

# Final Project Submission

Please fill out:

- Student name: Shayan Abdul Karim Khan
- Student pace: self paced
- Scheduled project review date/time: Monday 21st November 2022
- Instructor name: Abhineet Kulkarni
- Blog post URL:

## Problem Overview

Paragon Real Estate is a real estate agency with licensed operations in the King County Area. They help homeowners buy/sell homes.

Paragon wants to be able to advise it's clients about how home renovations might increase the estimated value of their homes, and approximately by what amount.

They want to help homeowners make smarter choices about investing in their properties so that homeowners can understand whether a renovation will be helpful for their property valuable.

## Business Questions

The business questions that will be explored in this analysis are as follows:

1. What kind of renovations increases the value?
  - This will give homeowners an understanding of which renovations to focus on to increase the value of their homes and which renovations to avoid.
2. What is the impact to the value of the kind of renovations identified in the first question?
  - This will provide insight to a homeowner to understand what kind of renovations to prioritise and what kind of change can they expect.

## Data Sources

We will be using the official King County House Sales dataset to conduct analysis and answer Paragon's **Business Questions**.

## Data Understanding

## Data Understanding

The dataset listed above has following characteristics discussed in this section:

- Contents and Features of the dataset
- Relevance of the Features to the Business Questions
- Relevant features of the datasets that will be used for analysis
- Limitations of the dataset
- Avenues of analysis that will be pursued

We will start by importing the appropriate python libraries to explore the datasets.

```
In [616]: 1 import pandas as pd #imports the pandas library as pd to work on databa
2 import sqlite3 as sql # imports the sqlite3 library to leverage sql wit
3 from pandasql import sqldf # imports pandas sql library
4 import matplotlib.pyplot as plt # importing matplotlib for visualizatio
5 %matplotlib inline
6 import numpy as np # imports the numpy library
7 import datetime as dt #import datetime module
8 import seaborn as sns #import seaborn
9 from collections import Counter #import Counter
10 import statsmodels.api as sm #import stats models
11 from statsmodels.stats.outliers_influence import variance_inflation_fac
12
13 #import scikit library functions
14 from sklearn.preprocessing import OneHotEncoder, StandardScaler
15 from sklearn.datasets import make_regression
16 from sklearn.linear_model import LinearRegression
17 from sklearn.metrics import mean_squared_error
18 from sklearn.model_selection import train_test_split, cross_validate, S
19 from sklearn.feature_selection import RFECV
20
21
22 #import scipy libraries
23 from scipy import stats as stats
24
25 #import plotly
26 import plotly.express as px
27 import plotly.graph_objects as go
```

Let's import the dataset and take an initial look at it. The King County House Sales dataset is stored as `kc_house_data.csv` in the `data` folder.

```
In [617]: 1 #import the dataset
          2 init_data = pd.read_csv('data/kc_house_data.csv')
          3
          4 #take an initial look at the dataset
          5 init_data.head()
```

Out[617]:

	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront
0	7129300520	10/13/2014	221900.0	3	1.00	1180	5650	1.0	NaN
1	6414100192	12/9/2014	538000.0	3	2.25	2570	7242	2.0	NO
2	5631500400	2/25/2015	180000.0	2	1.00	770	10000	1.0	NO
3	2487200875	12/9/2014	604000.0	4	3.00	1960	5000	1.0	NO
4	1954400510	2/18/2015	510000.0	3	2.00	1680	8080	1.0	NO

5 rows × 21 columns

Lets look at the overview of the data frame using the `.info()` function

```
In [618]: 1 init_data.info() # getting the overview info of the dataframe records
```

```
<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
```

There are **twenty-one** columns in this dataset.

We will use the `columns_names.md` file in the `data` folder to get insights into what information these columns provide us.

We will go through the columns and try to deduce whether the information is useful enough to include as a feature in our analysis.

- `id` : These are the unique identifiers of the houses. We can use this as the index to our dataframe.
- `date` : This column contains the date that the specific house was sold on. This column might be important depending on how much variation it has. Lets take a look at the range of years that we have available.

```
In [619]: 1 #convert string object to datetime
2 dt = pd.to_datetime(init_data['date'])
3
4 #convert to an array of only years
5 dt_yr = pd.DatetimeIndex(dt).year
6
7 #extract min and max year values
8 print('The oldest data we have available is for', dt_yr.min())
9 print(' ')
10 print('The newest data we have available is for', dt_yr.max())
```

The oldest data we have available is for 2014

The newest data we have available is for 2015

Looks like we have data only for **2014** and **2015**.

This means that our predictions will be lacking data from recent years and predictions will not be taking into account inflation and recent real estate market impacts.

We do not have ample data from varying years to give us a confidence on how house prices vary by months or years. Therefore, we will drop this column also. Also, This doesn't tell us anything about what kind of renovations to pursue therefore we will **ignore** this feature.

- **price** : This is the sale price of the house. This is the most important data for analysis. This data is **float** type.
- **bedrooms** : This is the number of bedrooms in a house. Intuitively, we know that bedrooms are an important factor in the value of the house. We will use this data to understand how changing the number of bedrooms impacts the house value.
- **bathrooms** : Similar to bedrooms, bathrooms are also an integral part of a house and intuitively, we know that they can play a part in the buying of a house therefore we will use this column for analysis. Interestingly, this dataset is a float. We know that there are 4 kinds of bathrooms: quarter bath, half bath, three-quarter bath, and a full bath. Let's make sure that the float data we have here is only in increments of 0.25.

```
In [620]: 1 #get the unique values and value counts in the bathroom columns
          2 init_data['bathrooms'].value_counts()
```

```
Out[620]: 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     155
          4.00     136
          4.50     100
          4.25      79
          0.75      71
          4.75      23
          5.00      21
          5.25      13
          5.50      10
          1.25       9
          6.00       6
          0.50       4
          5.75       4
          6.75       2
          8.00       2
          6.25       2
          6.50       2
          7.50       1
          7.75       1
          Name: bathrooms, dtype: int64
```

Through visual inspection, it is clear that the float data meets our knowledge of the 4 types of bathrooms. Let's move onto the other columns.

- `sqft_living` : This is the square footage of living space in the house. The size of the habitable area can play a big part in the price of a house. It will be interesting to develop insights into how changing the living space square footage can impact the value of a house.
- `sqft_lot` : This column contains the total size of the property in square footage. Similar to the habitable area, the total size of the property is also an important factor to consider when deducing house prices. Nonetheless, it is not something that can be changed through renovations. Therefore, we will **ignore** this column.
- `floors` : This is data on how many levels a house has. Generally, floors are also an important factor in a house price. Intuitively, this feature is tied together with `sqft_living` but it will be important to investigate how changes in the number of floors impacts the value of the house. Interestingly, this data is **float** type. We know from domain knowledge that there are full floors and half floors. Let's investigate to make sure that the data meets our domain knowledge assumptions.

```
In [621]: 1 #extract unique value counts
          2 init_data['floors'].value_counts()
```

```
Out[621]: 1.0    10673
          2.0     8235
          1.5     1910
          3.0      611
          2.5      161
          3.5         7
          Name: floors, dtype: int64
```

Looks like we either have full floors and half floors therefore this data should be good. We can move onto investigating the remaining columns.

- **waterfront** : This column contains information on whether a house is on a waterfront. Watrefront properties are generally known to be more valuable but this isn't something that can be changed through renovations therefore it won't help us answer the business questions. We will **ignore** this column in our analysis.
- **view** : This is the quality of view from the house. The quality of view is highly subjective to individual preferences and we don't have information on the basis for the quality. Also, the quality of view isn't something that can be renovated therefore we will **ignore** this column.
- **condition** : This data shows how good the overall condition of the house is. It is related to maintenance of the house. Intuitively, we know that the condition of the house plays a big part in the house value. This is also something that can be improved through renovations. Lets investigate what kind of data this columns contains.

```
In [622]: 1 #extract unique value counts
          2 init_data['condition'].value_counts()
```

```
Out[622]: Average    14020
          Good       5677
          Very Good   1701
          Fair        170
          Poor        29
          Name: condition, dtype: int64
```

Looks like these are categories stored in string format. We can process these later and convert them to integers to make them easire to work with. The codes respective to these descriptions listed on the King Coounty Assessor website is as follows:

1 = Poor- Worn out. Repair and overhaul needed on painted surfaces, roofing, plumbing, heating and numerous functional inadequacies. Excessive deferred maintenance and abuse, limited value-in-use, approaching abandonment or major reconstruction; reuse or change in occupancy is imminent. Effective age is near the end of the scale regardless of the actual chronological age.

2 = Fair- Badly worn. Much repair needed. Many items need refinishing or overhauling, deferred maintenance obvious, inadequate building utility and systems all shortening the life expectancy and increasing the effective age.

3 = Average- Some evidence of deferred maintenance and normal obsolescence with age in that a few minor repairs are needed, along with some refinishing. All major components still functional and contributing toward an extended life expectancy. Effective age and utility is standard for like properties of its class and usage.

4 = Good- No obvious maintenance required but neither is everything new. Appearance and utility are above the standard and the overall effective age will be lower than the typical property.

5= Very Good- All items well maintained, many having been overhauled and repaired as they have shown signs of wear, increasing the life expectancy and lowering the effective age with little deterioration or obsolescence evident with a high degree of utility.

- `grade` : This is the overall grade of the house which is related to the construction and design of the house. The Kings County Assessor provides official information on these grades which is listed below:
  - Grades run from grade 1 to 13. Generally defined as:
    - 1-3 Falls short of minimum building standards. Normally cabin or inferior structure.
    - 4 Generally older, low quality construction. Does not meet code.
    - 5 Low construction costs and workmanship. Small, simple design.
    - 6 Lowest grade currently meeting building code. Low quality materials and simple designs.
    - 7 Average grade of construction and design. Commonly seen in plats and older subdivisions.
    - 8 Just above average in construction and design. Usually better materials in both the exterior and interior finish work.
    - 9 Better architectural design with extra interior and exterior design and quality.
    - 10 Homes of this quality generally have high quality features. Finish work is better and more design quality is seen in the floor plans. Generally have a larger square footage.
    - 11 Custom design and higher quality finish work with added amenities of solid woods, bathroom fixtures and more luxurious options.
    - 12 Custom design and excellent builders. All materials are of the highest quality and all conveniences are present.
    - 13 Generally custom designed and built. Mansion level. Large amount of highest quality cabinet work, wood trim, marble, entry ways etc.

