

# How can you improve the value of your home?

The goal of this project is to examine how an existing homeowner can improve the value of their home. Using King County, WA sales data from 2014-15, I will determine which features are most valuable.

```
In [970]: 1 import pandas as pd
          2 import numpy as np
          3 import matplotlib.pyplot as plt
          4 import seaborn as sns
          5 from scipy import stats
          6 import statsmodels.api as sm
          7 from statsmodels.formula.api import ols
          8 from sklearn.model_selection import train_test_split
          9 from sklearn.model_selection import cross_validate, ShuffleSplit
         10 from sklearn.linear_model import LinearRegression
         11 from sklearn.metrics import mean_squared_error
         12 df = pd.read_csv('kc_house_data.csv')
         13 df
```

```
Out[970]:
```

	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	wat
0	7129300520	10/13/2014	221900.0	3	1.00	1180	5650	1.0	
1	6414100192	12/9/2014	538000.0	3	2.25	2570	7242	2.0	
2	5631500400	2/25/2015	180000.0	2	1.00	770	10000	1.0	
3	2487200875	12/9/2014	604000.0	4	3.00	1960	5000	1.0	
4	1954400510	2/18/2015	510000.0	3	2.00	1680	8080	1.0	
...	...	...	...	...	...	...	...	...	
21592	263000018	5/21/2014	360000.0	3	2.50	1530	1131	3.0	
21593	6600060120	2/23/2015	400000.0	4	2.50	2310	5813	2.0	
21594	1523300141	6/23/2014	402101.0	2	0.75	1020	1350	2.0	
21595	291310100	1/16/2015	400000.0	3	2.50	1600	2388	2.0	
21596	1523300157	10/15/2014	325000.0	2	0.75	1020	1076	2.0	

21597 rows × 21 columns

Let's examine the columns:

In [971]:

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

Let's create a 'Month' columns:

In [972]:

```
1 df['Month'] = pd.DatetimeIndex(df['date']).month
```

Out[972]:

```
0      10
1      12
2       2
3      12
4       2
..
21592   5
21593   2
21594   6
21595   1
21596  10
Name: Month, Length: 21597, dtype: int64
```

In [973]:

Out[973]:

```

5      2414
4      2229
7      2211
6      2178
8      1939
10     1876
3      1875
9      1771
--      ---

```

Eliminating irrelevant columns:

```

In [974]: 1 relevant_columns = ['price', 'bedrooms', 'bathrooms', 'sqft_living', 'floor
2
3 df = df.loc[:, relevant_columns]
4 df

```

Out[974]:

	price	bedrooms	bathrooms	sqft_living	floors	condition	yr_renovated	Month	gr
0	221900.0	3	1.00	1180	1.0	Average	0.0	10	Aver
1	538000.0	3	2.25	2570	2.0	Average	1991.0	12	Aver
2	180000.0	2	1.00	770	1.0	Average	NaN	2	6 I Aver
3	604000.0	4	3.00	1960	1.0	Very Good	0.0	12	Aver
4	510000.0	3	2.00	1680	1.0	Average	0.0	2	8 G
...	...	...	...	...	...	...	...	...	
21592	360000.0	3	2.50	1530	3.0	Average	0.0	5	8 G
21593	400000.0	4	2.50	2310	2.0	Average	0.0	2	8 G
21594	402101.0	2	0.75	1020	2.0	Average	0.0	6	Aver
21595	400000.0	3	2.50	1600	2.0	Average	0.0	1	8 G
21596	325000.0	2	0.75	1020	2.0	Average	0.0	10	Aver

21597 rows × 9 columns

Correlation Heat Map:

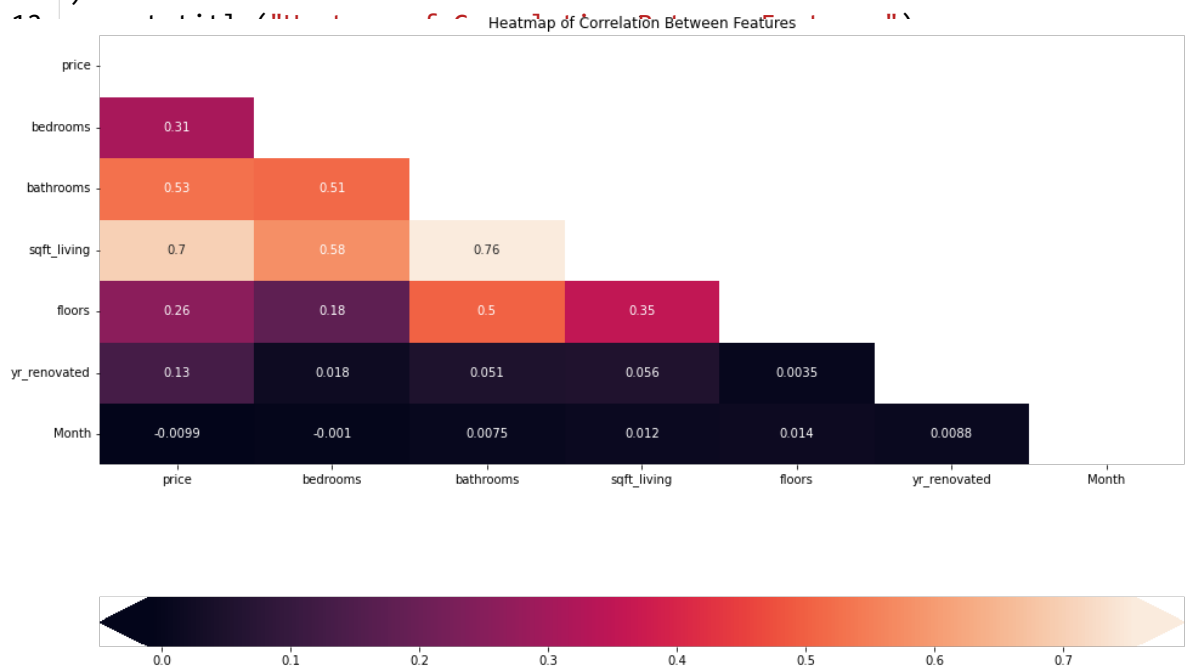
```

In [975]: 1 df_corr = df.corr()
2
3 fig, ax = plt.subplots(figsize=(16, 10))
4
5 sns.heatmap(
6     data=df_corr,
7     mask=np.triu(np.ones_like(df_corr, dtype=bool)),
8     ax=ax,
9     annot=True,

```

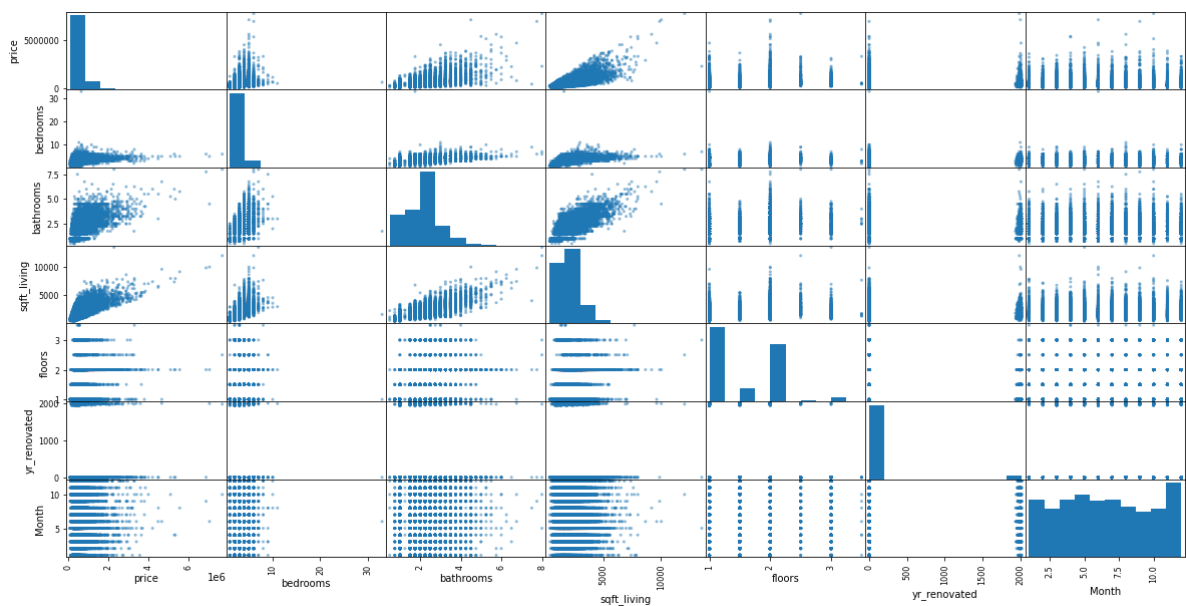
```
10 cbar_kws={"label": "Correlation", "orientation": "horizontal", "pad":
```

```
11 )
```



'Sqft\_Living' has the highest correlation with 'price'. Let's look at scatter plots and histograms to check the linearity and whether each variable is uniformly distributed.

