

Final Project Submission

Please fill out:

- Student name:
 - Student pace: self paced / part time / full time
 - Scheduled project review date/time:
 - Instructor name:
 - Blog post URL:
-
- **id** - unique identified for a house
 - **date** - date house was sold
 - **price** - price is prediction target
 - **bedrooms** - Number - of Bedrooms/House
 - **bathrooms** - of bathrooms/bedrooms
 - **sqft_living** - square footage of the home
 - **sqft_lots** - square footage of the lot
 - **floors** - floors (levels) in house
 - **waterfront** - House which has a view to a waterfront
 - **view** - has been viewed
 - **condition** - How good the condition is (Overall)
 - **grade** - Overall grade given to the housing unit, based on King County grading system
 - **sqft_above** - square footage of house apart from basement
 - **sqft_basement** - square footage of the basement
 - **yr_built** - Built Year
 - **yr_renovated** - Year when house was renovated
 - **zipcode** - zip
 - **lat** - Latitude coordinate
 - **long** - Longitude coordinate
 - **sqft_living15** - The square footage of interior housing living space for the nearest 15 neighbors
 - **sqft_lot15** - The square footage of the land lots of the nearest 15 neighbors

```
In [1]: import pandas as pd
import numpy as np

import matplotlib.pyplot as plt
import seaborn as sns
import plotly.graph_objects as go
import plotly.express as px

from statsmodels.formula.api import ols
import statsmodels.stats.outliers_influence import variance_inflation_factor
import statsmodels.api as sm
import scipy.stats as stats
from sklearn.model_selection import train_test_split
```

Overview of Data

```
In [2]: house_data = pd.read_csv('data/kc_house_data.csv')
house_data
```

	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view	...	grade	sqft_above	sqft_basement
Out[2]:	0	1/19/2005	219900	3	1.00	1180	5650	1.0	NaN	0.0	...	7	1180	
	1	6/14/100192	12/9/2014	538000.0	3	2.25	2570	7242	2.0	0.0	0.0	7	2170	
	2	563150400	2/25/2015	180000.0	2	1.00	770	10000	1.0	0.0	0.0	6	770	
	3	248720875	12/9/2014	604000.0	4	3.00	1960	5000	1.0	0.0	0.0	7	1050	
	4	1954400510	2/18/2015	510000.0	3	2.00	1680	8080	1.0	0.0	0.0	8	1680	

	21592	263000018	5/21/2014	360000.0	3	2.50	1530	1131	3.0	0.0	0.0	8	1530	
	21593	6600560120	2/23/2015	400000.0	3	2.50	2310	5813	2.0	0.0	0.0	8	2310	
	21594	1523300141	6/23/2014	402101.0	2	0.75	1020	1350	2.0	0.0	0.0	7	1020	
	21595	191310100	1/16/2015	400000.0	3	2.50	1600	2388	2.0	NaN	0.0	8	1600	
	21596	2523300157	10/15/2014	325000.0	2	0.75	1020	1076	2.0	0.0	0.0	7	1020	

21597 rows x 21 columns

```
In [3]: house_data.info()
house_data.shape
```

<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 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 int64
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: 3.8+ MB

```
Out[3]: (21597, 21)
```

```
In [4]: # columns to take a closer look at
# look at the unique values for price
house_data.price.unique()[1:5]
```

```
Out[4]: array([221900., 538000., 180000., 604000., 510000.]
```

```
In [5]: # sanity check missing values
house_data.isna().sum()
```

```
Out[5]: 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_renovated 3842
yr_built 0
zipcode 0
lat 0
long 0
sqft_living15 0
sqft_lot15 0
dtype: int64
```

```
In [6]: house_data.describe()
```

```
Out[6]:
```

	id	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
count	21597	21597	21597	21597	21597	21597	21597	21597	21597	21597	21597	21597	21597	21597	21597	21597	21597	21597	21597	21597
mean	580474.09	540296.66	3.373200	2.115826	2080.321850	1509941.64	1.494096	0.007596	0.233863	3.41										
std	2.876736e+09	3.673661e+05	1.000000	0.768984	918.106125	4.141264e+04	0.539683	0.086825	0.765686	1.0										
min	1.000102e+06	7.800000e+04	2.000000	0.500000	370.000000	5.200000e+02	1.000000	0.000000	0.000000	0.0										
25%	2.123049e+09	3.200000e+05	3.000000	1.750000	1430.000000	5.048000e+03	1.500000	0.000000	0.000000	3.0										
50%	3.904930e+09	4.500000e+05	3.000000	2.250000	1910.000000	7.618000e+03	1.500000	0.000000	0.000000	3.0										
75%	7.308900e+09	6.450000e+05	4.000000	2.500000	2550.000000	1.068500e+04	2.000000	0.000000	0.000000	4.0										
max	5.900000e+09	7.700000e+06	33.000000	8.000000	13540.000000	1.651359e+06	3.500000	1.000000	4.000000	5.0										