These grade descriptions tell us that it is possible to jump into higher grades through renovations and house improvements. This would be an interesting feature to investigate to develop insights for the user.

Lets investigate whether the values in this column match up with the official information.



```
In [623]: 1 #extract unique value counts
          2 init_data['grade'].value_counts()
```

```
Out[623]: 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
```

The data matches up with the official information. We can do some pre-processing to make it easier to work with this data which we'll explore later.

One thing to note is that from the King County assessor website, we know that `condition` is relative to `grade`. Lets use the `groupby` function to inspect the descriptive values of the two features to understand how they are related. Since the `condition` column is more generalized, we'll use that as the primary grouper and see what range of `grade` values does it cover.

```
In [624]: 1 #setting up the pivot table
          2 init_data.groupby(['condition', 'grade'])['grade'].count()
          3
```

```
Out[624]: condition grade
Average    10 Very Good    921
           11 Excellent    332
           12 Luxury       73
           13 Mansion      11
           4 Low          12
           5 Fair         100
           6 Low Average  1035
           7 Average     5229
           8 Good        4266
           9 Better     2041
Fair       10 Very Good     2
           4 Low           4
           5 Fair         15
           6 Low Average   59
           7 Average      75
           8 Good         13
           9 Better        2
Good       10 Very Good   156
           11 Excellent   56
           12 Luxury     13
           13 Mansion      2
           4 Low         10
           5 Fair        84
           6 Low Average  685
           7 Average    2831
           8 Good     1394
           9 Better    446
Poor       4 Low          1
           5 Fair         9
           6 Low Average  11
           7 Average      6
           8 Good         2
Very Good  10 Very Good   55
           11 Excellent   11
           12 Luxury      3
           3 Poor         1
           5 Fair        34
           6 Low Average  248
           7 Average    833
           8 Good       390
           9 Better    126
Name: grade, dtype: int64
```

We can see that every level of `condition` has multiple `grade` types. This tells us that a higher grade doesn't necessarily mean a better condition. For instance, we see that there are **mansions** that are in **Average** conditions while there are also buildings like cabins that are categorized as **poor** grade but are in **Very Good** condition. Therefore, we have to look at both of these variables during analysis.

As a general understanding, `condition` is related to the maintenance of the house while `grade` is more concerned with the architectural aspects of a house. Renovations to improve on the lack of maintenance is a lot easier financially as compared to architectural changes.

- `sqft_above` : This is square footage of house apart from basement. With the total square feet already available, it would not be extra insightful to look at the separate square footage. Therefore, we will **ignore** this column.
- `sqft_basement` : This is square footage of the basement. Similiar to the previous column, this won't add major insights for us with the total square footage already available. Therefore, we will **ignore** this column.
- `yr_built` : This column contains the year when the house was built. This is important information for estimating the house value because the age of a house can both be a pro or a con depending on the type of buyer and type of property. Nonetheless, this is not something that we can change or renovate therefore it will not providee valuable insights to answer the business questions. Therefore, we will **ignore** this column.
- `yr_renovated` : This is the year that the house was renovated in. Without more information on what kind of renovation took place, this data can be highly misleading because reenovations can be major overhaul or smaller improvements. For this reason, we will **ignore** this data.
- `zipcode` : This is the Zip Code that the property is located in which is used by the United States Postal Service. Zipcodes play a big part in establishing the value of a house because of different factors like rate of crime, school quality, etc. Nonetheless, this isn't something that can be changed. This can be something that is looked at separately in conjunction with the other factors at play in a neighbourhood to identify areas of most return on investment. But for the purpose of answering our business questions, wee will **ignore** this feature.
- `lat` : This is the latitude coordinates of the house. This is the same case as zipcode therefore we will **ignore** this feature also.
- `long` : This is the longitude coordinates of the house. Same casee as latitude, we will **ignore** this feature.
- `sqft_living15` : This is the square footage of interior housing living space for the nearest 15 neighbors. This isn't something that a homeowner looking to renovate has control over therefore we will **ignore** this data.
- `sqft_lot15` : This is the square footage of the land lots of the nearest 15 neighbors. Similar to the previous column, we will **ignore** this data because a homeowner doesn't have control over this feature.

## Summary

Out of the **21 features** that this dataset has, we were able to identify **7 features** which can be helpful in answering the business questions of our client. Before we review the features that we will be using, let's look at the features we will be ignoring.

### Columns/Features to be ignored

Although these features can be important in determining the value of the house, they are being ignored because of one or both of the reasons listed below:

1. They can't be changed through renovation
2. They are linked to other more insightful features which we have chosen instead.

The features to be ignored are:

- date
- sqft\_lot
- waterfront
- view
- sqft\_above
- sqft\_basement
- yr\_built
- yr\_renovated
- zipcode
- lat
- long
- sqft\_living15
- sqft\_lot15

### Columns/Features carried forward

The features we are carrying forward and how they will help us answer the business questions are the following:

- **price** : This is the main target that we will be investigating and the only source to find out the house price
- **bedrooms** : This will help us identify whether increasing the number of bedrooms increases the house value. This will also inform a homeowner of how important the number of bedrooms are for their house price.
- **bathrooms** : Similar to the number of bedrooms, this will help homeowners identify whether adding bathrooms increases their house value.
- **sqft\_living** : This will provide insights to the homeowner if increasing the living space by adding another floor or building an extension is a valuable renovation or not.

- `conditions` : This provides a qualitative and quantitative assessment of categorically grouping the condition of the house. This can provide valuable insights into how changing the condition can affect the value of the house. This can help homeowners identify exactly what kind of requirements they need to meet to jump to a better condition and help plan for more impactful renovations.
- `grade` : This is the more thorough categories of the quality of construction and the house. The condition data is related to it and as shown above, it will be important to investigate the relation of this in conjunction to other variables in determining types and impacts of renovations.
- `floor` : This is the number of levels in a house. This is important to take forward since large architectural changes such as adding more floors can improve the grade of the house and make way for more bedrooms, bathrooms and living square footage.

## Data Limitations

no way of identifying the cost effectiveness, detailed kind of renovation, difficulty of renovation, equity of improvement, location specific improvements. inflation impact, housing market trends, lack of more recent data

Our dataset has many limitations that are not being solved as a part of this analysis. Some of them are listed below:

- There is no data on how much it cost to make a new floor, bedroom, bathroom, etc. We also don't have information on how lavish the interior of the house is. Therefore there is no way to do a cost-benefit analysis to understand the return on investment and how long it would take for the increase in value to actualize.
- We also don't have definite information on how difficult the different types of renovations would be and how long they would take to accomplish. A quicker renovation can have a quick turn-around for customers who are looking to flip or sell houses quick.
- We also don't have detailed information on the zipcodes. Zipcodes and neighbourhoods play a big role during appraisals, and even minor renovations can have a major impact because of location. Also, bad neighbourhoods with high crime can have a ceiling to how much a house's value can go up.
- We also don't have housing market and inflation trends to gauge how renovations are impacted by features that are outside of the customer's control.
- There is also a lack of data on the type of houses that we are looking at. Renovations like more bedrooms on a rental multi-family property in a university area would have a different impact on the house price as compared to other types of properties. The lack of this kind of data also severely hampers our ability to confidently explain the trend and correlation.

## Avenues of Analysis

We are going to use multiple linear regression to understand how the different features impact house prices. The reason we want to use linear regression is because it allows us to reliably predict values and understand the impact of multiple variables on the target variable.

Linear Regression also gives us the power of inferring the relationships of multiple variables to the target variable. We will use that insight to infer which kind of renovations would improve the house prices. Since we don't have a plethora of relevant features, we will avoid building a predictive model.

## Data Preparation

The following steps will be followed in preparing the data:

- Data Cleaning
- Data Processing for Regression

### Data Cleaning

Lets start with selecting the columns we will be using.

```
In [625]: 1 #copy dataframe to avoid changing the original data
2 data_cln = init_data.copy()
3
4 #note the columns to keep
5 col_kp = ['price', 'bedrooms', 'bathrooms', 'floors', 'sqft_living', 'condition']
6
7 #select the column to use
8 data_cln = data_cln[col_kp]
9
10 #preview the info
11 print(data_cln.info())
12
13 #preview the data
14 data_cln.head()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21597 entries, 0 to 21596
Data columns (total 7 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   floors           21597 non-null  float64
4   sqft_living      21597 non-null  int64
5   condition        21597 non-null  object
6   grade            21597 non-null  object
dtypes: float64(3), int64(2), object(2)
memory usage: 1.2+ MB
None
```

Out[625]:

	price	bedrooms	bathrooms	floors	sqft_living	condition	grade
0	221900.0	3	1.00	1.0	1180	Average	7 Average
1	538000.0	3	2.25	2.0	2570	Average	7 Average
2	180000.0	2	1.00	1.0	770	Average	6 Low Average
3	604000.0	4	3.00	1.0	1960	Very Good	7 Average
4	510000.0	3	2.00	1.0	1680	Average	8 Good

Lets seperate grade descriptions and numbers so that we can process them seperately if we need to.

```
In [626]: 1 #split the grade column into 2 new columns
2 data_cln[['grade#', 'grade_desc']] = data_cln['grade'].str.split(' ', 1)
3
4 #drop the grade column
5 data_cln.drop(['grade'],axis=1, inplace=True)
6
7 #convert grade# column to int64
8 data_cln['grade#'] = data_cln['grade#'].astype('int64')
9 #preview the data
10 print(data_cln.info())
11 data_cln.head()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21597 entries, 0 to 21596
Data columns (total 8 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   floors          21597 non-null  float64
 4   sqft_living     21597 non-null  int64
 5   condition       21597 non-null  object
 6   grade#          21597 non-null  int64
 7   grade_desc      21597 non-null  object
dtypes: float64(3), int64(3), object(2)
memory usage: 1.3+ MB
None
```

Out[626]:

price	bedrooms	bathrooms	floors	sqft_living	condition	grade#	grade_desc
-------	----------	-----------	--------	-------------	-----------	--------	------------

Using the `column_names.md` file and the King County Assessor website, let's add a column with the respective condition numbers for the condition descriptions in the `condition` column.

