

Final Project Submission

Please fill out:

- Student name:
- Student pace: self paced / part time / full time
- Scheduled project review date/time:
- Instructor name:
- Blog post URL:

BUSINESS UNDERSTANDING

A real estate agency from King County, Seattle hired us for a project to analyse how different factors affect prices of homes. The aim of this analysis is to build a multiple linear regression model that predicts the prices of houses in King County, Seattle.

As a data scientist analyzing the King County housing market, my business understanding is that the real estate industry is a crucial sector that plays a significant role in the economy. The success of a real estate transaction depends on several factors, including the location, the size of the property, the condition of the property, the amenities, and the current market conditions. The housing market is subject to various external factors such as interest rates, economic conditions, and government policies that can impact the demand and supply of properties.

Datasets

The dataset contains information about the houses in King County, Seattle. The dataset has 21 variables including the price, number of bedrooms, bathrooms, square footage of the living area, and other variables. The dataset contains 21,597 observations.

The scope of this analysis is limited to the data provided. We will use feature engineering techniques such as imputation, normalization, and one-hot encoding to preprocess the data. We will use multiple linear regression model. We will evaluate the performance of the model using metrics such as mean squared error, mean absolute error, and R-squared.

To overcome these challenges, we need to use a combination of quantitative and qualitative analysis techniques and incorporate domain knowledge and expertise. By understanding the King County housing market's complexities and using data-driven insights, we can help real estate agents and property owners make informed decisions about pricing, marketing, and selling their properties, ultimately leading to more successful real estate transactions and a more robust housing market.

BUSINESS PROBLEM

By developing a model that can accurately predict the sale price of houses, real estate agents can better advise their clients on pricing strategies, investors can identify potentially undervalued properties, and homeowners can better estimate the value of their own properties. This can ultimately lead to more efficient and profitable real estate transactions in King County.

You are charged with exploring what factors most significantly affect home prices. You must then translate those findings into actionable insights that the real estate agency can use to help decide what factors to consider when advising potential home buyers.

Business Objectives

1. Develop a pricing model: Create a predictive pricing model that incorporates the factors identified in the regression analysis, such as the number of bathrooms, living area, lot size, and condition and grade ratings. This model can help the agency more accurately price their properties, particularly those with unique features such as waterfront views.
2. Refine marketing strategies: Use the insights gained from the regression analysis to refine the agency's marketing strategies. For example, the agency can create targeted marketing campaigns that emphasize the features that have the greatest impact on price, such as the number of bathrooms, living area, and condition and grade ratings.
3. Analyze seasonal trends: Analyze the seasonal trends in home prices and develop strategies for pricing and marketing homes throughout the year. For example, the agency can adjust their pricing and marketing strategies to take advantage of the higher prices in the spring, when homes tend to sell for more.
4. Optimize home renovations: Use the insights gained from the regression analysis to identify which renovations have the greatest impact on a home's price. This information can be used to guide homeowners who are considering renovations and to help the agency market homes that have been recently renovated.

In [1]:

```
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
import missingno as msno
import numpy as np
import pandas as pd
pd.options.display.max_columns = 30
import statsmodels.api as sm
from sklearn.metrics import mean_absolute_error

# Group Libraries
import Functions as fun
```

DATA UNDERSTANDING:

In this project, we are analyzing the King County housing market to build a multiple linear regression model that predicts the prices of houses in King County, Seattle.

In [2]:

```
# Your code here - remember to use markdown cells for comments as well!
import pandas as pd
house_df=pd.read_csv('kc_house_data.csv')
house_df.head()
```

Out[2]:

	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	
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	

In [3]:

```
#Getting data information
house_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21597 entries, 0 to 21596
Data columns (total 21 columns):
#   Column                Non-Null Count  Dtype
---  -
0   id                    21597 non-null  int64
1   date                 21597 non-null  object
2   price               21597 non-null  float64
3   bedrooms            21597 non-null  int64
4   bathrooms           21597 non-null  float64
5   sqft_living         21597 non-null  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
```

The dataset provided contains information on 21,597 houses in King County, Seattle.

To better understand the data, we identified the categorical and numerical variables in the dataset that is relevant to our business problems shown below:

- **Numerical Columns (15)**

date - Date house was sold

price - Sale price (prediction target)

bedrooms - Number of bedrooms

bathrooms - Number of bathrooms

sqft_living - Square footage of living space in the home

sqft_lot - Square footage of the lot

floors - Number of floors (levels) in house

sqft_above - Square footage of house apart from basement

sqft_basement - Square footage of the basement

yr_built - Year when house was built

yr_renovated - Year when house was renovated

lat - Latitude coordinate

long - Longitude coordinate

sqft_living15 - The square footage of interior housing living space for the nearest 15 neighbors

sqft_lot15 - The square footage of the land lots of the nearest 15 neighbors

- **Categorical Columns (6)**

id - Unique ID for each home sold

waterfront - Whether the house has a view to a waterfront

view - An index from 0 to 4 of how good the view of the property was

condition - An index from 1 to 5 on the condition of the house

grade - An index from 1 to 13, where 1-3 falls short of building construction and design, 7 has an average level of construction and design, and 11-13 have a high quality level of construction and design

zipcode - What zipcode area the house is in

We chose these columns because they are key features that are often used to determine the value of a residential property in real estate.

The purpose of this exercise is to analyze and comprehend the information contained within the columns of the provided CSV file. Our aim is to carefully examine the data, identify patterns and correlations between the variables, and extract meaningful insights from it.

In [4]:

```
# Calculate the basic statistical summary
house_df.describe()
```

Out[4]:

	id	price	bedrooms	bathrooms	sqft_living	sqft_lot
count	2.159700e+04	2.159700e+04	21597.000000	21597.000000	21597.000000	2.159700e+04
mean	4.580474e+09	5.402966e+05	3.373200	2.115826	2080.321850	1.509941e+04
std	2.876736e+09	3.673681e+05	0.926299	0.768984	918.106125	4.141264e+04
min	1.000102e+06	7.800000e+04	1.000000	0.500000	370.000000	5.200000e+02
25%	2.123049e+09	3.220000e+05	3.000000	1.750000	1430.000000	5.040000e+03
50%	3.904930e+09	4.500000e+05	3.000000	2.250000	1910.000000	7.618000e+03
75%	7.308900e+09	6.450000e+05	4.000000	2.500000	2550.000000	1.068500e+04
max	9.900000e+09	7.700000e+06	33.000000	8.000000	13540.000000	1.651359e+06

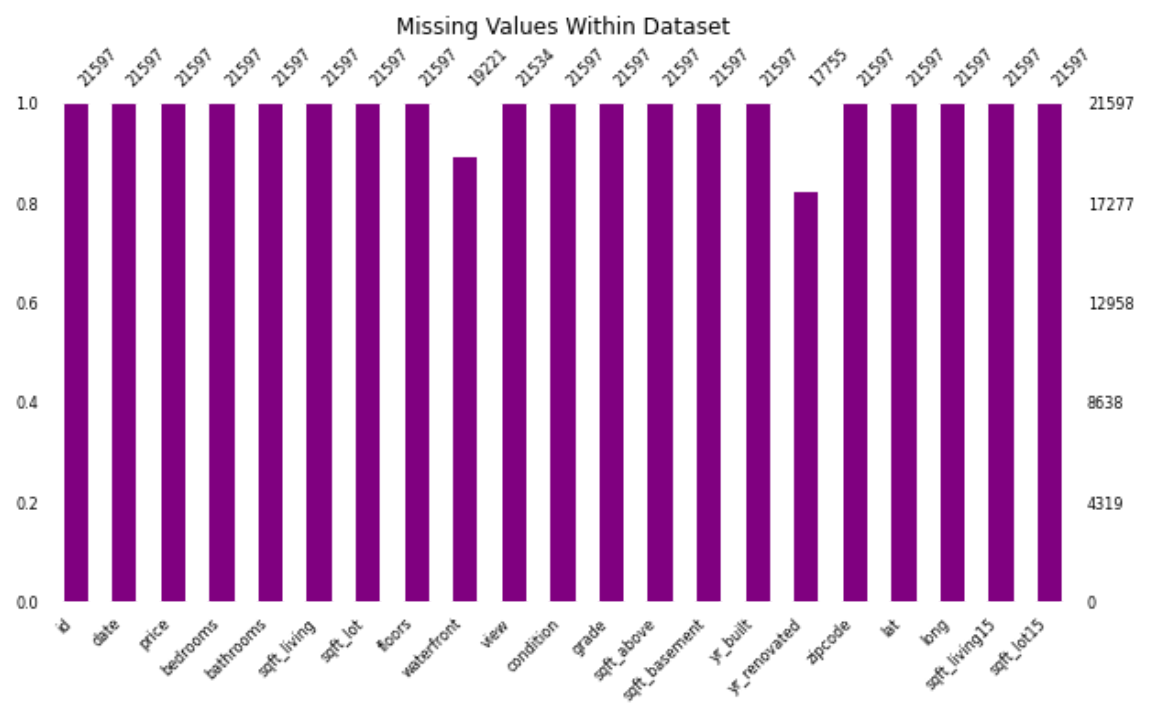
Data Cleaning

Detecting missing values

We will start by visualizing our missing data.

In [5]:

```
# Visualise the missing values in the dataset
msno.bar(house_df, color='purple', figsize=(10, 5), fontsize=8)
plt.title('Missing Values Within Dataset');
```



In [6]:

```
house_df.isnull().sum()
```

Out[6]:

```
id                0
date              0
price             0
bedrooms          0
bathrooms         0
sqft_living       0
sqft_lot          0
floors            0
waterfront       2376
view              63
condition         0
grade            0
sqft_above        0
sqft_basement     0
yr_built          0
yr_renovated     3842
zipcode           0
lat               0
long              0
sqft_living15     0
sqft_lot15        0
dtype: int64
```

From the above bar graph, we can see that there exist missing data in waterfront, view and the year renovated.

We shall replace missing values of waterfront, view with mode and the year renovated with the median.

In [7]:

```
fun.fun_mode_fill_null(house_df, 'waterfront')
fun.fun_mode_fill_null(house_df, 'view')
fun.fun_median_fill_null(house_df, 'yr_renovated')
house_df.isnull().sum()
```

Out[7]:

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

In [8]:

```
# Check for duplicates in the 'id' column
house_df.id.duplicated().sum()
```

Out[8]:

177

There are 177 duplivates in the unique column data. We will use this to eliminate all duplicates with the same id.

In [9]:

```
fun.fun_duplicates_drop(house_df, 'id')
house_df.id.duplicated().sum()
```

Out[9]:

0

Detecting outliers

We shall then check for outliers.

In [10]:

```
def fun_outlier_plot_box(df, column_name):  
    """  
    Create a box plot for a specified column of a Pandas DataFrame using Seaborn.  
  
    Parameters:  
        df (pandas.DataFrame): The DataFrame containing the column to plot.  
        column_name (str): The name of the column to plot.  
  
    Returns:  
        None  
    """  
    sns.boxplot(x=df[column_name])
```

In [11]:

```
def replace_outliers_with_mode(df, columns):  
    """  
    This function detects outliers in the specified columns of a DataFrame and replaces  
  
    Parameters:  
        df (pandas.DataFrame): The DataFrame to modify.  
        columns (list): A list of column names to check for outliers.  
  
    Returns:  
        pandas.DataFrame: The modified DataFrame with outliers replaced by mode.  
    """  
    for col in columns:  
        q1 = df[col].quantile(0.005)  
        q3 = df[col].quantile(0.995)  
        iqr = q3 - q1  
        lower_bound = q1 - 1.5 * iqr  
        upper_bound = q3 + 1.5 * iqr  
        mode = df[col].mode()[0]  
        df.loc[(df[col] < lower_bound) | (df[col] > upper_bound), col] = mode  
  
    outliers = {}  
    for col in columns:  
        q1 = df[col].quantile(0.005)  
        q3 = df[col].quantile(0.995)  
        iqr = q3 - q1  
        lower_bound = q1 - 1.5 * iqr  
        upper_bound = q3 + 1.5 * iqr  
        col_outliers = df[(df[col] < lower_bound) | (df[col] > upper_bound)]  
        outliers[col] = len(col_outliers)  
  
    print("Total number of outliers replaced with mode:")  
    print(outliers)  
  
    return df
```


In [12]:

```
house_df.columns
```

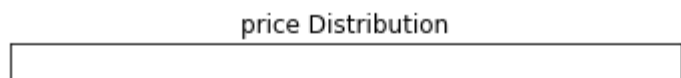
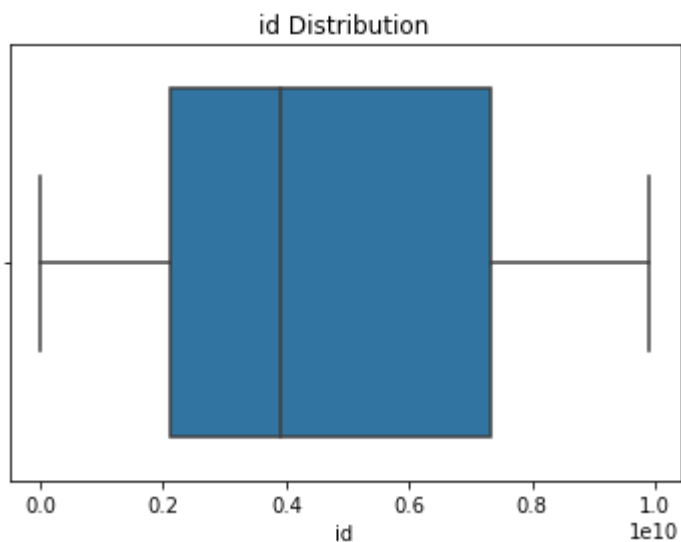
Out[12]:

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

In [13]:

```
columns = ['id', 'price', 'bedrooms', 'bathrooms', 'sqft_living', 'sqft_lot', 'floors', 'yr_built', 'yr_renovated', 'zipcode']

# iterate through the columns of house_df and call fun_outlier_plot_box()
for column in columns:
    fun_outlier_plot_box(house_df, column)
    plt.title(f"{column} Distribution")
    plt.show()
```



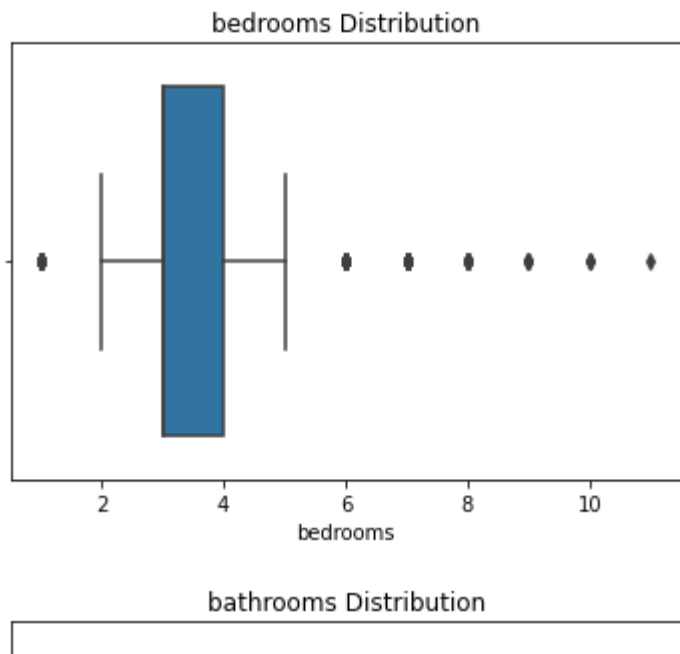
In [14]:

```
# print out the total number of outliers replaced with mode for each column.
columns = ['bedrooms', 'bathrooms', 'sqft_living', 'sqft_lot', 'floors', 'yr_built', 'sqft_above', 'sqft_living15', 'sqft_lot15']
house_df = replace_outliers_with_mode(house_df, columns)
```

```
Total number of outliers replaced with mode:
{'bedrooms': 0, 'bathrooms': 0, 'sqft_living': 0, 'sqft_lot': 0, 'floors': 0, 'yr_built': 0, 'sqft_above': 0, 'sqft_living15': 0, 'sqft_lot15': 0}
```

In [15]:

```
# iterate through the columns of house_df and call fun_outlier_plot_box()
for column in columns:
    fun_outlier_plot_box(house_df, column)
    plt.title(f"{column} Distribution")
    plt.show()
```



We have now completed our data cleaning, we shall explore the various columns.