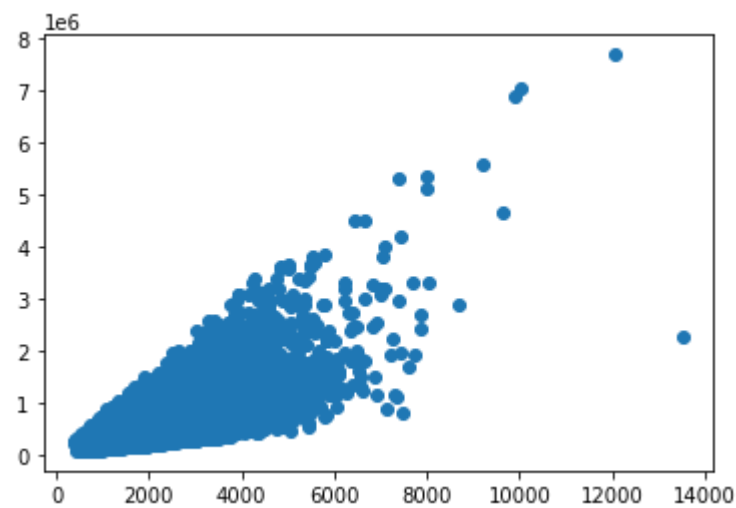
In [976]:



Let's take a closer look at 'sqft\_living' and price:

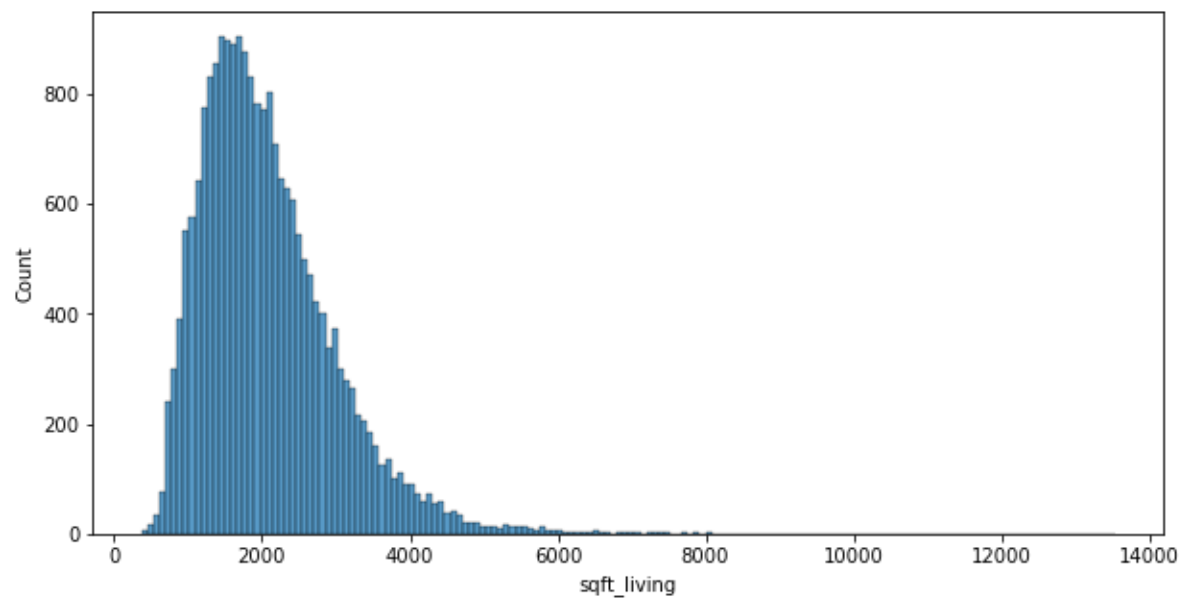
In [977]:

Out[977]: &lt;matplotlib.collections.PathCollection at 0x2873cc999a0&gt;



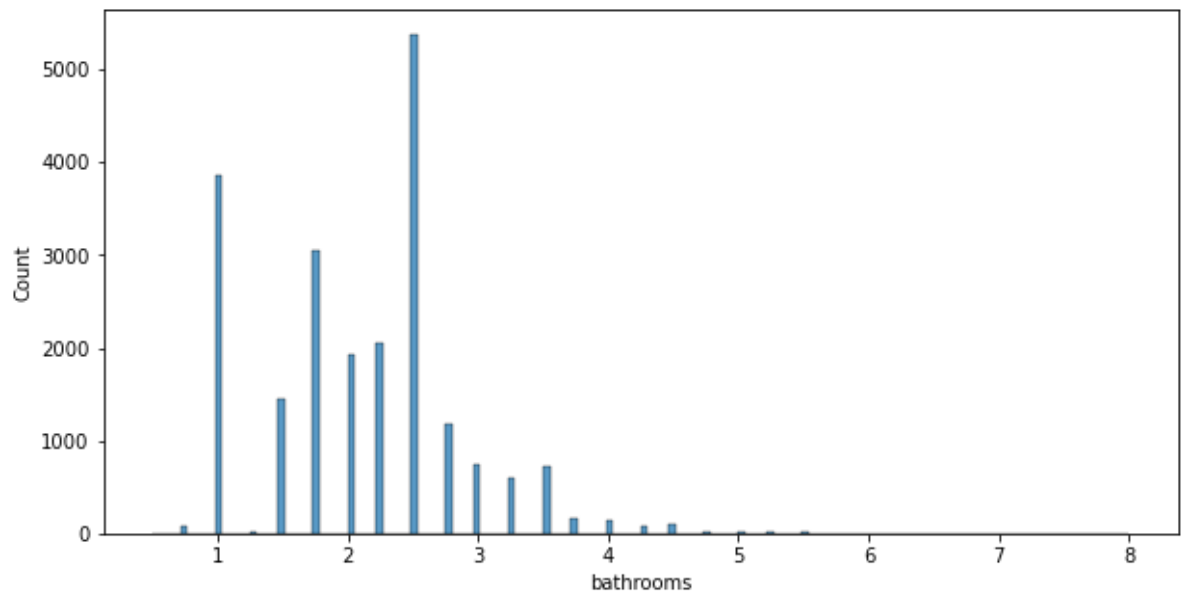
In [978]:

1 fig, ax = plt.subplots(figsize=(10, 5))

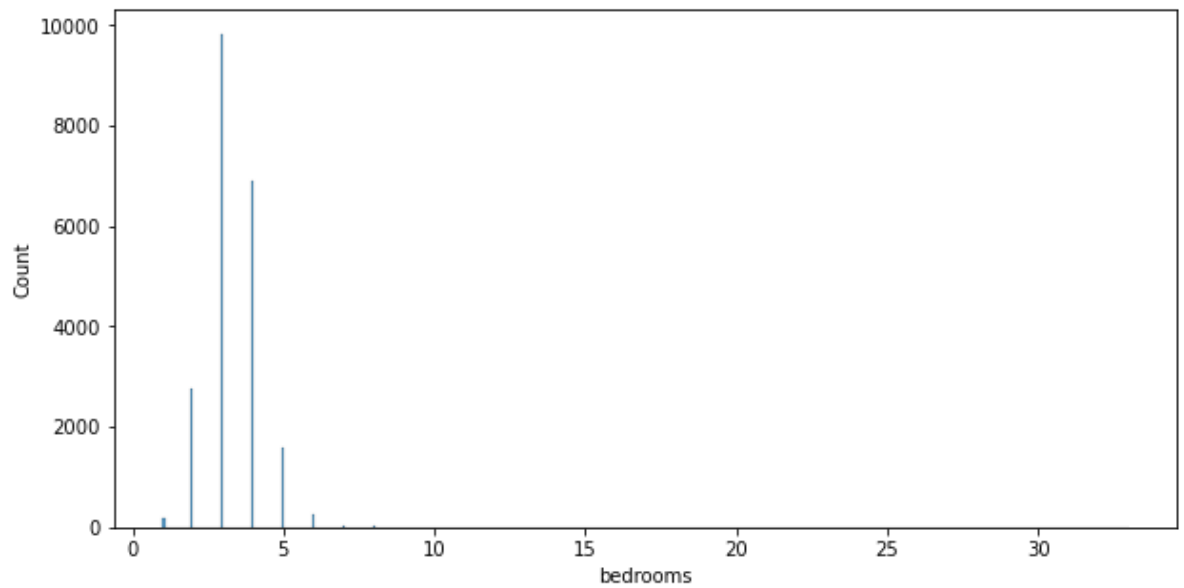


The relationship is certainly linear but it isn't normally distributed. Let's look at the distributions of some of the other variables.

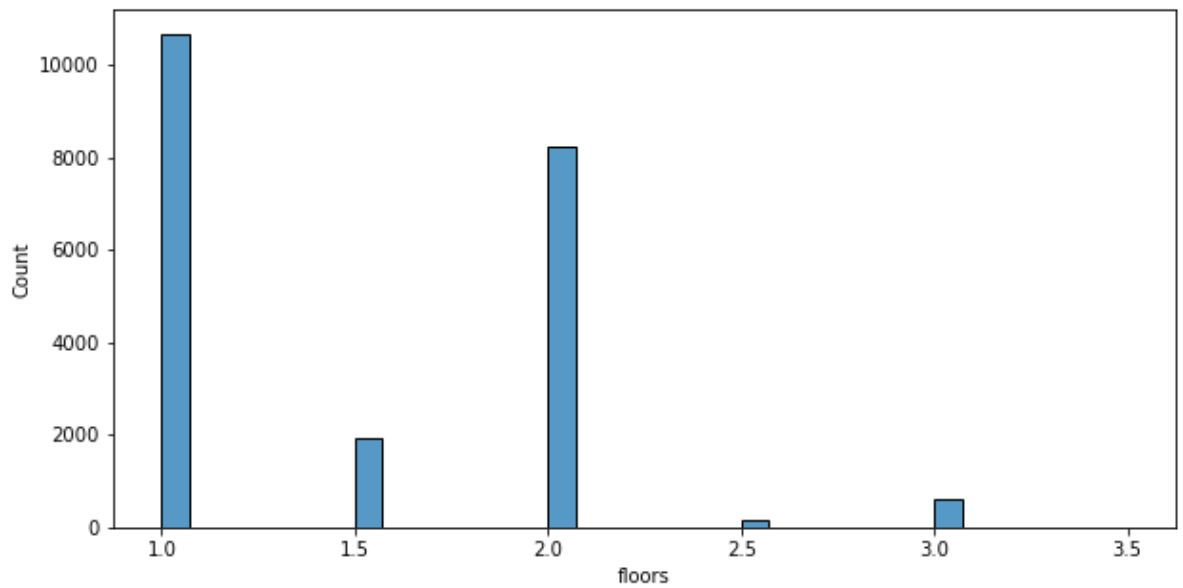
```
In [979]: 1 fig, ax = plt.subplots(figsize=(10, 5))
```