The building conditions are coded from 1-5. They are coded as follows:

1 = Poor

2 = Fair

3 = Average

4 = Good

5 = Very Good



```
In [627]: 1 #check the values in the condition column
          2 data_cln['condition'].value_counts()
```

```
Out[627]: Average      14020
          Good        5677
          Very Good   1701
          Fair         170
          Poor         29
          Name: condition, dtype: int64
```

```
In [628]: 1 #create a dictionary of the description and codes
          2 cond_cd = {'Poor':1, 'Fair':2, 'Average':3, 'Good':4, 'Very Good':5}
          3
          4 #map every condition record and compare with the the dictionary to creat
          5 data_cln['cond_code'] = data_cln['condition'].map(lambda x: cond_cd[x])
          6
          7 #preview the data
          8 print(data_cln.info())
          9 data_cln.head()
```

```
<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   floors           21597 non-null  float64
4   sqft_living      21597 non-null  int64
5   condition        21597 non-null  object
6   grade#           21597 non-null  int64
7   grade_desc       21597 non-null  object
8   cond_code        21597 non-null  int64
dtypes: float64(3), int64(4), object(2)
memory usage: 1.5+ MB
None
```

```
Out[628]:
```

	price	bedrooms	bathrooms	floors	sqft_living	condition	grade#	grade_desc	cond_code
0	221900.0	3	1.00	1.0	1180	Average	7	Average	3
1	538000.0	3	2.25	2.0	2570	Average	7	Average	3
2	180000.0	2	1.00	1.0	770	Average	6	Low Average	3
3	604000.0	4	3.00	1.0	1960	Very Good	7	Average	5
4	510000.0	3	2.00	1.0	1680	Average	8	Good	3

We will make the `price` column easier to work with. Lets change the scale to **\$100,000**.

```
In [629]: 1 #divide the price column by 100000
          2 data_cln['price'] = data_cln['price']/100000
          3
```

```
In [630]: 1 #rename the price column
          2 data_cln.rename(columns={'price':'price($100,000)'},inplace=True)
          3
          4 #preview the dataframe
          5 data_cln.head()
```

Out[630]:

	price(\$100,000)	bedrooms	bathrooms	floors	sqft_living	condition	grade#	grade_desc	cond_c
0	2.219	3	1.00	1.0	1180	Average	7	Average	
1	5.380	3	2.25	2.0	2570	Average	7	Average	
2	1.800	2	1.00	1.0	770	Average	6	Low Average	
3	6.040	4	3.00	1.0	1960	Very Good	7	Average	
4	5.100	3	2.00	1.0	1680	Average	8	Good	

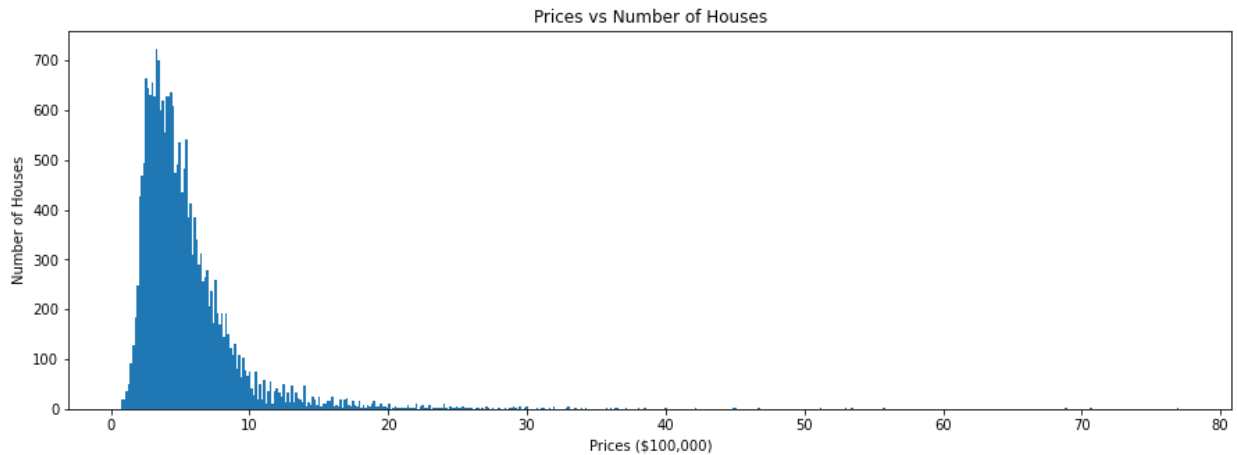
That should end the data cleaning process. We will now process this data to prepare it for regression analysis.

## Data Processing for Regression Analysis

Lets start with separating the target variable, i.e price .

Lets take a look at the distribution of the target variable to get an understanding of how the data is distributed.

```
In [631]: 1 #defining the skeleton of the graph
2 fig, ax = plt.subplots(figsize=(15,5))
3
4 #plotting the histogram
5 ax.hist(data_cln['price($100,000)'],bins=500);
6
7 #setting titles and axis labels
8 ax.set_title('Prices vs Number of Houses');
9 ax.set_xlabel('Prices ($100,000)');
10 ax.set_ylabel('Number of Houses');
```



We can see a right-skewed normal distribution with very few houses above **\$1 million**.

Lets take a quick look at how many these are.

```
In [632]: 1 #filter for >1 million and measure the length
2 print('Number of houses sold for more than $1 million:', len(data_cln[d
3
4 print('The percentage of houses sold for than $1 million:',
5       len(data_cln[data_cln['price($100,000)'] > 10])*100/len(data_cln)
```

Number of houses sold for more than \$1 million: 1458

The percentage of houses sold for than \$1 million: 6.75093763022642

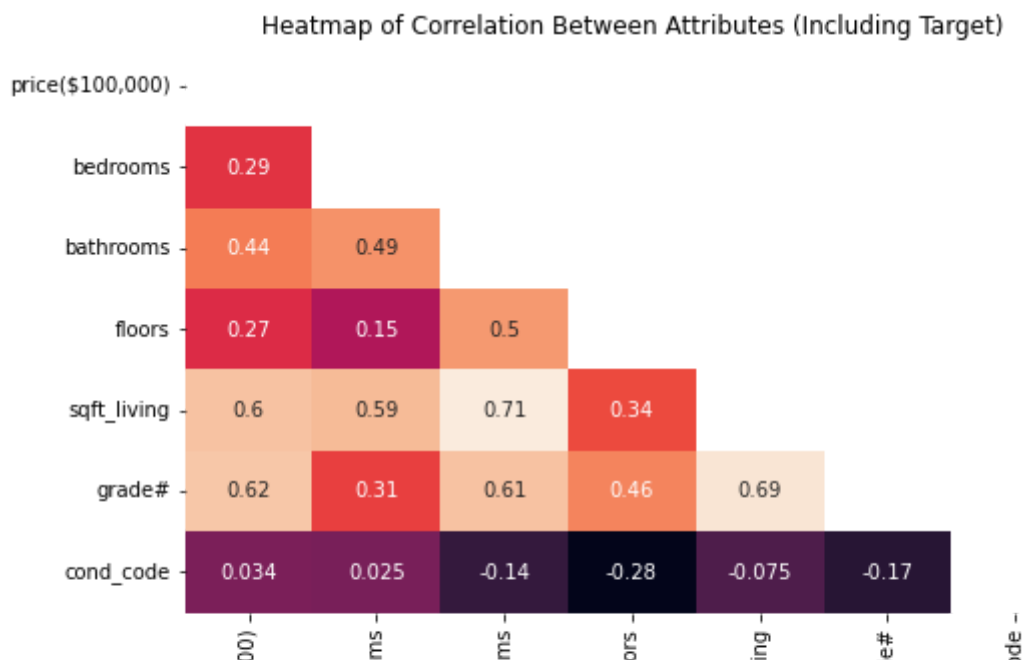
Approximately 6.75% is a considerably small number which we will drop so as not to skew our model.

Lets take a quick look at how the correlation between the different variables and the target variable of price changes as we go over **\$1 million** as compared to if we stay under that.

```

In [633]: 1
2 #compute the correlation matrix
3 corr = data_cln[data_cln['price($100,000)']<=10].corr()
4
5 # Set up figure and axes
6 fig, ax = plt.subplots(figsize=(8, 8))
7
8 # Plot a heatmap of the correlation matrix, with both
9 # numbers and colors indicating the correlations
10 sns.heatmap(
11     # Specifies the data to be plotted
12     data=corr,
13     # The mask means we only show half the values,
14     # instead of showing duplicates. It's optional.
15     mask=np.triu(np.ones_like(corr, dtype=bool)),
16     # Specifies that we should use the existing axes
17     ax=ax,
18     # Specifies that we want labels, not just colors
19     annot=True,
20     # Customizes colorbar appearance
21     cbar_kws={"label": "Correlation", "orientation": "horizontal", "pad
22 )
23
24 # Customize the plot appearance
25 ax.set_title("Heatmap of Correlation Between Attributes (Including Targ

```



Interestingly, there is not a lot of difference except for correlatin of price and grade# .

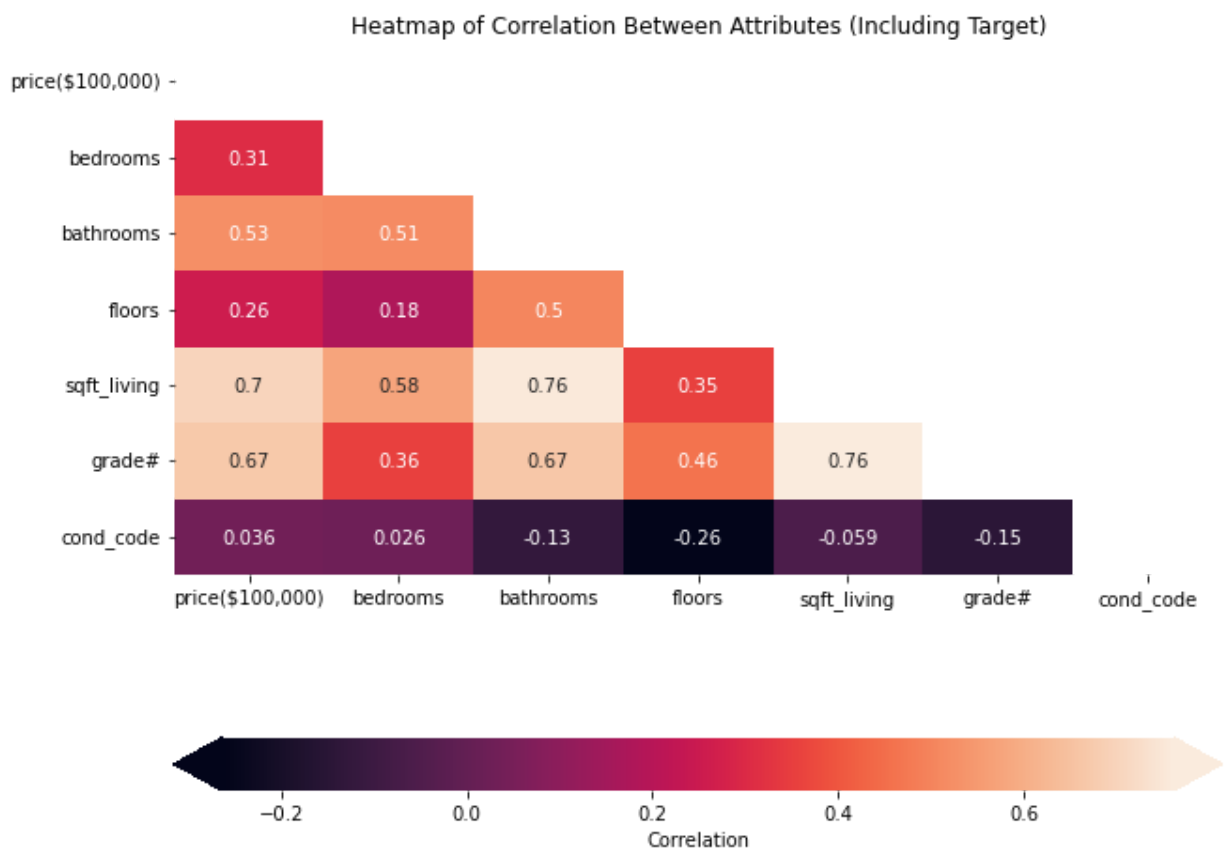
The other trends remain the same. We can't accept to create a prefect model for all scenarios and we will try to fit our regression model so that it has the minimum error when evaluating prices for all price scenarios.

Now lets take a look at the correlation of our complete dataset.

```

In [634]: 1 # heatmap_data = pd.concat([y_train, X_train], axis=1)
2 corr = data_cln.corr()
3
4 # Set up figure and axes
5 fig, ax = plt.subplots(figsize=(10, 8))
6
7 # Plot a heatmap of the correlation matrix, with both
8 # numbers and colors indicating the correlations
9 sns.heatmap(
10     # Specifies the data to be plotted
11     data=corr,
12     # The mask means we only show half the values,
13     # instead of showing duplicates. It's optional.
14     mask=np.triu(np.ones_like(corr, dtype=bool)),
15     # Specifies that we should use the existing axes
16     ax=ax,
17     # Specifies that we want labels, not just colors
18     annot=True,
19     # Customizes colorbar appearance
20     cbar_kws={"label": "Correlation", "orientation": "horizontal", "pad
21 )
22
23 # Customize the plot appearance
24 ax.set_title("Heatmap of Correlation Between Attributes (Including Target)

```



Based on the plot above, `sqft_living` is the most correlated feature with `price`. We will use this feature to create the baseline model before iterating on it.

`grade#` is the second highest correlated feature.

`cond_code` has very little correlation with the `price` of a house which shows that this feature should be ignored.

Another interesting observation is the high correlation between `grade#`, `sqft_living`, and `bathrooms`. To avoid having a highly collinear model, we will drop `bathrooms`, `bedrooms`.

We will investigate models with combinations of all these variables to generate the most appropriate regression model.

Before we move onto creating the model, we still need to investigate these variables for other properties and transform them if needed.

Let's start with checking which variables are **categorical** variables after dropping `cond_code` and `grade#`. These were two columns we added earlier but even if these two features are categorical, we can work with them using the descriptive analyses.

```
In [635]: 1 #drop cond_code and condition
          2 data_cln = data_cln.drop(['cond_code', 'grade#', 'bathrooms', 'bedrooms'],
```

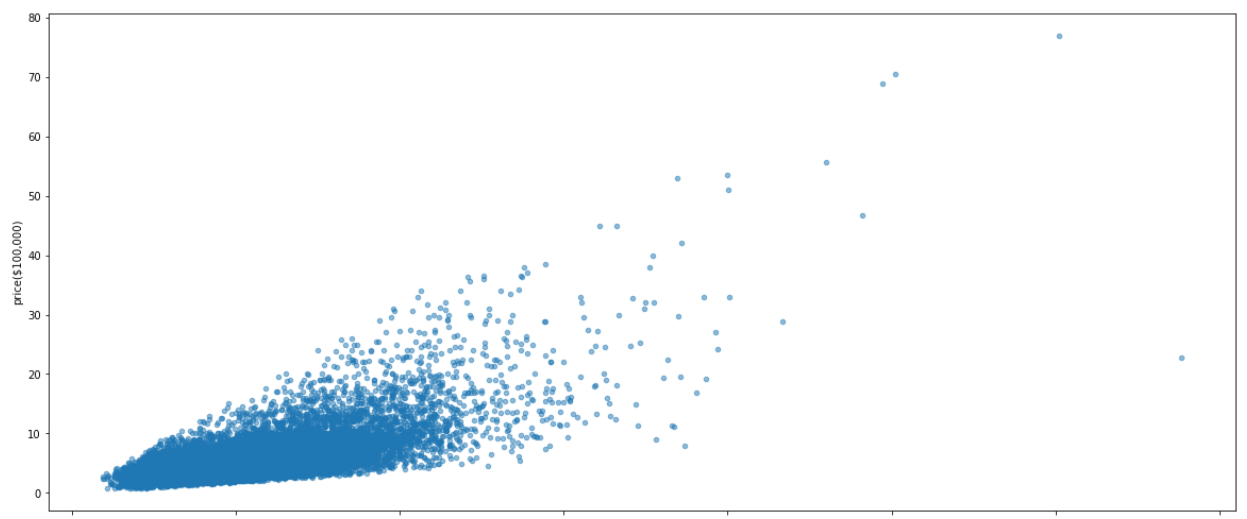
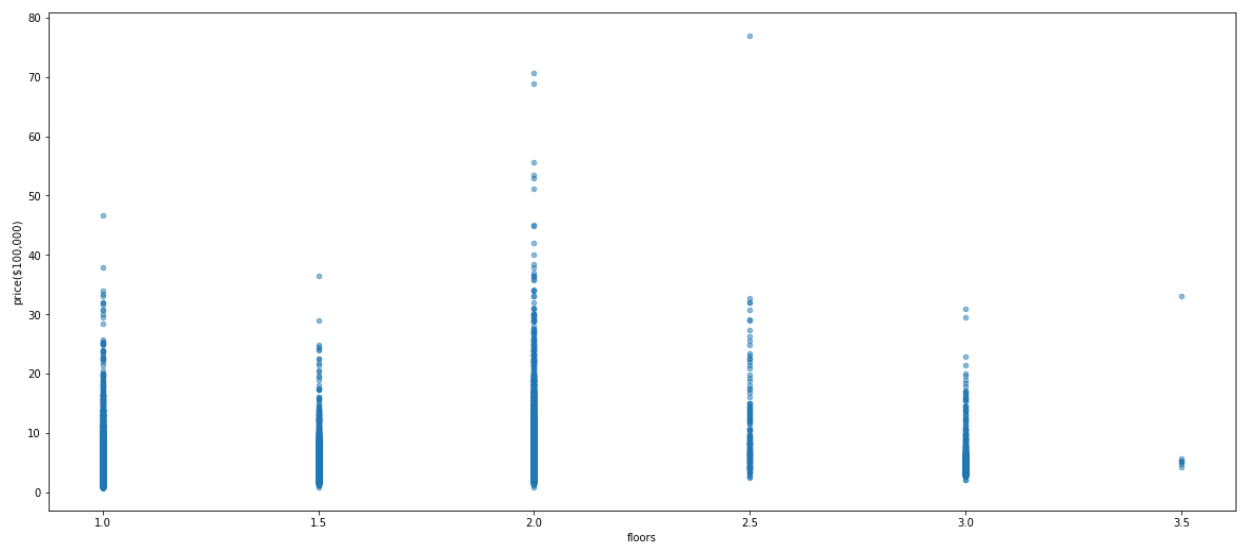
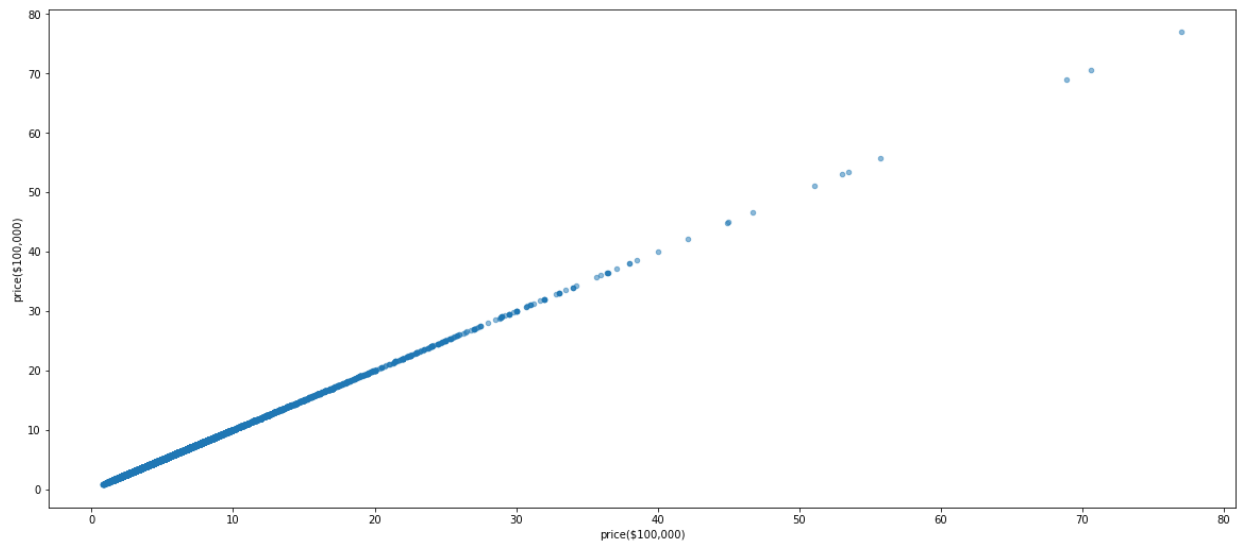
```
In [636]: 1 #preview data
          2 data_cln
```

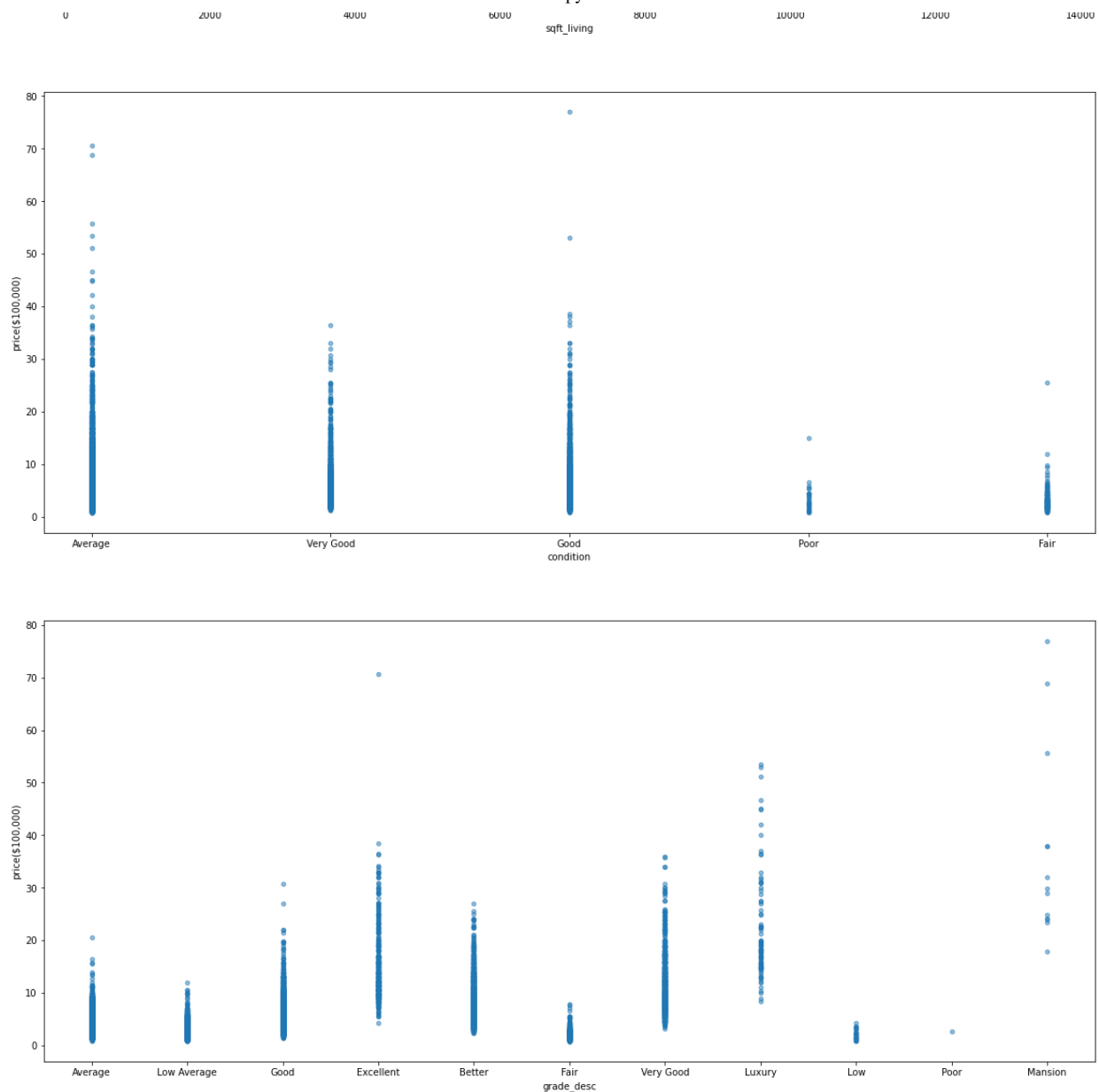
Out[636]:

	price(\$100,000)	floors	sqft_living	condition	grade_desc
0	2.21900	1.0	1180	Average	Average
1	5.38000	2.0	2570	Average	Average
2	1.80000	1.0	770	Average	Low Average
3	6.04000	1.0	1960	Very Good	Average
4	5.10000	1.0	1680	Average	Good
...	...	...	...	...	...
21592	3.60000	3.0	1530	Average	Good
21593	4.00000	2.0	2310	Average	Good
21594	4.02101	2.0	1020	Average	Average
21595	4.00000	2.0	1600	Average	Good
21596	3.25000	2.0	1020	Average	Average

21597 rows × 5 columns

```
In [637]: 1 #Set up figure and axes
2 fig, axes = plt.subplots(nrows=5, ncols=1, figsize=(20, 50))
3
4 #plot the variables to check which ones are categoricals
5 for xcol, ax in zip(np.array(data_cln.columns), axes):
6     data_cln.plot(kind='scatter', x=xcol, y='price($100,000)', ax=ax, a
```





From visual inspection, we see that the features forming vertical plots are categorical variables. They are:

1. floors
2. grade
3. condition

The only continuous feature is `sqft_living`.

Since the values for `floors` is numeric, we won't have to encode them for our regression analysis but `grade` and `condition` has be to be encoded to use in our regression beause reegression modeels can only work with numeeric valuees.

Lets take a look at `sqft_living` to check if it will need any transformations.

Before we move onto investigating transformations, we should take a closer look at the graphs above. We can see that their are a lot of outliers in the `sqft_living` and `price` columns. We deduced earlier that we can ignore houses with prices above \$1 million.



We can also see that there is very little data for sqft\_living greater than 6000 sqft. Including the houses with sqft areas bigger than 6000 sqft can potentially skew our data. Therefore we will ignore houses with sqft\_living area greater than 6000 sqft.

```
In [638]: 1 data_x = data_cln[(data_cln['sqft_living']<=6000)]
          2
          3 data_x = data_cln[(data_cln['price($100,000)']<=10)]
          4
```

```
In [639]: 1 #preview the data
          2 data_x
```

Out[639]:

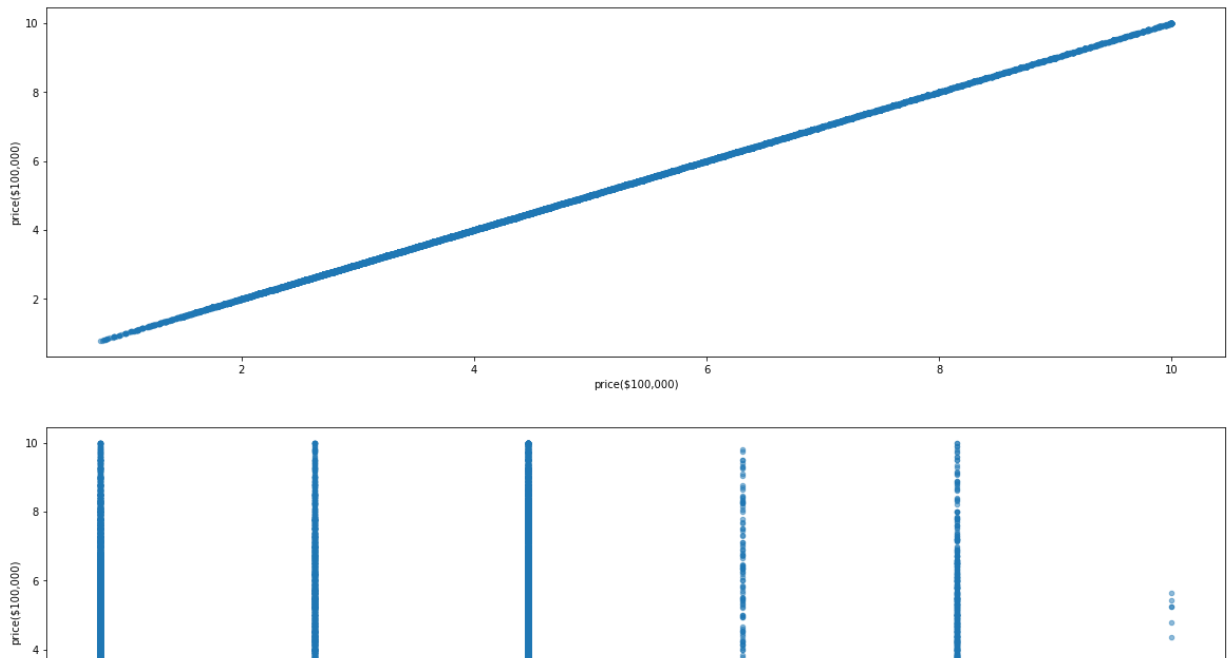
	price(\$100,000)	floors	sqft_living	condition	grade_desc
0	2.21900	1.0	1180	Average	Average
1	5.38000	2.0	2570	Average	Average
2	1.80000	1.0	770	Average	Low Average
3	6.04000	1.0	1960	Very Good	Average
4	5.10000	1.0	1680	Average	Good
...	...	...	...	...	...
21592	3.60000	3.0	1530	Average	Good
21593	4.00000	2.0	2310	Average	Good
21594	4.02101	2.0	1020	Average	Average
21595	4.00000	2.0	1600	Average	Good
21596	3.25000	2.0	1020	Average	Average

20139 rows × 5 columns

Filtering the two variables above basically means that according to our data source, it is rare to have houses that have living space areas greater than 6000 sqft or prices greater than a million dollars.

Lets recheck the graphs we created from earlier to see if things have improved.

```
In [640]: 1 #Set up figure and axes
2 fig, axes = plt.subplots(nrows=7, ncols=1, figsize=(20, 50))
3
4 #plot the variables
5 for xcol, ax in zip(np.array(data_x.columns), axes):
6     data_x.plot(kind='scatter', x=xcol, y='price($100,000)', ax=ax, alp
```



Lets investigate thee distribution of our features now. We, idally want our features to be normally distributed without skeewness. If they are not, we will transform that feature for regression analysis.

```
In [641]: 1 #scattr plot matrix
2 pd.plotting.scatter_matrix(data_x, figsize=[12, 12]);
```

We will have to transform `sqft_living` . We can see that it is a skewed normal curve.

Linear regression models perform better with variables distributed normally. Therefore, it will be in our benefit to transform this variable.

We can use some transformation techniques to improve the normal distribution characteristics for this variable but before can start transforming our data we need to split it into test and train datasets to ensure that we are not contaminating our test and train tests, we will split the data first with 25% of the dataset assigned for testing purposes.

Before we split our data, we need to check if the descriptive categorical variables at least have 2 records for each category. This is to ensure that we can stratify without issues.

```
In [642]: 1 #check grade_desc values
          2 data_x['grade_desc'].value_counts()
```

```
Out[642]: Average      8951
          Good        5873
          Better      2224
          Low Average  2033
          Very Good    694
          Fair         242
          Excellent    92
          Low          27
          Luxury        2
          Poor          1
          Name: grade_desc, dtype: int64
```

Looks like we have 2 categories with only 1 record each. We will drop these records so that we can stratify properly.

```
In [643]: 1 data_x = data_x[(data_x['grade_desc'] != 'Luxury')] #drop luxury
          2 data_x = data_x[(data_x['grade_desc'] != 'Poor')] # drop poor
```

```
In [644]: 1 data_x['grade_desc'].value_counts() #recheck values
```

```
Out[644]: Average      8951
          Good        5873
          Better      2224
          Low Average  2033
          Very Good    694
          Fair         242
          Excellent    92
          Low          27
          Name: grade_desc, dtype: int64
```

Lets check condition for the same thing and make changes accordingly.

```
In [645]: 1 data_x['condition'].value_counts() #check condition column
```

```
Out[645]: Average      13074
          Good         5347
          Very Good    1519
          Fair         168
          Poor         28
          Name: condition, dtype: int64
```

The condition column should be good. Lets move on to splitting our data. We will have to separate the target variable and the features before splitting them up.

```
In [646]: 1 pr_x = data_x['price($100,000)'] # create target variable dataset
          2 data_prep = data_x.drop(['price($100,000)'],axis=1) #create feature dataset
          3
          4 data_prep #preview
          5
```

```
Out[646]:
```

	floors	sqft_living	condition	grade_desc
0	1.0	1180	Average	Average
1	2.0	2570	Average	Average
2	1.0	770	Average	Low Average
3	1.0	1960	Very Good	Average
4	1.0	1680	Average	Good
...	...	...	...	...
21592	3.0	1530	Average	Good
21593	2.0	2310	Average	Good
21594	2.0	1020	Average	Average
21595	2.0	1600	Average	Good
21596	2.0	1020	Average	Average

Lets split our data into train and test sets.

```
In [647]: 1 #use train test split to separate the data
          2 X_train, X_test, pr_train, pr_test = train_test_split(data_prep,
          3                                                         pr_x,
          4                                                         random_state=42,
          5                                                         test_size=0.25,
          6                                                         stratify=data_prep
```

Lets apply log transformation to the sqft\_living data and check the distribution again.

```
In [648]: 1
2 #copy the datasets
3 X_train_log = X_train.copy()
4 X_test_log = X_test.copy()
5
6 #transform the X_train data
7 X_train_log['sqft_living_log'] = X_train['sqft_living'].map(lambda x: np.
8
9 #transform the X_test data
10 X_test_log['sqft_living_log'] = X_test['sqft_living'].map(lambda x: np.
11
```

```
In [649]: 1 X_train_log #preview the data
```

Out[649]:

	floors	sqft_living	condition	grade_desc	sqft_living_log
<b>12455</b>	2.0	2650	Average	Better	7.882315
<b>6914</b>	1.0	3040	Good	Good	8.019613
<b>12510</b>	2.0	2040	Average	Average	7.620705
<b>3518</b>	1.5	1550	Average	Average	7.346010
<b>6472</b>	1.0	2060	Good	Average	7.630461
...	...	...	...	...	...
<b>16348</b>	2.0	3260	Average	Better	8.089482
<b>19047</b>	1.0	1240	Good	Average	7.122867
<b>15157</b>	2.0	2570	Average	Good	7.851661
<b>14481</b>	1.0	1150	Average	Average	7.047517
<b>17902</b>	1.0	2260	Average	Better	7.723120

15102 rows × 5 columns

```
In [650]: 1 data_tr_log = X_train_log.drop(['sqft_living'],axis=1) #drop the sqft_l
2 data_test_log = X_test_log.drop(['sqft_living'],axis=1) #drop the sqft_l
```

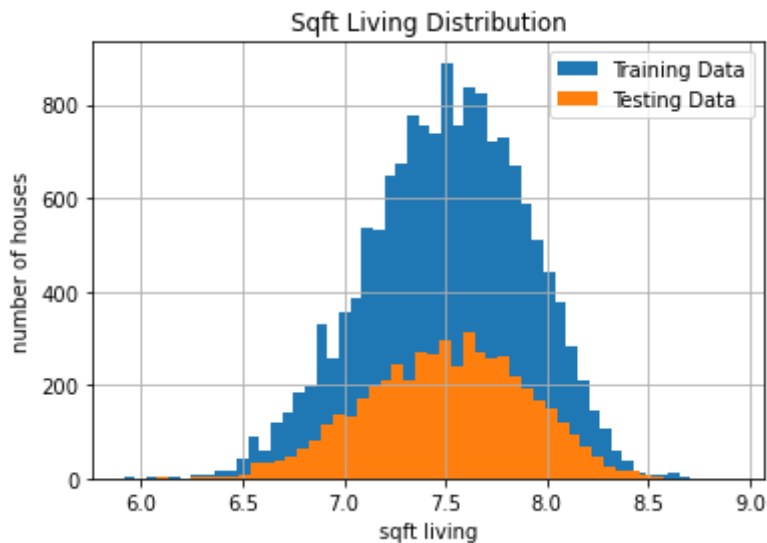
After log transformation lets scale the data also.

```
In [651]: 1 # scaler = StandardScaler() #set up scaler
2
3 # #scale train data
4 # data_tr_log['sqft_living_log_scaled'] = scaler.fit_transform(np.array
5
6 # #drop the log column
7 # data_tr_log = data_tr_log.drop(['sqft_living_log'],axis=1)
8
9
```

```
In [652]: 1 # #scale test data
2 # data_test_log['sqft_living_log_scaled'] = scaler.fit_transform(np.arr
3
4 # #drop the log column
5 # data_test_log = data_test_log.drop(['sqft_living_log'],axis=1)
6
```

Now lets check to make sure that the variable is normally distributed

```
In [653]: 1 #plot the training and testing data
2
3 #set up figure and axis
4 fig,ax = plt.subplots()
5
6
7 #plot the data
8 data_tr_log['sqft_living_log'].hist(bins=50,label='Training Data');
9
10
11 data_test_log['sqft_living_log'].hist(bins=50,label='Testing Data');
12
13 #labels and titles
14 ax.set_title('Sqft Living Distribution');
15 ax.set_xlabel('sqft living');
16 ax.set_ylabel('number of houses');
17 ax.legend();
```



Now we need to encode the descriptive categorical variables

```
In [654]: 1 X_train_ohe = pd.get_dummies(data_tr_log[['condition','grade_desc']], d
2 X_test_ohe = pd.get_dummies(data_test_log[['condition','grade_desc']],
```

```
In [655]: 1 X_train_pro = pd.concat([data_tr_log, X_train_ohe], axis=1) #merge the
          2 X_test_pro = pd.concat([data_test_log, X_test_ohe], axis=1) #merge the
          3 X_train_pro #preview
```

Out[655]:

	floors	condition	grade_desc	sqft_living_log	condition_Fair	condition_Good	condition_Poor
<b>12455</b>	2.0	Average	Better	7.882315	0	0	0
<b>6914</b>	1.0	Good	Good	8.019613	0	1	0
<b>12510</b>	2.0	Average	Average	7.620705	0	0	0
<b>3518</b>	1.5	Average	Average	7.346010	0	0	0
<b>6472</b>	1.0	Good	Average	7.630461	0	1	0
...	...	...	...	...	...	...	...
<b>16348</b>	2.0	Average	Better	8.089482	0	0	0
<b>19047</b>	1.0	Good	Average	7.122867	0	1	0
<b>15157</b>	2.0	Average	Good	7.851661	0	0	0
<b>14481</b>	1.0	Average	Average	7.047517	0	0	0
<b>17902</b>	1.0	Average	Better	7.723120	0	0	0

15102 rows × 15 columns

```
In [656]: 1 X_train_mod = X_train_pro.drop(['condition', 'grade_desc'], axis=1) #drop
          2 X_test_mod = X_test_pro.drop(['condition', 'grade_desc'], axis=1) #drop t
```

```
In [657]: 1 X_train_mod #preview
```

Out[657]:

	floors	sqft_living_log	condition_Fair	condition_Good	condition_Poor	condition_Very Good	grade_
<b>12455</b>	2.0	7.882315	0	0	0	0	
<b>6914</b>	1.0	8.019613	0	1	0	0	
<b>12510</b>	2.0	7.620705	0	0	0	0	
<b>3518</b>	1.5	7.346010	0	0	0	0	
<b>6472</b>	1.0	7.630461	0	1	0	0	
...	...	...	...	...	...	...	...
<b>16348</b>	2.0	8.089482	0	0	0	0	
<b>19047</b>	1.0	7.122867	0	1	0	0	
<b>15157</b>	2.0	7.851661	0	0	0	0	
<b>14481</b>	1.0	7.047517	0	0	0	0	
<b>17902</b>	1.0	7.723120	0	0	0	0	

15102 rows × 13 columns

This brings our Data Preparation to an end. We will now move on to modelling the variables.

## Modelling

We are going to start with a baseline model using the variable with the highest correlation to price .

During data processing, we found that `sqft_living` had the highest correlation. We will solely use that variable to develop a baseline linear regression model.

A baseline model is used to compare against models that have multiple variables. We will use the most correlated feature to the price to build it.

```
In [658]: 1 base_model = LinearRegression()
```

We will be using cross validation to perform multiple separate train-test splits within the training datasets.

This ensures that our Train dataset itself is spliced into different test sets. The r-squared value we get at the end is the mean therefore there is less likelihood of getting an anomalous r-squared value.

```
In [659]: 1 #create a splitter function to slice the training set into test samples
2 splitter = ShuffleSplit(n_splits=10, test_size=0.2, random_state=0)
3
4 #Setting baseline scores
5 baseline_scores = cross_validate(
6     estimator=base_model,
7     X=np.array(X_train_mod['sqft_living_log
8     y=pr_train,
9     return_train_score=True,
10    cv=splitter
11    )
12
13 #print the etrain and validation scores
14 print('Train score: ', baseline_scores['train_score'].mean())
15 print('Validation score: ', baseline_scores['test_score'].mean())
```

```
Train score:      0.33965784294596224
Validation score:  0.34116912874681377
```

The values above are R-squared values. R-squared value is the proportion of the variation in the dependent variable that is predictable from the independent variable.

The Validation score is the more important of the two r-squared values because it tells us how our model actually performs.

Nonetheless, the Train and Validation Score are very similar for the baseline model but both of them are low.

This is a very weak model because of the low r-squared value. A value of 0.34 means that the model can only predict 34% of the variation in the dependent variable.



Therefore we have to try adding other features to the model to investigate if the preformance improves.

We will add all of the relevant features we prepared and compare against the baseline model.

```
In [660]: 1 #create the second model
2 s_model=LinearRegression()
3
4 #create th splitter function
5 splitter = ShuffleSplit(n_splits=10, test_size=0.2, random_state=0)
6
7 #run the model and record the scores
8 second_model_scores = cross_validate(
9     estimator=s_model,
10    X=np.array(X_train_mod),
11    y=pr_train,
12    return_train_score=True,
13    cv=splitter
14 )
15 print('Second Model including grade#')
16 print('Train score: ', second_model_scores['train_score'].mean())
17 print('Validation score: ', second_model_scores['test_score'].mean())
18 print()
19 print('Baseline Model with sqft_living_log')
20 print('Train score: ', baseline_scores['train_score'].mean())
21 print('Validation score: ', baseline_scores['test_score'].mean())
```

```
Second Model including grade#
Train score:    0.45429709503946525
Validation score:    0.45665411022858293
```

```
Baseline Model with sqft_living_log
Train score:    0.33965784294596224
Validation score:    0.34116912874681377
```

With a r-squared value of 0.45, this model can predict atleast 45% of the variation in the housing prices.

This is a much better result than the baseline model using all the features that we prepared. Nonetheless, the problem is that we are not optimizing our model to the most relevant features. In order to have the best results, we have the option of using the Recursive Feature Elimination process to improve hone in on the base dfeatures to use.

```

In [661]: 1  #RFECV specific dataset
          2  X_train_RFECV = StandardScaler().fit_transform(X_train_mod)
          3
          4  #set up model for RFECV
          5  model_RFECV = LinearRegression()
          6
          7  #set up the selector and fit the model
          8  selector = RFECV(model_RFECV, cv=splitter)
          9  selector.fit(X_train_RFECV, pr_train)
         10
         11  # print the rankings
         12  # print(ft.ranking_)
         13
         14  #print which features does the model select
         15  for index, col in enumerate(X_train_mod.columns):
         16      print(f"{col}:{selector.support_[index]}")
         17
         18  #setup an array with the best features
         19  best_features = []
         20  for index, col in enumerate(X_train_mod.columns):
         21      if selector.support_[index] == True:
         22          best_features.append(col)
         23  print(best_features)
         24

```

```

floors:False
sqft_living_log:True
condition_Fair:False
condition_Good:True
condition_Poor:False
condition_Very Good:True
grade_desc_Better:True
grade_desc_Excellent:True
grade_desc_Fair:True
grade_desc_Good:True
grade_desc_Low:False
grade_desc_Low Average:True
grade_desc_Very Good:True
['sqft_living_log', 'condition_Good', 'condition_Very Good', 'grade_desc_Better', 'grade_desc_Excellent', 'grade_desc_Fair', 'grade_desc_Good', 'grade_desc_Low Average', 'grade_desc_Very Good']

```

```

In [662]: 1  #set up the final model
          2  final_model = LinearRegression()
          3
          4  #fit the model using thee best features from RFECV analysis
          5  final_model.fit(X_train_mod[best_features], pr_train)
          6
          7  #get the r-squared value
          8  final_model.score(X_test_mod[best_features], pr_test)

```

Out[662]: 0.4639812796704821

Looks like if we ignore floors and a few other variables, we get a slight improvement in our model where we move from predicting 45% of the variation in the prices to 46%.

We will use this model as our final because it is giving us the best r-squared results. Ideally, we want the r-squared values to be as high as possible so that the confidence in the model predictions is high. Generally, an r-squared value of 0.5 and higher is considered to be a strong model. With the value of 0.46, we can say that our model is close to being a strong model.

Let's set up a OLS regression model and go through the summary of the model to understand more about it.

```
In [663]: 1 sm.OLS(pr_train, sm.add_constant(X_train_mod[best_features])).fit().sum
```

Out[663]: OLS Regression Results

<b>Dep. Variable:</b>	price(\$100,000)	<b>R-squared:</b>	0.455
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.454
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	1398.
<b>Date:</b>	Fri, 18 Nov 2022	<b>Prob (F-statistic):</b>	0.00
<b>Time:</b>	16:14:35	<b>Log-Likelihood:</b>	-27029.
<b>No. Observations:</b>	15102	<b>AIC:</b>	5.408e+04
<b>Df Residuals:</b>	15092	<b>BIC:</b>	5.416e+04
<b>Df Model:</b>	9		
<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>	-7.2455	0.307	-23.620 0.000 -7.847 -6.644
<b>sqft_living_log</b>	1.4976	0.042	36.048 0.000 1.416 1.579
<b>condition_Good</b>	0.3430	0.028	12.440 0.000 0.289 0.397
<b>condition_Very Good</b>	0.9390	0.046	20.383 0.000 0.849 1.029
<b>grade_desc_Better</b>	2.0459	0.046	44.621 0.000 1.956 2.136
<b>grade_desc_Excellent</b>	3.6396	0.179	20.308 0.000 3.288 3.991
<b>grade_desc_Fair</b>	-0.7042	0.111	-6.323 0.000 -0.923 -0.486
<b>grade_desc_Good</b>	0.8641	0.030	28.498 0.000 0.805 0.924
<b>grade_desc_Low Average</b>	-0.4958	0.044	-11.346 0.000 -0.581 -0.410
<b>grade_desc_Very Good</b>	2.8722	0.073	39.520 0.000 2.730 3.015
<b>Omnibus:</b>	590.217	<b>Durbin-Watson:</b>	1.993
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	662.094
<b>Skew:</b>	0.497	<b>Prob(JB):</b>	1.69e-144
<b>Kurtosis:</b>	3.253	<b>Cond. No.</b>	203.

Notes:

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

These are the important observations to extract from the above data:

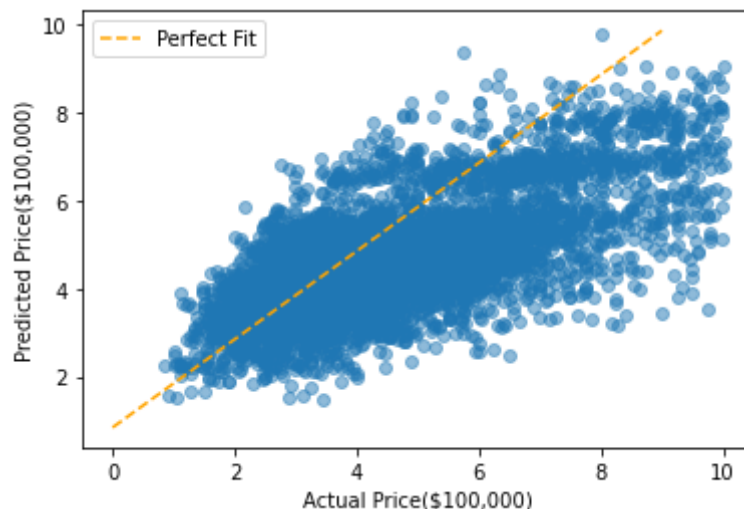
- All the **P>|t|** values are well under 0.05 therefore we can be confident that the features selected are statistically significant.

- The **Cond. No.** is an indication of multi-collinearity. This number should ideally be below 3. We have to investigate the multi-collinearity assumption to see whether our model violates it or not.
- The **Skew** is 0 and **Kurtosis** is 3 for perfect Normal Distributions but we see a slight deviation so we will investigate the Normality Assumption too.

Lets check our assumptions for the model.

## Linearity Assumption

```
In [664]: 1 #predict the values for the final model
2 preds = final_model.predict(X_test_mod[best_features])
3
4 #set up figure and axis
5 fig, ax = plt.subplots()
6
7 #set perfect line
8 perfect_line = np.arange(pr_test.min(), pr_test.max())
9
10 #plot the perfect line
11 ax.plot(perfect_line, linestyle="--", color="orange", label="Perfect Fit")
12
13 #plot the real values and the predicted values
14 ax.scatter(pr_test, preds, alpha=0.5)
15
16 #set the label
17 ax.set_xlabel("Actual Price($100,000)")
18 ax.set_ylabel("Predicted Price($100,000)")
19 ax.legend();
```

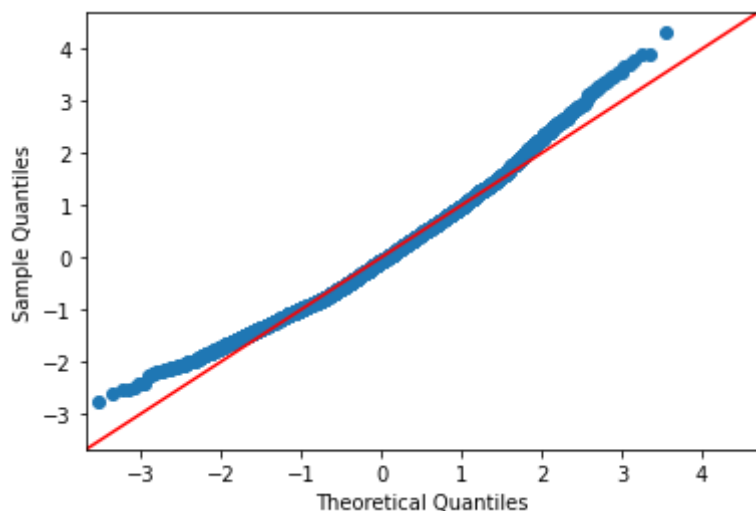


While our residuals are not perfectly linear, we can visually inspect that we are close to linearity. As the actual price tends to increase, our values are biased towards the lower end. Which means that the model is under-predicting.

## Normality Assumption:

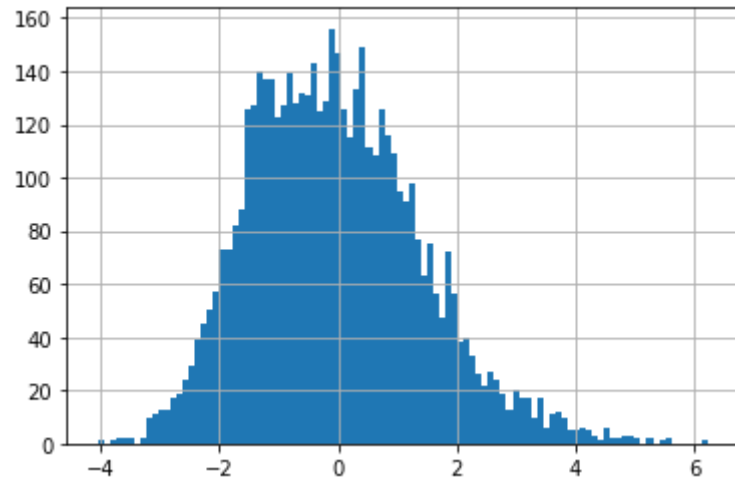
## Linearity Assumption

```
In [665]: 1 import scipy.stats as stats
2 residuals = (pr_test - preds)
3 sm.graphics.qqplot(residuals, dist=stats.norm, line='45', fit=True);
4
```



Lets check a histogram distribution for better understanding.

```
In [666]: 1 residuals.hist(bins=100);
2
```



We have a bell shaped-curve which means the residuals are normally distributed with a slight right skewness as we saw from our OLS results. Wee know from this graph, that our model is under-predicting the prices, which is something we saw in the Linearity Assumption Test also.

## Multi-Colinearity Assumption:

```
In [667]: 1 #multi-collinearity
2 X_f=X_train_log.copy()
3 a = X_train_mod[best_features]
4 vif = [variance_inflation_factor(a.values,i) for i in range(a.shape[1])]
5 pd.Series(vif, index=a.columns, name='VIF')
```

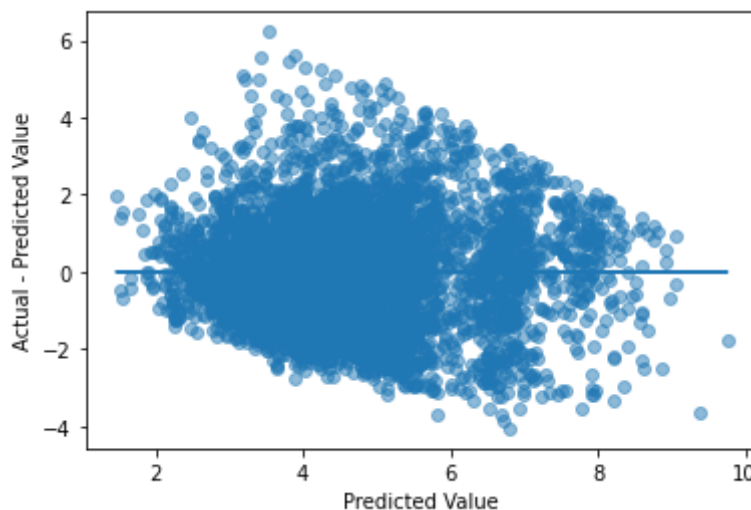
```
Out[667]: sqft_living_log      3.123886
condition_Good      1.454836
condition_Very Good  1.138377
grade_desc_Better    1.319111
grade_desc_Excellent 1.016561
grade_desc_Fair      1.023719
grade_desc_Good      1.724535
grade_desc_Low Average 1.207339
grade_desc_Very Good  1.109697
Name: VIF, dtype: float64
```

Multi-Collinearity is important to test because it causes issues with the interpretation of the coefficients. Specifically, you can interpret a coefficient as “an increase of 1 in this predictor results in a change of (coefficient) in the response variable, holding all other predictors constant.” This becomes problematic when multicollinearity is present because we can’t hold correlated predictors constant. Additionally, it increases the standard error of the coefficients, which results in them potentially showing as statistically insignificant when they might actually be significant.

Generally VIF values lower than 5 indicate that there is low or no multi-collinearity.

## Homoscedasticity:

```
In [668]: 1 fig, ax = plt.subplots()
2
3 ax.scatter(preds, residuals, alpha=0.5)
4 ax.plot(preds, [0 for i in range(len(X_test))])
5 ax.set_xlabel("Predicted Value")
6 ax.set_ylabel("Actual - Predicted Value");
```



Homoscedasticity is the is the same variance within our error terms. Homoscedasticity can impact significance tests for coefficients due to the standard errors being biased. Additionally, the confidence intervals can be either too wide or too narrow.

While, we do not have a perfect homoscedasticity, we can see that we have the same trend that we saw before which is that it is biased towards under-prediction.

## Interpreting the Final Model:

Lets first take a quick look at the mean squared error, which tells us how off our model is as compared to the real values.

```
In [669]: 1 from sklearn.metrics import mean_squared_error
          2
          3 mean_squared_error(pr_test, final_model.predict(X_test_mod[best_feature
```

Out[669]: 1.4436625399896426

This means that our model is off by almost \$140,000 from the true prediction. This shows that this model should be refined a lot more before using it for preedictions. Nonetheless, our mmodel's main goal was for inference purposes and we can achieve that by evaluating the coefficients of our final model's variables.

```
In [670]: 1 #print cefficients
          2 print(pd.Series(final_model.coef_, index=X_train_mod[best_features].col
          3 print()
          4 print("Intercept:", final_model.intercept_) #print intercept
```

```
sqft_living_log      1.497633
condition_Good       0.342965
condition_Very Good  0.938995
grade_desc_Better    2.045854
grade_desc_Excellent 3.639608
grade_desc_Fair      -0.704242
grade_desc_Good      0.864142
grade_desc_Low Average -0.495792
grade_desc_Very Good  2.872165
Name: Coefficients, dtype: float64
```

Intercept: -7.245529345435795

The intercept being negative shows that out model's basic house price without any variables is in the negative which is not plausible in reality. Nonetheless, we are still able to use it for inference purposes.



```
In [671]: 1 #reference for undeerstanding grade categories
          2 init_data['grade'].unique()
```

```
Out[671]: array(['7 Average', '6 Low Average', '8 Good', '11 Excellent', '9 Bette
r',
        '5 Fair', '10 Very Good', '12 Luxury', '4 Low', '3 Poor',
        '13 Mansion'], dtype=object)
```

Our model is biasing towards under-predicting the prices towards the higher end.

Using the coefficients, the performance of certain renovations and how they impact the housing prices is given below. These can be used to inform homeowners what kind of renovations for these features would help them improve their house's value.

- **sqft living** : For every percentage increase in square footage of living space, the house price goes up by 1.5%
  - This means that improvements, such as increasing floors or extending the house to add extra rooms and living quarters, can increase the house prices significantly.
- **condition** :
  - Intuitively, it is clear that as the condition of the house improves, the house price also increases.
  - The modell showcases that as the condition shifts to 'Good' the house price starts to improves by a mmultiple of 0.34, i.e. No obvious maintenance required but neither is everything new. Also, appearance and utility are above the standard.
  - This means that renovating houses to make sure that everything that the current house has is properly functional along with regular maintainenace can help a lot too.
  - Also, moving into the Very Good condition increases thee house value by a multiple Of 0.94 but it also reequires that All items well maintained, many having been overhauled and repaired as they have shown signs of wear, increasing the life expectancy and lowering the effective age with little deterioration or obsolescence evident with a high degree of utility.
  - Complete remodelling project and flipping homes would see this kind of a change.
- **grade** :
  - We can see that having a Fair or Low Average grade for a house poses a negative penalty to a house's value but moving into just the Good Category starts having a positive impact on the house price. the Assessor describes Good (Grade 8) as 'Just above average in construction and design. Usually better materials in both the exterior and interior finish work'.
  - As the grades keep increasing, the price starts going up substantially and at Excellent, described as 'Custom design and higher quality finish work with added amenities of solid woods, bathroom fixtures and more luxurious options' increases the price the most by a factor of 3.64 which is very substantial.
  - This means that renovations geared towards adding extra luxurious amenities can increase the price of the house significantly.

