```

MultiCollinearity Check

```
In [517]: df=X.corr().abs().stack().reset_index().sort_values(0, ascending=False)
df['pairs'] = list(zip(df.level_0, df.level_1))
df.set_index('pairs', inplace = True)
df.drop_duplicates(inplace=True)
df.columns = ['cc']
df.drop_duplicates(inplace=True)
df
abs(df.corr()) > 0.75
df.corr().abs().stack().reset_index().sort_values(0, ascending=False)
df[(df[0]>.75) & (df[0] <1)]

Out[517]:
level_0 level_1 0
```

Model Score

```
In [518]: def model_score(model, X_train, X_test, y_train, y_test):
    train_score = model.score(X_train, y_train)
    test_score = model.score(X_test, y_test)
    return f'Train Score = {train_score}, Test Score = {test_score}'

In [519]: model_score(second_model, X_train, X_test, y_train, y_test)

Out[519]: 'Train Score = 0.5054460285367008, Test Score = 0.49330562646352305'
```

Residual Helpers

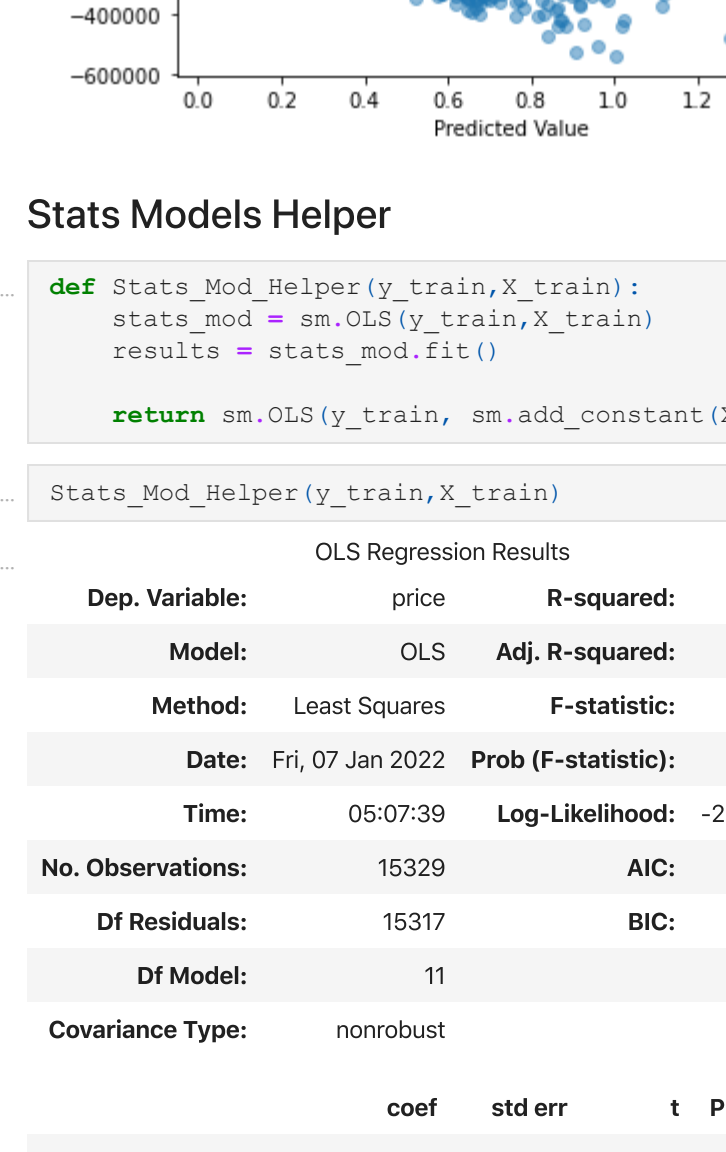
```
In [520]: def residual_helper(model,X_train,y_train):
    prediction = model.predict(X_train)
    residuals = (y_train - prediction)
    return prediction, residuals

In [521]: prediction, residuals = residual_helper(second_model,X_train,y_train)
```

Residual Plots

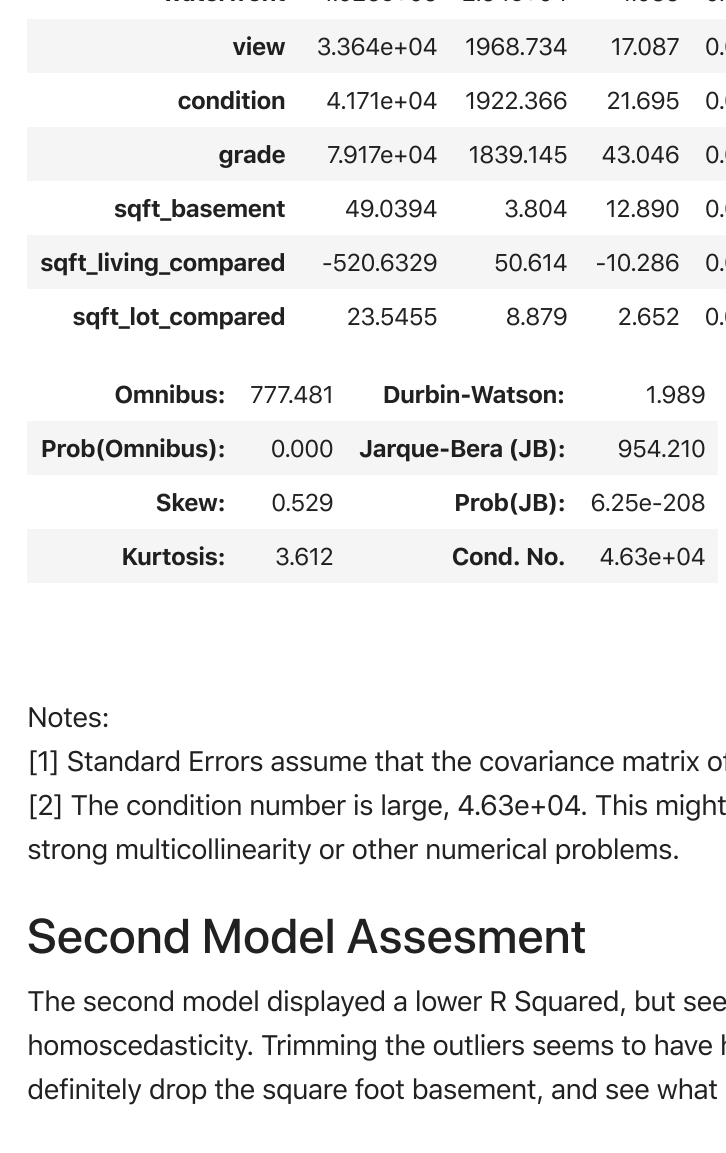
```
In [522]: sm.graphics.qqplot(residuals, dist=stats.norm, line='45', fit=True);

Out[522]:
```



```
In [523]: # Run this cell without changes
fig, ax = plt.subplots()
ax.scatter(prediction, residuals, alpha=0.5)
ax.plot(prediction, [0 for i in range(len(X_train))])
ax.set_xlabel("Predicted Value")
ax.set_ylabel("Residuals")
ax.set_title("Homoscedasticity Test")

Out[523]:
```



Stats Models Helper

```
In [526]: def Stats_Mod_Helper(y_train,X_train):
    stats_mod = sm.OLS(y_train,X_train)
    results= stats_mod.fit()
    return sm.OLS(y_train, sm.add_constant(X_train)).fit().summary()

In [527]: Stats_Mod_Helper(y_train,X_train)
```

OLS Regression Results				
Dep. Variable:	price	R-squared:	0.506	
Model:	OLS	Adj. R-squared:	0.506	
Method:	Least Squares	F-statistic:	1423.	
Date:	Thu, 07 Jan 2022	Prob (F-statistic):	0.00	
Time:	05:07:09	Log-Likelihood:	-2.040e+05	
No. Observations:	16329	AIC:	4.082e+05	
Df Residuals:	15317	BIC:	4.083e+05	
Df Model:	11			
Covariance Type:	nonrobust			
	coef	std err	t	P> t
const	-4.14e+05	1.44e+04	-28.786	0.000
bedrooms	-713919e+04	16483.788	-4.330	0.000
bathrooms	-1.9599e+04	26893.76	-0.7283	0.000
sqft_living	102.4033	3.1392	32.082	0.000
floors	3.153e+04	30363.078	10.385	0.000
waterfront	1.025e+05	2.54e+04	4.035	0.000
view	3.864e+04	19683.74	17.087	0.000
condition	4.171e+04	19222.366	21.695	0.000
grade	19717e+04	1839145	43.046	0.000
sqft_basement	409.3934	12.890	0.000	1.582
sqft_living_compared	-520.6329	50.614	-10.286	0.000
sqft_lot_compared	23.5455	83.79	0.2852	0.008
Omnibus:	777.481	Durbin-Watson:	1.989	
Prob(Omnibus):	0.000	Jarque-Bera (JB):	954.210	
Skew:	0.529	Prob(JB):	0.25e-208	
Kurtosis:	3.612	Cond. No.	4.62e+04	

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 4.63e+04. This might indicate that there are strong multicollinearity or other numerical problems.

Second Model Assessment

The second model displayed a lower R Squared, but seemed to show a better adherence to the assumptions of normality and homoscedasticity. Trimming the outliers results to have helped the model, but maybe some of the features were not helpful. I will definitely drop the square foot basement, and see what other features can be dropped.

Feature Selection

```
In [528]: from sklearn.feature_selection import RFECV
from sklearn.model_selection import cross_validate, ShuffleSplit
splitter = ShuffleSplit(n_splits=5, test_size=0.25, random_state=0)
```

```
In [529]: # Importances are based on coefficient magnitude, so
# we need to scale the data to normalize the coefficients
X_train_for_RFECV = StandardScaler().fit_transform(X_train)
model_for_RFECV = LinearRegression()
```

```
# Instantiate and fit the selector
selector = RFECV(model_for_RFECV, cv=splitter)
selector.fit(X_train_for_RFECV, y_train)
```

```
# Print the results
print("Was the column selected?")
for index, col in enumerate(X_train.columns):
    print(f"{col}: {selector.support_((index))}")
```

Was the column selected?
bedrooms: True
bathrooms: True
sqft_living: True
floors: True
waterfront: True
view: True
condition: True
grade: True
sqft_basement: True
sqft_living_compared: True
sqft_lot_compared: False

I will keep the square foot Basement column, the Square Foot Lot Compared feature.

```
In [530]: to_drop = ['sqft_lot_compared']
```

```
In [531]: # Run this cell without changes
print(f"X_train is a DataFrame with {X_train.shape[0]} rows and {X_train.shape[1]} columns")
print(f"y_train is a Series with {y_train.shape[0]} values")
assert X_train.shape[0] == y_train.shape[0]
```

X_train is a DataFrame with 15329 rows and 11 columns
y_train is a Series with 15329 values

Third Model Iteration

Categorical Feature Engineering

Location Information

The first two iterations of the model did not utilize location information, but location most likely plays a role in the valuation of homes. Currently the data features 70 zip codes. In order to reduce the number of features and make them interpretable I will add join the zip codes to city names.

I will also add some information about seasonality of sale to give some additional information to the model to improve the R-Squared. I hope that dropping the insignificant features and adding this categorical information will show improvement.

```
In [532]: len(data['zipcode'].value_counts())

Out[532]: 70

In [533]: zip_df = pd.read_csv('data/uszipis.csv')

In [534]: zip_df

Out[534]:
```

98042:	'Kent',
98033:	'Kirkland',
98034:	'Kirkland',
98038:	'Maple Valley',
98039:	'Medina',
98040:	'Mercer Island',
98045:	'North Bend',
98052:	'Redmond',
98053:	'Redmond',
98055:	'Renton',
98056:	'Renton',
98058:	'Renton',
98059:	'Renton',
98074:	'Sammamish',
98075:	'Sammamish',
98125:	'Seattle',
98126:	'Seattle',
98133:	'Seattle',
98136:	'Seattle',
98144:	'Seattle',
98146:	'Seattle',
98168:	'Seattle',
98155:	'Seattle',
98166:	'Seattle',
98122:	'Seattle',
98177:	'Seattle',
98178:	'Seattle',
98188:	'Seattle',

The zip code Data Frame contains a lot of information I probably won't use, so I will only make a Data Frame of columns that seem relevant.

```
In [535]: zip_df = zip_df[['zip','city','state_name','population','density','county_name']]

This is a large data set and I will narrow it down to only Kings County zip codes.
```

```
In [536]: data['zipcode'].value_counts().index

Out[536]: Int64Index([98013, 98038, 98052, 98115, 98042, 98117, 98034, 98118, 98023, 98113, 98058, 98059, 98155, 98074, 98066, 98066, 98125, 98027, 98053, 98033, 98031, 98126, 98029, 98106, 98075, 98144, 98029, 98116, 98055, 98146, 98028, 98003, 98198, 98031, 98199, 98122, 98159, 98169, 98055, 98039, 98072, 98107, 98178, 98136, 98030, 98166, 98022, 98177, 98045, 98002, 98011, 98019, 98077, 98108, 98105, 98112, 98040, 98034, 98119, 98009, 98007, 98188, 98032, 98014, 98076, 98010, 98035, 98039, 98024, 98149, 98091], dtype='int64')
```

```
In [537]: #Creating a list of zips to filter from
kings_zips = list(data['zipcode'].value_counts().index)
```

```
In [538]: zip_df = zip_df[zip_df['zip'].isin(kings_zips)]

In [539]: zip_df

Out[539]:
```

```

19 sqft_living>
20 sqft_lot>
21 sqft_living_compared
22 sqft_lot_compared
23 age
24 renovated
25 city
dtypes: datetime64[ns](1), float64(11), int64(13), object(1)
memory usage: 4.2+ MB

```

One Hot Encoding "City"

```

len(data['city'].value_counts())

```

24

For this iteration I will be one hot encoding the cities as separate features. There will be 23 features that I will be modeled along with house features.

```

city_dummies = pd.get_dummies(data['city'], drop_first=True)

city_cols = list(city_dummies.columns)

data = pd.concat([data, city_dummies], axis=1)
data.info()

```

70 rows x 6 columns

```
In [540]: #Creating a dataframe of zipcodes and cities
zip_city = zip_df[['zip','city']].sort_values(by='city')
```

```
In [541]: zip_city

Out[541]:
```

14	yr_built	20439	non-null	int64
15	yr_renovated	20439	non-null	float64
16	zipcode	20439	non-null	int64
17	lat	20439	non-null	float64
18	long	20439	non-null	float64
19	sqft_living15	20439	non-null	int64
20	sqft_lot15	20439	non-null	int64
21	sqft_living_compared	20439	non-null	float64
22	sqft_lot_compared	20439	non-null	float64
23	age	20439	non-null	int64
24	renovated	20439	non-null	int64
25	city	20439	non-null	object
26	Bellevue	20439	non-null	uint8
27	Black Diamond	20439	non-null	uint8
28	Bochell	20439	non-null	uint8
29	Cararnation	20439	non-null	uint8
30	Duwali	20439	non-null	uint8
31	Enchancel	20439	non-null	uint8
32	Fall City	20439	non-null	uint8
33	Federal Way	20439	non-null	uint8
34	Issaquah	20439	non-null	uint8
35	Kennmore	20439	non-null	uint8
36	Kent	20439	non-null	uint8
37	Kirkland	20439	non-null	uint8
38	Maple Valley	20439	non-null	uint8
39	Medina	20439	non-null	uint8
40	Merrest Island	20439	non-null	uint8
41	North Bend	20439	non-null	uint8
42	Redmond	20439	non-null	uint8

70 rows x 2 columns

```
In [542]: #Making a dictionary of zip codes and cities
zip_dict = dict(zip(zip_city.zip, zip_city.city))
zip_dict

Out[542]:
```

```
{98001: 'Auburn',
 98002: 'Auburn',
 98029: 'Auburn',
 98004: 'Bellevue',
 98005: 'Bellevue',
 98006: 'Bellevue',
 98007: 'Bellevue',
 98008: 'Bellevue',
 98010: 'Black Diamond',
 98011: 'Bothell',
 98014: 'Carnation',
 98179: 'Duvali',
 98022: 'Enumclaw',
 98024: 'Fall City',
 98193: 'Federal Way',
 98003: 'Federal Way',
 98027: 'Issaquah',
 98029: 'Issaquah',
 98028: 'Kennesaw',
 98031: 'Kent',
 98033: 'Kent',
 98032: 'Kent',
 98039: 'Kent',
 98033: 'Kirkland',
 98034: 'Kirkland',
 98038: 'Maple Valley',
 98039: 'Medina',
 98040: 'Mercer Island',
 98041: 'North Bend',
 98052: 'Redmond',
 98151: 'Redmond',
 98055: 'Renton',
 98056: 'Renton',
 98058: 'Renton',
 98059: 'Renton',
 98074: 'Sammamish',
 98075: 'Sammamish',
 98125: 'Seattle',
 98126: 'Seattle',
 98133: 'Seattle',
 98136: 'Seattle',
 98144: 'Seattle',
 98146: 'Seattle',
 98169: 'Seattle',
 98155: 'Seattle',
 98166: 'Seattle',
 98122: 'Seattle',
 98177: 'Seattle',
 98178: 'Seattle',
 98188: 'Seattle',
 98149: 'Seattle',
 98119: 'Seattle',
 98199: 'Seattle',
 98117: 'Seattle',
 98198: 'Seattle',
 98102: 'Seattle',
 98103: 'Seattle',
 98118: 'Seattle',
 98106: 'Seattle',
 98108: 'Seattle',
 98109: 'Seattle',
 98112: 'Seattle',
 98115: 'Seattle',
 98116: 'Seattle',
 98107: 'Seattle',
 98071: 'Shoquahmie',
 98072: 'Woodinville',
 98072: 'Woodinville'}
```

The DataSet know contains the additional features for the a second iteration of modeling. Having now droppe the outliers, and one-hot encoded the cities, it may be time for another pass at modeling.

I will be making some helper functions for the second iteration of the model.

Month of Sale

```
In [548]: data['month_sale']=data['date'].dt.month
data['month_sale']

Out[548]:
```

```
0 10
1 12
2 12
3 12
4 2
5 2
21592 1
21593 2
21594 6
21595 1
21596 10
Name: month_sale, Length: 20439, dtype: int64
```

```
In [549]: def season_of_date(date.UTC):
    year = str(date.UTC.year)
    seasons = ['spring': pd.date_range(start=year+'-03-21 00:00:00', end=year+'-06-20 00:00:00'),
               'summer': pd.date_range(start=year+'-06-21 00:00:00', end=year+'-09-22 00:00:00'),
               'autumn': pd.date_range(start=year+'-09-23 00:00:00', end=year+'-12-20 00:00:00')]
    if date.UTC in seasons['spring']:
        return 'spring'
    if date.UTC in seasons['summer']:
        return 'summer'
    if date.UTC in seasons['autumn']:
        return 'autumn'
    else:
        return 'winter'

data['season'] = data.date.map(season_of_date)
```

One Hot Encoding Season

```
In [550]: season_dummies = pd.get_dummies(data['season'],drop_first=True)

In [551]: season_cols = list(season_dummies.columns)
```

```
In [552]: data = pd.concat((data, season_dummies), axis=1)
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 20439 entries, 0 to 21596
Data columns (total 49 columns):
# Column Non-Null Count Dtype
---  ---
0 id 20439 non-null int64
1 date 20439 non-null datetime64[ns]
2 price 20439 non-null float64
3 bedrooms 20439 non-null int64
4 bathrooms 20439 non-null float64
5 sqft_living 20439 non-null int64
6 sqft_lot 20439 non-null int64
7 floors 20439 non-null float64
8 waterfront 20439 non-null float64
9 view 20439 non-null float64
10 condition 20439 non-null int64
11 grade 20439 non-null int64
12 sqft_above 20439 non-null int64
13 sqft_basement 20439 non-null float64
14 yr_built 20439 non-null int64
15 yr_renovated 20439 non-null int64
16 zipcode 20439 non-null int64
17 lat 20439 non-null float64
18 long 20439 non-null float64
19 sqft_living15 20439 non-null int64
20 sqft_lot15 20439 non-null int64
21 sqft_living_compared 20439 non-null float64
22 sqft_lot_compared 20439 non-null float64
23 age 20439 non-null int64
24 renovated 20439 non-null int64
25 city 20439 non-null object
26 Bellevue 20439 non-null uint8
27 Black Diamond 20439 non-null uint8
28 Bothell 20439 non-null uint8
29 Carnation 20439 non-null uint8
30 Duvali 20439 non-null uint8
31 Enumclaw 20439 non-null uint8
32 Fall City 20439 non-null uint8
33 Federal Way 20439 non-null uint8
34 Issaquah 20439 non-null uint8
35 Kenmore 20439 non-null uint8
36 Kent 20439 non-null uint8
37 Kirkland 20439 non-null uint8
38 Maple Valley 20439 non-null uint8
39 Medina 20439 non-null uint8
40 Mercer Island 20439 non-null uint8
41 North Bend 20439 non-null uint8
42 Redmond 20439 non-null uint8
43 Renton 20439 non-null uint8
44 Sammamish 20439 non-null uint8
45 Snoqualmie 20439 non-null uint8
46 Vashon 20439 non-null uint8
47 Woodinville 20439 non-null uint8
48 Woodinville 20439 non-null object
49 season 20439 non-null object
dtypes: datetime64[ns](1), float64(11), int64(14), object(2), uint8(26)
memory usage: 3.0+ MB
```

Third Model Iteration

```
Out[573]: house_features

('bedrooms',
 'bathrooms',
 'sqft_living',
 'floors',
 'waterfront',
 'view',
 'condition',
 'sqft_basement',
 'sqft_living_compared')
```

```
In [574]: house_features = [element for element in house_features if element not in to_drop]

X = data[house_features+city_cols+season_cols]
y = data['price']

In [580]: third_model, X_train, X_test, y_train, y_test = model_iteration(X,y)
```

Scoring Third Model

```
In [581]: model_score(third_model, X_train, X_test, y_train, y_test)

Out[581]: 'Train Score = 0.696033221702934, Test Score = 0.7063994618859418'
```

```
In [582]: # Run this cell without changes
from sklearn.model_selection import cross_val_score
cross_val_score = cross_val_score(third_model, X_train, y_train, cv=5)
cross_val_score.mean()

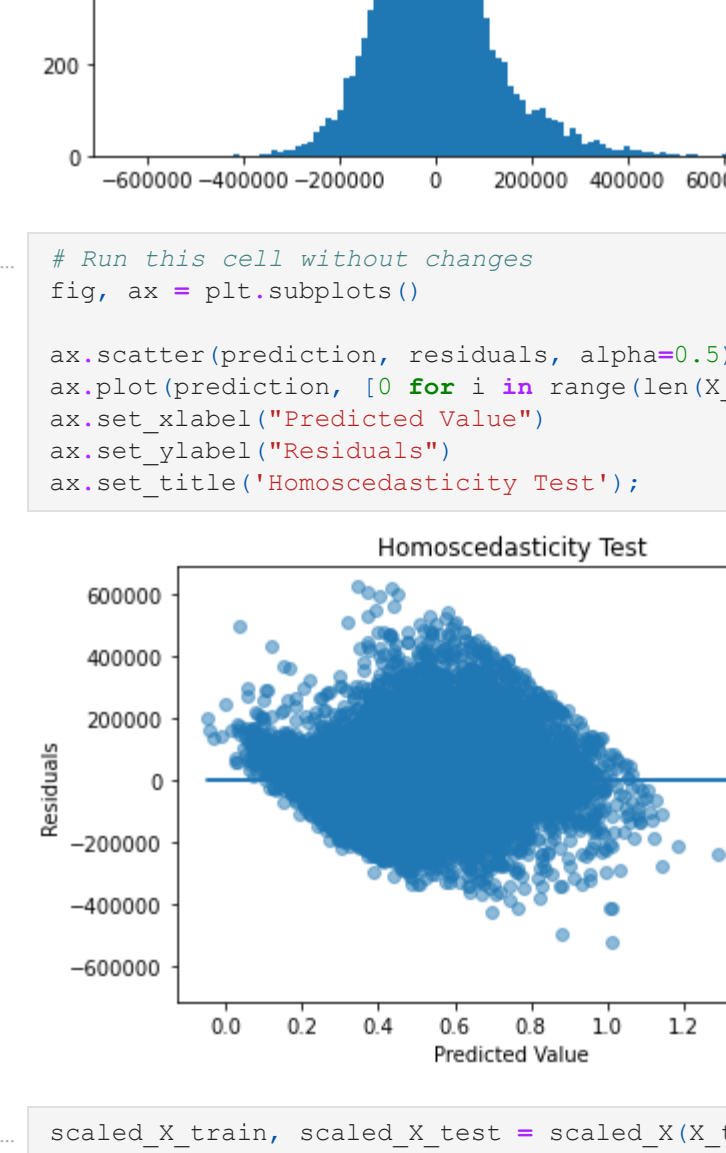
Out[582]: 0.694307715930957
```

Residual Tests

```
In [583]: prediction, residuals = residual_helper(third_model,X_train,y_train)

In [584]: sm.graphics.qqplot(residuals, dist=stats.norm, line='45', fit=True);

Out[584]:
```



```
In [585]: plt.hist(residuals, bins=100)

Out[585]:
```


OLS Regression Results

Dep. Variable:

price

R-squared:

0.695

Method:

Least Squares

Adj. R-squared:

0.972

Date:

Fri, 07-Jan-2022

Prob (F-statistic):

0.00

Time:

05:17:39

Log-Likelihood:

-2.0034e+05

No. Observations:

15329

AIC:

4.008e+05

DF Residuals:

15292

BIC:

4.010e+05

DF Model:

36

Covariance Type:

nonrobust

	coef	std err	t	P> t	[0.025	0.975]
const	4.778e+05	927.793	515.001	0.000	4.76e+05	4.8e+05
bedrooms	-7621.3444	1219.395	-6.250	0.000	-1e+04	-5231.186
bathrooms	2976.6549	1520.429	1.958	0.050	-3.567	5956.877
sqft_living	9.74e+04	2075.070	46.939	0.000	9.33e+04	1.01e+05
floors	6297.1748	1323.422	4.758	0.000	3703.111	8891.239
waterfront	5151.6933	988.907	5.209	0.000	3213.317	7090.069
view	1857e+04	1027.991	18.064	0.000	1.66e+04	2.06e+04
condition	2.541e+04	994.092	25.558	0.000	2.35e+04	2.74e+04
sqft_basement	-5926.4246	1283.988	-4.616	0.000	-8443.194	-3409.655
sqft_living_compared	-1518e+04	1272.797	-11.907	0.000	-1.77e+04	-1.27e+04
grade	6.046e+04	1539.644	39.268	0.000	5.47e+04	6.35e+04
Bellevue	6.759e+04	1363.330	49.580	0.000	6.49e+04	7.03e+04
Black Diamond	6630.0476	974.964	6.800	0.000	4719.002	8541.093
Bothell	1.408e+04	1017.004	13.843	0.000	1.21e+04	1.61e+04
Carnation	1.031e+04	980.426	10.512	0.000	8384.630	122e+04
Duvall	1.035e+04	1020.658	10.145	0.000	8353.791	124e+04
Enumclaw	1910.0521	1044.886	1.732	0.083	-238.049	3858.154
Fall City	1.156e+04	965.063	11.974	0.000	9663.733	1.34e+04
Federal Way	-3366.1730	1238.572	-2.718	0.007	-5793.921	-938.425
Issaquah	3.54e+04	1226.183	28.869	0.000	3.3e+04	3.78e+04
Kenmore	1.521e+04	1047.461	14.520	0.000	1.32e+04	1.73e+04
Kent	1199.5403	1365.811	0.878	0.380	-1477.813	3876.693
Kirkland	4.984e+04	1278.005	39.000	0.000	4.73e+04	5.23e+04
Maple Valley	6602.7026	1179.804	5.598	0.000	4290.146	8951.259
Medina	8660.6464	930.377	9.309	0.000	6836.997	1.05e+04
Mercer Island	3.828e+04	1021.955	37.457	0.000	3.63e+04	4.03e+04
North Bend	1.23e+04	1032.270	11.916	0.000	1.03e+04	1.41e+04
Redmond	5.154e+04	1310.080	39.340	0.000	4.9e+04	5.41e+04
Renton	1.987e+04	1462.436	13.586	0.000	1.7e+04	2.27e+04
Sammamish	4.075e+04	1265.375	32.202	0.000	3.83e+04	4.32e+04
Seattle	1.102e+05	2313.182	47.655	0.000	1.06e+05	1.15e+05
Snoqualmie	1.495e+04	1078.946	13.854	0.000	1.29e+04	1.71e+04
Vashon	1.024e+04	1009.192	10.142	0.000	8257.040	1.22e+04
Woodinville	2.701e+04	1137.784	23.735	0.000	2.48e+04	2.92e+04
spring	8520.3226	1189.738	7.162	0.000	6188.294	1.09e+04
summer	-556.8421	1181.079	-0.471	0.637	-2871.899	1758.214
winter	3376.8585	1122.481	3.008	0.003	1176.662	5577.055

Omnibus:

1549.817

Durbin-Watson:

2.003

Prob(Omnibus):

0.000

Jarque-Bera (JB):

3321.279

Skew:

0.640

Prob(JB):

0.00

Kurtosis:

4.888

Cond. No.

8.03

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

Model Assessment

The final model showed a marked improvement on the R-Squared score. The residuals look normally distributed, with still some skew. Heteroscedasticity remains an issue so I believe that calls the model into question. The p-values for "summer" in the categorical season is above the desired level, as is the cities of Kent and Enumclaw. I would weigh these against the low p-values of the features in the same category.

```
In [589]: results_summary = Stats_Mod_Helper(y_train,scaled_X_train)
         results_as_html = results_summary.tables[1].as_html()
         Results_df = pd.read_html(results_as_html, headers=0, index_col=0)[0]
```

```
In [590]: Results_df
```

	coef	std err	t	P> t	[0.025	0.975]
const	477800.00000	927.79300	515.00000	0.00000	476000.00000	480000.00000
bedrooms	-7621.34444	1219.39500	-6.25000	0.00000	-10000.00000	-5231.18600
bathrooms	2976.65490	1520.42900	1.95800	0.05000	-3.56700	5956.87700
sqft_living	97400.00000	2075.07000	46.93900	0.00000	93300.00000	101000.00000
floors	6297.17480	1323.42200	4.75800	0.00000	3703.11100	8891.23900
waterfront	5151.69330	988.90700	5.20900	0.00000	3213.31700	7090.06900
view	18570.00000	1027.99100	18.06400	0.00000	16600.00000	20600.00000
condition	25410.00000	994.09200	25.55800	0.00000	23500.00000	27400.00000
sqft_basement	-5926.42460	1283.98800	-4.61600	0.00000	-8443.19400	-3409.65500
sqft_living_compared	-15180.00000	1272.79700	-11.90700	0.00000	-17700.00000	-12700.00000
grade	60460.00000	1539.64400	39.26800	0.00000	57400.00000	63500.00000
Bellevue	67590.00000	1363.33000	49.58000	0.00000	64900.00000	70300.00000
Black Diamond	6630.04760	974.96400	6.80000	0.00000	4719.00200	8541.09300
Bothell	14080.00000	1017.00400	13.84300	0.00000	12100.00000	16100.00000
Carnation	10310.00000	980.42600	10.51200	0.00000	8384.63000	12200.00000
Duvall	10350.00000	1020.65800	10.14500	0.00000	8353.79100	12400.00000
Enumclaw	1910.05210	1044.88600	1.73200	0.08300	-238.04900	3858.15400
Fall City	11560.00000	965.06300	11.97400	0.00000	9663.73300	13400.00000
Federal Way	-3366.17300	1238.57200	-2.71800	0.00700	-5793.92100	-938.42500
Issaquah	35400.00000	1226.18300	28.86900	0.00000	33000.00000	37800.00000
Kenmore	15210.00000	1047.46100	14.52000	0.00000	13200.00000	17300.00000
Kent	1199.54030	1365.81100	0.87800	0.38000	-1477.81300	3876.69300
Kirkland	49840.00000	1278.00500	39.00000	0.00000	47900.00000	52300.00000
Maple Valley	6602.70260	1179.80400	5.59800	0.00000	4290.14600	8951.25900
Medina	8660.64640	930.37700	9.30900	0.00000	6836.99700	10500.00000
Mercer Island	38280.00000	1021.95500	37.45700	0.00000	36300.00000	40300.00000
North Bend	12300.00000	1032.27000	11.91600	0.00000	10300.00000	14300.00000
Redmond	51540.00000	1310.08000	39.34000	0.00000	49000.00000	54100.00000
Renton	19870.00000	1462.43600	13.58600	0.00000	17000.00000	22700.00000
Sammamish	40750.00000	1265.37500	32.20200	0.00000	38300.00000	43200.00000
Seattle	110200.00000	2313.18200	47.65500	0.00000	106000.00000	115000.00000
Snoqualmie	14950.00000	1078.94600	13.85400	0.00000	12800.00000	17100.00000
Vashon	10240.00000	1009.19200	10.14200	0.00000	8257.04000	12200.00000
Woodinville	27010.00000	1137.78400	23.73500	0.00000	24800.00000	29200.00000
spring	8520.32260	1189.73900	7.16200	0.00000	6188.29400	10900.00000
summer	-556.84210	1181.07900	-0.47100	0.63700	-2871.89900	1758.21400
winter	3376.85850	1122.48100	3.00800	0.00300	1176.66200	5577.05500

House Features Assessment

Scaled Comparison

```
In [591]: House_Sum = Results_df.loc[house_features,:]
```

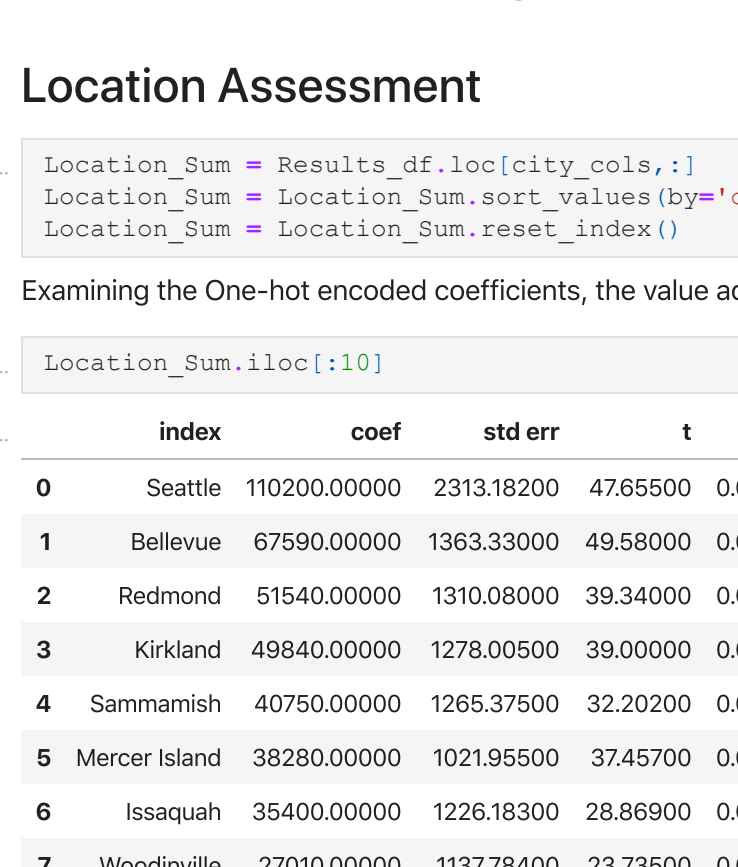
```
In [610]: House_Sum = House_Sum.sort_values(by='coef',ascending=False)
         House_Sum = House_Sum.reset_index()
         House_Sum
```

```
Out[610]:
```

	index	coef	std err	t	P> t	[0.025	0.975]
0	sqft_living	97400.00000	2075.07000	46.93900	0.00000	93300.00000	101000.00000
1	grade	60460.00000	1539.64400	39.26800	0.00000	57400.00000	63500.00000
2	condition	25410.00000	994.09200	25.55800	0.00000	23500.00000	27400.00000
3	view	18570.00000	1027.99100	18.06400	0.00000	16600.00000	20600.00000
4	floors	6297.17480	1323.42200	4.75800	0.00000	3703.11100	8891.23900
5	waterfront	5151.69330	988.90700	5.20900	0.00000	3213.31700	7090.06900
6	bathrooms	2976.65490	1520.42900	1.95800	0.05000	-3.56700	5956.87700
7	sqft_basement	-5926.42460	1283.98800	-4.61600	0.00000	-8443.19400	-3409.65500
8	bedrooms	-7621.34440	1219.39500	-6.25000	0.00000	-10000.00000	-5231.18600
9	sqft_living_compared	-15180.00000	1272.79700	-11.90700	0.00000	-17700.00000	-12700.00000

```
In [611]: sns.barplot(data=House_Sum.iloc[:10],x='index',y='coef', color='Green')
         plt.xlabel('City', fontsize=18)
         plt.ylabel('Value Added Dollars',fontsize=15)
         plt.xticks(rotation=45)
         plt.title('Top Ten Coefficients Per House Feature Scaled',size=15)
```

```
Out[611]: Text(0.5, 1.0, 'Top Ten Coefficients Per House Feature Scaled')
```



Based on the scaled data, an improvement of one standard deviation of the square-foot living space would add 97,400 dollars to the value of a home. The grade and condition column are the next ranked coefficients when scaled. View also ranks as another feature, but certainly not one that cannot be altered.

I would think it best to use these coefficients to look at their real dollar value.

Real Dollar Coefficients

```
In [608]: real_results_summary = Stats_Mod_Helper(y_train,X_train)
         real_results_as_html = real_results_summary.tables[1].as_html()
         real_results_df = pd.read_html(real_results_as_html, headers=0, index_col=0)[0]
```

```
In [594]: order = House_Sum.index
```

```
In [602]: real_results_df
         real_House_Sum = real_results_df.loc[house_features,:]
         real_House_Sum = real_House_Sum.reindex(index=order).reset_index()
```

In the unscaled version of the model, ranked by the scaled values, An additional square footage of living space would add 126 to a property with only a standard error of \$2.66. An additional floor would add 11,700. Certainly an improvement of Grade or Condition would yield value, but these may be harder to practically achieve.

```
In [604]: sns.barplot(data=real_House_Sum.iloc[:10],x='index',y='coef', color='Green')
         plt.xlabel('City', fontsize=18)
         plt.ylabel('Value Added Dollars',fontsize=15)
         plt.xticks(rotation=45)
         plt.title('Top Ten Coefficients Per House Feature Unscaled',size=15)
```

```
Out[604]: Text(0.5, 1.0, 'Top Ten Coefficients Per House Feature Unscaled')
```



Location Assessment

```
In [598]: Location_Sum = Results_df.loc[city_cols,:]
         Location_Sum = Location_Sum.sort_values(by='coef',ascending=False)
         Location_Sum = Location_Sum.reset_index()
```

```
In [597]: Location_Sum.iloc[:10]
```

	index	coef	std err	t	P> t	[0.025	0.975]
0	Seattle	110200.00000	2313.18200	47.65500	0.00000	106000.00000	115000.00000
1	Bellevue	67590.00000	1363.33000	49.58000	0.00000	64800.00000	70300.00000
2	Redmond	51540.00000	1310.08000	39.34000	0.00000	49000.00000	54100.00000
3	Kirkland	49840.00000	1278.00500	39.00000	0.00000	47900.00000	52300.00000
4	Sammamish	40750.00000	1265.37500	32.20200	0.00000	38300.00000	43200.00000
5	Mercer Island	38280.00000	1021.95500	37.45700	0.00000	36300.00000	40300.00000
6	Issaquah	35400.00000	1226.18300	28.86900	0.00000	33000.00000	37800.00000
7	Woodinville	27010.00000	1137.78400	23.73500	0.00000	24800.00000	29200.00000
8	Renton	19870.00000	1462.43600	13.58600	0.00000	17000.00000	22700.00000
9	Kenmore	15210.00000	1047.46100	14.52000	0.00000	13200.00000	17300.00000

```
In [598]: sns.barplot(data=Location_Sum.iloc[:10],x='index',y='coef', color='Green')
         plt.xlabel('City', fontsize=18)
         plt.ylabel('Value Added Dollars',fontsize=15)
         plt.xticks(rotation=45)
         plt.title('Top Ten King County City Coefficients',size=15)
```

```
Out[598]: Text(0.5, 1.0, 'Top Ten King County City Coefficients')
```



Conclusions

- The model showed an R-Squared of nearly 70%, with what appeared to be close to normality of residuals and some problems with Heteroskedasticity. Outliers were dropped from the dataset and I believe the model struggled with higher priced homes.
- I superior model could have dealt with location data differently, as its my belief location plays a strong role in the price. I think there might be a city interaction that could have been built into the model.
- As always more recent data, especially during and after Covid would be helpful.

```
In [ ]:
```