Cleaning & Exploring

- Plan
- Set the columns(waterfront, view, yr_renovated) with missing values to '0'
 - Replace column dtypes for necessary columns.

```
In [7]: # check the value counts for each column
for value in house_data.columns:
    print(value, house_data[value].value_counts(ascending=False)) # distribution of frequently occurring values
    print('')
```

```
id
195500620 3
1825069031 2
2019202220 2
112530456 2
1781500435 2
7812801125 1
4364700875 1
3021059726 1
88000205 1
1777500160 1
Name: id, Length: 21420, dtype: int64

date
6/23/2014 142
6/26/2014 131
7/10/2015 1
7/8/2014 127
4/27/2015 126
...
5/17/2014 1
5/24/2015 1
7/10/2015 1
8/30/2014 1
7/20/2014 1
Name: date, Length: 372, dtype: int64

price
350000.0 172
450000.0 172
500000.0 159
500000.0 152
425000.0 150
...
870515.0 1
336950.0 1
386100.0 1
176250.0 1
884746.0 1
Name: price, Length: 3622, dtype: int64

bedrooms
3 9824
4 6882
2 2760
5 1801
6 272
1 196
7 38
8 13
9 6
10 3
11 1
12 1
13 1
Name: bedrooms, dtype: int64

bathrooms
2.50 5377
1.00 3851
1.75 3048
2.25 2047
2.00 1930
1.50 1445
2.75 1185
3.00 753
3.50 731
3.25 589
3.75 185
4.00 136
4.25 79
0.75 71
4.75 23
5.00 21
5.25 13
5.50 9
1.25 9
6.00 6
5.75 4
0.50 4
7.00 2
6.25 2
6.75 2
6.50 2
7.50 1
7.75 1
Name: bathrooms, dtype: int64

sqft_living
1300 138
1400 135
1440 133
1660 129
1010 129
...
4970 1
2505 1
2793 1
4810 1
1975 1
Name: sqft_living, Length: 1034, dtype: int64

sqft_lot
5000 358
6000 290
4000 251
7200 220
7500 118
...
1448 1
38884 1
17313 1
884746 1
315374 1
Name: sqft_lot, Length: 9776, dtype: int64

floors
1.0 10473
2.0 8235
1.5 1910
3.0 611
2.5 161
3.5 7
Name: floors, dtype: int64

waterfront
0.0 18075
1.0 146
Name: waterfront, dtype: int64

view
0.0 19422
2.0 957
3.0 508
1.0 330
4.0 317
Name: view, dtype: int64

condition
3 14020
4 5677
5 1701
2 170
1 29
Name: condition, dtype: int64

grade
9 8974
8 6065
9 2615
6 3038
10 1134
11 399
5 242
12 89
4 27
13 13
3 1
Name: grade, dtype: int64

sqft_above
1300 212
1010 210
200 206
1220 192
1140 184
...
2601 1
440 1
2473 1
2441 1
1975 1
Name: sqft_above, Length: 942, dtype: int64

sqft_basement
0.0 12826
0 454
600.0 217
500.0 209
700.0 208
1248.0 1
588.0 1
1275.0 1
2300.0 1
374.0 1
Name: sqft_basement, Length: 304, dtype: int64

yr_built
2014 559
2006 453
2005 450
2004 423
2003 420
...
1901 29
1902 27
1955 24
1934 21
Name: yr_built, Length: 116, dtype: int64

yr_renovated
0.0 17011
2014.0 73
2005.0 31
2013.0 31
2007.0 30
...
1946.0 1
1959.0 1
1971.0 1
1951.0 1
1954.0 1
Name: yr_renovated, Length: 70, dtype: int64

zipcode
98103 602
98038 589
98115 583
98052 574
98117 553
...
98102 104
98010 100
98024 80
98148 57
98039 50
Name: zipcode, Length: 70, dtype: int64

lat
47.6624 17
47.6493 17
47.5322 17
47.6846 17
47.6713 16
...
47.6785 1
47.4162 1
47.3870 1
47.2213 1
47.2715 1
Name: lat, Length: 5033, dtype: int64

long
-122.290 115
-122.350 111
-122.362 104
-122.291 100
-122.372 99
...
-121.403 1
-121.804 1
-121.726 1
-121.835 1
-121.893 1
Name: long, Length: 751, dtype: int64

sqft_living15
1540 197
1440 195
1560 192
1500 180
1460 169
...
4890 1
2873 1
492 1
3193 1
2049 1
Name: sqft_living15, Length: 777, dtype: int64

sqft_lot15
5000 427
4000 356
6000 286
7200 210
4800 145
...
11036 1
8883 1
871200 1
809 1
6147 1
Name: sqft_lot15, Length: 8682, dtype: int64
```

```
In [8]: # id has multiple counts that means there are duplicate housing info on here
# is it safe to drop the duplicated data? check the #
house_data.duplicated('id', value_counts)
```

```
Out[8]: False    21420
True      177
dtype: int64
```

```
In [9]: # Fill missing values and fix the dtypes
house_data['waterfront'].fillna(0, inplace=True)
house_data['waterfront'] = house_data['waterfront'].astype('int64')
house_data['view'].fillna(0, inplace=True)
house_data['view'] = house_data['view'].astype('int64')
house_data['yr_renovated'].fillna(0, inplace=True)
house_data['yr_renovated'] = house_data['yr_renovated'].astype('int64')
house_data['sqft_basement'] = house_data['sqft_basement'].replace(' ', value=0)
house_data['sqft_basement'] = pd.to_numeric(house_data['sqft_basement'])
house_data['sqft_basement'] = house_data['sqft_basement'].astype('int64')
house_data['date'] = pd.to_datetime(house_data['date'])

# double check data summary
house_data.info()
house_data.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 datetime64[ns]
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 int64
9 view 21597 non-null int64
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 int64
14 yr_built 21597 non-null int64
15 yr_renovated 21597 non-null int64
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: datetime64[ns](1), float64(5), int64(15)
memory usage: 3.5 MB

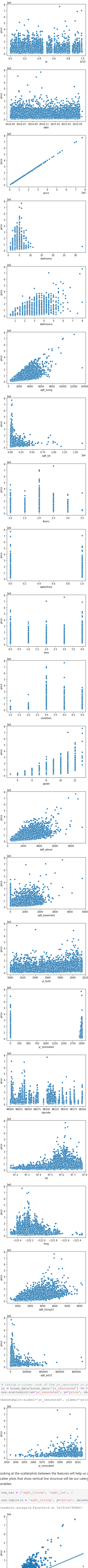
```
Out[9]: id 0
date 0
price 0
bedrooms 0
bathrooms 0
sqft_living 0
sqft_lot 0
floors 0
waterfront 0
view 0
condition 0
grade 0
sqft_above 0
sqft_basement 0
yr_built 0
yr_renovated 0
zipcode 0
lat 0
long 0
sqft_living15 0
sqft_lot15 0
dtype: int64
```