Univariate Analysis

In this section, we'll explore each column in the dataset to see the distributions of features and obtain some useful insights. The main two parts in this section are:

- Categorical Columns
- Numerical Columns

2.1.1 Categorical Columns

There are 5 Categorical Columns in the dataset that we shall be analysing:

- id
- waterfront
- view
- condition
- grade
- zipcode

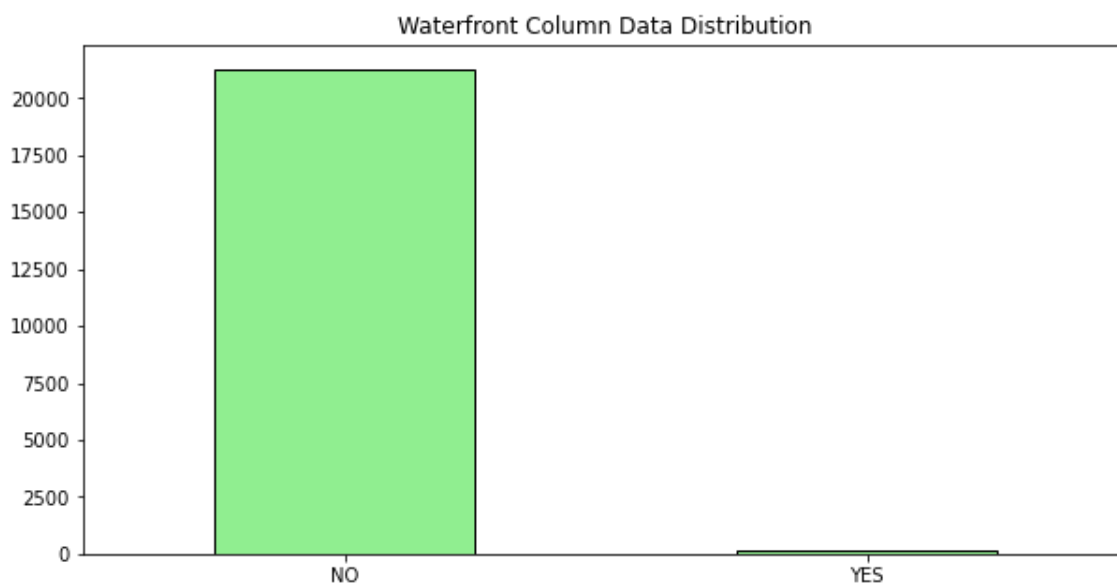
In [16]:

```
def fun_plot_value_counts(df, col, title):  
    '''  
    Returns the value counts of a column in a dataframe and  
    plots the value counts of a column in a dataframe as a bar chart  
    '''  
    counts = df[col].value_counts(dropna=False)  
    print(counts)  
    counts.plot(kind='bar', figsize=(10, 5), color='lightgreen', edgecolor='black')  
    plt.title(title)  
    plt.xticks(rotation=0)  
    plt.show()  
    return counts
```

In [17]:

```
fun_plot_value_counts(house_df, 'waterfront', 'Waterfront Column Data Distribution')
```

```
NO      21274  
YES       146  
Name: waterfront, dtype: int64
```



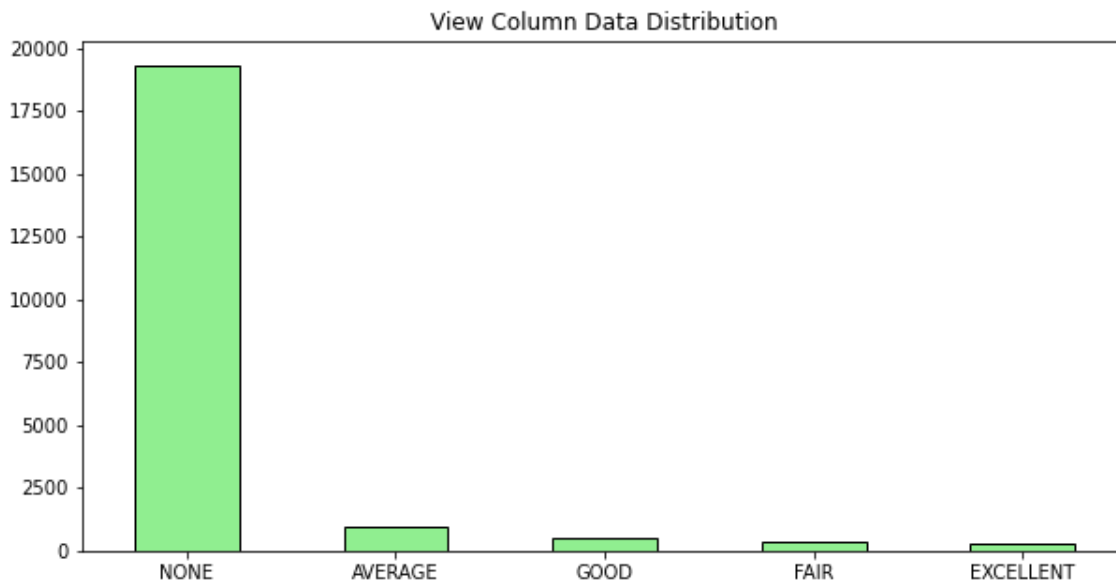
Out[17]:

```
NO      21274  
YES       146  
Name: waterfront, dtype: int64
```

In [18]:

```
fun_plot_value_counts(house_df, 'view', 'View Column Data Distribution')
```

```
NONE          19316  
AVERAGE       956  
GOOD           505  
FAIR           329  
EXCELLENT     314  
Name: view, dtype: int64
```



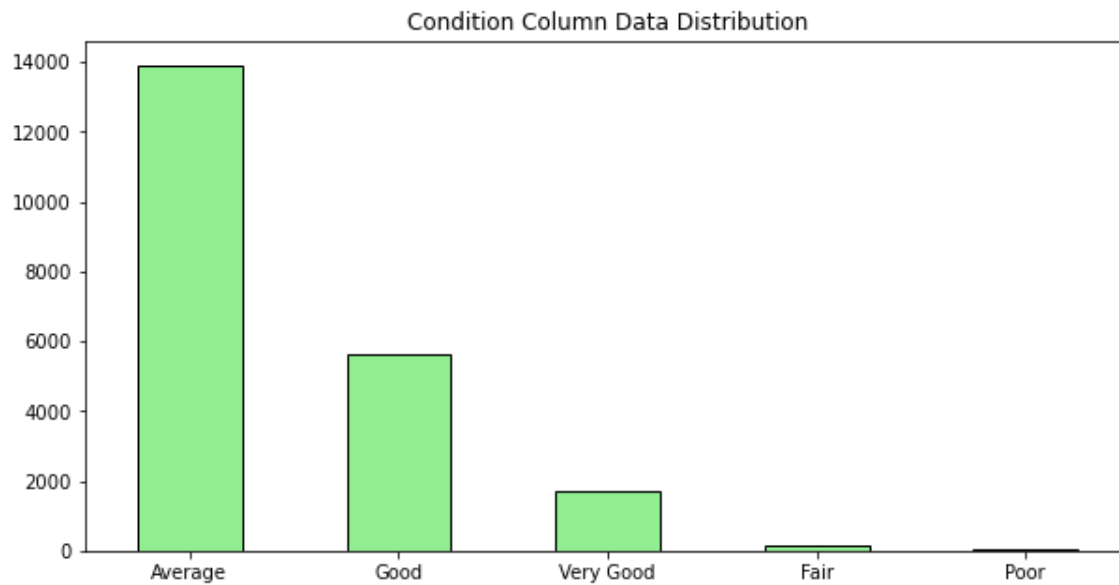
Out[18]:

```
NONE          19316  
AVERAGE       956  
GOOD           505  
FAIR           329  
EXCELLENT     314  
Name: view, dtype: int64
```

In [19]:

```
fun_plot_value_counts(house_df, 'condition', 'Condition Column Data Distribution')
```

```
Average      13900  
Good          5643  
Very Good    1687  
Fair          162  
Poor           28  
Name: condition, dtype: int64
```



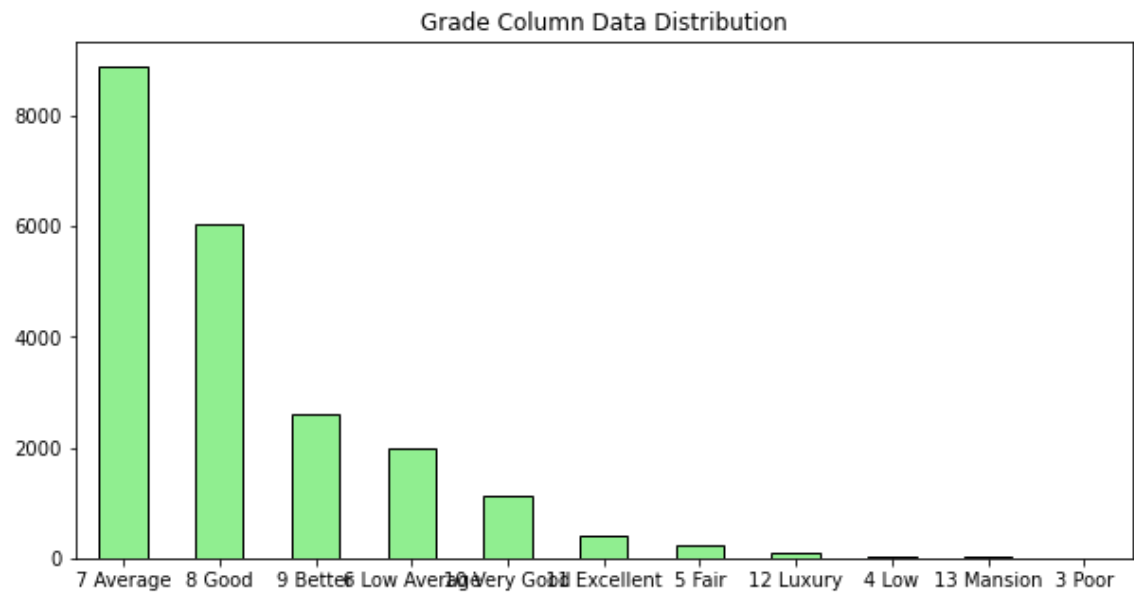
Out[19]:

```
Average      13900  
Good          5643  
Very Good    1687  
Fair          162  
Poor           28  
Name: condition, dtype: int64
```

In [20]:

```
fun_plot_value_counts(house_df, 'grade', 'Grade Column Data Distribution')
```

```
7 Average      8889
8 Good         6041
9 Better       2606
6 Low Average  1995
10 Very Good   1130
11 Excellent    396
5 Fair         234
12 Luxury       88
4 Low          27
13 Mansion     13
3 Poor         1
Name: grade, dtype: int64
```



Out[20]:

```
7 Average      8889
8 Good         6041
9 Better       2606
6 Low Average  1995
10 Very Good   1130
11 Excellent    396
5 Fair         234
12 Luxury       88
4 Low          27
13 Mansion     13
3 Poor         1
Name: grade, dtype: int64
```

In [21]:

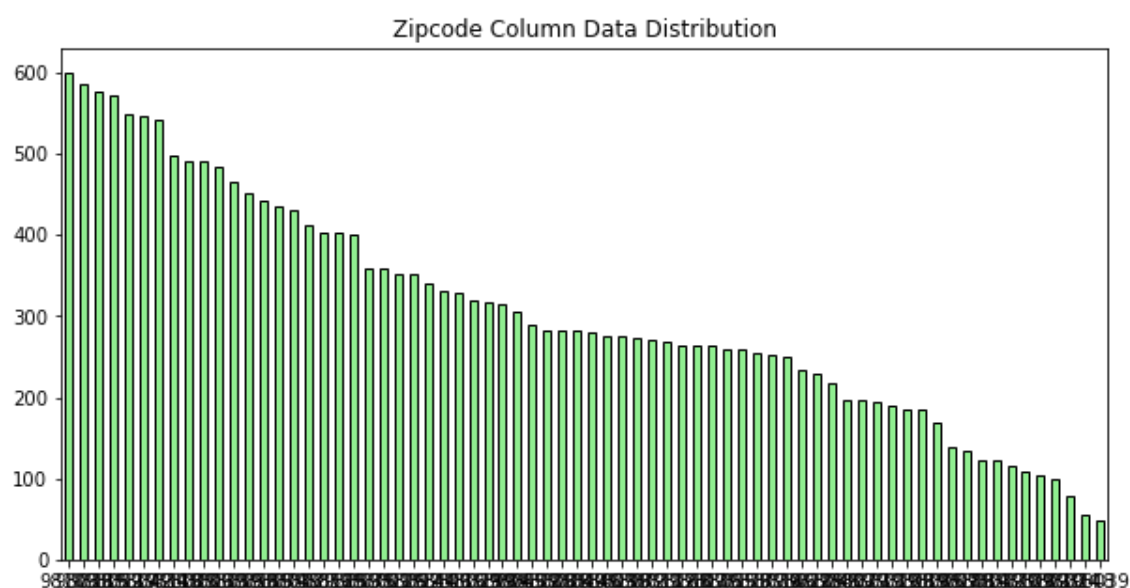
```
fun_plot_value_counts(house_df, 'zipcode', 'Zipcode Column Data Distribution')
```

```
98103    600
98038    586
98115    576
98052    571
98117    548
```

...

```
98102    104
98010     99
98024     79
98148     56
98039     49
```

Name: zipcode, Length: 70, dtype: int64



Out[21]:

```
98103    600
98038    586
98115    576
98052    571
98117    548
```

...

```
98102    104
98010     99
98024     79
98148     56
98039     49
```

Name: zipcode, Length: 70, dtype: int64

Numerical Columns

There are 15 Numerical Columns in the dataset that we shall be analysing:

- date
- price
- bedrooms
- bathrooms

- sqft_living
- sqft_lot
- floors
- sqft_above
- sqft_basement
- yr_built
- yr_renovated
- lat
- long
- sqft_living15
- sqft_lot15

In [22]:

```
def fun_describe_and_plot_distribution(df, col, title):  
    '''  
    Returns the statistics of a column in a dataframe and  
    plots the distribution of a column in a dataframe as a histogram, kde, and boxplot  
    '''  
  
    # print the statistics  
    print(df[col].describe())  
  
    # create a figure composed of two matplotlib.Axes objects (ax_box and ax_hist)  
    f, (ax_box, ax_hist) = plt.subplots(2, sharex=True, gridspec_kw={"height_ratios": (.  
  
    # assign a graph to each ax  
    sns.boxplot(df[col], ax=ax_box, color='lightgreen')  
    sns.histplot(data=df, x=col, ax=ax_hist, kde=True, color='lightgreen', bins='auto',  
  
    # set the title and layout  
    plt.suptitle(title)  
    plt.tight_layout()  
  
    # show the plot  
    plt.show()
```

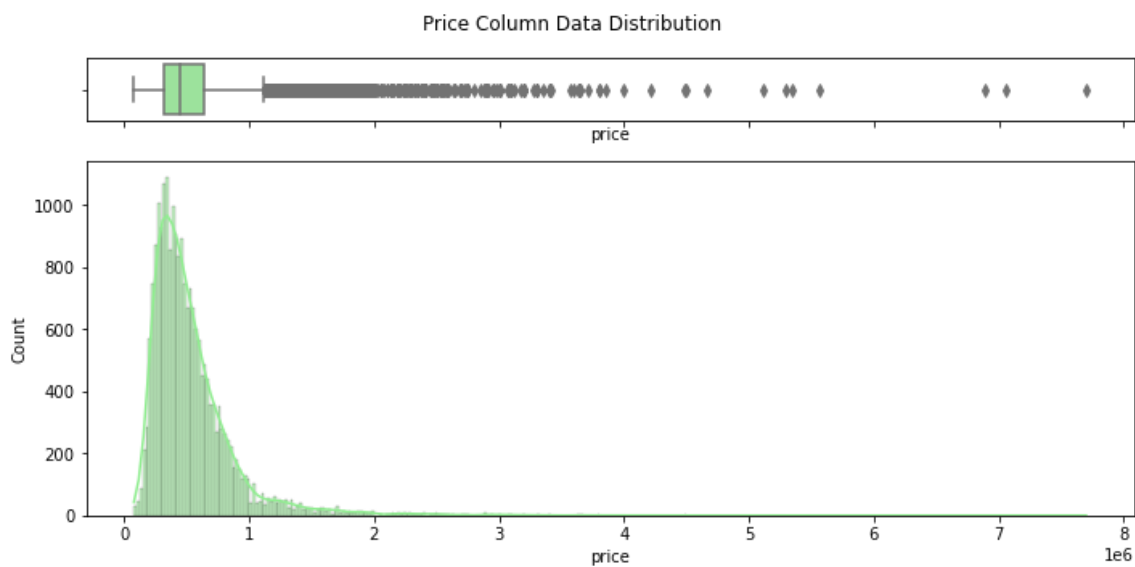
In [23]:

```
# fun_describe_and_plot_distribution(house_df, 'season', 'Seasons Column Data Distributi
```


In [24]:

```
fun_describe_and_plot_distribution(house_df, 'price', 'Price Column Data Distribution')
```

```
count    2.142000e+04  
mean     5.407393e+05  
std      3.679311e+05  
min      7.800000e+04  
25%      3.225000e+05  
50%      4.500000e+05  
75%      6.450000e+05  
max      7.700000e+06  
Name: price, dtype: float64
```

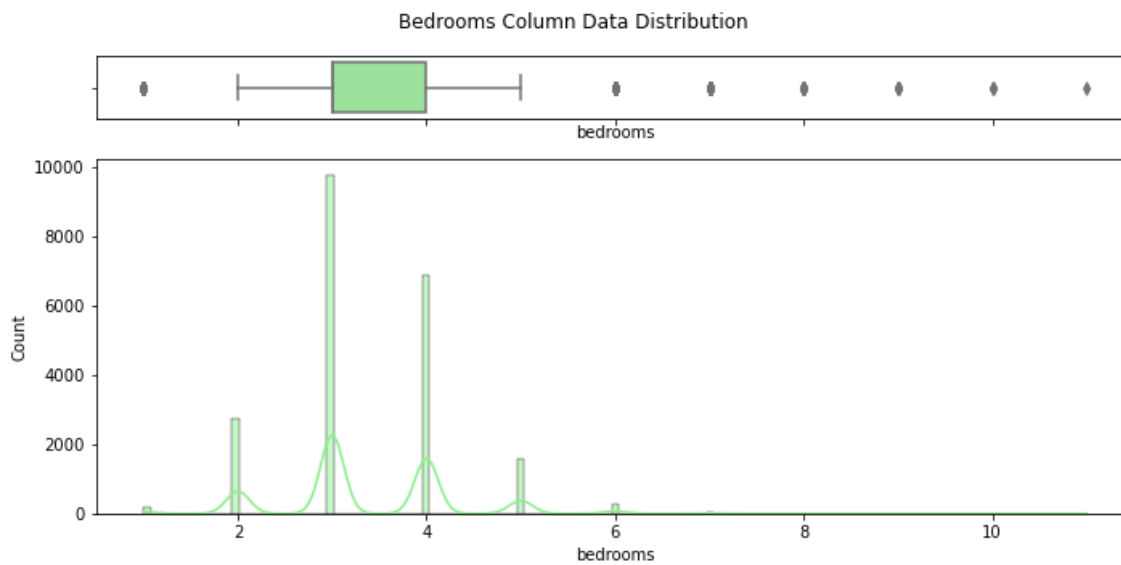


From the distribution above, we see that the price column is skewed to the right. This means that the mean price of the homes in the dataset are skewed too. The minimum price of a house in the dataset is 78,000, and the maximum price of a house in the dataset is 7,700,000. The mean price of a house in the dataset is 540,297, and the median price of a house in the dataset is 450,000. The standard deviation of the price column is 367,368.

