

Final Project Submission

Please fill out:

- Student name: Margaret Milley
- Student pace: Full time
- Scheduled project review date/time: 9/7/2023
- Instructor name: Veronica Islaho
- Blog post URL:https://github.com/Chemami12/Phase2Project/blob/main/student_ipynb

BUSINESS UNDERSTANDING

The real estate market is a critical driver of economic growth in many countries globally. It is a dynamic and ever-changing industry where key stakeholders (buyers and sellers) have to make accurate prediction of home prices in order to make informed decisions through use of reliable and comprehensive data.

Real estate is dramatically affected by its location and factors such as employment rates, economy, crime rates, transport facilities and land rates.

The King County House Sale dataset is a valuable source of information for analysing house sales in the Northwestern County. The dataset has several variables that can act as key indicators in price of houses in this area such as Zip code, number of bedrooms, number of bathrooms, square feet, numeric grade and the year the house was built. The dataset will be used for analysis and modelling to better understand the factor that affects house pricing.

BUSINESS PROBLEM A real estate agency that helps homeowners buy and/or sell homes is looking into how certain features could affect the pricing of houses in King County and by what amount so as to make appropriate recommendations/ predictions to their clients.

SPECIFIC OBJECTIVES For the real estate agency to make the appropriate recommendations. The following objectives need to be answered:

1. Determine key factors that impact housing prices in King County.
2. How number of bedrooms, bathrooms, grade, square footage, year built and the zipcode of the house correlate with its sale price.
3. Predict house sale price given certain house specification.

DATA UNDERSTANDING The project will use the King County House Sales dataset

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

import seaborn as sns

import statsmodels.api as sm
import numpy as np
import scipy.stats as stats

In [ ]: data = pd.read_csv('data/kc_hc_house_data.csv')
data
```

Out []:	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view	...	grade	sqft_above	sqft_basement
	0	7129300520	10/13/2014	221900.0	3	1.00	1180	5650	1.0	NaN	NONE	...	7 Average	1180
	1	6414100192	12/9/2014	538000.0	3	2.25	2570	7242	2.0	NO	NONE	...	7 Average	2170
	2	5631500400	2/25/2015	180000.0	2	1.00	770	10000	1.0	NO	NONE	...	6 Low Average	770
	3	2487200875	12/9/2014	604000.0	4	3.00	1960	5000	1.0	NO	NONE	...	7 Average	1050
	4	1954400510	2/18/2015	510000.0	3	2.00	1680	8080	1.0	NO	NONE	...	8 Good	1680

	21592	263000018	5/21/2014	360000.0	3	2.50	1530	1131	3.0	NO	NONE	...	8 Good	1530
	21593	6600601020	2/23/2015	400000.0	4	2.50	2310	5613	2.0	NO	NONE	...	8 Good	2310
	21594	1523300141	6/23/2014	402101.0	2	0.75	1020	1380	2.0	NO	NONE	...	7 Average	1020
	21595	291310100	1/16/2015	400000.0	3	2.50	1600	2388	2.0	NaN	NONE	...	8 Good	1600
	21596	1523300157	10/15/2014	325000.0	2	0.75	1020	1076	2.0	NO	NONE	...	7 Average	1020

21597 rows x 21 columns

```
In [ ]: data.columns
Out [ ]: Index(['id', 'date', 'price', 'bedrooms', 'bathrooms', 'sqft_living', 'sqft_lot', 'floors', 'waterfront', 'view', '...', 'grade', 'sqft_above', 'sqft_basement', 'yr_built', 'zipcode', 'lat', 'long', 'sqft_living15', 'sqft_lot15'],
      dtype='object')
```

```
In [ ]: data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21597 entries, 0 to 21596
Data columns (total 21 columns):
#   Column              Non-Null Count  Dtype
---  --
0   id                   21597 non-null   int64
1   date                 21597 non-null   object
2   price                21597 non-null   float64
3   bedrooms             21597 non-null   int64
4   bathrooms            21597 non-null   float64
5   sqft_living          21597 non-null   float64
6   sqft_lot             21597 non-null   int64
7   floors               21597 non-null   float64
8   waterfront           19223 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: 3.4+ MB
```

DATA PREPARATION This is to select the data that will be used, make necessary conversions from categorical to numerical and perform data cleaning

```
In [ ]: # Display the first few rows of the dataset
pd.DataFrame(data.head())

Out [ ]:
```

	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view	...	grade	sqft_above	sqft_basement
0	7129300520	10/13/2014	221900.0	3	1.00	1180	5650	1.0	NaN	NONE	...	7 Average	1180	0.0
1	6414100192	12/9/2014	538000.0	3	2.25	2570	7242	2.0	NO	NONE	...	7 Average	2170	400.0
2	5631500400	2/25/2015	180000.0	2	1.00	770	10000	1.0	NO	NONE	...	6 Low Average	770	0.0
3	2487200875	12/9/2014	604000.0	4	3.00	1960	5000	1.0	NO	NONE	...	7 Average	1050	910.0
4	1954400510	2/18/2015	510000.0	3	2.00	1680	8080	1.0	NO	NONE	...	8 Good	1680	0.0

5 rows x 21 columns

```
In [ ]: # Check for the data type of each column in the dataframe
data.dtypes