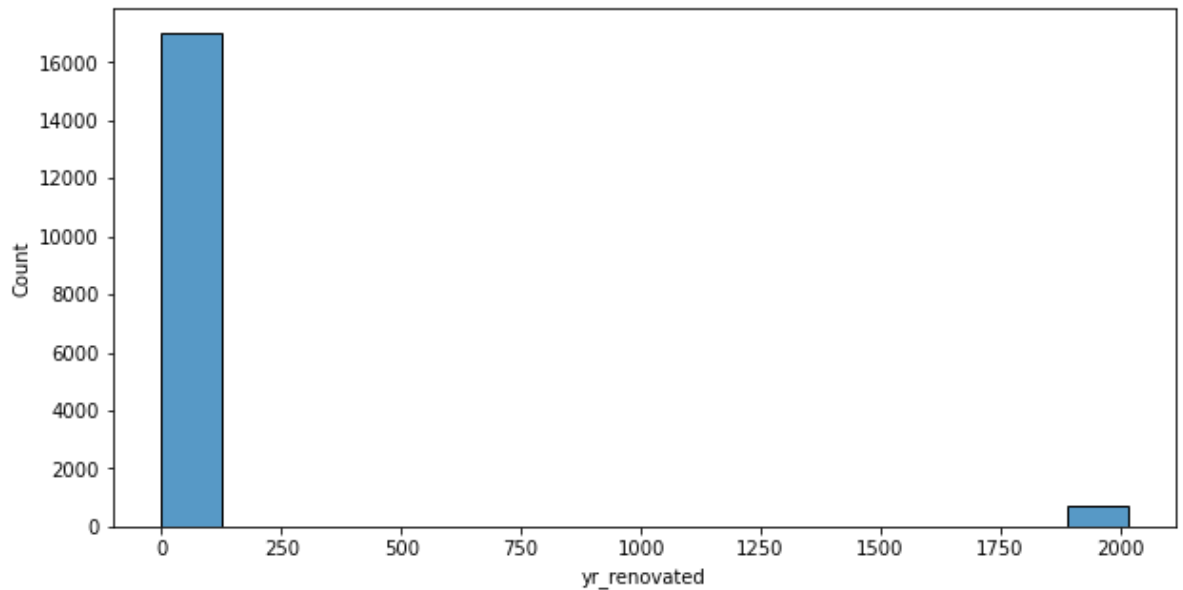
```
In [980]: 1 fig, ax = plt.subplots(figsize=(10, 5))
```



```
In [981]: 1 fig, ax = plt.subplots(figsize=(10, 5))
```



```
In [982]: 1 fig, ax = plt.subplots(figsize=(10, 5))
```



Ok, so it looks like it's best to use 'sqft\_living' for our independent variable. So, let's do it:

```
In [983]: 1 X_baseline= df['sqft_living']  
2 y_baseline= df['price']
```

Out[983]:

```

0      1180
1      2570
2       770
3      1960

```

```
In [984]: 1 X_baseline= sm.add_constant(X_baseline)
          2
```

C:\Users\isaia\anaconda3\lib\site-packages\statsmodels\tsa\tsatools.py:142: FutureWarning: In a future version of pandas all arguments of concat except for the argument 'objs' will be keyword-only

```
x = pd.concat(x[::order], 1)
```

```
In [985]:
```

Out[985]: OLS Regression Results

<b>Dep. Variable:</b>	price	<b>R-squared:</b>	0.493
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.493
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	2.097e+04
<b>Date:</b>	Thu, 19 May 2022	<b>Prob (F-statistic):</b>	0.00
<b>Time:</b>	13:02:17	<b>Log-Likelihood:</b>	-3.0006e+05
<b>No. Observations:</b>	21597	<b>AIC:</b>	6.001e+05
<b>Df Residuals:</b>	21595	<b>BIC:</b>	6.001e+05
<b>Df Model:</b>	1		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>const</b>	-4.399e+04	4410.023	-9.975	0.000	-5.26e+04	-3.53e+04
<b>sqft_living</b>	280.8630	1.939	144.819	0.000	277.062	284.664

<b>Omnibus:</b>	14801.942	<b>Durbin-Watson:</b>	1.982
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	542662.604
<b>Skew:</b>	2.820	<b>Prob(JB):</b>	0.00
<b>Kurtosis:</b>	26.901	<b>Cond. No.</b>	5.63e+03

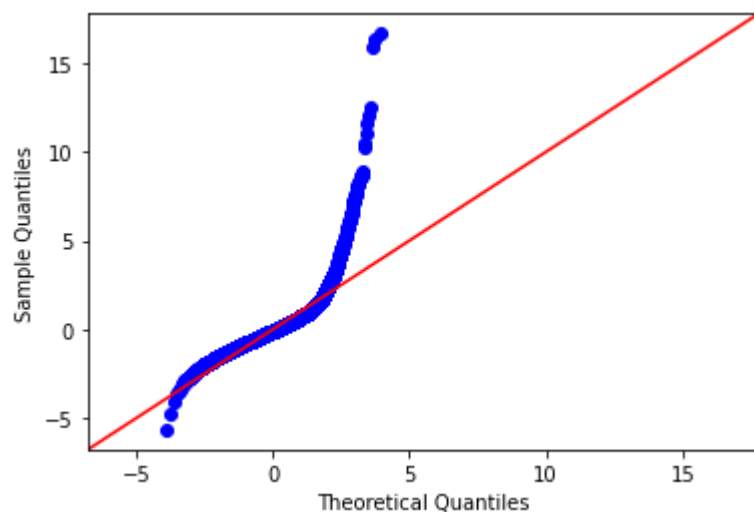
Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 5.63e+03. This might indicate that there are strong multicollinearity or other numerical problems.

```
In [986]: 1 residuals = baseline_model.resid
          2 fig = sm.graphics.qqplot(residuals, dist=stats.norm, fit=True, line='45',)
```

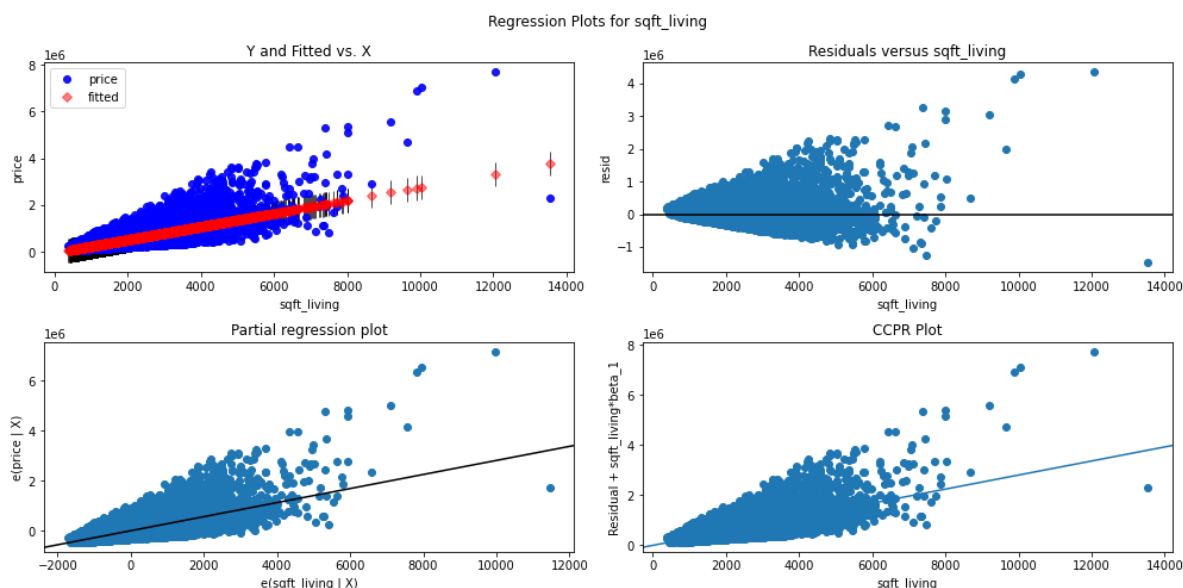


C:\Users\isaia\anaconda3\lib\site-packages\statsmodels\graphics\gofplots.py:93: UserWarning: marker is redundantly defined by the 'marker' keyword argument and the fmt string "bo" (-> marker='o'). The keyword argument will take pr



In [987]:

```
plt.figure(figsize=(10, 10))
plt.subplot(2, 2, 1)
plt.scatter(price, sqft_living)
plt.title('Y and Fitted vs. X')
plt.legend(['price', 'fitted'])
plt.subplot(2, 2, 2)
plt.scatter(resid, sqft_living)
plt.title('Residuals versus sqft_living')
plt.subplot(2, 2, 3)
plt.scatter(e(price | X), e(sqft_living | X))
plt.title('Partial regression plot')
plt.subplot(2, 2, 4)
plt.scatter(resid + sqft_living*beta_1, sqft_living)
plt.title('CCPR Plot')
```



Ok, it really isn't a great fit at all and the residuals are not homoscedastic. ***Let's improve by adding more features.***

## Data Preparation ,Cleaning, and Preprocessing

So, before we add more features, let's get it ready for processing.