```
In [10]: # checking the raw distribution for each feature
house_data.hist(figsize=(15,15));
```



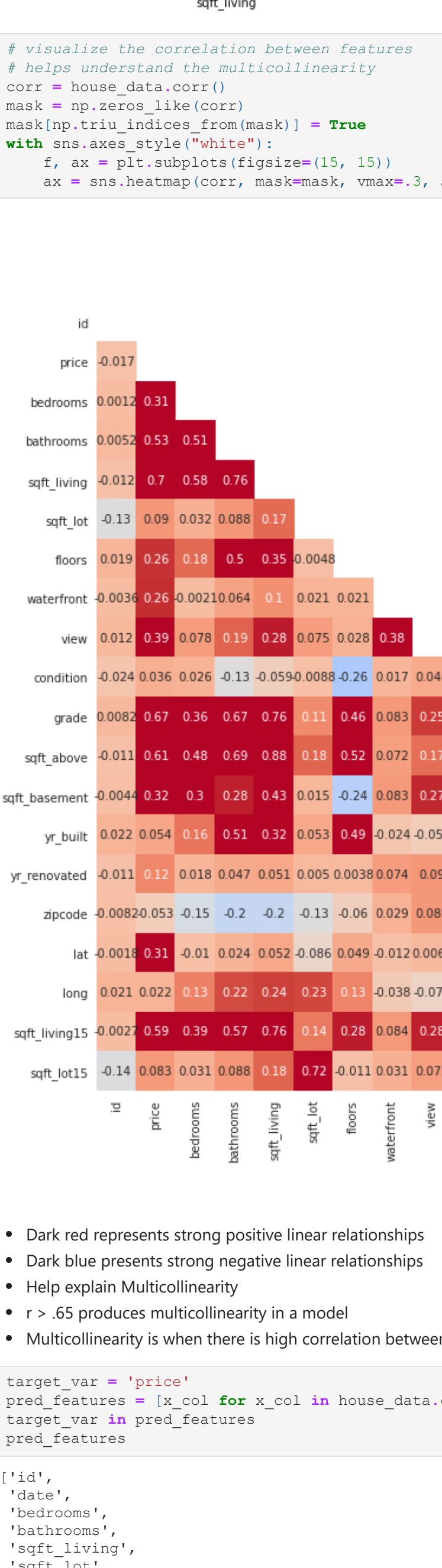
The distribution for each feature is either left or right skewed. This will not pass the normality assumption test. We will work on this when preparing the data.