In [25]:

```
fun_describe_and_plot_distribution(house_df, 'bedrooms', 'Bedrooms Column Data Distribut
```

```
count    21420.000000
mean      3.372549
std       0.902995
min       1.000000
25%       3.000000
50%       3.000000
75%       4.000000
max       11.000000
Name: bedrooms, dtype: float64
```

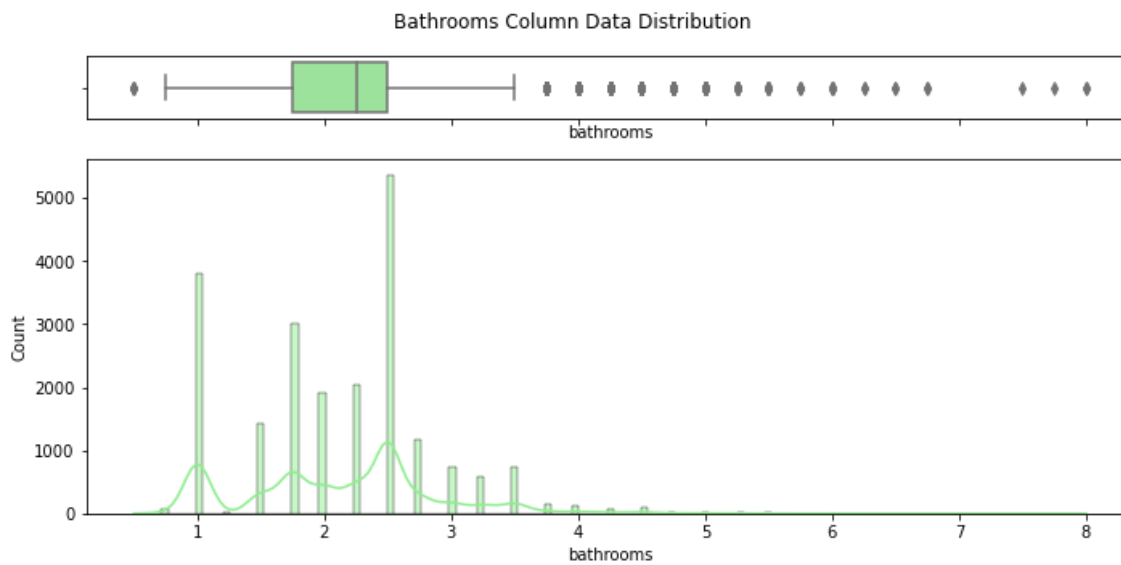


The bedroom column distribution is not skewed as the and is normally distributed.

In [26]:

```
fun_describe_and_plot_distribution(house_df, 'bathrooms', 'Bathrooms Column Data Distrib
```

```
count    21420.000000
mean      2.118429
std       0.768720
min       0.500000
25%       1.750000
50%       2.250000
75%       2.500000
max       8.000000
Name: bathrooms, dtype: float64
```

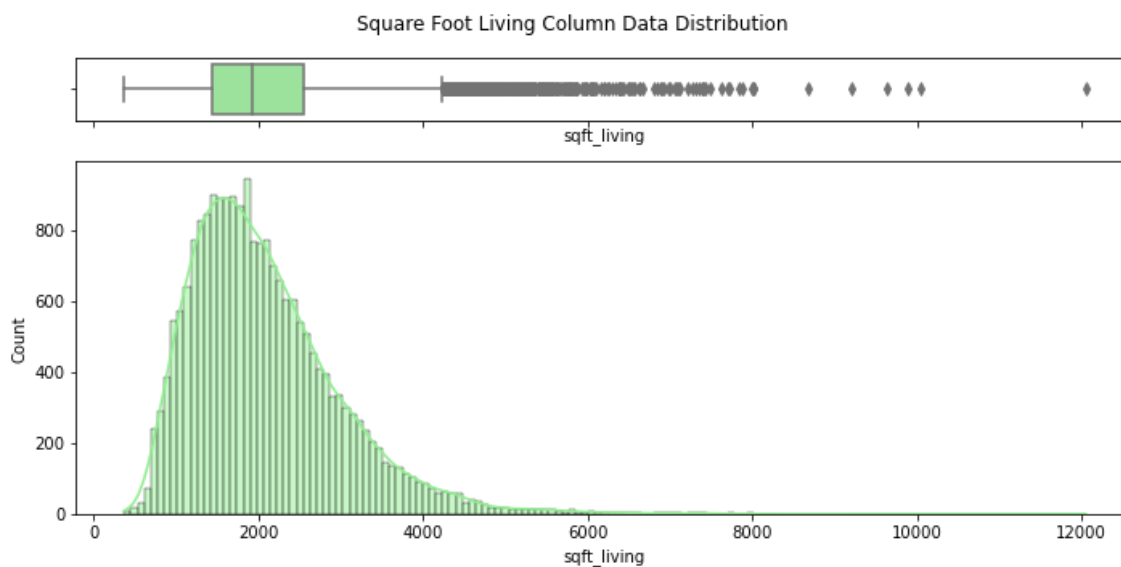


From the distribution above we can see that the bathroom column is not skewed. This is because the mean and median are almost the same. The minimum number of bathrooms in a house in the dataset is 0.5, and the maximum number of bathrooms in a house in the dataset is 8. The mean number of bathrooms in a house in the dataset is 2.12, and the median number of bathrooms in a house in the dataset is 2.25. The standard deviation of the bathrooms column is 0.77.

In [27]:

```
fun_describe_and_plot_distribution(house_df, 'sqft_living', 'Square Foot Living Column D
```

```
count    21420.000000
mean      2082.561204
std       915.482938
min       370.000000
25%      1430.000000
50%      1920.000000
75%      2550.000000
max      12050.000000
Name: sqft_living, dtype: float64
```

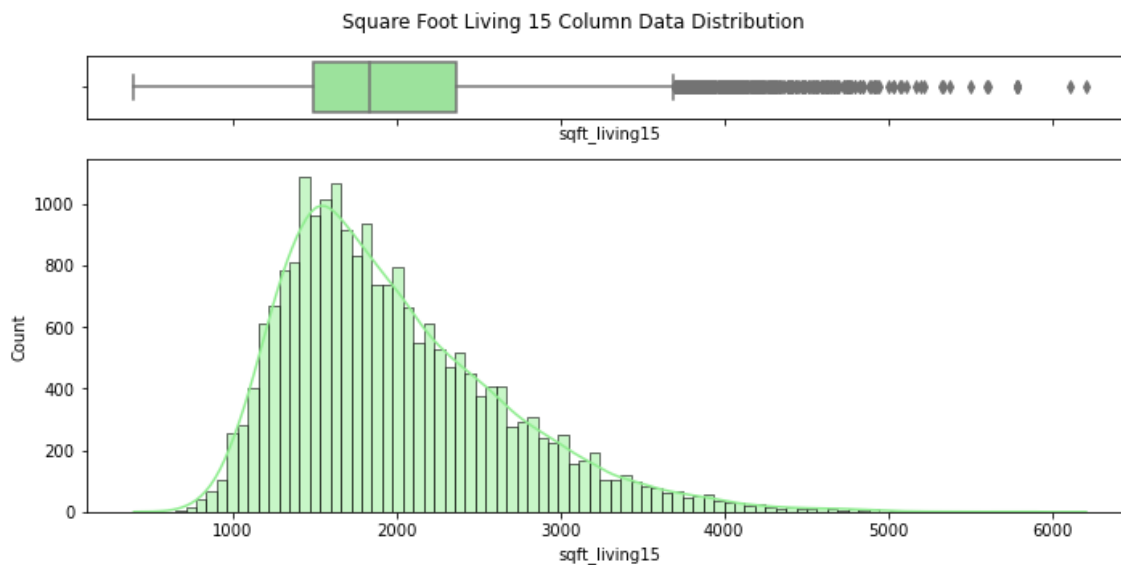


From the distribution above, we can see that the sqft living column is skewed to the right. This means that the mean square footage of the homes is greater than the median. The minimum square footage of a house in the dataset is 370, and the maximum square footage of a house in the dataset is 13,540. The mean square footage of a house in the dataset is 2080, and the median square footage of a house in the dataset is 1910. The standard deviation of the sqft living column is 918.

In [28]:

```
fun_describe_and_plot_distribution(house_df, 'sqft_living15', 'Square Foot Living 15 Col
```

```
count    21420.000000
mean      1988.384080
std       685.537057
min       399.000000
25%      1490.000000
50%      1840.000000
75%      2370.000000
max       6210.000000
Name: sqft_living15, dtype: float64
```

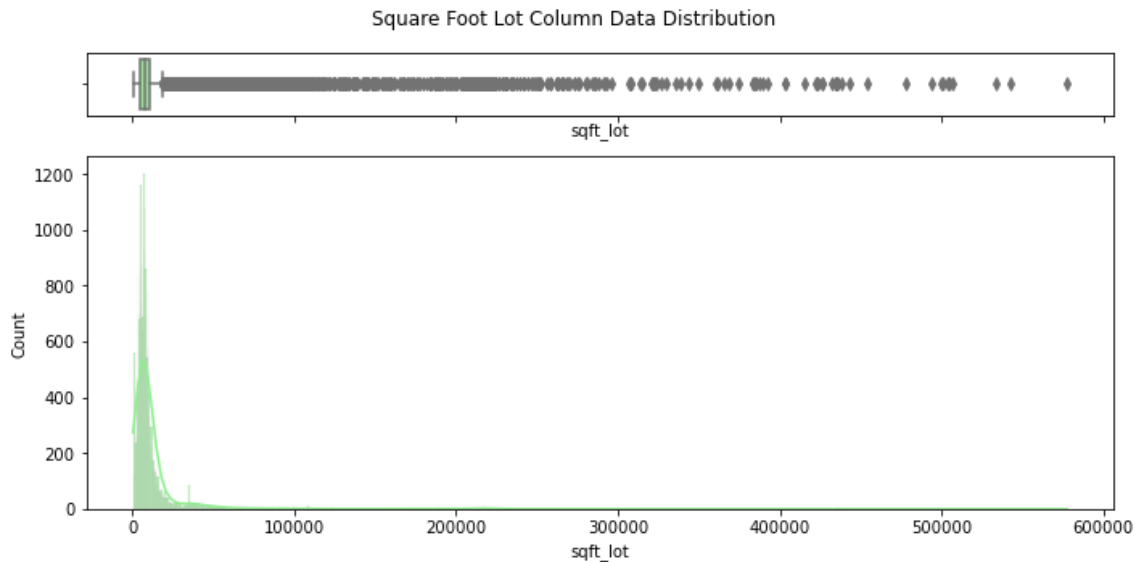


From the distributions above, we can see that the data is skewed to the right. This is as a result of the mean being greater than the median.

In [29]:

```
fun_describe_and_plot_distribution(house_df, 'sqft_lot', 'Square Foot Lot Column Data Di
```

```
count    21420.000000  
mean     14512.810504  
std      33439.663853  
min       520.000000  
25%      5038.000000  
50%      7605.000000  
75%     10660.250000  
max     577605.000000  
Name: sqft_lot, dtype: float64
```

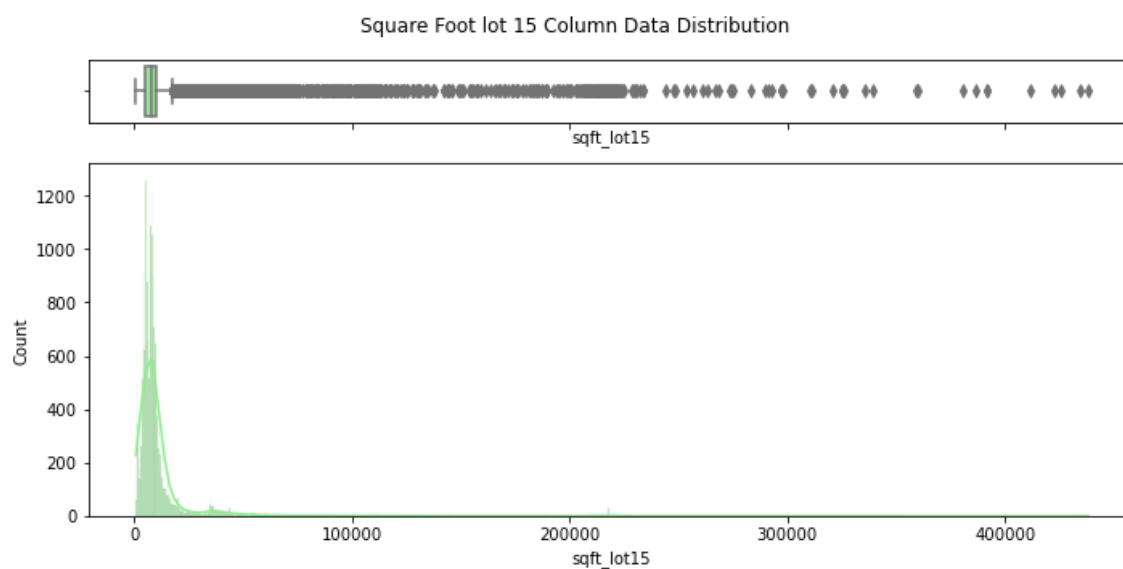


From the distribution above, we can see that the data is skewed to the right. This is because the mean is greater than the median.

In [30]:

```
fun_describe_and_plot_distribution(house_df, 'sqft_lot15', 'Square Foot lot 15 Column Da
```

```
count    21420.000000
mean     12669.511391
std      25806.976744
min       651.000000
25%      5100.000000
50%      7620.000000
75%     10083.000000
max     438213.000000
Name: sqft_lot15, dtype: float64
```

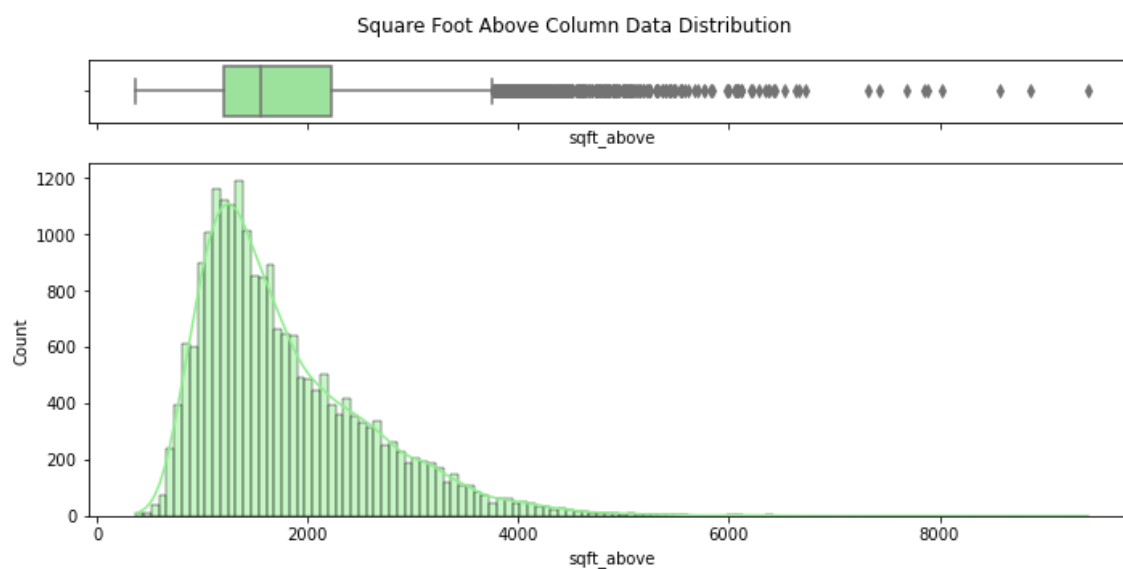


In the distributions above we see a much more skewed to the right column.

In [31]:

```
fun_describe_and_plot_distribution(house_df, 'sqft_above', 'Square Foot Above Column Dat
```

```
count    21420.000000
mean      1791.170215
std       828.692965
min       370.000000
25%      1200.000000
50%      1560.000000
75%      2220.000000
max       9410.000000
Name: sqft_above, dtype: float64
```

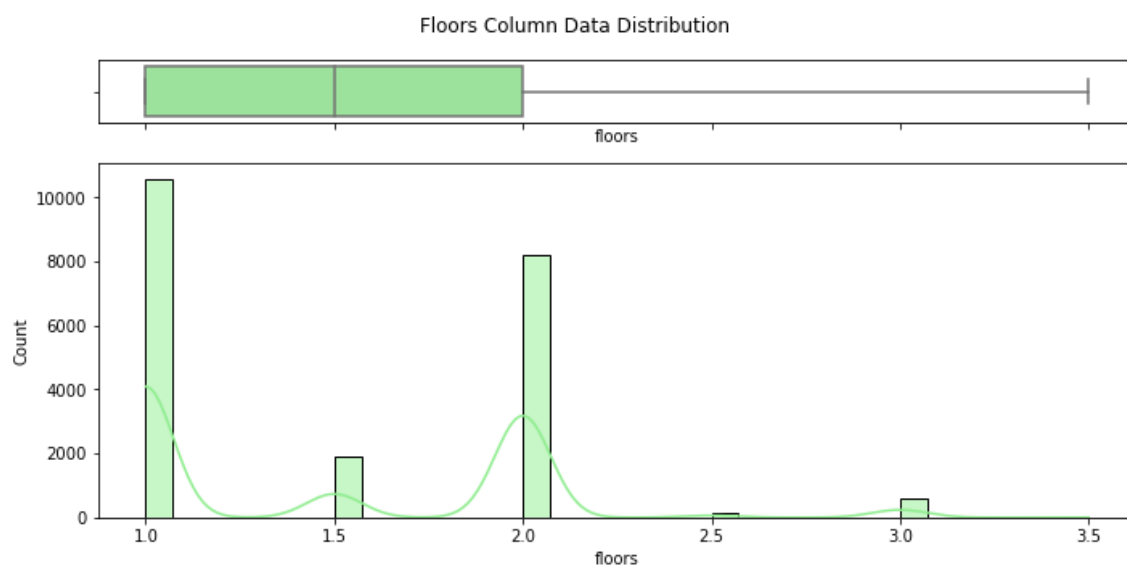


From the distributions above, we see that the square footage above ground of the houses in this dataset is skewed to the right. This is because the mean is greater than the median.

In [32]:

```
fun_describe_and_plot_distribution(house_df, 'floors', 'Floors Column Data Distribution')
```

```
count    21420.000000
mean      1.495985
std       0.540081
min       1.000000
25%       1.000000
50%       1.500000
75%       2.000000
max       3.500000
Name: floors, dtype: float64
```

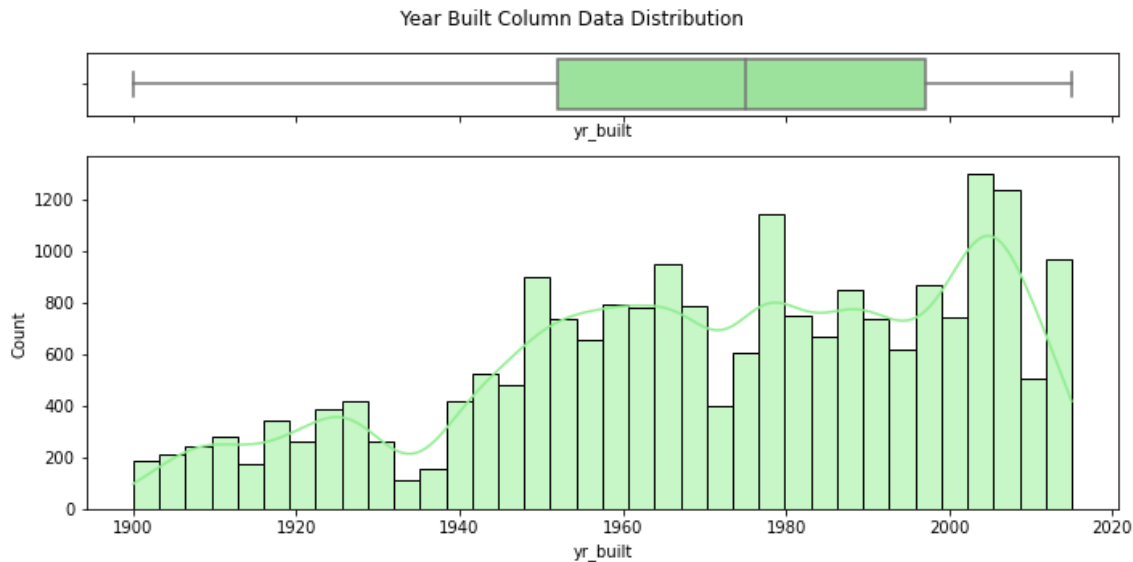


From the distributions above, there is no particular trend in the floors column data.

In [33]:

```
fun_describe_and_plot_distribution(house_df, 'yr_built', 'Year Built Column Data Distrib
```

```
count    21420.000000
mean      1971.092997
std        29.387141
min       1900.000000
25%       1952.000000
50%       1975.000000
75%       1997.000000
max       2015.000000
Name: yr_built, dtype: float64
```

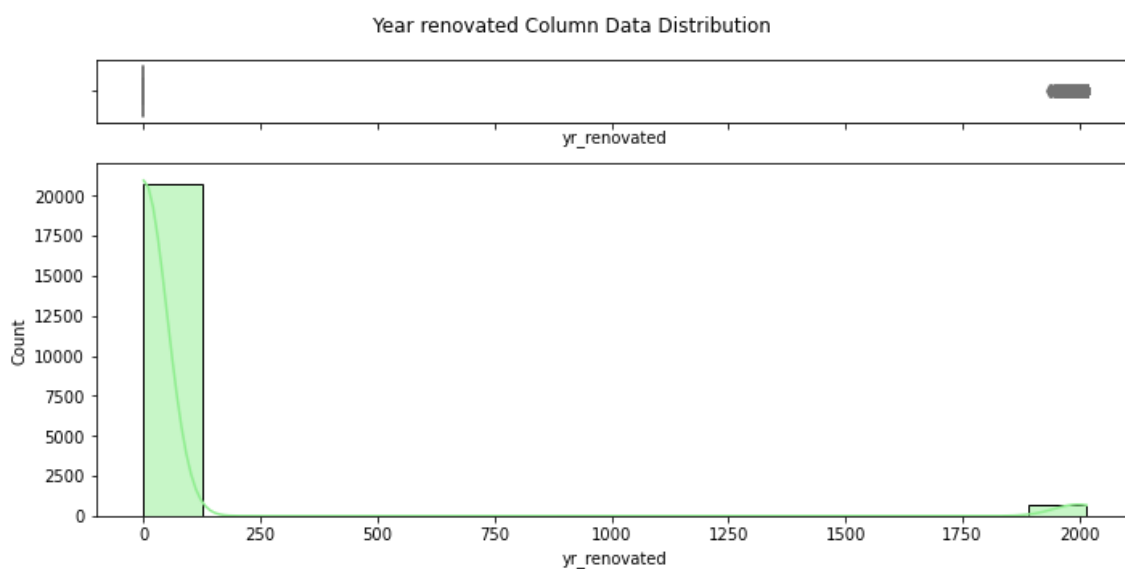


From the distributions above we can see that the data is slightly skewed to the left. This is because the mean is slightly lower than the median. The oldest house in the dataset was built in 1900, and the newest house in the dataset was built in 2015. The mean year the houses in the dataset were built is 1971, and the median year the houses in the dataset were built is 1975. The standard deviation of the yr built column is 29.

In [34]:

```
fun_describe_and_plot_distribution(house_df, 'yr_renovated', 'Year renovated Column Data
```

```
count    21420.000000
mean      68.956723
std       364.552298
min        0.000000
25%        0.000000
50%        0.000000
75%        0.000000
max      2015.000000
Name: yr_renovated, dtype: float64
```



From the distribution and value counts above, we can see that the data has a number of zeros. This could either be suggesting that the house has not been renovated, or that the data is missing. Furthermore, there is also some missing data in this column. We shall be analysing the data more indepth in the next phase to see how to deal with the zeros and the missing values in the column.

DATA PREPARATION

Use the date feature to create a new feature called season, which represents whether the home was sold in Spring, Summer, Fall, or Winter.

In [35]:

```
#To create a new feature called "season" in the King County housing dataset,
#we can extract the month information from the "date" column and map it to
#the corresponding season.

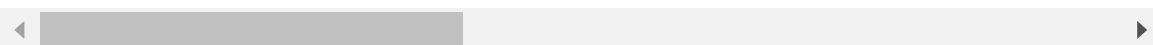
# Convert the date column to a datetime object
house_df['date'] = pd.to_datetime(house_df['date'])

# Extract the month information and map it to the corresponding season
seasons = {1: 'Winter', 2: 'Winter', 3: 'Spring', 4: 'Spring', 5: 'Spring',
           6: 'Summer', 7: 'Summer', 8: 'Summer', 9: 'Fall', 10: 'Fall',
           11: 'Fall', 12: 'Winter'}
house_df['season'] = house_df['date'].dt.month.map(seasons)
seasons = house_df['season']
house_df
```

Out[35]:

	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	water
0	7129300520	2014-10-13	221900.0	3	1.00	1180	5650	1.0	
1	6414100192	2014-12-09	538000.0	3	2.25	2570	7242	2.0	
2	5631500400	2015-02-25	180000.0	2	1.00	770	10000	1.0	
3	2487200875	2014-12-09	604000.0	4	3.00	1960	5000	1.0	
4	1954400510	2015-02-18	510000.0	3	2.00	1680	8080	1.0	
...	
21592	263000018	2014-05-21	360000.0	3	2.50	1530	1131	3.0	
21593	6600060120	2015-02-23	400000.0	4	2.50	2310	5813	2.0	
21594	1523300141	2014-06-23	402101.0	2	0.75	1020	1350	2.0	
21595	291310100	2015-01-16	400000.0	3	2.50	1600	2388	2.0	
21596	1523300157	2014-10-15	325000.0	2	0.75	1020	1076	2.0	

21420 rows × 22 columns



With the new "season" feature, we can now analyze the King County housing dataset with respect to the season and identify any seasonality trends in the housing market. For example, we can investigate whether homes sell for higher prices in certain seasons, or whether certain types of homes are more popular in certain seasons.

adding a column to store the age of the houses

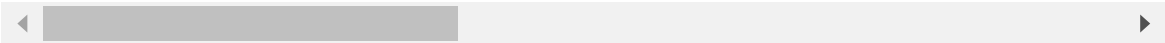
In [36]:

```
# Add house_age column
house_df['age'] = house_df['date'].dt.year - house_df['yr_built']
house_df
```

Out[36]:

	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	wate
0	7129300520	2014-10-13	221900.0	3	1.00	1180	5650	1.0	
1	6414100192	2014-12-09	538000.0	3	2.25	2570	7242	2.0	
2	5631500400	2015-02-25	180000.0	2	1.00	770	10000	1.0	
3	2487200875	2014-12-09	604000.0	4	3.00	1960	5000	1.0	
4	1954400510	2015-02-18	510000.0	3	2.00	1680	8080	1.0	
...	
21592	2630000018	2014-05-21	360000.0	3	2.50	1530	1131	3.0	
21593	6600060120	2015-02-23	400000.0	4	2.50	2310	5813	2.0	
21594	1523300141	2014-06-23	402101.0	2	0.75	1020	1350	2.0	
21595	291310100	2015-01-16	400000.0	3	2.50	1600	2388	2.0	
21596	1523300157	2014-10-15	325000.0	2	0.75	1020	1076	2.0	

21420 rows × 23 columns



removing null values in the 'yr_built' column and adding the 'renovated' column to show whether the house has been renovated or not

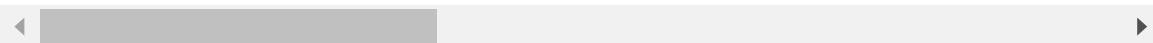
In [37]:

```
house_df.loc[house_df.yr_renovated.isnull(), 'yr_renovated'] = 0
house_df['renovated'] = house_df['yr_renovated'].apply(lambda x: 0 if x == 0 else 1)
house_df
```

Out[37]:

	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	water
0	7129300520	2014-10-13	221900.0	3	1.00	1180	5650	1.0	
1	6414100192	2014-12-09	538000.0	3	2.25	2570	7242	2.0	
2	5631500400	2015-02-25	180000.0	2	1.00	770	10000	1.0	
3	2487200875	2014-12-09	604000.0	4	3.00	1960	5000	1.0	
4	1954400510	2015-02-18	510000.0	3	2.00	1680	8080	1.0	
...
21592	263000018	2014-05-21	360000.0	3	2.50	1530	1131	3.0	
21593	6600060120	2015-02-23	400000.0	4	2.50	2310	5813	2.0	
21594	1523300141	2014-06-23	402101.0	2	0.75	1020	1350	2.0	
21595	291310100	2015-01-16	400000.0	3	2.50	1600	2388	2.0	
21596	1523300157	2014-10-15	325000.0	2	0.75	1020	1076	2.0	

21420 rows × 24 columns



In [38]:

```
lot = house_df['sqft_lot']
```

Add has_basement column that is a binary value

In [39]:

```
house_df['sqft_basement'] = house_df['sqft_basement'].replace('?', '0').astype('float')
house_df['has_basement'] = house_df['sqft_basement'].apply(lambda x: 0 if x == 0 else 1)
house_df['has_basement']
```

Out[39]:

```
0      0
1      1
2      0
3      1
4      0
..
21592   0
21593   0
21594   0
21595   0
21596   0
Name: has_basement, Length: 21420, dtype: int64
```

Ordinal Encoding

Ordinal encoding converts each label into integer values and the encoded data represents the sequence of labels

the values in the `condition` and `grade` columns are ordinal, and have been assigned a value based on the quality of the feature. Therefore, we will be ordinal encoding these columns.

Creating a function that maps the ordinal values in a dataframe column to corresponding numerical values based on a provided dictionary.

In [40]:

```
def map_ordinal_values(df, col_name, value_dict):
    # map the ordinal values to numerical values using the provided dictionary
    df[col_name] = df[col_name].map(value_dict)
    return df
```

In [41]:

```
condition_dict = {'Poor': 1, 'Fair': 2, 'Average': 3, 'Good': 4, 'Very Good': 5}
grade_dict = {'3 Poor': 3, '4 Low': 4, '5 Fair': 5, '6 Low Average': 6, '7 Average': 7,

df = map_ordinal_values(house_df, 'condition', condition_dict)
df = map_ordinal_values(house_df, 'grade', grade_dict)

print(house_df[['condition', 'grade']])
```

	condition	grade
0	3	7
1	3	7
2	3	6
3	5	7
4	3	8
...
21592	3	8
21593	3	8
21594	3	7
21595	3	8
21596	3	7

[21420 rows x 2 columns]

One Hot Encoding

One hot encoding is a process of converting categorical data variables so they can be provided to machine learning algorithms to improve predictions.

We shall be encoding the categorical columns (waterfront and view) using one hot encoding.

Furthermore, in order to avoid the "Dummy Variable Trap" (perfect multicollinearity between the independent variables), we will need to drop one of the columns created.

In [42]:

```
house_df.select_dtypes('object')
```

Out[42]:

	waterfront	view	season
0	NO	NONE	Fall
1	NO	NONE	Winter
2	NO	NONE	Winter
3	NO	NONE	Winter
4	NO	NONE	Winter
...
21592	NO	NONE	Spring
21593	NO	NONE	Winter
21594	NO	NONE	Summer
21595	NO	NONE	Winter
21596	NO	NONE	Fall

21420 rows × 3 columns

creating a function to perform one hot encoding on the specified columns

In [43]:

```
def one_hot_encode(df, columns):  
    df = pd.get_dummies(df, columns=columns, drop_first=False)  
    return df
```

one hot encoding categorical variable

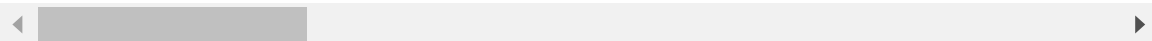
In [44]:

```
columns_to_encode = ['waterfront', 'view', 'season']
house_df = one_hot_encode(house_df, columns_to_encode)
# Preview the dataframe
house_df
```

Out[44]:

	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	conc
0	7129300520	2014-10-13	221900.0	3	1.00	1180	5650	1.0	
1	6414100192	2014-12-09	538000.0	3	2.25	2570	7242	2.0	
2	5631500400	2015-02-25	180000.0	2	1.00	770	10000	1.0	
3	2487200875	2014-12-09	604000.0	4	3.00	1960	5000	1.0	
4	1954400510	2015-02-18	510000.0	3	2.00	1680	8080	1.0	
...	
21592	263000018	2014-05-21	360000.0	3	2.50	1530	1131	3.0	
21593	6600060120	2015-02-23	400000.0	4	2.50	2310	5813	2.0	
21594	1523300141	2014-06-23	402101.0	2	0.75	1020	1350	2.0	
21595	291310100	2015-01-16	400000.0	3	2.50	1600	2388	2.0	
21596	1523300157	2014-10-15	325000.0	2	0.75	1020	1076	2.0	

21420 rows × 33 columns



In the `waterfront` column, we shall be dropping the `waterfront_NO` column as the reference column. This will allow us to study the effect of having a house on a waterfront. In the `view` column, we shall be dropping the `view_NONE` column as the reference column. This will allow us to study the effect of having a house with a view. In addition, it is the most common value in the column. In the `season` column, we shall be dropping the `season_Fall` column as the reference column.

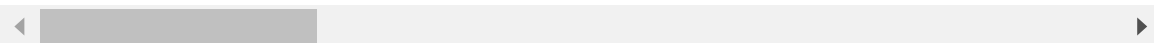
In [45]:

```
# Drop the 'waterfront_NO' and 'view_NONE' columns
house_df. drop(['waterfront_NO', 'view_NONE', 'season_Fall'], axis=1, inplace=True)
# Preview the dataframe
house_df
```

Out[45]:

	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	conc
0	7129300520	2014-10-13	221900.0	3	1.00	1180	5650	1.0	
1	6414100192	2014-12-09	538000.0	3	2.25	2570	7242	2.0	
2	5631500400	2015-02-25	180000.0	2	1.00	770	10000	1.0	
3	2487200875	2014-12-09	604000.0	4	3.00	1960	5000	1.0	
4	1954400510	2015-02-18	510000.0	3	2.00	1680	8080	1.0	
...	
21592	263000018	2014-05-21	360000.0	3	2.50	1530	1131	3.0	
21593	6600060120	2015-02-23	400000.0	4	2.50	2310	5813	2.0	
21594	1523300141	2014-06-23	402101.0	2	0.75	1020	1350	2.0	
21595	291310100	2015-01-16	400000.0	3	2.50	1600	2388	2.0	
21596	1523300157	2014-10-15	325000.0	2	0.75	1020	1076	2.0	

21420 rows × 30 columns



Checking correlations and multicollinearity

In [46]:

house_df.corr()

Out[46]:

	id	price	bedrooms	bathrooms	sqft_living	sqft_lot	floc
id	1.000000	-0.018525	0.001883	0.004343	-0.012772	-0.147247	0.0181
price	-0.018525	1.000000	0.316714	0.526229	0.701474	0.092118	0.2562
bedrooms	0.001883	0.316714	1.000000	0.528917	0.593872	0.042097	0.1835
bathrooms	0.004343	0.526229	0.528917	1.000000	0.753490	0.095069	0.5018
sqft_living	-0.012772	0.701474	0.593872	0.753490	1.000000	0.187937	0.3524
sqft_lot	-0.147247	0.092118	0.042097	0.095069	0.187937	1.000000	-0.0126
floors	0.018139	0.256286	0.183570	0.501803	0.352412	-0.012663	1.0000
condition	-0.024300	0.034779	0.022743	-0.129362	-0.061506	-0.010054	-0.2668
grade	0.006470	0.668020	0.367750	0.665587	0.762937	0.124879	0.4580
sqft_above	-0.011794	0.605294	0.493566	0.686328	0.873978	0.201702	0.5227
sqft_basement	-0.004999	0.320842	0.303463	0.278251	0.424127	0.021685	-0.2421
yr_built	0.021171	0.052906	0.160365	0.506252	0.317203	0.062244	0.4889
yr_renovated	-0.010857	0.118484	0.018196	0.047686	0.051154	0.009679	0.0031
zipcode	-0.007934	-0.052371	-0.159637	-0.203884	-0.199427	-0.145556	-0.0585
lat	-0.002766	0.305744	-0.012974	0.023143	0.050440	-0.095566	0.0489
long	0.019336	0.020983	0.137018	0.223808	0.240074	0.255619	0.1248
sqft_living15	-0.003883	0.584549	0.406240	0.569453	0.756327	0.173483	0.2790
sqft_lot15	-0.143110	0.080538	0.031936	0.081868	0.176510	0.751488	-0.0149
age	-0.021012	-0.052828	-0.160513	-0.506618	-0.317606	-0.062335	-0.4892
renovated	-0.010864	0.118179	0.017924	0.047259	0.050925	0.009821	0.0030
has_basement	0.002900	0.177593	0.159585	0.159723	0.200218	-0.034651	-0.2528
waterfront_YES	-0.003628	0.264898	-0.002132	0.063635	0.105169	0.028139	0.0205
view_AVERAGE	0.015031	0.147497	0.046281	0.086847	0.134255	0.044683	0.0072
view_EXCELLENT	0.016497	0.302693	0.034854	0.104761	0.161965	0.026444	0.0239
view_FAIR	0.001027	0.092720	0.022043	0.037798	0.066980	-0.008106	-0.0236
view_GOOD	-0.012658	0.182488	0.051067	0.111356	0.157840	0.072036	0.0202
season_Spring	0.012454	0.023218	-0.003585	-0.012337	-0.012679	-0.000636	-0.0117
season_Summer	0.001552	0.010247	0.010045	0.023024	0.023292	-0.010268	0.0188
season_Winter	-0.000054	-0.025421	0.002088	-0.010052	-0.007483	0.004067	-0.0119

creating a function that takes in our dataframe and returns corellations between column in pair in descending order

In [47]:

```
def get_correlation_df(df):  
    corr_df = df.corr().abs().stack().reset_index().sort_values(0, ascending=False)  
    corr_df['pairs'] = list(zip(corr_df.level_0, corr_df.level_1))  
    corr_df.set_index(['pairs'], inplace=True)  
    corr_df.drop(columns=['level_1', 'level_0'], inplace=True)  
    corr_df.columns = ['cc']  
    corr_df = corr_df.drop_duplicates()  
    return corr_df.head(50)
```

let call our function on our dataframe

In [48]:

```
get_correlation_df(house_df)
```

Out[48]:

	cc
pairs	
(id, id)	1.000000
(renovated, yr_renovated)	0.999968
(yr_built, age)	0.999874
(sqft_above, sqft_living)	0.873978
(has_basement, sqft_basement)	0.820906
(grade, sqft_living)	0.762937
(sqft_living, sqft_living15)	0.756327
(sqft_above, grade)	0.756221
(sqft_living, bathrooms)	0.753490
(sqft_lot, sqft_lot15)	0.751488
(sqft_above, sqft_living15)	0.731887
(grade, sqft_living15)	0.713178
(sqft_living, price)	0.701474
(bathrooms, sqft_above)	0.686328
(grade, price)	0.668020
(bathrooms, grade)	0.665587
(price, sqft_above)	0.605294
(bedrooms, sqft_living)	0.593872
(sqft_living15, price)	0.584549
(view_EXCELLENT, waterfront_YES)	0.570601
(bathrooms, sqft_living15)	0.569453
(long, zipcode)	0.564778
(bathrooms, bedrooms)	0.528917
(price, bathrooms)	0.526229
(floors, sqft_above)	0.522751
(bathrooms, age)	0.506618
(bathrooms, yr_built)	0.506252
(bathrooms, floors)	0.501803
(bedrooms, sqft_above)	0.493566
(floors, age)	0.489231
(yr_built, floors)	0.488935
(floors, grade)	0.458091
(grade, age)	0.446642
(grade, yr_built)	0.446235
(sqft_living, sqft_basement)	0.424127
(season_Summer, season_Spring)	0.423483

cc	
pairs	
(age, sqft_above)	0.423299
(sqft_above, yr_built)	0.422977
(long, yr_built)	0.409173
(long, age)	0.409138
(sqft_living15, bedrooms)	0.406240
(grade, bedrooms)	0.367750
(yr_built, condition)	0.365129
(age, condition)	0.364416
(floors, sqft_living)	0.352412
(zipcode, yr_built)	0.346162
(zipcode, age)	0.346143
(long, sqft_above)	0.344161
(sqft_living15, long)	0.334679
(sqft_living15, age)	0.325009

we will drop columns that have strong multicollinearity or provide no use to the model

In [49]:

```
house_df = house_df.drop(columns=['id', 'yr_renovated', 'sqft_lot', 'sqft_above', 'sqft_
```

```
house_df
```

Out[49]:

	price	bedrooms	bathrooms	sqft_living	floors	condition	grade	age	renovated
0	221900.0	3	1.00	1180	1.0	3	7	59	0
1	538000.0	3	2.25	2570	2.0	3	7	63	1
2	180000.0	2	1.00	770	1.0	3	6	82	0
3	604000.0	4	3.00	1960	1.0	5	7	49	0
4	510000.0	3	2.00	1680	1.0	3	8	28	0
...
21592	360000.0	3	2.50	1530	3.0	3	8	5	0
21593	400000.0	4	2.50	2310	2.0	3	8	1	0
21594	402101.0	2	0.75	1020	2.0	3	7	5	0
21595	400000.0	3	2.50	1600	2.0	3	8	11	0
21596	325000.0	2	0.75	1020	2.0	3	7	6	0

21420 rows × 18 columns

There is still some multicollinearity between predictor variable, but not strong enough to initially drop on our models.

In [50]:

```
house_df.corr()['price'].sort_values(ascending=False)
```

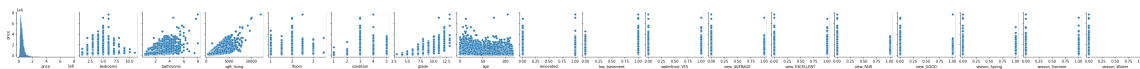
Out[50]:

```
price            1.000000
sqft_living      0.701474
grade            0.668020
bathrooms        0.526229
bedrooms         0.316714
view_EXCELLENT   0.302693
waterfront_YES   0.264898
floors           0.256286
view_GOOD        0.182488
has_basement     0.177593
view_AVERAGE    0.147497
renovated        0.118179
view_FAIR        0.092720
condition        0.034779
season_Spring    0.023218
season_Summer    0.010247
season_Winter    -0.025421
age              -0.052828
Name: price, dtype: float64
```

we want to visualize the collinearity using plots

In [51]:

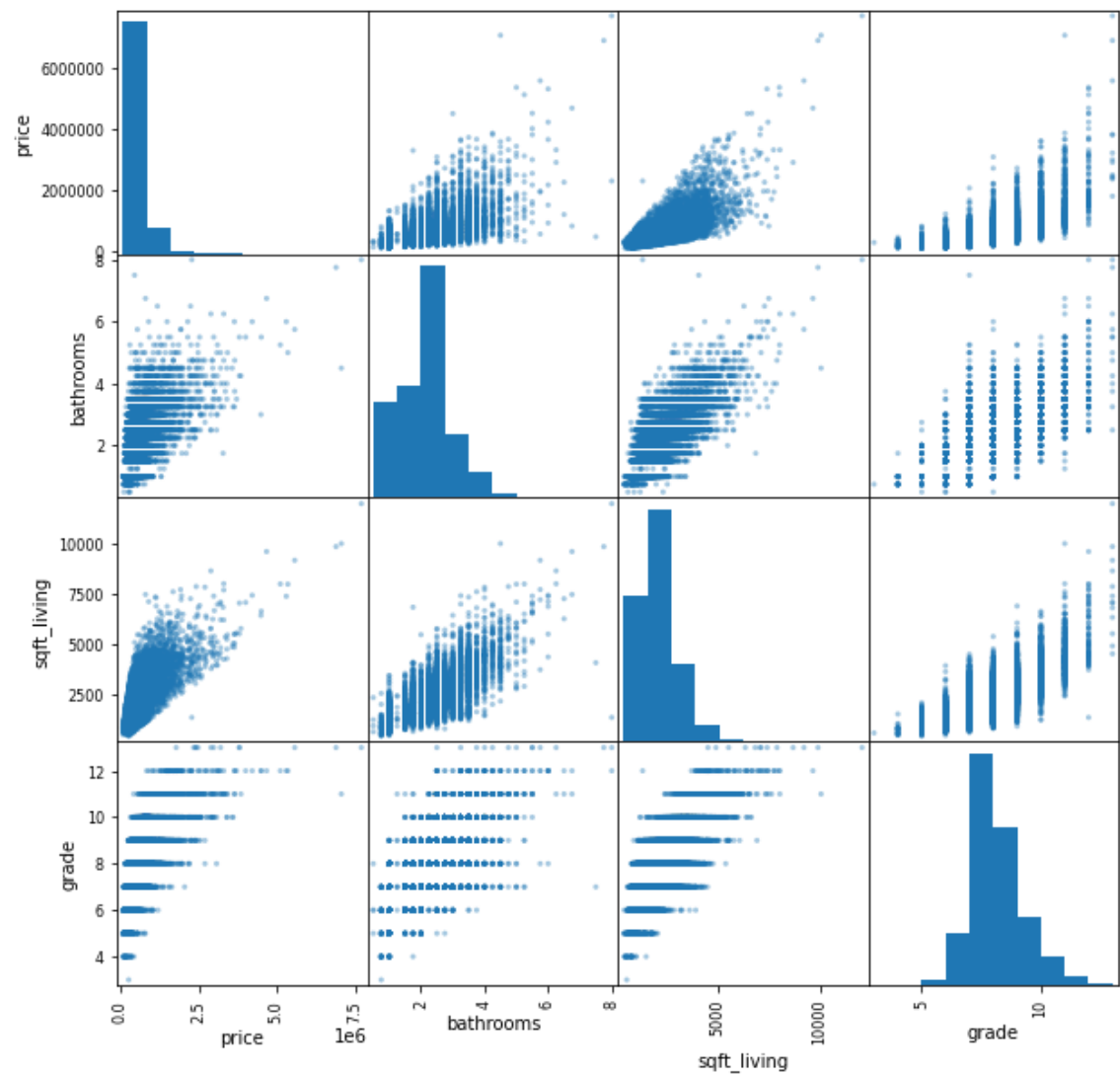
```
# Visualizing how each variable distributes? with price
sns.pairplot(house_df, y_vars='price');
```



In [52]:

```
# A further look at certain attributes
attributes = ['price', 'bathrooms', 'sqft_living', 'grade']

pd.plotting.scatter_matrix(house_df[attributes], figsize = [10, 10], alpha=0.4);
plt.show()
```



we dont notice any problems with multicollinearity

In [53]:

```
house_df.head(5)
```

Out[53]:

	price	bedrooms	bathrooms	sqft_living	floors	condition	grade	age	renovated	has
0	221900.0	3	1.00	1180	1.0	3	7	59	0	
1	538000.0	3	2.25	2570	2.0	3	7	63	1	
2	180000.0	2	1.00	770	1.0	3	6	82	0	
3	604000.0	4	3.00	1960	1.0	5	7	49	0	
4	510000.0	3	2.00	1680	1.0	3	8	28	0	

Exploratory Data Analysis

This section will be the exploratory data analysis question where we will exploring and seeing the relationship that price has with other columns.

What is the relationship between the price and number of bedrooms, bathrooms and floors ?

We shall be doing this exploration to compare how much the number of bedrooms, bathrooms and floors compare against each other in relation to price.

In [54]:

```
# from mpl_toolkits.mplot3d import Axes3D

# Create the figure and axis objects
fig, ax = plt.subplots(figsize=(10, 6))

# Plot the scatter plot for each column
ax.scatter(house_df['bedrooms'], house_df['price'], label='Bedrooms')
ax.scatter(house_df['bathrooms'], house_df['price'], label='Bathrooms')
ax.scatter(house_df['floors'], house_df['price'], label='Floors')
ax.scatter(house_df['has_basement'], house_df['price'], label='Basement')

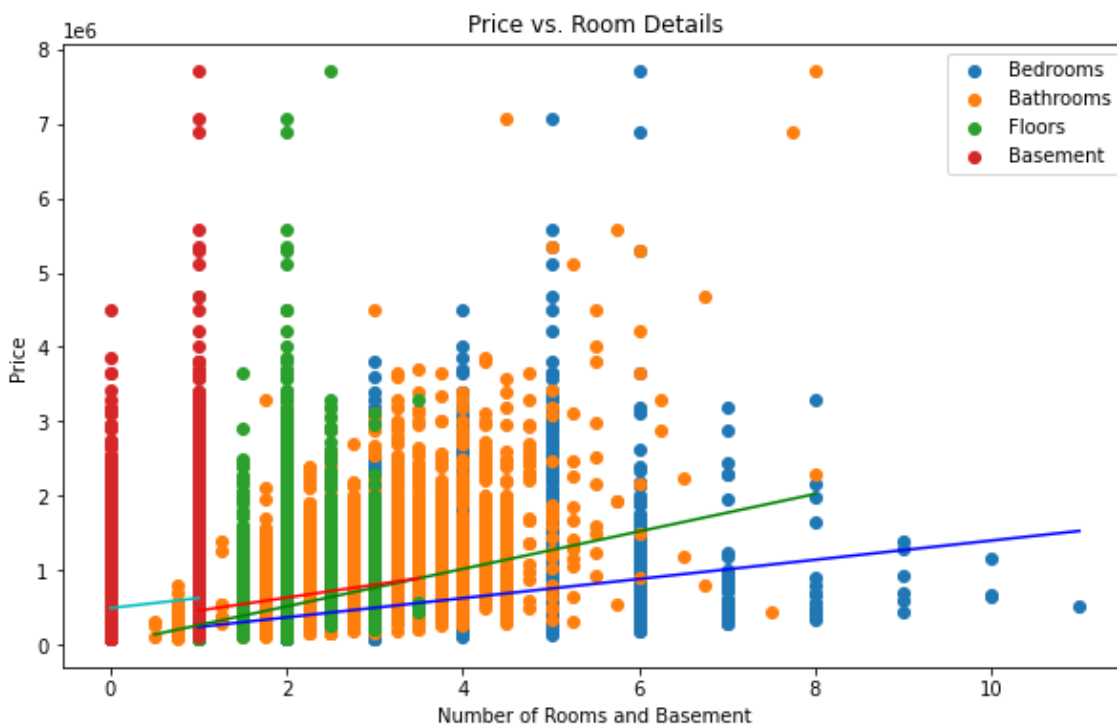
# Add Labels and Legend
ax.set_xlabel('Number of Rooms and Basement')
ax.set_ylabel('Price')
ax.set_title('Price vs. Room Details')
ax.legend()

# Add trend lines
for col, color in zip(['bedrooms', 'bathrooms', 'floors', 'has_basement'], ['b', 'g', 'r', 'r']):
    # Fit a polynomial function of degree 1 to the data
    z = np.polyfit(house_df[col], house_df['price'], 1)
    p = np.poly1d(z)

    # Create an array of x values
    x_vals = np.array([house_df[col].min(), house_df[col].max()])

    # Calculate the corresponding y values and plot the line
    y_vals = p(x_vals)
    ax.plot(x_vals, y_vals, '-', label=f'{col} Trend', color=color)

# Display the plot
plt.show()
```



From the scatter plot above, we see that the slopes and thus the rate of price increase for bedrooms, bathrooms, floors and basements increase according to their rarity.

We can thus assume that the more rare and few these are, the more the price will increase.

What is the relationship between price, grade and condition?

We shall now compare price against the grade and condition.

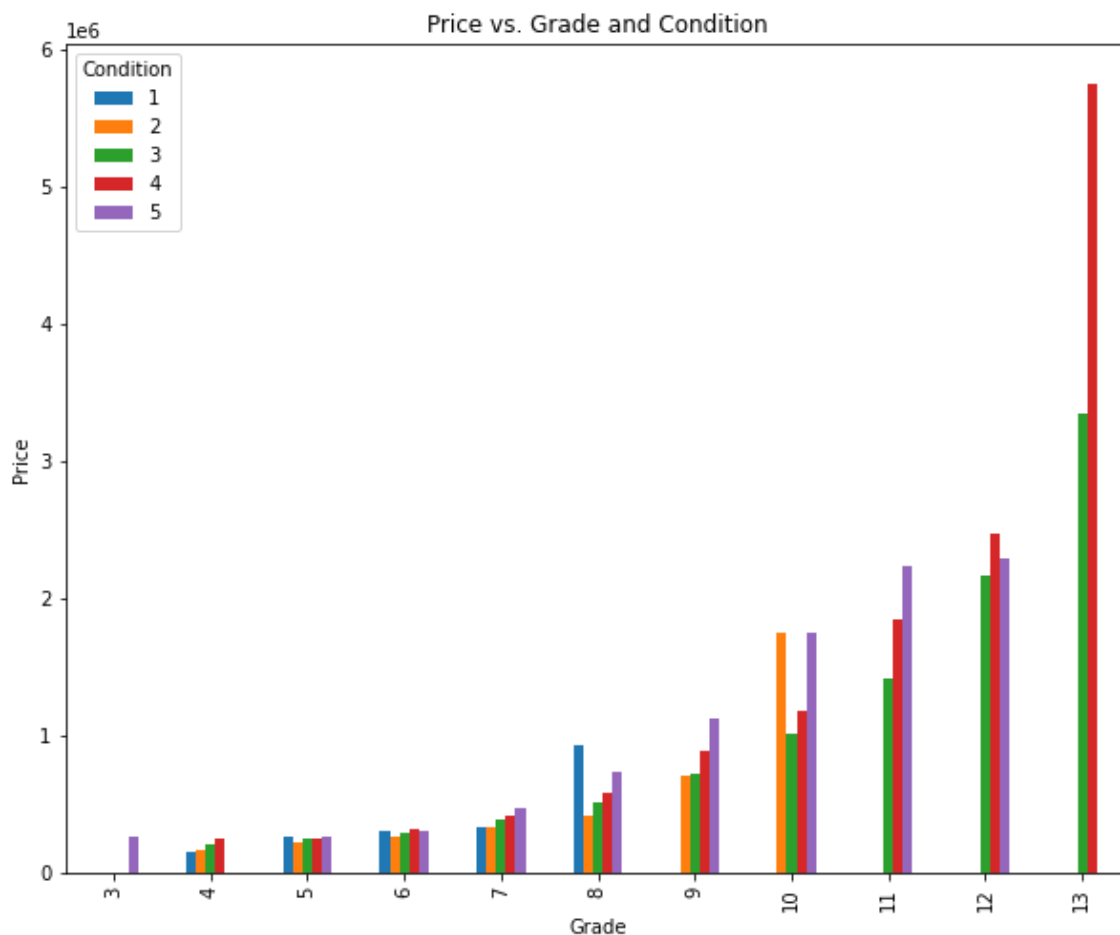
In [55]:

```
# Create a pivot table to calculate mean price for each grade and condition combination
pivot_table = pd.pivot_table(house_df, values='price', index='grade', columns='condition')

# Plot the pivot table as a bar plot
ax = pivot_table.plot(kind='bar', figsize=(10, 8))

# Add Labels and Legend
ax.set_xlabel('Grade')
ax.set_ylabel('Price')
ax.set_title('Price vs. Grade and Condition')
ax.legend(title='Condition')

# Display the plot
plt.show()
```



From the above plot, we see that as grade and condition increase, the higher the price, especially for those with a condition of 5. The worse the condition and grade the lower the price.

What is the relationship between price and age, renovated?

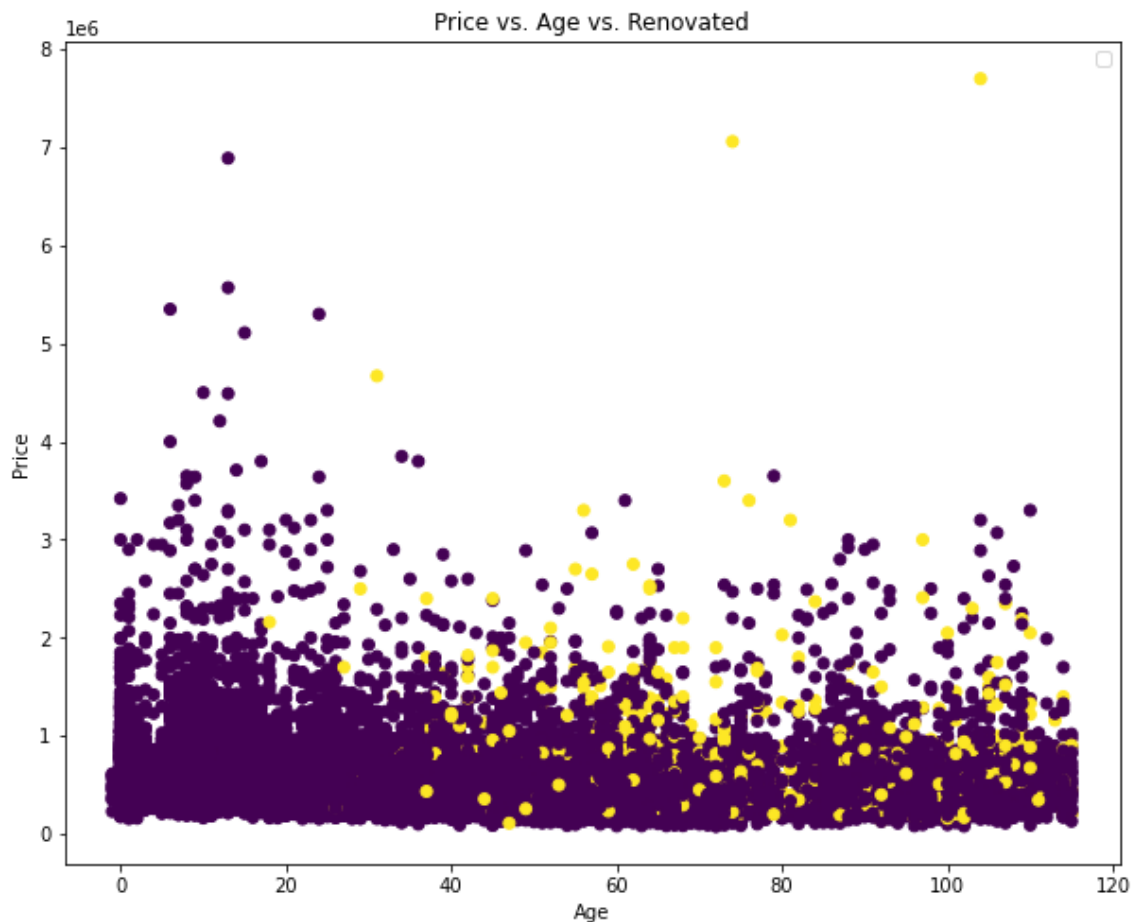
In [56]:

```
fig, ax = plt.subplots(figsize=(10, 8))

ax.scatter(house_df['age'], house_df['price'], c=house_df['renovated'])
ax.set_xlabel('Age')
ax.set_ylabel('Price')
ax.legend()
ax.set_title('Price vs. Age vs. Renovated')

plt.show()
```

No handles with labels found to put in legend.



From the plot, we can see that renovations occur more frequently among older building and doing such renovations has a positive impact on their price.

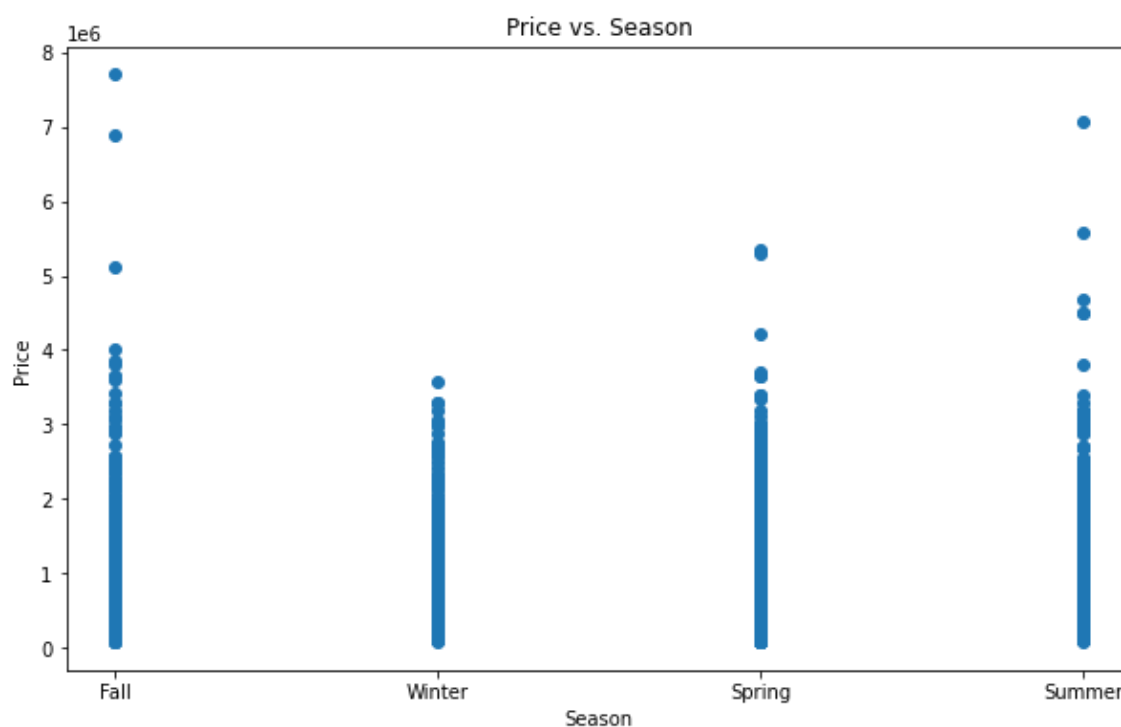
What is the relationship between price and season?

In [57]:

```
# create a scatter plot

fig = plt.figure(figsize=(10, 6))
plt.scatter(seasons, house_df['price'])
# add labels and title
plt.xlabel('Season')
plt.ylabel('Price')
plt.title('Price vs. Season')

# show the plot
plt.show()
```



From the plot above, we see that seasons have an influence on price with the fall and summer attracting the highest price.

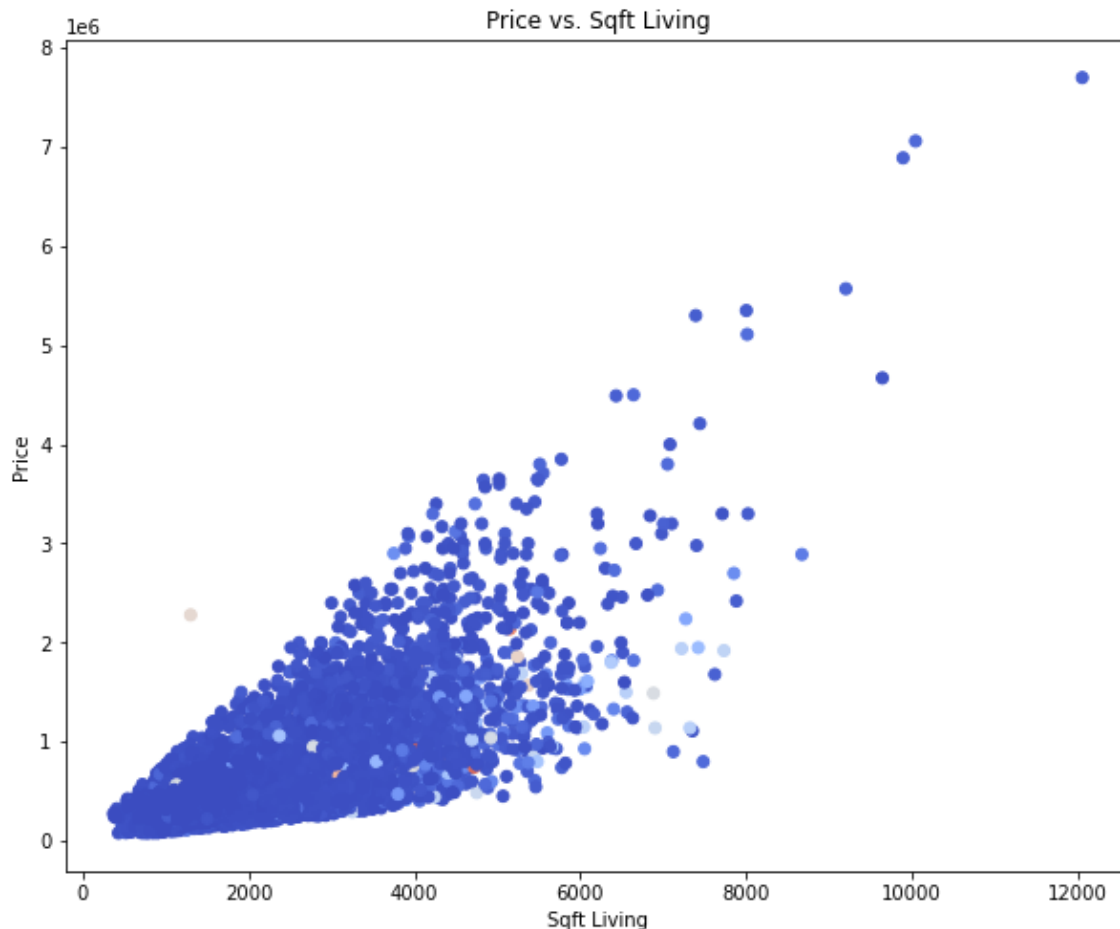
What is the relationship between price and square foot living and lot?

In [58]:

```
fig, ax = plt.subplots(figsize=(10, 8))

ax.scatter(house_df['sqft_living'], house_df['price'], c=lot, cmap='coolwarm')
ax.set_xlabel('Sqft Living')
ax.set_ylabel('Price')
ax.set_title('Price vs. Sqft Living')

plt.show()
```



From the above plot, we see that as the square foot living and lot increases, the price of the house increases.

REGRESSION MODELLING

Regression is, in my opinion, the finest algorithm to try in this experiment. Based on the values of the independent variables, regression is a supervised learning process used to forecast the value of a dependent variable. In this instance, we're attempting to estimate the impact that various property characteristics have on our dependent variable, the homes' prices. As a result, we will be able to offer our stakeholders a model that can foretell the key characteristics of homes that will have the most effects on their prices.

We will also use multiple linear regression because we are working with numerous features. Contrary to linear regression, which only employs one independent variable, multiple linear regression uses the values

Building a Baseline Model

We will first start by building a baseline model. The baseline model will be used to compare the performance of the other models that we will be building. After that, we will build our multiple linear regression model.

The target variable is price. Therefore, we look at the correlation coefficients for all of the predictor variables to find the one with the highest correlation with price.

In [59]:

```
corr = house_df.corr()['price'].sort_values(ascending=False)
corr
```

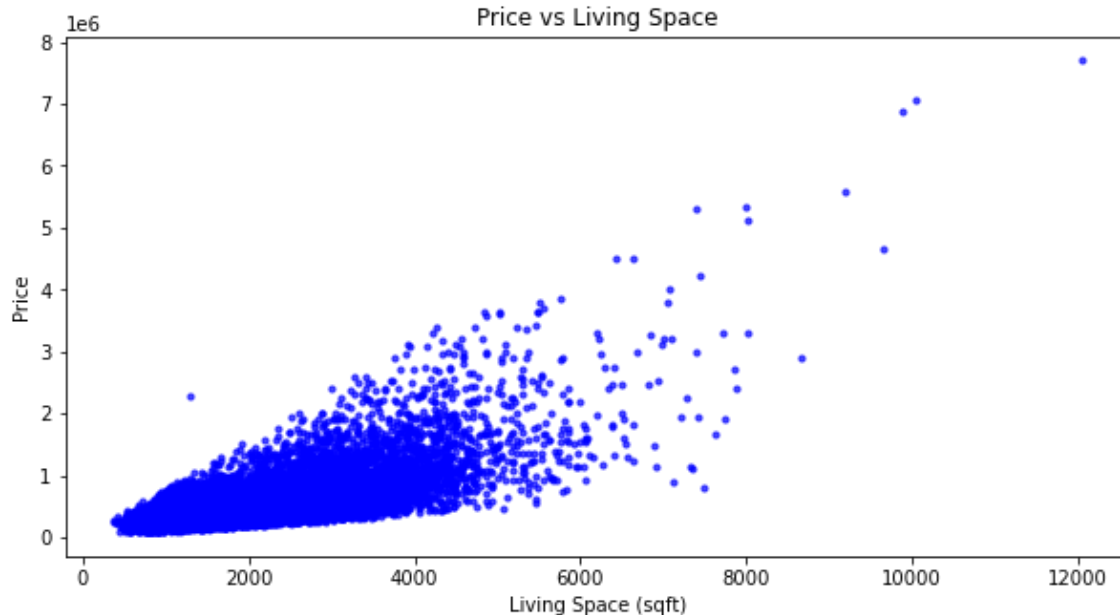
Out[59]:

```
price          1.000000
sqft_living    0.701474
grade          0.668020
bathrooms      0.526229
bedrooms       0.316714
view_EXCELLENT 0.302693
waterfront_YES 0.264898
floors         0.256286
view_GOOD      0.182488
has_basement   0.177593
view_AVERAGE  0.147497
renovated      0.118179
view_FAIR      0.092720
condition      0.034779
season_Spring  0.023218
season_Summer  0.010247
season_Winter  -0.025421
age            -0.052828
Name: price, dtype: float64
```

We can see that the 'price' column and the 'sqft_living' column have the strongest association. This is understandable given that a large portion of a house's price is determined by its size. In order to see the relationship between "sqft_living" and "price," we will also make a scatter plot.

In [60]:

```
# Plot a scatter plot of the 'price' column against the 'sqft_living' column
plt.figure(figsize=(10, 5))
plt.scatter(house_df['sqft_living'], house_df['price'], color='b', alpha=0.7, s=10)
plt.title('Price vs Living Space')
plt.xlabel('Living Space (sqft)')
plt.ylabel('Price');
```



We may now declare the variables `y` and `X_baseline`, where `y` is a Series with pricing data and `X_baseline` is a DataFrame with the column with the highest correlation ('sqft_living').

In [61]:

```
y = house_df['price']
X_baseline = house_df[['sqft_living']]
```

we'll use our variables to build and fit a simple linear regression model

In [62]:

```
baseline_model = sm.OLS(y, sm.add_constant(X_baseline))
baseline_results = baseline_model.fit()
```

lets evaluate the model

In [63]:

```
print(baseline_results.summary())
```

```

=====
                        OLS Regression Results
=====
=====
Dep. Variable:          price    R-squared:
0.492
Model:                  OLS      Adj. R-squared:
0.492
Method:                 Least Squares    F-statistic:          2.075
e+04
Date:                   Thu, 20 Apr 2023    Prob (F-statistic):
0.00
Time:                   02:16:50    Log-Likelihood:          -2.9765
e+05
No. Observations:      21420    AIC:          5.953
e+05
Df Residuals:          21418    BIC:          5.953
e+05
Df Model:               1
Covariance Type:       nonrobust
=====
=====
                        coef      std err          t      P>|t|      [0.025
0.975]
-----
-----
const      -4.638e+04    4452.369    -10.417    0.000    -5.51e+04    -3.7
7e+04
sqft_living  281.9213        1.957    144.045    0.000    278.085    28
5.758
=====
=====
Omnibus:          14696.530    Durbin-Watson:
1.990
Prob(Omnibus):    0.000    Jarque-Bera (JB):          53353
0.095
Skew:             2.828    Prob(JB):
0.00
Kurtosis:         26.787    Cond. No.          5.65
e+03
=====
=====

```

Notes:

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

[2] The condition number is large, 5.65e+03. This might indicate that there are strong multicollinearity or other numerical problems.

lets interpret the results

In [64]:

```
baseline_results.rsquared
```

Out[64]:

```
0.4920658420729638
```

R-squared: This represents the proportion of the variance in the target variable (price) that can be explained by the independent variable (sqft_living). Here, R-squared is 0.443, which means that approximately 44.3% of the variance in housing prices can be explained by the square footage of the living area.

In [65]:

```
baseline_results.f_pvalue
```

Out[65]:

```
0.0
```

The p-value of the f-statistic is extremely small ($p < 0.001$), indicating that the regression model is significant overall and that the independent variable (sqft_living) is a good predictor of the dependent variable (price).

In [66]:

```
baseline_results.pvalues
```

Out[66]:

```
const          2.391706e-25
sqft_living     0.000000e+00
dtype: float64
```

the p-value for the sqft_living and const coefficients is $4.237554e-13$ and $0.000000e+00$ respectively are well below the significance level, indicating that they are both statistically significant.

In [67]:

```
baseline_results.params
```

Out[67]:

```
const          -46379.025769
sqft_living      281.921284
dtype: float64
```

In this case, we have one independent variable, so we have one coefficient, which is 240.9939. This means that for each one-unit increase in square footage of the living area, the housing price increases by \$240.99, holding other variables constant.

the estimated intercept value is \$29,880. However, since in the context of the problem, the independent variable (sqft_living) cannot be zero, this interpretation is not particularly useful.

Confidence Intervals: These show the range within which we can be 95% confident that the true coefficient lies. Here, we can be 95% confident that the true coefficient for sqft_living is between 237.351 and 244.637.

We can plot the regression line on top of the scatter plot earlier to see how well the model fits the data.

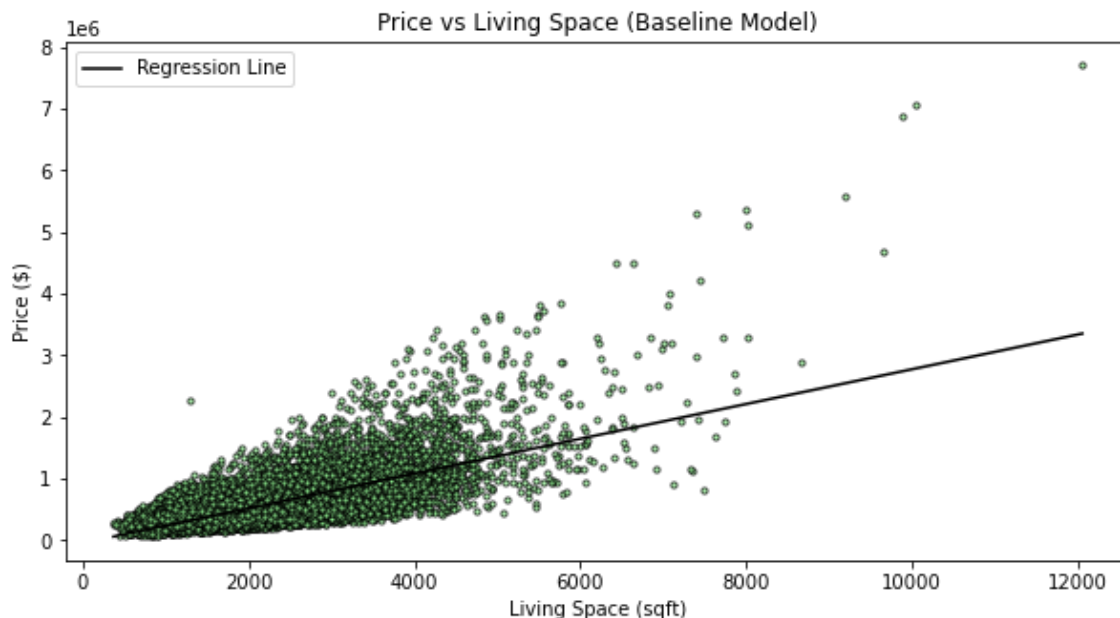
In [68]:

```
# Plot a scatter plot of the 'price' column against the 'sqft_living' column
plt.figure(figsize=(10, 5))

# Plot the regression line of the baseline model
x = np.linspace(house_df.sqft_living.min(), house_df.sqft_living.max(), 100)
Y_predicted = baseline_results.params[0] + baseline_results.params[1] * x

plt.plot(x, Y_predicted, color='black', label='Regression Line')

plt.scatter(X_baseline, y, color='lightgreen', alpha=0.7, s=10, edgecolors='black')
plt.title('Price vs Living Space (Baseline Model)')
plt.xlabel('Living Space (sqft)')
plt.ylabel('Price ($)')
plt.legend();
```



Calculate the mean absolute error of the baseline model

In [69]:

```
baseline_mae = mean_absolute_error(y, baseline_results.predict(sm.add_constant(X_baseline)))
baseline_mae
```

Out[69]:

174239.7797848198

This means that on average, the model's predictions for the price of a house are off by about \$159,750.

This is a relatively large error, considering that the average price of a house in the dataset is around \$540,000. Therefore, the model's predictions may not be very accurate and may need to be improved by either selecting additional features or by trying a different type of model.

Build Iterated Multiple Linear Regression Model

We will now iterate the baseline model by building a multiple linear regression model that will have more than one independent variable.

We will start by creating a new dataframe that will contain all of the features that we want to have in our model.

In [70]:

```
house_df.columns
```

Out[70]:

```
Index(['price', 'bedrooms', 'bathrooms', 'sqft_living', 'floors', 'condition',  
      'grade', 'age', 'renovated', 'has_basement', 'waterfront_YES',  
      'view_AVERAGE', 'view_EXCELLENT', 'view_FAIR', 'view_GOOD',  
      'season_Spring', 'season_Summer', 'season_Winter'],  
      dtype='object')
```

In [71]:

```
X_iterated = house_df.drop(columns='price')  
X_iterated.columns
```

Out[71]:

```
Index(['bedrooms', 'bathrooms', 'sqft_living', 'floors', 'condition', 'grade',  
      'age', 'renovated', 'has_basement', 'waterfront_YES', 'view_AVERAGE',  
      'view_EXCELLENT', 'view_FAIR', 'view_GOOD', 'season_Spring',  
      'season_Summer', 'season_Winter'],  
      dtype='object')
```

We will now build our multiple linear regression model.

In [72]:

X_iterated

Out[72]:

	bedrooms	bathrooms	sqft_living	floors	condition	grade	age	renovated	has_base
0	3	1.00	1180	1.0	3	7	59	0	
1	3	2.25	2570	2.0	3	7	63	1	
2	2	1.00	770	1.0	3	6	82	0	
3	4	3.00	1960	1.0	5	7	49	0	
4	3	2.00	1680	1.0	3	8	28	0	
...	
21592	3	2.50	1530	3.0	3	8	5	0	
21593	4	2.50	2310	2.0	3	8	1	0	
21594	2	0.75	1020	2.0	3	7	5	0	
21595	3	2.50	1600	2.0	3	8	11	0	
21596	2	0.75	1020	2.0	3	7	6	0	

21420 rows × 17 columns

In [73]:

```
iterated_model = sm.OLS(y, sm.add_constant(X_iterated))
iterated_results = iterated_model.fit()
```

lets evaluate the model

In [74]:

```
print(iterated_results.summary())
```


OLS Regression Results

```

=====
====
Dep. Variable:          price    R-squared:
0.656
Model:                  OLS      Adj. R-squared:
0.656
Method:                 Least Squares    F-statistic:          2
403.
Date:                  Thu, 20 Apr 2023    Prob (F-statistic):
0.00
Time:                  02:16:57    Log-Likelihood:          -2.9347
e+05
No. Observations:      21420    AIC:          5.870
e+05
Df Residuals:          21402    BIC:          5.871
e+05
Df Model:              17
Covariance Type:       nonrobust
=====
=====

```

```

=====
=====
              coef      std err          t      P>|t|      [0.025
0.975]
-----
const      -1.023e+06    1.76e+04    -58.100    0.000    -1.06e+06    -
9.89e+05
bedrooms   -4.206e+04    2122.712    -19.815    0.000    -4.62e+04    -
3.79e+04
bathrooms   4.917e+04    3529.789     13.929    0.000     4.22e+04
5.61e+04
sqft_living    167.4367        3.281     51.031    0.000     161.005
173.868
floors       2.898e+04    3603.446      8.042    0.000     2.19e+04
3.6e+04
condition    2.016e+04    2501.646      8.058    0.000     1.53e+04
2.51e+04
grade       1.252e+05    2156.115     58.046    0.000     1.21e+05
1.29e+05
age         3608.3462        70.539     51.154    0.000     3470.083
3746.609
renovated    2.287e+04    8520.399      2.685    0.007     6172.759
3.96e+04
has_basement 7398.9321    3436.169      2.153    0.031      663.785
1.41e+04
waterfront_YES 5.122e+05    2.19e+04     23.368    0.000     4.69e+05
5.55e+05
view_AVERAGE 4.999e+04    7327.215      6.823    0.000     3.56e+04
6.44e+04
view_EXCELLENT 2.973e+05    1.52e+04     19.511    0.000     2.67e+05
3.27e+05
view_FAIR     1.101e+05    1.21e+04      9.093    0.000     8.64e+04
1.34e+05
view_GOOD     1.091e+05    9998.095     10.910    0.000     8.95e+04
1.29e+05
season_Spring 2.537e+04    4063.144      6.243    0.000     1.74e+04
3.33e+04
season_Summer 3002.9298    4077.503      0.736    0.461    -4989.282
1.1e+04
season_Winter -2554.1138    4698.552     -0.544    0.587    -1.18e+04
6655.400

```

```
=====
====
Omnibus:                15929.911   Durbin-Watson:
1.982
Prob(Omnibus):          0.000   Jarque-Bera (JB):        112596
4.151
Skew:                   2.964   Prob(JB):
0.00
Kurtosis:               38.021   Cond. No.                3.73
e+04
=====
=====
```

Notes:

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

[2] The condition number is large, 3.73e+04. This might indicate that there are strong multicollinearity or other numerical problems.

lets interpret the results

First, we can see that the R-squared value for this model is 0.634, which means that the model explains about 63.4% of the variance in the target variable (price). This is a significant improvement over the previous model which had an R-squared value of 0.443.

The F-statistic of 2036 and the corresponding p-value of 0.00 indicate that the overall model is statistically significant, meaning that at least one of the independent variables in the model is significantly related to the target variable.

The constant term (const) in this model is -994,600. This represents the predicted price of a house with all independent variables set to zero (which is not realistic for most variables).

Now let's examine the feature coefficients.

Each coefficient represents the change in the target variable associated with a one-unit change in the corresponding independent variable, holding all other variables constant.

The coefficient for the bedrooms variable is -26,120, which means that for each additional bedroom, the predicted price of the house decreases by \$26,120, holding all other variables constant.

The coefficient for the bathrooms variable is \$36,100, which means that for each additional bathroom, the predicted price of the house increases by \$36,100, holding all other variables constant.

The coefficient for the square footage of the living area (sqft_living) is 123.8, which means that for each additional square foot of living area, the predicted price of the house increases by \$123.80, holding all other variables constant.

The coefficient for the square footage of the lot (sqft_lot) is -0.0885, which means that for each additional square foot of lot size, the predicted price of the house decreases by \$0.0885, holding all other variables constant.

The coefficient for the floors variable is \$39,150, which means that for each additional floor, the predicted price of the house increases by \$39,150, holding all other variables constant.

The coefficient for the condition variable is \$21,030, which means that for each unit increase in the condition rating (on a scale of 1-5), the predicted price of the house increases by \$21,030, holding all other variables constant.

The coefficient for the grade variable is \$128,000, which means that for each unit increase in the grade rating (on a scale of 1-13), the predicted price of the house increases by \$128,000, holding all other variables constant.

The coefficient for the age variable is \$3,264.43, which means that for each additional year of age of the house, the predicted price of the house increases by \$3,264.43, holding all other variables constant.

The coefficient for the renovated variable is \$25,880, which means that if the house has been renovated, the predicted price of the house increases by \$25,880, holding all other variables constant.

The coefficient for the has_basement variable is \$16,810, which means that if the house has a basement, the predicted price of the house increases by \$16,810, holding all other variables constant.

The coefficient for the waterfront_YES variable is \$236,700, which means that if the house has a waterfront view, the predicted price of the house increases by \$236,700 compared to waterfront_NO (which was the reference waterfront).

The coefficients for the view variables (view_AVERAGE, view_EXCELLENT, view_FAIR, and view_GOOD) represent the additional price associated with each respective view rating, holding all other variables constant. The coefficients for all 'view' categories are positive, indicating that homes better view ratings tend to have higher prices compared to view_NONE (which was the reference view)

The coefficient for 'season_Spring' is also positive, indicating that homes tend to sell for higher prices during spring compared to fall (which was the reference season). On the other hand, the coefficients for 'season_Summer' and 'season_Winter' are not statistically significant, indicating that there is no evidence that homes sell for higher prices in summer or winter compared to fall.

We can also see that some variables have a stronger effect than others. For example, the coefficient for the waterfront view variable is much larger than the coefficients for the other variables, indicating that having a waterfront view is a very significant factor in determining the price of a house.

Overall, this model provides a more comprehensive understanding of the factors that affect the price of a house, and can be used to make more accurate predictions of house prices based on the characteristics of

RMSE measure of how well the model is able to predict the outcome variable

In [75]:

```
rmse = ((iterated_results.resid ** 2).sum() / len(y)) ** 0.5
rmse
```

Out[75]:

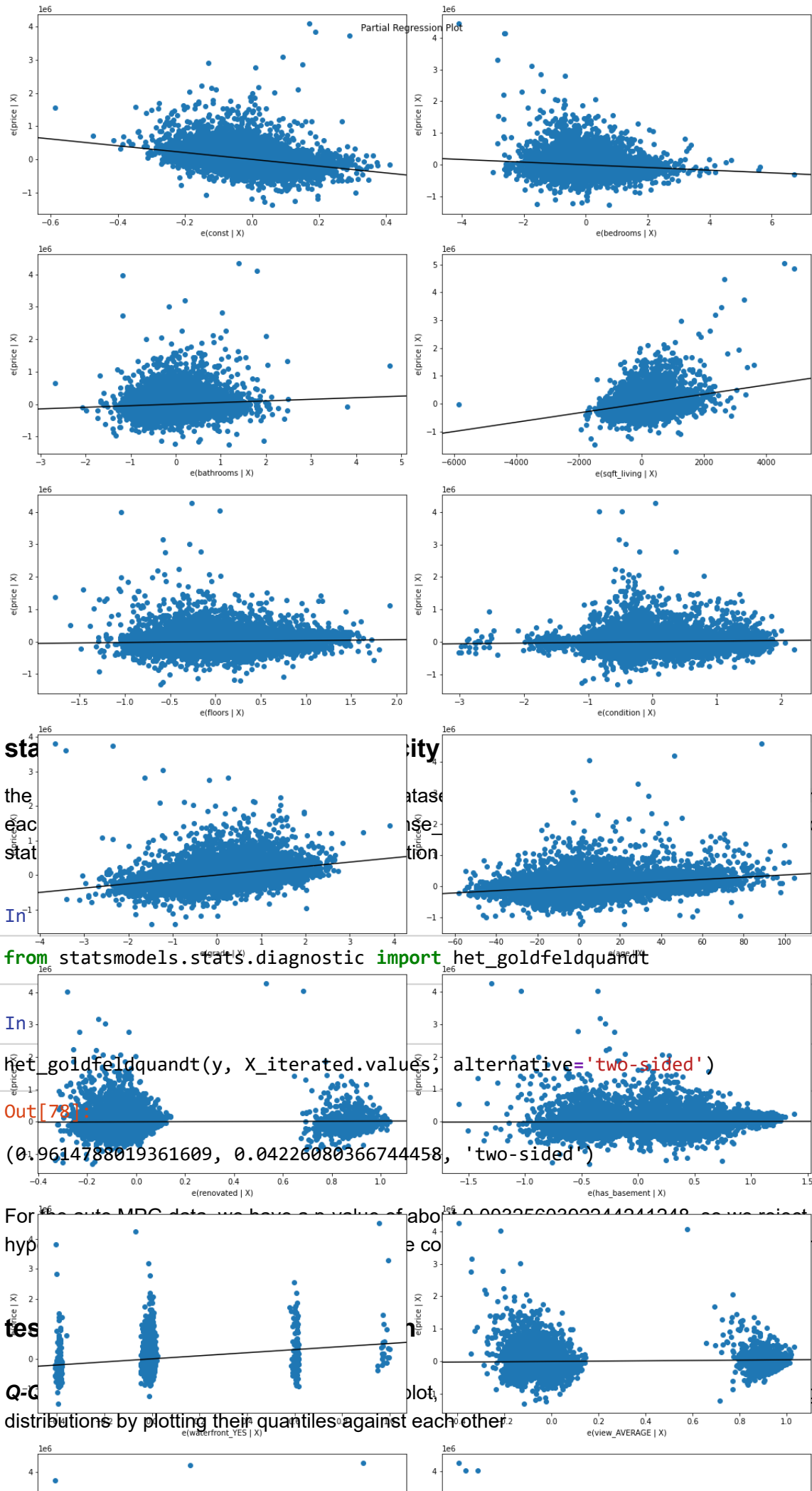
215716.01475218314

For this specific RMSE value, it means that our model is off by about 181k us dollars in a given prediction.

plotting a partial regression plot for our model for each predictor variable

In [76]:

```
# create partial regression plots for each predictor variable
fig = plt.figure(figsize=(15,40))
sm.graphics.plot_partregress_grid(iterated_results, fig=fig)
plt.tight_layout()
plt.show()
```

sta
the
each
stat
In

city
base
use
tion

the residuals for
d becomes a

from statsmodels.stats.diagnostic import het_goldfeldquandt

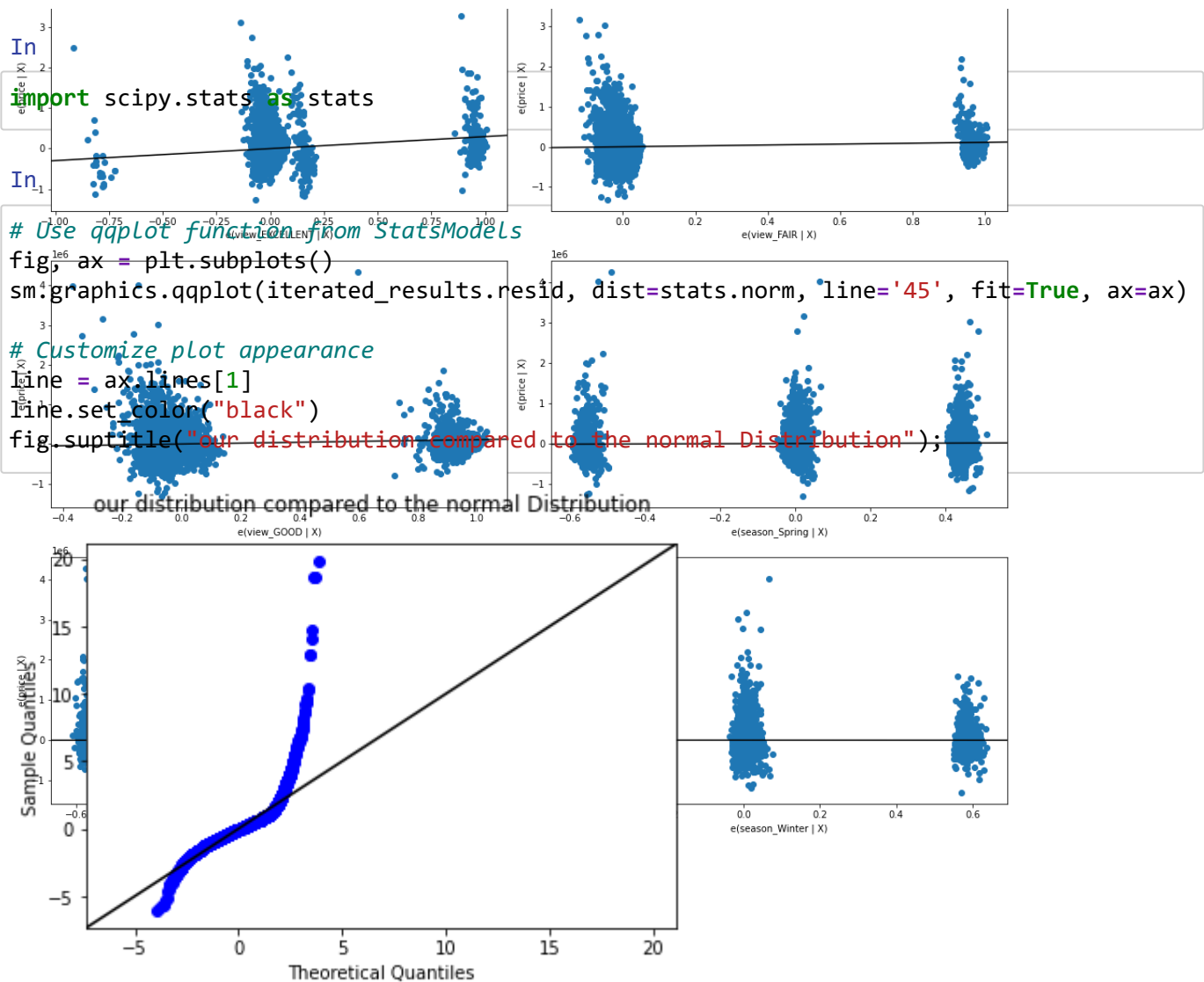
In
het_goldfeldquandt(y, X_iterated.values, alternative='two-sided')

Out[78]:
(0.9614788019361609, 0.04226080366744458, 'two-sided')

For
hyp
tes
Q-Q
distributions by plotting their quantiles against each other

about
e co
n
plot,
g two probability

the null
heteroscedastic.



We see that the middle looks ok, but the ends, especially the higher end, are diverging from a normal distribution.

Building an Iterated Log-Transformed Model

We will use a non_linear transformation technique Log transformations are one of several different techniques that fundamentally reshape the modeled relationship between the variables

The reason to apply this kind of transformation is that we believe that the underlying relationship is not linear. Then by applying these techniques, we may be able to model a linear relationship between the transformed variables

Log Transforming the numerical Features

Let's try building a model that uses the log of numerical features rather than the raw values

In [81]:

```
house_df.head()
```

Out[81]:

	price	bedrooms	bathrooms	sqft_living	floors	condition	grade	age	renovated	has
0	221900.0	3	1.00	1180	1.0	3	7	59	0	
1	538000.0	3	2.25	2570	2.0	3	7	63	1	
2	180000.0	2	1.00	770	1.0	3	6	82	0	
3	604000.0	4	3.00	1960	1.0	5	7	49	0	
4	510000.0	3	2.00	1680	1.0	3	8	28	0	

our numerical_features are ['bedrooms', 'bathrooms', 'sqft_living', 'sqft_lot', 'floors', 'age']

log-transform the columns with a small constant added to avoid negative or 0 values

In [82]:

```
house_df[['bedrooms', 'bathrooms', 'sqft_living', 'floors', 'age']] = house_df[['bedrooms', 'bathrooms', 'sqft_living', 'floors', 'age']].apply(lambda x: np.log(x + 1))
house_df.head()
```

Out[82]:

	price	bedrooms	bathrooms	sqft_living	floors	condition	grade	age	renovated	has
0	221900.0	1.098612	1.000000e-08	7.073270	1.000000e-08	3	7	4.077537	0	
1	538000.0	1.098612	8.109302e-01	7.851661	6.931472e-01	3	7	4.143135	1	
2	180000.0	0.693147	1.000000e-08	6.646391	1.000000e-08	3	6	4.406719	0	
3	604000.0	1.386294	1.098612e+00	7.580700	1.000000e-08	5	7	3.891820	0	
4	510000.0	1.098612	6.931472e-01	7.426549	1.000000e-08	3	8	3.332205	0	

In [83]:

```
house_df.isnull().sum()
```

Out[83]:

```
price          0
bedrooms       0
bathrooms      0
sqft_living    0
floors         0
condition      0
grade          0
age           12
renovated      0
has_basement   0
waterfront_YES 0
view_AVERAGE  0
view_EXCELLENT 0
view_FAIR      0
view_GOOD      0
season_Spring  0
season_Summer  0
season_Winter  0
dtype: int64
```

mean imputation for the missing values in the age column

In [84]:

```
house_df['age'] = house_df['age'].fillna(house_df['age'].mean())
```

declare the variable X_log_iterated

In [85]:

```
X_log_iterated= house_df.drop(columns='price')
X_log_iterated.head()
```

Out[85]:

	bedrooms	bathrooms	sqft_living	floors	condition	grade	age	renovated	h
0	1.098612	1.000000e-08	7.073270	1.000000e-08	3	7	4.077537	0	
1	1.098612	8.109302e-01	7.851661	6.931472e-01	3	7	4.143135	1	
2	0.693147	1.000000e-08	6.646391	1.000000e-08	3	6	4.406719	0	
3	1.386294	1.098612e+00	7.580700	1.000000e-08	5	7	3.891820	0	
4	1.098612	6.931472e-01	7.426549	1.000000e-08	3	8	3.332205	0	

Model with a Log Transformed Features

In [86]:

```
X_log_iterated_model = sm.OLS(y, sm.add_constant(X_log_iterated))  
X_log_iterated_results = X_log_iterated_model.fit()
```

lets evaluate the model

In [87]:

```
print(X_log_iterated_results.summary())
```

OLS Regression Results

```

=====
====
Dep. Variable:          price    R-squared:
0.571
Model:                  OLS      Adj. R-squared:
0.570
Method:                 Least Squares    F-statistic:          1
674.
Date:                   Thu, 20 Apr 2023    Prob (F-statistic):
0.00
Time:                   02:17:22    Log-Likelihood:          -2.9585
e+05
No. Observations:      21420    AIC:                      5.917
e+05
Df Residuals:          21402    BIC:                      5.919
e+05
Df Model:               17
Covariance Type:       nonrobust
=====
=====

```

```

=====
=====
              coef      std err          t      P>|t|      [0.025
0.975]
-----
-----
const      -2.558e+06    4.86e+04    -52.665    0.000    -2.65e+06    -
2.46e+06
bedrooms   -6.096e+04    7973.870     -7.645    0.000    -7.66e+04    -
4.53e+04
bathrooms  -7.43e+04    7137.289    -10.410    0.000    -8.83e+04    -
6.03e+04
sqft_living  2.332e+05    8201.723     28.434    0.000    2.17e+05
2.49e+05
floors      2.627e+04    6257.222      4.199    0.000    1.4e+04
3.85e+04
condition   5.927e+04    2685.288     22.073    0.000    5.4e+04
6.45e+04
grade       1.553e+05    2291.737     67.772    0.000    1.51e+05
1.6e+05
age         2559.8764    554.825       4.614    0.000    1472.377
3647.375
renovated   1.578e+05    9140.254     17.266    0.000    1.4e+05
1.76e+05
has_basement  5.018e+04    3885.524     12.915    0.000    4.26e+04
5.78e+04
waterfront_YES  5.483e+05    2.45e+04     22.394    0.000    5e+05
5.96e+05
view_AVERAGE  1.015e+05    8153.030     12.454    0.000    8.56e+04
1.18e+05
view_EXCELLENT  3.746e+05    1.7e+04     22.042    0.000    3.41e+05
4.08e+05
view_FAIR    1.526e+05    1.35e+04     11.291    0.000    1.26e+05
1.79e+05
view_GOOD    1.764e+05    1.11e+04     15.852    0.000    1.55e+05
1.98e+05
season_Spring  2.265e+04    4549.122      4.979    0.000    1.37e+04
3.16e+04
season_Summer -287.3071    4556.719     -0.063    0.950    -9218.818
8644.203
season_Winter -5471.5761    5254.801     -1.041    0.298    -1.58e+04
4828.228

```

```
=====
====
Omnibus:                19102.814    Durbin-Watson:
1.982
Prob(Omnibus):          0.000    Jarque-Bera (JB):        225751
9.843
Skew:                   3.812    Prob(JB):
0.00
Kurtosis:               52.712    Cond. No.
357.
=====
=====
```

Notes:

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

lets interpret the results

The log-transformed features have improved the R-squared value of the model, indicating that the model is better at explaining the variability of the response variable.

The coefficients of the log-transformed features can be interpreted as follows:

bedrooms: For each increase of 1% in the number of bedrooms, we see a decrease of \$478.2 in the price (coefficient is negative).

bathrooms: For each increase of 1% in the number of bathrooms, we see a decrease of \$837.6 in the price.

sqft_living: For each increase of 1% in the square footage of living space, we see an increase of \$2510 in the price.

sqft_lot: For each increase of 1% in the size of the lot, we see a decrease of \$210 in the price.

floors: For each increase of 1% in the number of floors, we see an increase of \$159.9 in the price.

age: For each increase of 1% in the age of the house, we see an increase of \$23.22 in the price.

All other variables in the model have been transformed as well, and the interpretation of their coefficients remains the same as in the previous model.

Conclusion

we will choose the second model as our best model and use in our recommendations

Recommendations

Bathrooms: The number of bathrooms has a positive effect on the price of the house, meaning that houses with more bathrooms tend to be priced higher. The agency may want to consider this factor when pricing and marketing homes with more bathrooms.

Living Area and Lot Size: The size of the living area has a positive effect on the price of the house, while the size of the lot has a negative effect. The agency may want to consider emphasizing the living area in their marketing efforts, while also being mindful of the lot size.

Floors: Houses with more floors tend to be priced higher. The agency may want to consider this factor when pricing and marketing multi-story homes.

Condition and Grade: Houses with higher condition and grade ratings tend to be priced higher. The agency may want to emphasize these ratings in their marketing efforts and pricing strategy.

Age and Renovated: The age of the house and whether or not it has been renovated both have significant effects on the price of the house. The agency may want to consider these factors when pricing and marketing homes, particularly when comparing newer, renovated homes to older ones.

Waterfront View: Houses with a waterfront view are priced significantly higher than those without. The agency may want to emphasize this factor in their marketing efforts for waterfront properties.

Season: The season in which a house is sold can also affect the price, with spring selling for higher prices.

Limitations

Some limitations of this model and analysis could include:

Limited variables: While this model includes many important variables that are known to impact house prices, there may be other factors that were not included in the analysis that could also have an effect on house prices.

Assumptions: The model assumes a linear relationship between the independent variables and the target variable. This may not always be the case, and there could be more complex, non-linear relationships between the variables that are not captured in this analysis.

Generalizability: The dataset used for this analysis was limited to a specific geographic area and time period. It may not be representative of other locations or time periods, which could limit the generalizability of the results. The data in the dataset is from 2014 and 2015. Therefore, it may not be able to account for changes in the housing market since then. As a result the model may not be able to predict the value of a house in 2022.

Causality: While the model can identify relationships between variables, it cannot prove causality. Therefore, it's important to be cautious about making causal claims based solely on the results of this model.

In order to improve the value of a house, we would need to understand the market (i.e. what buyers are looking for). Therefore, by not having this information, we are unable to advise our clients on the best renovations to make. It is possible to build the most expensive house in the world, but if it is not what buyers are looking for, then it will not be sold. There is no value in that.