In [988]:

```
from sklearn.preprocessing import StandardScaler
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21597 entries, 0 to 21596
Data columns (total 9 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   price           21597 non-null  float64
1   bedrooms        21597 non-null  int64
2   bathrooms       21597 non-null  float64
3   sqft_living     21597 non-null  int64
4   floors          21597 non-null  float64
5   condition       21597 non-null  object
6   yr_renovated    17755 non-null  float64
7   Month           21597 non-null  int64
8   grade           21597 non-null  object
```

Fill in null values for 'yr\_renovated'. 0 is actually the mode of this column so it's a good pick.

```
In [989]: 1 df['yr_renovated'].fillna(0.0, inplace=True)
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21597 entries, 0 to 21596
Data columns (total 9 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   price           21597 non-null  float64
1   bedrooms        21597 non-null  int64
2   bathrooms       21597 non-null  float64
3   sqft_living     21597 non-null  int64
4   floors          21597 non-null  float64
5   condition       21597 non-null  object
6   yr_renovated    21597 non-null  float64
7   Month           21597 non-null  int64
8   grade           21597 non-null  object
dtypes: float64(4), int64(3), object(2)
memory usage: 1.5+ MB
```

```
In [990]: 1 df['yr_renovated'].isna().value_counts()
```

```
Out[990]: False      21597
Name: yr_renovated, dtype: int64
```

```
In [991]:
```

```
Out[991]:
```

	price	bedrooms	bathrooms	sqft_living	floors	condition	yr_renovated	Month	gr
0	221900.0	3	1.00	1180	1.0	Average	0.0	10	Aver
1	538000.0	3	2.25	2570	2.0	Average	1991.0	12	Aver
2	180000.0	2	1.00	770	1.0	Average	0.0	2	61 Aver
3	604000.0	4	3.00	1960	1.0	Very Good	0.0	12	Aver

	price	bedrooms	bathrooms	sqft_living	floors	condition	yr_renovated	Month	gr
4	510000.0	3	2.00	1680	1.0	Average	0.0	2	8 G
...	...	...	...	...	...	...	...	...	...
21592	360000.0	3	2.50	1530	3.0	Average	0.0	5	8 G

Let's make a model with only numeric features:

In [992]:

In [993]:

Out[993]:

	price	bedrooms	bathrooms	sqft_living	floors	yr_renovated
0	221900.0	3	1.00	1180	1.0	0.0
1	538000.0	3	2.25	2570	2.0	1991.0
2	180000.0	2	1.00	770	1.0	0.0
3	604000.0	4	3.00	1960	1.0	0.0
4	510000.0	3	2.00	1680	1.0	0.0
...	...	...	...	...	...	...
21592	360000.0	3	2.50	1530	3.0	0.0
21593	400000.0	4	2.50	2310	2.0	0.0
21594	402101.0	2	0.75	1020	2.0	0.0
21595	400000.0	3	2.50	1600	2.0	0.0
21596	325000.0	2	0.75	1020	2.0	0.0

21597 rows × 6 columns

In [994]:

```
1 X_numeric= Numeric.drop('price', axis=1)
```

In [995]:

```
1 X_numeric = sm.add_constant(X_numeric)
```

C:\Users\isaia\anaconda3\lib\site-packages\statsmodels\tsa\tsatools.py:142: FutureWarning: In a future version of pandas all arguments of concat except for the argument 'objs' will be keyword-only  
x = pd.concat(x[::order], 1)

In [996]:

Out[996]:

OLS Regression Results

<b>Dep. Variable:</b>	price	<b>R-squared:</b>	0.513
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.513
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	4556.
<b>Date:</b>	Thu, 19 May 2022	<b>Prob (F-statistic):</b>	0.00

<b>Time:</b>	13:02:22	<b>Log-Likelihood:</b>	-2.9961e+05				
<b>No. Observations:</b>	21597	<b>AIC:</b>	5.992e+05				
<b>Df Residuals:</b>	21591	<b>BIC:</b>	5.993e+05				
<b>Df Model:</b>	5						
<b>Covariance Type:</b>	nonrobust						
	<b>coef</b>	<b>std err</b>	<b>t</b>	<b>P&gt; t </b>	<b>[0.025</b>	<b>0.975]</b>	
<b>const</b>	7.171e+04	7661.047	9.361	0.000	5.67e+04	8.67e+04	
<b>bedrooms</b>	-5.78e+04	2344.310	-24.654	0.000	-6.24e+04	-5.32e+04	
<b>bathrooms</b>	5910.5602	3808.216	1.552	0.121	-1553.825	1.34e+04	
<b>sqft_living</b>	308.7572	3.083	100.162	0.000	302.715	314.799	
<b>floors</b>	2104.8528	3761.480	0.560	0.576	-5267.927	9477.632	
<b>yr_renovated</b>	81.2053	4.799	16.920	0.000	71.798	90.612	
<b>Omnibus:</b>	14310.293	<b>Durbin-Watson:</b>	1.987				
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	479360.010				
<b>Skew:</b>	2.709	<b>Prob(JB):</b>	0.00				
<b>Kurtosis:</b>	25.435	<b>Cond. No.</b>	1.04e+04				

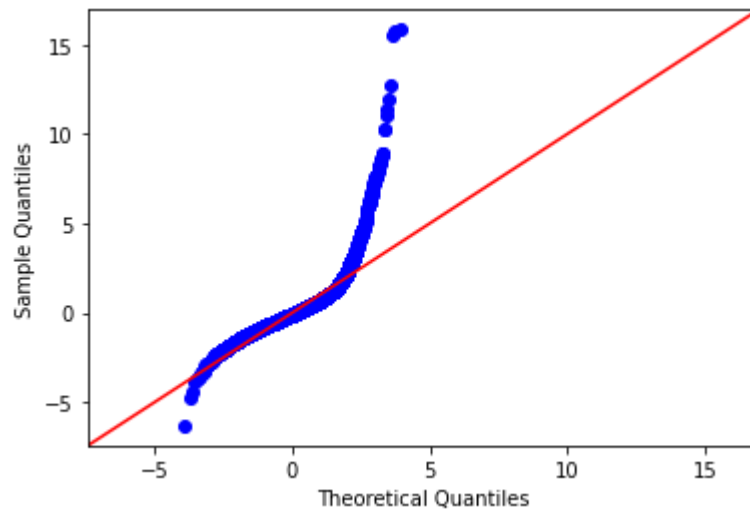
## Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 1.04e+04. This might indicate that there are strong multicollinearity or other numerical problems.

In [997]:

```
C:\Users\isaia\anaconda3\lib\site-packages\statsmodels\graphics\gofplots.py:9
93: UserWarning: marker is redundantly defined by the 'marker' keyword argume
nt and the fmt string "bo" (-> marker='o'). The keyword argument will take pr
ecedence.
```

```
ax.plot(x, y, fmt, **plot_style)
```



So, this looks ok but there are problems as the quantiles increase too dramatically and skew the line. This is sort of similar to the fat tails of the sqft\_living histogram from earlier. Log transforming might help it out. But, let's look at the categorical values first.

## One Hot Encoding categories

In [998]:

```
Out[998]: 7 Average      8974
8 Good      6065
9 Better    2615
6 Low Average 2038
10 Very Good 1134
11 Excellent 399
5 Fair      242
12 Luxury   89
4 Low       27
13 Mansion  13
3 Poor       1
Name: grade, dtype: int64
```

In [999]:

```
1 new_grades = {'3 Poor':3, '4 Low':4, '5 Fair':5, '6 Low Average':6, '7 Ave
2 '10 Very Good':10, '11 Excellent':11, '12 Luxury':12, '13 Mar
```

In [1000]:

Out[1000]:

```
price bedrooms bathrooms sqft_living floors condition yr_renovated Month grad
```

	price	bedrooms	bathrooms	sqft_living	floors	condition	yr_renovated	Month	grad
0	221900.0	3	1.00	1180	1.0	Average	0.0	10	
1	538000.0	3	2.25	2570	2.0	Average	1991.0	12	
2	180000.0	2	1.00	770	1.0	Average	0.0	2	
3	604000.0	4	3.00	1960	1.0	Very Good	0.0	12	
4	510000.0	3	2.00	1680	1.0	Average	0.0	2	
...	...	...	...	...	...	...	...	...	...
21592	360000.0	3	2.50	1530	3.0	Average	0.0	5	
21593	400000.0	4	2.50	2310	2.0	Average	0.0	2	
21594	402101.0	2	0.75	1020	2.0	Average	0.0	6	
21595	400000.0	3	2.50	1600	2.0	Average	0.0	1	
21596	325000.0	2	0.75	1020	2.0	Average	0.0	10	

```
In [1001]: 1 from sklearn.preprocessing import OneHotEncoder
2 Conditions = df[['condition']]
3 ohe = OneHotEncoder(categories="auto", sparse=False, handle_unknown="ignore")
4 ohe.fit(Conditions)
```

```
Out[1001]: [array(['Average', 'Fair', 'Good', 'Poor', 'Very Good'], dtype=object)]
```

```
In [1002]: 1 Conditions_encoded = ohe.transform(Conditions)
```

```
Out[1002]: array([[1., 0., 0., 0., 0.],
 [1., 0., 0., 0., 0.],
 [1., 0., 0., 0., 0.],
 ...,
 [1., 0., 0., 0., 0.],
 [1., 0., 0., 0., 0.],
 [1., 0., 0., 0., 0.]])
```

```
In [1003]: 1 Conditions_encoded = pd.DataFrame(Conditions_encoded, columns=ohe.categories_)
```