```
In [11]: # taking a closer look at our features to closer identify which columns are categorical or numerical
# picture an imagination regression line
for x in house_data.columns:
    sns.scatterplot(x=x, y='price', data=house_data)
    plt.show()
```

In [12]: # taking a closer look of the yr_renovated to price
yr = house_data[house_data['yr_renovated'] != 0]
sns.scatterplot(x=yr_renovated, y=price, data=yr)

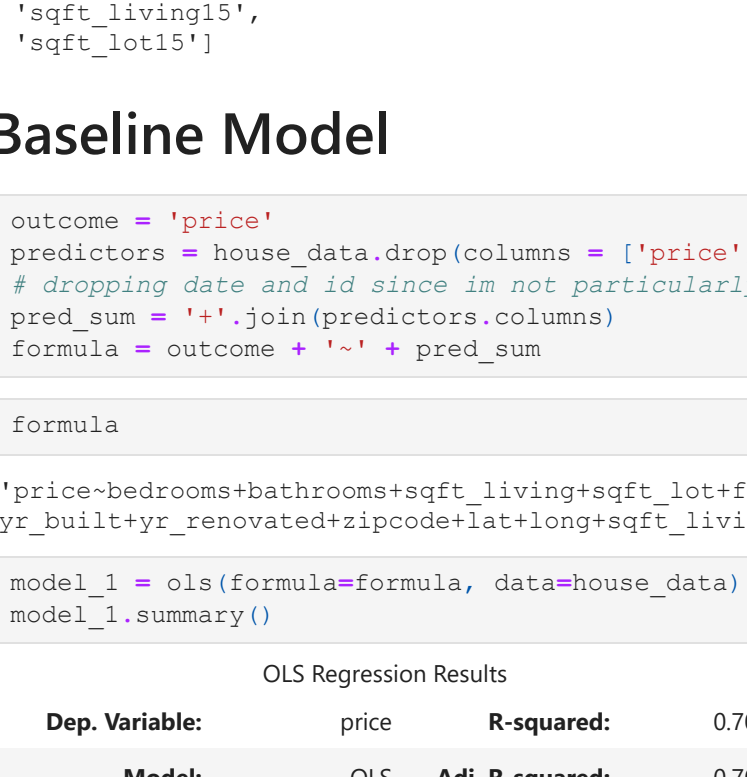
Out[12]: <AxesSubplot: xlabel='yr_renovated', ylabel='price'>



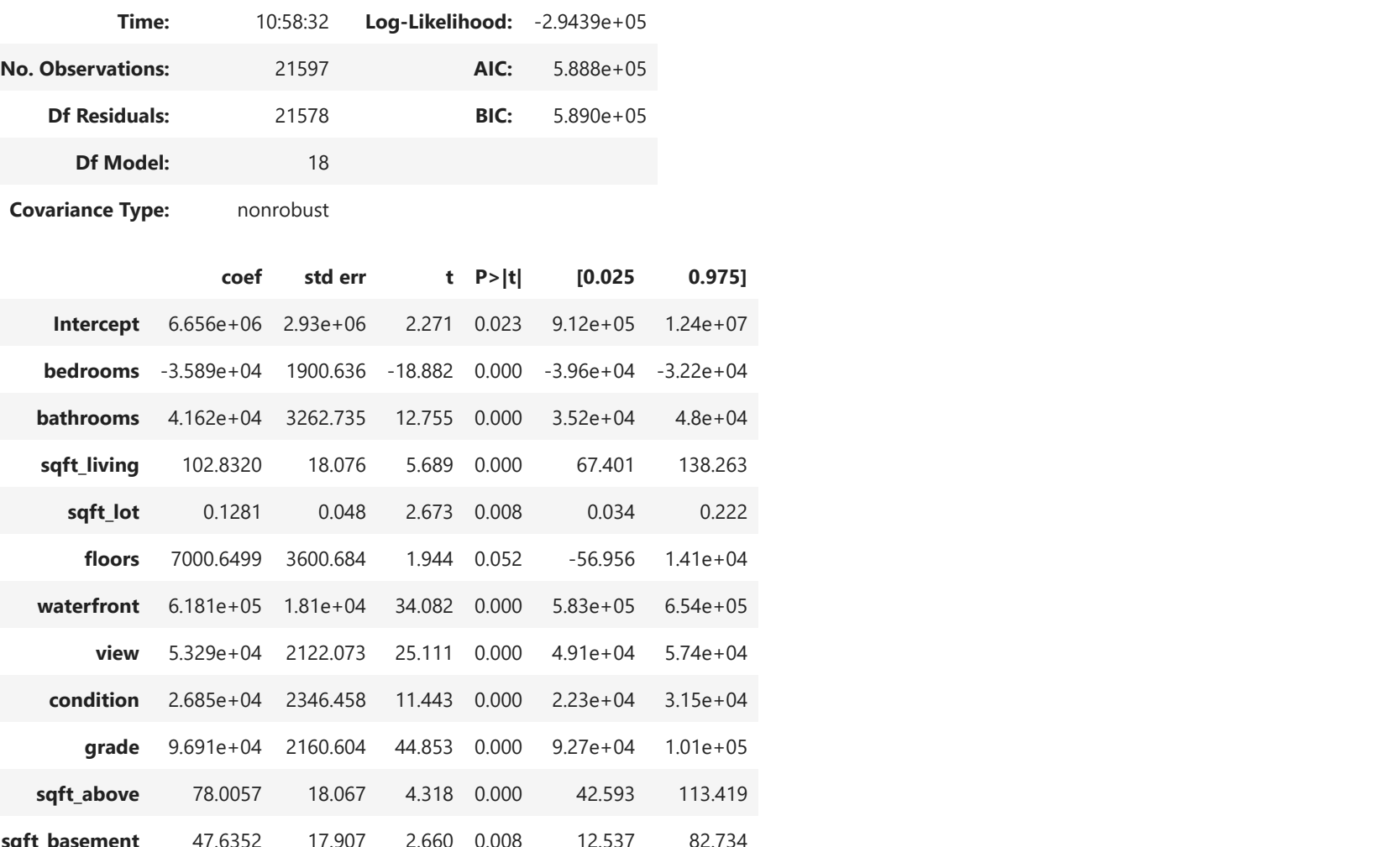
Looking at the scatterplots between the features will help us determine the numerical and categorical variables in regards to price. The scatter plots that show vertical line structure will be our categorical variables while the graphs with a cloud structure will be our numerical variables.

In [13]: num_var = ['sqft_living', 'sqft_lot',]
sns.lmplot(x = 'sqft_living', y='price', data=house_data)

Out[13]: <seaborn.axisgrid.FacetGrid at 0x910e7958b0>



In [14]: # visualize the correlation between features
helps understand the multicollinearity
corr = house_data.corr()
mask = np.zeros_like(corr)
mask[np.triu_indices_from(mask)] = True
with sns.axes_style("white"):
f, ax = plt.subplots(figsize=(15, 15))
ax = sns.heatmap(corr, mask=mask, vmax=.3, square=True, cmap="coolwarm", annot=True)



- Dark red represents strong positive linear relationships
- Dark blue represents strong negative linear relationships