Out [ ]:
```

	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view	...	grade	sqft_above	sqft_basement	yr_built	zipcode	lat	long	sqft_living15	sqft_lot15
id	int64	object	float64	int64	float64	float64	int64	float64	object	object	...	object	int64	object	int64	int64	float64	float64	int64	int64
price	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	...	float64	float64	float64	float64	float64	float64	float64	float64	float64
bedrooms	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	...	int64	int64	int64	int64	int64	int64	int64	int64	int64
bathrooms	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	...	float64	float64	float64	float64	float64	float64	float64	float64	float64
sqft_living	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	...	float64	float64	float64	float64	float64	float64	float64	float64	float64
sqft_lot	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	...	int64	int64	int64	int64	int64	int64	int64	int64	int64
floors	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	...	float64	float64	float64	float64	float64	float64	float64	float64	float64
waterfront	object	object	object	object	object	object	object	object	object	object	...	object	object	object	object	object	object	object	object	object
view	object	object	object	object	object	object	object	object	object	object	...	object	object	object	object	object	object	object	object	object
condition	object	object	object	object	object	object	object	object	object	object	...	object	object	object	object	object	object	object	object	object
grade	object	object	object	object	object	object	object	object	object	object	...	object	object	object	object	object	object	object	object	object
sqft_above	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	...	int64	int64	int64	int64	int64	int64	int64	int64	int64
sqft_basement	object	object	object	object	object	object	object	object	object	object	...	object	object	object	object	object	object	object	object	object
yr_built	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	...	int64	int64	int64	int64	int64	int64	int64	int64	int64
yr_renovated	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	...	float64	float64	float64	float64	float64	float64	float64	float64	float64
zipcode	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	...	int64	int64	int64	int64	int64	int64	int64	int64	int64
lat	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	...	float64	float64	float64	float64	float64	float64	float64	float64	float64
long	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	...	float64	float64	float64	float64	float64	float64	float64	float64	float64
sqft_living15	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	...	int64	int64	int64	int64	int64	int64	int64	int64	int64
sqft_lot15	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	...	int64	int64	int64	int64	int64	int64	int64	int64	int64
dtype:	int64	object	float64	int64	float64	float64	int64	float64	object	object	...	object	int64	object	int64	int64	float64	float64	int64	int64

```
In [ ]: # Check for missing values
data.isnull().sum()

Out [ ]:
```

	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view	...	grade	sqft_above	sqft_basement	yr_built	zipcode	lat	long	sqft_living15	sqft_lot15
id	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
price	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
bedrooms	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
bathrooms	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
sqft_living	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
sqft_lot	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
waterfront	2376	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
view	63	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
condition	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
grade	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
sqft_above	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
sqft_basement	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
yr_built	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
yr_renovated	3842	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
zipcode	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
lat	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
long	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
sqft_living15	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
sqft_lot15	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
dtype:	int64	object	float64	int64	float64	float64	int64	float64	object	object	...	object	int64	object	int64	int64	float64	float64	int64	int64

```
In [ ]: # Drop columns with missing values
print(data.dropna(axis=1, inplace=True))

In [ ]: # Check for duplicated in the row
print(data.duplicated().sum())

In [ ]: data.head()
```

Out []:	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	condition	grade	sqft_above	sqft_basement	yr_built	zi
	0	7129300520	10/13/2014	221900.0	3	1.00	1180	5650	1.0	Average	7 Average	1180	0.0	1955
	1	6414100192	12/9/2014	538000.0	3	2.25	2570	7242	2.0	Average	7 Average	2170	400.0	1951
	2	5631500400	2/25/2015	180000.0	2	1.00	770	10000	1.0	Average	6 Low Average	770	0.0	1933
	3	2487200875	12/9/2014	604000.0	4	3.00	1960	5000	1.0	Very Good	7 Average	1050	910.0	1965
	4	1954400510	2/18/2015	510000.0	3	2.00	1680	8080	1.0	Average	8 Good	1680	0.0	1987

```
In [ ]: # Splitting the grade cell into numerical and descriptive
# Define a function to extract the numeric grade from the grade column
def extract_numeric_grade(s):
    numeric_str = ''
    for char in s:
        if char.isdigit():
            numeric_str += char
    return int(numeric_str)

# Define a function to extract the grade description from the grade column
def extract_grade_desc(s):
    grade_desc = ''
    for char in s:
        if not char.isdigit() and char != '.':
            grade_desc += char
    return grade_desc.strip()

# Extract the numeric grade and grade description using the custom functions
data['numeric_grade'] = data['grade'].apply(lambda x: extract_numeric_grade(x))
data['grade_desc'] = data['grade'].apply(lambda x: extract_grade_desc(x))

data.head()
```

Out []:	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	condition	grade	sqft_above	sqft_basement	yr_built	zi
	0	7129300520	10/13/2014	221900.0	3	1.00	1180	5650	1.0	Average	7 Average	1180	0.0	1955
	1	6414100192	12/9/2014	538000.0	3	2.25	2570	7242	2.0	Average	7 Average	2170	400.0	1951
	2	5631500400	2/25/2015	180000.0	2	1.00	770	10000	1.0	Average	6 Low Average	770	0.0	1933
	3	2487200875	12/9/2014	604000.0	4	3.00	1960	5000	1.0	Very Good	7 Average	1050	910.0	1965
	4	1954400510	2/18/2015	510000.0	3	2.00	1680	8080	1.0	Average	8 Good	1680	0.0	1987

```
In [ ]: # Dropping id, date, condition and
```