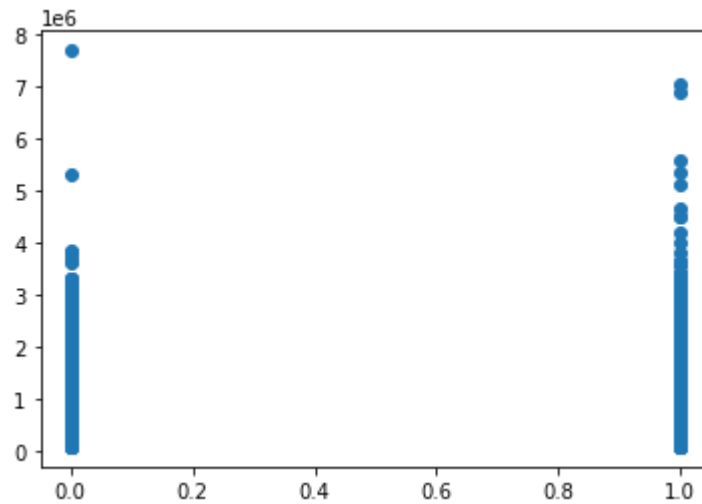
```
Out[1003]:
```

	Average	Fair	Good	Poor	Very Good
0	1.0	0.0	0.0	0.0	0.0
1	1.0	0.0	0.0	0.0	0.0
2	1.0	0.0	0.0	0.0	0.0
3	0.0	0.0	0.0	0.0	1.0
4	1.0	0.0	0.0	0.0	0.0
...	...	...	...	...	...
21592	1.0	0.0	0.0	0.0	0.0
21593	1.0	0.0	0.0	0.0	0.0
21594	1.0	0.0	0.0	0.0	0.0

	Average	Fair	Good	Poor	Very Good
21595	1.0	0.0	0.0	0.0	0.0
21596	1.0	0.0	0.0	0.0	0.0

In [1004]:

Out[1004]: &lt;matplotlib.collections.PathCollection at 0x2873a9c4be0&gt;



In [1005]:

```
1 Month = df[['Month']]
2 ohe = OneHotEncoder(categories="auto", sparse=False, handle_unknown="ignore")
3 ohe.fit(Month)
```

Out[1005]: [array([ 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12], dtype=int64)]

In [1006]:

```
1 Month_encoded = ohe.transform(Month)
```

Out[1006]: array([[0., 0., 0., ..., 1., 0., 0.],
 [0., 0., 0., ..., 0., 0., 1.],
 [0., 1., 0., ..., 0., 0., 0.],
 ...,
 [0., 0., 0., ..., 0., 0., 0.],
 [1., 0., 0., ..., 0., 0., 0.],
 [0., 0., 0., ..., 1., 0., 0.]])

In [1007]:

```
1 Month_encoded = pd.DataFrame(Month_encoded, columns=ohe.categories_[0], ir
```

Out[1007]:

	1	2	3	4	5	6	7	8	9	10	11	12
0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0
1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0
2	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0
4	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
...	...	...	...	...	...	...	...	...	...	...	...	...

	1	2	3	4	5	6	7	8	9	10	11	12
<b>21592</b>	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
<b>21593</b>	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
<b>21594</b>	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0
<b>21595</b>	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
<b>21596</b>	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0

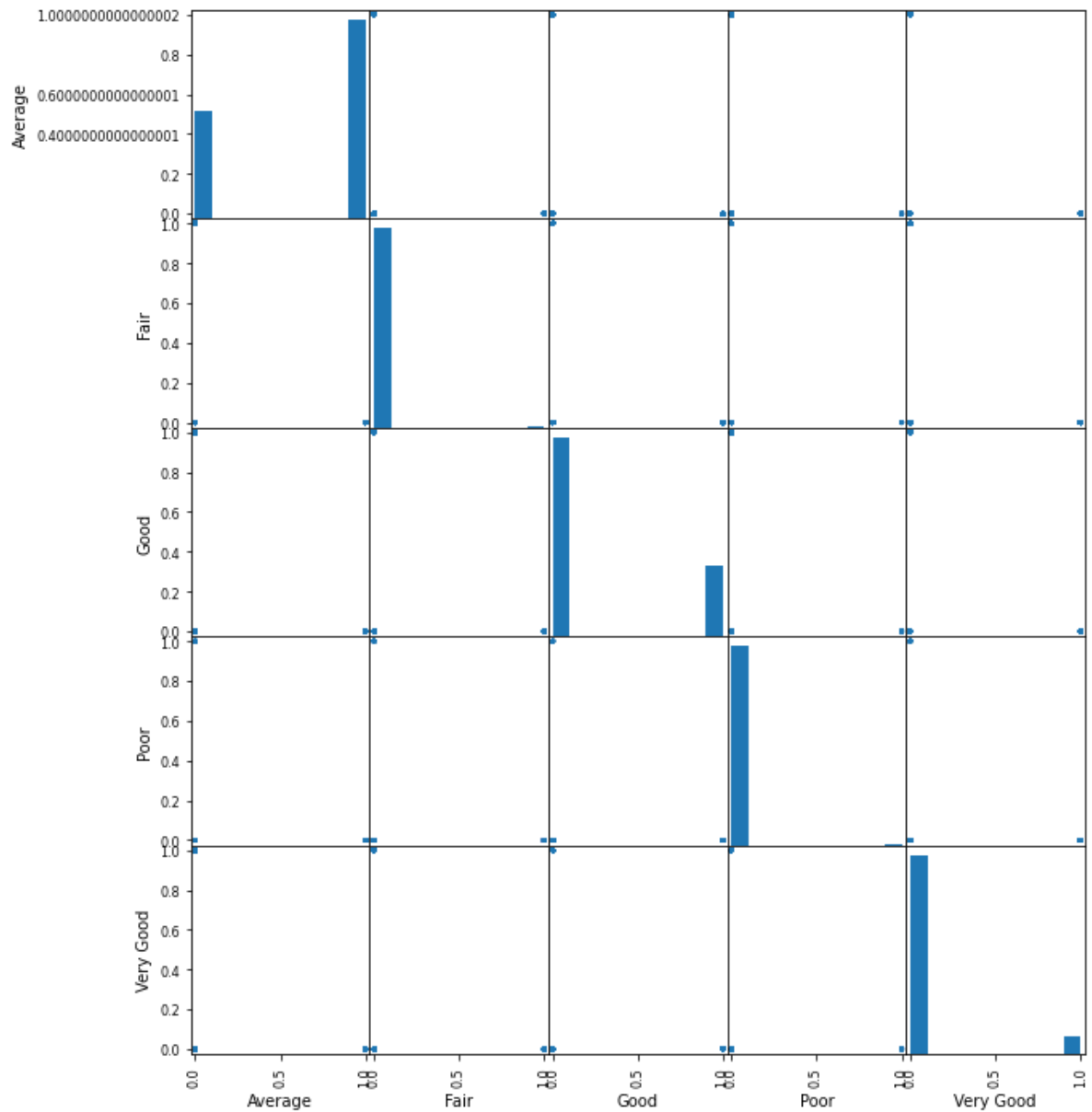
```
In [1008]: 1 Month_encoded.columns = ['January', 'February', 'March', 'April', 'May', '
2           'September', 'October', 'November', 'December']
```

Out[1008]:

	January	February	March	April	May	June	July	August	September	October	November	December
<b>0</b>	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0
<b>1</b>	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
<b>2</b>	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
<b>3</b>	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
<b>4</b>	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
...	...	...	...	...	...	...	...	...	...	...	...	...
<b>21592</b>	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
<b>21593</b>	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
<b>21594</b>	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0
<b>21595</b>	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
<b>21596</b>	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0



In [1009]:



Creating a Completed Renovations Column and a Recent renovations Column:

```
In [1010]: 1 df['Completed_renovation']=df['yr_renovated'] > 1
```

```
In [1011]: 1 df['Recent_renovation']=df['yr_renovated'] > 1
```

```
In [1012]: 1 df
```

Out[1012]:

	price	bedrooms	bathrooms	sqft_living	floors	condition	yr_renovated	Month	grad
0	221900.0	3	1.00	1180	1.0	Average	0.0	10	
1	538000.0	3	2.25	2570	2.0	Average	1991.0	12	
2	180000.0	2	1.00	770	1.0	Average	0.0	2	

	price	bedrooms	bathrooms	sqft_living	floors	condition	yr_renovated	Month	grad
3	604000.0	4	3.00	1960	1.0	Very Good	0.0	12	
4	510000.0	3	2.00	1680	1.0	Average	0.0	2	
...	...	...	...	...	...	...	...	...	...
21592	360000.0	3	2.50	1530	3.0	Average	0.0	5	
21593	400000.0	4	2.50	2310	2.0	Average	0.0	2	
21594	402101.0	2	0.75	1020	2.0	Average	0.0	6	

```
In [1013]: 1 Finished_Reno = df[['Completed_renovation']]
           2 ohe = OneHotEncoder(categories="auto", sparse=False, handle_unknown="ignore")
           3 ohe.fit(Finished_Reno)
```

```
Out[1013]: array([False,  True])
```

```
In [1014]: 1 Finished_Reno_encoded = ohe.transform(Finished_Reno)
```

```
Out[1014]: array([[1., 0.],
                  [0., 1.],
                  [1., 0.],
                  ...,
                  [1., 0.],
                  [1., 0.],
                  [1., 0.]])
```

```
In [1015]: 1 Finished_Reno_encoded = pd.DataFrame(Finished_Reno_encoded, columns=ohe.categories_[0])
           2 Finished_Reno_encoded.columns=['Not_Renovated', 'Renovated']
```

```
Out[1015]:
```

	Not_Renovated	Renovated
0	1.0	0.0
1	0.0	1.0
2	1.0	0.0
3	1.0	0.0
4	1.0	0.0
...	...	...
21592	1.0	0.0
21593	1.0	0.0
21594	1.0	0.0
21595	1.0	0.0
21596	1.0	0.0

```
In [1016]: 1 Recent_Reno = df[['Recent_Renovation']]
           2 ohe = OneHotEncoder(categories="auto", sparse=False, handle_unknown="ignore")
           3 ohe.fit(Recent_Reno)
```

```
Out[1016]: [array([False,  True])]
```

```
In [1017]: 1 Recent_Reno_encoded = ohe.transform(Recent_Reno)
```

```
Out[1017]: array([[1., 0.],
                  [1., 0.],
                  [1., 0.],
                  ...,
                  [1., 0.],
                  [1., 0.],
                  [1., 0.]])
```

```
In [1018]: 1 Recent_Reno_encoded = pd.DataFrame(Recent_Reno_encoded, columns=ohe.categories_[0])
           2 Recent_Reno_encoded.columns=['Not_Renovated_Recently', 'Renovated_Recently']
```

```
Out[1018]:
```

	Not_Renovated_Recently	Renovated_Recently
0	1.0	0.0
1	1.0	0.0
2	1.0	0.0
3	1.0	0.0
4	1.0	0.0
...	...	...
21592	1.0	0.0
21593	1.0	0.0
21594	1.0	0.0
21595	1.0	0.0
21596	1.0	0.0

Combining numerical features:

```
In [1019]: 1 month_and_cond = pd.concat((Month_encoded, Conditions_encoded, Recent_Reno_
```

```
Out[1019]:
```

	January	February	March	April	May	June	July	August	September	October	...	Dec
0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	...	
1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	
2	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	
3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	
4	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	

	January	February	March	April	May	June	July	August	September	October	...	Dec
...	...	...	...	...	...	...	...	...	...	...	...	...
21592	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	...	...
21593	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	...
21594	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	...	...
21595	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	...

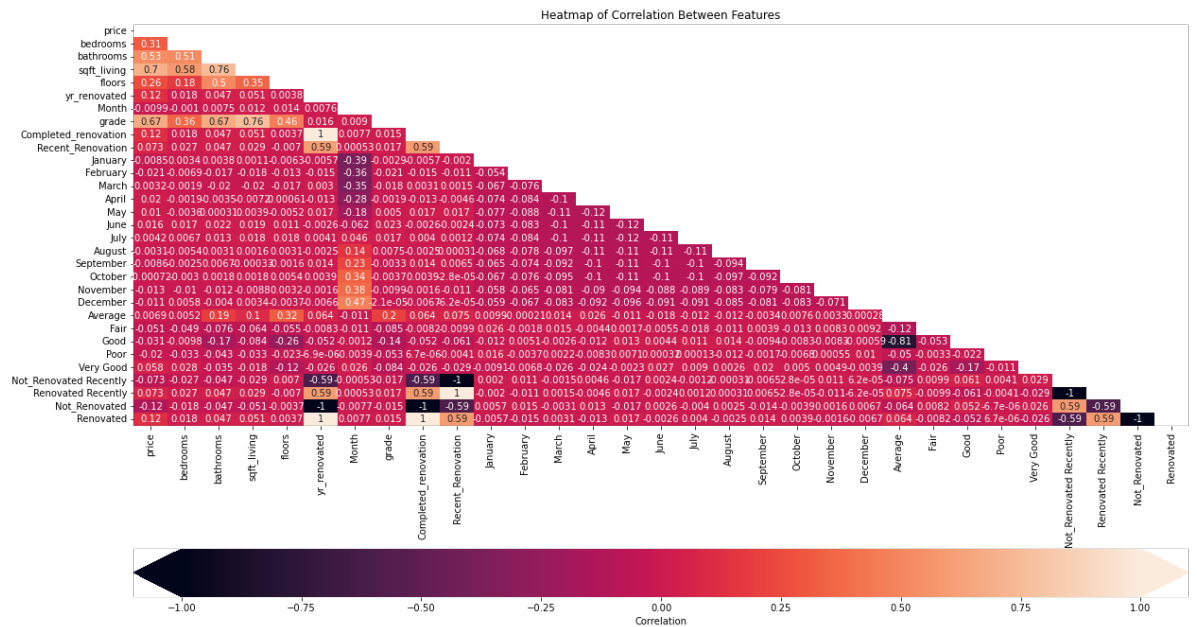
```
In [1020]: 1 df= pd.concat((df,month_and_cond), axis=1)
```

```
Out[1020]:
```

	price	bedrooms	bathrooms	sqft_living	floors	condition	yr_renovated	Month	grad
0	221900.0	3	1.00	1180	1.0	Average	0.0	10	
1	538000.0	3	2.25	2570	2.0	Average	1991.0	12	
2	180000.0	2	1.00	770	1.0	Average	0.0	2	
3	604000.0	4	3.00	1960	1.0	Very Good	0.0	12	
4	510000.0	3	2.00	1680	1.0	Average	0.0	2	
...	...	...	...	...	...	...	...	...	...
21592	360000.0	3	2.50	1530	3.0	Average	0.0	5	
21593	400000.0	4	2.50	2310	2.0	Average	0.0	2	
21594	402101.0	2	0.75	1020	2.0	Average	0.0	6	
21595	400000.0	3	2.50	1600	2.0	Average	0.0	1	
21596	325000.0	2	0.75	1020	2.0	Average	0.0	10	

21597 rows × 32 columns

```
In [1021]: 1 df_corr_2 = df.corr()
2
3 fig, ax = plt.subplots(figsize=(20, 12))
4
5 sns.heatmap(
6     data=df_corr_2,
7     mask=np.triu(np.ones_like(df_corr_2, dtype=bool)),
8     ax=ax,
9     annot=True,
10     cbar_kws={"label": "Correlation", "orientation": "horizontal", "pad":
11 )
12
```



Some of the correlations are on the Boolean columns will probably create noise in any model.

```
In [1022]: Out[1022]:
```

	price	bedrooms	bathrooms	sqft_living	floors	condition	yr_renovated	Month	grad
0	221900.0	3	1.00	1180	1.0	Average	0.0	10	
1	538000.0	3	2.25	2570	2.0	Average	1991.0	12	
2	180000.0	2	1.00	770	1.0	Average	0.0	2	
3	604000.0	4	3.00	1960	1.0	Very Good	0.0	12	
4	510000.0	3	2.00	1680	1.0	Average	0.0	2	
...	...	...	...	...	...	...	...	...	...
21592	360000.0	3	2.50	1530	3.0	Average	0.0	5	
21593	400000.0	4	2.50	2310	2.0	Average	0.0	2	
21594	402101.0	2	0.75	1020	2.0	Average	0.0	6	
21595	400000.0	3	2.50	1600	2.0	Average	0.0	1	

On to the model:

In [1023]:

In [1024]:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21597 entries, 0 to 21596
Data columns (total 29 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   price                                21597 non-null  float64
1   bedrooms                            21597 non-null  int64
2   bathrooms                           21597 non-null  float64
3   sqft_living                          21597 non-null  int64
4   floors                              21597 non-null  float64
5   yr_renovated                        21597 non-null  float64
6   Month                               21597 non-null  int64
7   grade                               21597 non-null  int64
8   January                             21597 non-null  float64
9   February                             21597 non-null  float64
10  March                               21597 non-null  float64
11  April                               21597 non-null  float64
12  May                                 21597 non-null  float64
13  June                               21597 non-null  float64
14  July                               21597 non-null  float64
15  August                             21597 non-null  float64
16  September                           21597 non-null  float64
17  October                             21597 non-null  float64
18  November                             21597 non-null  float64
19  December                             21597 non-null  float64
20  Average                             21597 non-null  float64
21  Fair                                21597 non-null  float64
22  Good                                21597 non-null  float64
23  Poor                                21597 non-null  float64
24  Very Good                           21597 non-null  float64
25  Not_Renovated_Recently              21597 non-null  float64
26  Renovated_Recently                  21597 non-null  float64
27  Not_Renovated                       21597 non-null  float64
28  Renovated                           21597 non-null  float64
dtypes: float64(25), int64(4)
memory usage: 4.8 MB
```

In [1025]:

```
1 X_model_4= df.drop('price', axis=1)
```

In [1026]:

```
1 X_model_4= sm.add_constant(X_model_4)
```

```
C:\Users\isaia\anaconda3\lib\site-packages\statsmodels\tsa\tsatools.py:142: FutureWarning: In a future version of pandas all arguments of concat except for the argument 'objs' will be keyword-only
  x = pd.concat(x[::order], 1)
```

In [1027]:

Out[1027]:

OLS Regression Results

<b>Dep. Variable:</b>	price	<b>R-squared:</b>	0.570
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.569
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	1242.
<b>Date:</b>	Thu, 19 May 2022	<b>Prob (F-statistic):</b>	0.00
<b>Time:</b>	13:02:34	<b>Log-Likelihood:</b>	-2.9828e+05
<b>No. Observations:</b>	21597	<b>AIC:</b>	5.966e+05
<b>Df Residuals:</b>	21573	<b>BIC:</b>	5.968e+05
<b>Df Model:</b>	23		
<b>Covariance Type:</b>	nonrobust		

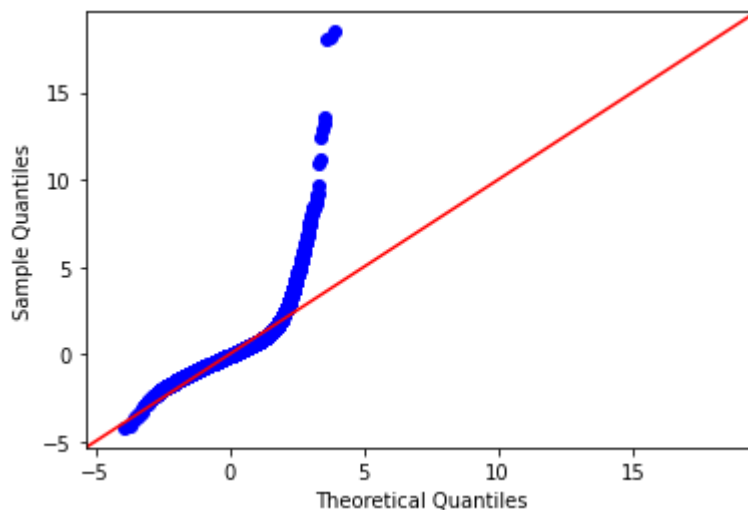
	coef	std err	t	P> t	[0.025	0.975]
const	-1.595e+06	3.07e+05	-5.193	0.000	-2.2e+06	-9.93e+05

In [1028]:

```
1 residuals = Fourth_model.resid
2 fig = sm.graphics.qqplot(residuals, dist=stats.norm, fit=True, line='45',)
```

C:\Users\isaia\anaconda3\lib\site-packages\statsmodels\graphics\gofplots.py:93: UserWarning: marker is redundantly defined by the 'marker' keyword argument and the fmt string "bo" (-> marker='o'). The keyword argument will take precedence.

ax.plot(x, y, fmt, \*\*plot\_style)



Again, this model violates the assumption of normality of errors so we need to do a few things to change the data. First, let's drop some boolean and redundant columns.

In [1029]:

```
1 #drop not recently renovated, not renovated, Month
```

In [1030]:

Out[1030]:

	price	bedrooms	bathrooms	sqft_living	floors	yr_renovated	grade	January	February
0	221900.0	3	1.00	1180	1.0	0.0	7	0.0	
1	538000.0	3	2.25	2570	2.0	1991.0	7	0.0	
2	180000.0	2	1.00	770	1.0	0.0	6	0.0	
3	604000.0	4	3.00	1960	1.0	0.0	7	0.0	
4	510000.0	3	2.00	1680	1.0	0.0	8	0.0	
...	...	...	...	...	...	...	...	...	
21592	360000.0	3	2.50	1530	3.0	0.0	8	0.0	
21593	400000.0	4	2.50	2310	2.0	0.0	8	0.0	
21594	402101.0	2	0.75	1020	2.0	0.0	7	0.0	
21595	400000.0	3	2.50	1600	2.0	0.0	8	1.0	
21596	325000.0	2	0.75	1020	2.0	0.0	7	0.0	

21597 rows × 26 columns

In [1031]: 1 X\_model\_5= df.drop('price', axis=1)

In [1032]: 1 X\_model\_5= sm.add\_constant(X\_model\_5)

C:\Users\isaia\anaconda3\lib\site-packages\statsmodels\tsa\tsatools.py:142: FutureWarning: In a future version of pandas all arguments of concat except for the argument 'objs' will be keyword-only

```
x = pd.concat(x[::order], 1)
```

In [1033]:

Out[1033]:

OLS Regression Results

<b>Dep. Variable:</b>	price	<b>R-squared:</b>	0.570
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.569
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	1242.
<b>Date:</b>	Thu, 19 May 2022	<b>Prob (F-statistic):</b>	0.00
<b>Time:</b>	13:02:35	<b>Log-Likelihood:</b>	-2.9828e+05
<b>No. Observations:</b>	21597	<b>AIC:</b>	5.966e+05
<b>Df Residuals:</b>	21573	<b>BIC:</b>	5.968e+05
<b>Df Model:</b>	23		
<b>Covariance Type:</b>	nonrobust		
	<b>coef</b>	<b>std err</b>	<b>t</b> <b>P&gt; t </b> <b>[0.025</b> <b>0.975]</b>



<b>const</b>	-3.908e+05	1.32e+04	-29.591	0.000	-4.17e+05	-3.65e+05
<b>bedrooms</b>	-4.305e+04	2247.409	-19.158	0.000	-4.75e+04	-3.86e+04
<b>bathrooms</b>	-1.356e+04	3631.989	-3.733	0.000	-2.07e+04	-6440.477
<b>sqft_living</b>	212.7252	3.544	60.016	0.000	205.778	219.673
<b>floors</b>	-1.663e+04	3752.016	-4.433	0.000	-2.4e+04	-9279.393
<b>yr_renovated</b>	3707.6403	787.462	4.708	0.000	2164.156	5251.124
<b>grade</b>	1.1e+05	2344.069	46.938	0.000	1.05e+05	1.15e+05
<b>January</b>	-4.107e+04	7328.167	-5.604	0.000	-5.54e+04	-2.67e+04
<b>February</b>	-3.701e+04	6556.983	-5.644	0.000	-4.99e+04	-2.42e+04
<b>March</b>	-5258.1896	5466.352	-0.962	0.336	-1.6e+04	5456.264
<b>April</b>	892.5514	5084.056	0.176	0.861	-9072.574	1.09e+04
<b>May</b>	-2.797e+04	4905.965	-5.701	0.000	-3.76e+04	-1.84e+04
<b>June</b>	-3.353e+04	5149.683	-6.511	0.000	-4.36e+04	-2.34e+04
<b>July</b>	-4.41e+04	5106.588	-8.636	0.000	-5.41e+04	-3.41e+04
<b>August</b>	-4.043e+04	5405.426	-7.479	0.000	-5.1e+04	-2.98e+04
<b>September</b>	-4.401e+04	5604.592	-7.853	0.000	-5.5e+04	-3.3e+04
<b>October</b>	-3.156e+04	5466.775	-5.772	0.000	-4.23e+04	-2.08e+04
<b>November</b>	-3.93e+04	6201.602	-6.338	0.000	-5.15e+04	-2.71e+04
<b>December</b>	-4.749e+04	6086.226	-7.802	0.000	-5.94e+04	-3.56e+04
<b>Average</b>	-1.607e+05	9076.636	-17.704	0.000	-1.78e+05	-1.43e+05
<b>Fair</b>	-1.126e+05	1.73e+04	-6.521	0.000	-1.46e+05	-7.87e+04
<b>Good</b>	-9.55e+04	9106.462	-10.487	0.000	-1.13e+05	-7.77e+04
<b>Poor</b>	-1.744e+04	3.74e+04	-0.467	0.641	-9.07e+04	5.58e+04
<b>Very Good</b>	-4600.8581	9903.028	-0.465	0.642	-2.4e+04	1.48e+04
<b>Renovated Recently</b>	-4.954e+04	2.56e+04	-1.936	0.053	-9.97e+04	622.763
<b>Renovated</b>	-7.19e+06	1.57e+06	-4.593	0.000	-1.03e+07	-4.12e+06

<b>Omnibus:</b>	16634.219	<b>Durbin-Watson:</b>	1.990
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	964690.215
<b>Skew:</b>	3.209	<b>Prob(JB):</b>	0.00
<b>Kurtosis:</b>	35.107	<b>Cond. No.</b>	3.09e+17

Notes:

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

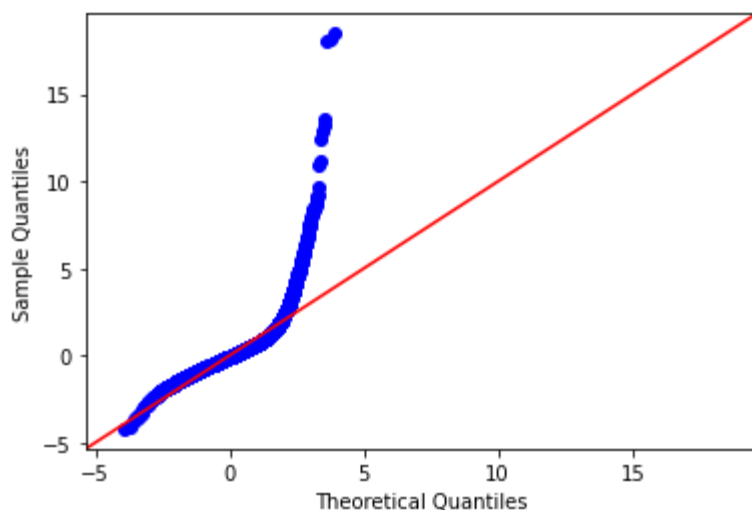
[2] The smallest eigenvalue is 1.17e-24. This might indicate that there are

In [ ]: 1

```
In [1034]: 1 residuals = Fifth_model.resid
2 fig = sm.graphics.qqplot(residuals, dist=stats.norm, fit=True, line='45',)
3 fig.show()
```

C:\Users\isaia\anaconda3\lib\site-packages\statsmodels\graphics\gofplots.py:93: UserWarning: marker is redundantly defined by the 'marker' keyword argument and the fmt string "bo" (-> marker='o'). The keyword argument will take precedence.

```
ax.plot(x, y, fmt, **plot_style)
```



Again, the assumption of normality is violated. The solution may lie in log transforming.

## Log transforming

```
In [1035]: 1 continuous = ['sqft_living', 'price']
2 df_conti= df[continuous]
3 log_names = [f'{column}_log' for column in df_conti.columns]
4 df_log = np.log(df_conti)
5 df_log.columns = log_names
6 def normalize(feature):
7     return (feature - feature.mean()) / feature.std()
8
```

```
In [1036]: 1 df_log_normal = df_log.apply(normalize, axis=1)
```

```
In [1037]: 1 X_log= df_log_normal.drop(['price', 'price_log', 'bathrooms', 'sqft_living',
2                                'Good', 'Very Good'], axis=1)
3 y_log= df_log_normal['price_log']
```

```
In [1038]: 1 X_log = sm.add_constant(X_log)
```

C:\Users\isaia\anaconda3\lib\site-packages\statsmodels\tsa\tsatools.py:142: FutureWarning: In a future version of pandas all arguments of concat except for the argument 'objs' will be keyword-only

```
x = pd.concat(x[::order], 1)
```

In [1039]:

Out[1039]:

OLS Regression Results

<b>Dep. Variable:</b>	price_log	<b>R-squared:</b>	0.558
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.557
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	1601.
<b>Date:</b>	Thu, 19 May 2022	<b>Prob (F-statistic):</b>	0.00
<b>Time:</b>	13:02:36	<b>Log-Likelihood:</b>	-21833.
<b>No. Observations:</b>	21597	<b>AIC:</b>	4.370e+04
<b>Df Residuals:</b>	21579	<b>BIC:</b>	4.385e+04
<b>Df Model:</b>	17		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>const</b>	-2.5342	0.050	-50.841	0.000	-2.632	-2.436
<b>sqft_living_log</b>	0.3725	0.008	45.081	0.000	0.356	0.389
<b>bedrooms</b>	-0.0510	0.006	-7.998	0.000	-0.063	-0.038
<b>floors</b>	-0.0492	0.009	-5.204	0.000	-0.068	-0.031
<b>yr_renovated</b>	0.0002	1.55e-05	13.887	0.000	0.000	0.000
<b>grade</b>	0.3875	0.006	62.933	0.000	0.375	0.400
<b>January</b>	-0.2673	0.020	-13.059	0.000	-0.307	-0.227
<b>February</b>	-0.2380	0.018	-13.004	0.000	-0.274	-0.202
<b>March</b>	-0.1508	0.015	-9.821	0.000	-0.181	-0.121
<b>April</b>	-0.1006	0.014	-7.039	0.000	-0.129	-0.073
<b>May</b>	-0.1822	0.014	-13.157	0.000	-0.209	-0.155
<b>June</b>	-0.1941	0.014	-13.410	0.000	-0.223	-0.166
<b>July</b>	-0.2239	0.014	-15.616	0.000	-0.252	-0.196
<b>August</b>	-0.2270	0.015	-15.002	0.000	-0.257	-0.197
<b>September</b>	-0.2288	0.016	-14.573	0.000	-0.260	-0.198
<b>October</b>	-0.2163	0.015	-14.119	0.000	-0.246	-0.186
<b>November</b>	-0.2406	0.017	-13.863	0.000	-0.275	-0.207
<b>December</b>	-0.2645	0.017	-15.523	0.000	-0.298	-0.231
<b>Renovated Recently</b>	0.0731	0.051	1.428	0.153	-0.027	0.173

<b>Omnibus:</b>	81.303	<b>Durbin-Watson:</b>	1.978
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	81.949
<b>Skew:</b>	0.148	<b>Prob(JB):</b>	1.60e-18
<b>Kurtosis:</b>	2.937	<b>Cond. No.</b>	2.28e+17

Notes:

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

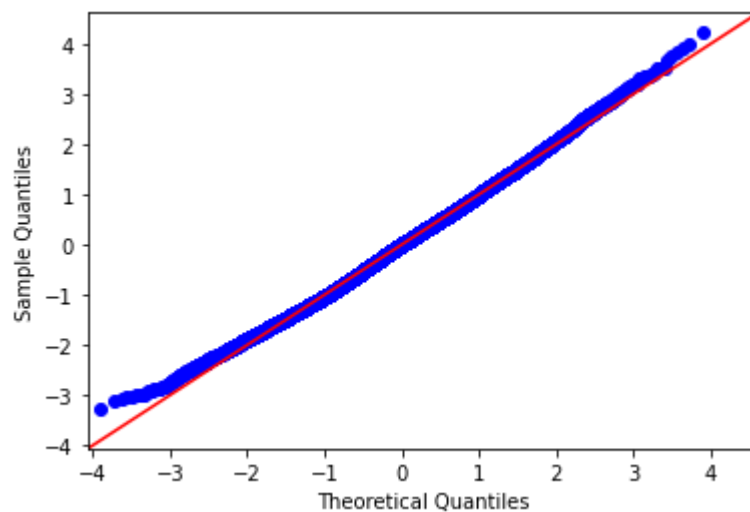
[2] The smallest eigenvalue is 5.72e-26. This might indicate that there are strong multicollinearity problems or that the design matrix is singular.

In [1040]:

```
1 resid = log_model.resid
2
```

C:\Users\isaia\anaconda3\lib\site-packages\statsmodels\graphics\gofplots.py:93: UserWarning: marker is redundantly defined by the 'marker' keyword argument and the fmt string "bo" (-> marker='o'). The keyword argument will take precedence.

```
ax.plot(x, y, fmt, **plot_style)
```



In [1041]:

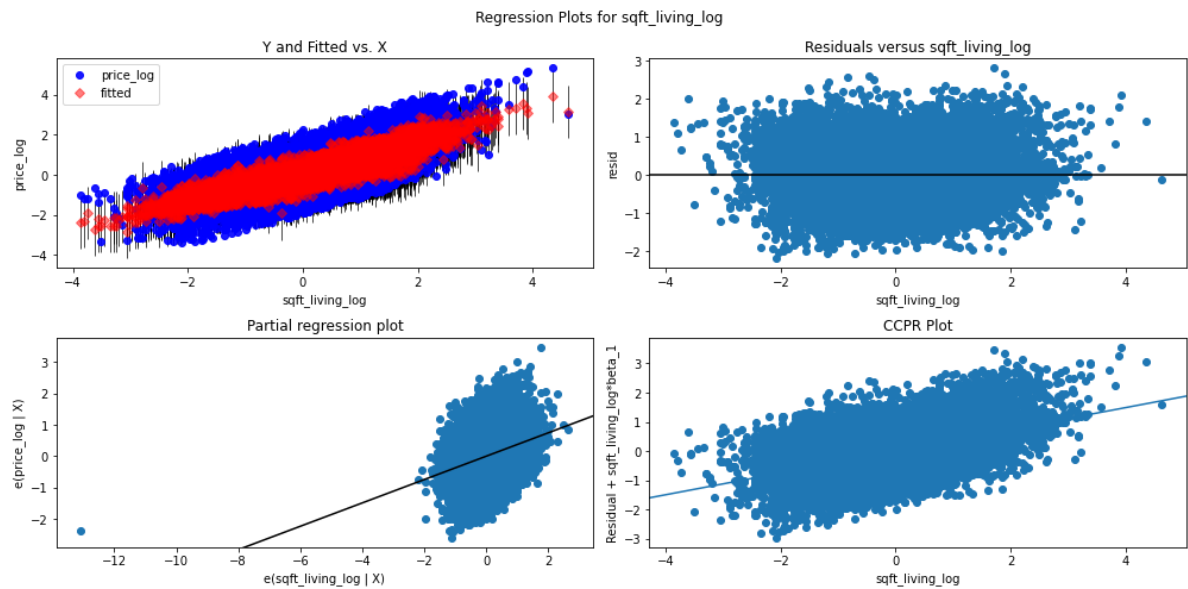
```
1
2 import statsmodels.stats.api as sms
3 name = ['Jarque-Bera', 'Prob', 'Skew', 'Kurtosis']
4 test = sms.jarque_bera(log_model.resid)
```

Out[1041]:

```
[('Jarque-Bera', 81.94898639194479),
 ('Prob', 1.6032590726147647e-18),
 ('Skew', 0.1476112765272201),
 ('Kurtosis', 2.9374642366512265)]
```

So, this looks alot better. Now, we have a much better fit. Also, the skewness and kurtosis have decreased to acceptable numbers. Now let's look at the scedasticity of 'sqft\_living\_log'.

In [1042]:



## Calculations for Square Footage, Grade, and Renovation.

For a 5% increase in square footage:

In [1043]:

Out[1043]: 1.0183901806917963

1 unit change in grade:

In [1044]:

Out[1044]: 38.75

Recently Renovated:

In [1045]:

Out[1045]: 7.31

## Recommendations

1. Improving the square footage by 5% leads to an increase in price of 2%.
2. Improving the grade by 1 unit will increase price by 39%.

### 3. Recently Renovated homes sell for 7% more.

In [1046]:

Out[1046]: 1.0183901806917963

Overall, it looks good but there are certainly some concerns. First, the model may be overfitted as it hews very closely to the line for most of its duration. Another concern is the R-squared value although the R-squared value is a reflection of the fact that the model does not contain many important features like zip code, waterfront, and view. Even so, those factors should be excluded because an existing homeowner cannot change them.

We also should be wary of generalizing too much with this data since it is from a very unique real estate market where prices have continued to rise over the last few decades. Comparing the Seattle area with other areas of comparable size may not work since Seattle is home to so much wealth. However, for homeowners in this area, the results are very valuable and can lead to greater potential profits.