- Help explain Multicollinearity
- > .45 produces multicollinearity in a model
- Multicollinearity is 'when there is high correlation between three or more variables'

In [15]: target_var = 'price'
pred_features = [x_col for x_col in house_data.columns if x_col != target_var]
target_var in pred_features
pred_features

Out[15]: ['id',
'date',
'bedrooms',
'bathrooms',
'sqft_living',
'sqft_lot',
'waterfront',
'view',
'condition',
'grade',
'sqft_above',
'sqft_baseament',
'yr_built',
'yr_renovated',
'zipcode',
'lat',
'long',
'sqft_living15',
'sqft_lot15']

Baseline Model

In [16]: outcome = 'price'
pred_features = [x_col for x_col in house_data.columns if x_col != target_var]
dropping date and id since in not particularly interested in the predictions of these features
pred_sum = '+'.join(pred_features.columns)
formula = outcome + ' <= > ' + pred_sum

In [17]: formula

Out[17]: 'price~bedrooms+bathrooms+sqft_living+sqft_lot+floors+waterfront+view+condition+grade+sqft_above+sqft_baseament+yr_built+yr_renovated+zipcode+lat+long+sqft_living15+sqft_lot15'

In [18]: model_1 = ols(formula=formula, data=house_data).fit()
model_1.summary()

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 2.15e+08. This might indicate that there are strong multicollinearity or other numerical problems.

In [19]: # use a qqplot to check for normality
ha clearly we can't use this model
fig = sm.graphics.qqplot(model_1.resid, dist='state.norm, line='45', fit=True)



In []:

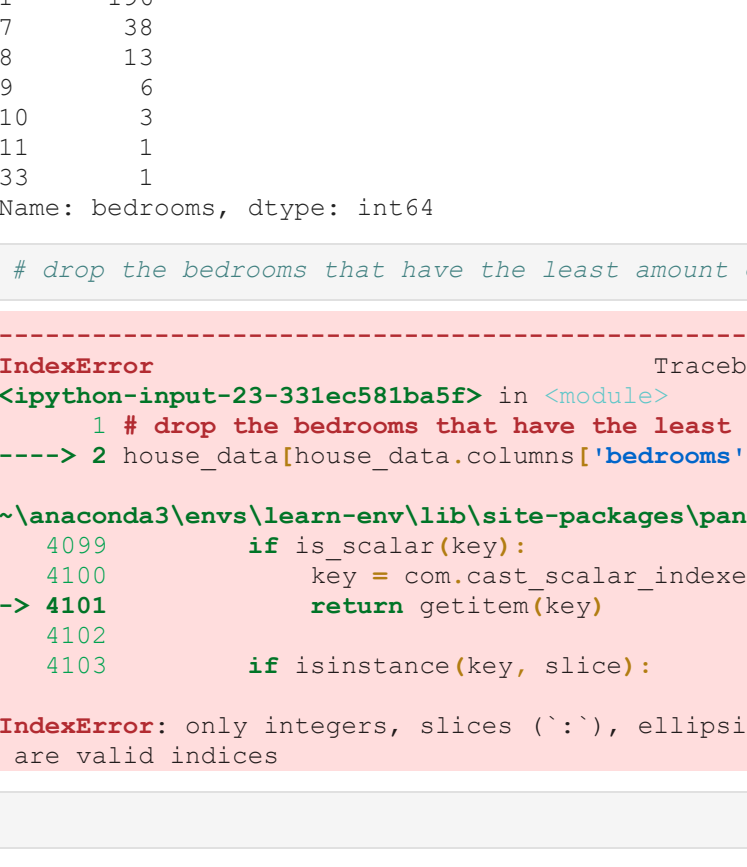
In [20]:

In []:

Dropping Outliers

Bedrooms

In [21]: sns.lmplot(x='bedrooms', y='price', data=house_data)
plt.show()



In [22]: house_data['bedrooms'].value_counts()

Out[22]: 3 2924
4 6882
5 1601
6 272
7 196
8 38
9 13
10 3
11 1
33 1
Name: bedrooms, dtype: int64

In [23]: # drop the bedrooms that have the least amount of information

```
----- Traceback (most recent call last)
<ipython-input-23-31ec581ba5> in <module>
      1 # drop the bedrooms that have the least amount of information
----> 2 house_data[house_data.columns['bedrooms'] == 33]

~\anaconda3\envs\learn-env\lib\site-packages\pandas\core\indexes\base.py in getitem(self, key)
    4099         if is_scalar(key):
    4100             key = com.cast_scalar_indexer(key, warn_float=True)
-> 4101             return getitem(key)
    4102         else:
    4103             if isinstance(key, slice):
IndexError: only integers, slices (':'), ellipsis ('...'), numpy.newaxis ('None') and integer or boolean arrays are valid indices
```

In []:

In []:

In []:

In []: corr = house_data[pred_features].corr()

- Dark red represents strong positive linear relationships
- Dark blue represents strong negative linear relationships
- Help explain Multicollinearity

In []: # check outliers and remove them
look at dates maybe add a months colm
decide which features to take a closer look at

Feature Engineering

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []: