

Contents

Active Group Members.....	2
Dataset overview:	2
Types of Data	2
Project Questions.....	3
Data Cleaning:.....	4
Irrelevant Data	4
Duplicate Data.....	6
Structural Errors.....	8
Missing Data.....	9
Data Outliers	10
Descriptive Statistics.....	11
Feature Engineering.....	11
Normality Tests	12
Conclusion on Normality Tests	13
Hypothesis Testing	14
Pattern Recognition	15
PO2, PO3, and PO4	15
PO1: Impact of Location on prices	17
Findings	18
Appendix A: Project Presentations	19
Project Week 1.....	19
Project week 2	25
Final Presentation	35
Appendix B: Tools Used	39

Active Group Members

Farrukh Nadeem

Iris Kontogiorgi

Matheus Pereira da Silva Florencio

Cynthia Anyanwu

Dataset overview:

We started our project by getting familiar with overall picture of our data. The dataset represents houses data with 21 columns and 21597 records.

Types of Data

Variable	Type of Variable	Explanation
Id	Nominal	They are unique identifiers randomly assigned without any order in mind
Date	Ordinal	Each date is either bigger or smaller than the other
Price	Ratio	They cannot be less than zero and have equal spacing.
Bedrooms	Ratio	They cannot be less than zero and have equal spacing.
Bathrooms	Ratio	They cannot be less than zero and have equal spacing.
sqft_living	Ratio	They cannot be less than zero and have equal spacing.
sqft_lot	Ratio	They cannot be less than zero and have equal spacing.
Floors	Ratio	They cannot be less than zero and have equal spacing.
Waterfront	Nominal	They only represent existence of waterfront with no inherent order.
View	Ordinal	They have an order with 0 being lowest and 4 highest
Condition	Ordinal	They have an order with 1 being lowest and 5 highest

Grade	Ordinal	They have an order with 1 being lowest and 13 highest
sqft_above	Ratio	They cannot be less than zero and have equal spacing.
sqft_basement	Ratio	They cannot be less than zero and have equal spacing.
yr_built	Ordinal	Each year is either bigger or smaller than the other
yr_renovated	Ordinal	Each year is either bigger or smaller than the other
Zipcode	Nominal	They are unique identifiers randomly assigned without any order in mind
Lat	Interval	They can be less than zero and have equal spacing.
Long	Interval	They can be less than zero and have equal spacing.
sqft_living15	Ratio	They cannot be less than zero and have equal spacing.
sqft_lot15	Ratio	They cannot be less than zero and have equal spacing.

Project Questions

What are the ideal locations (Zip codes) to invest in considering profit margins on resale in mind?

Does the ratio of bedroom to lot area have any impact on price? If yes, what is the suggested ratio?

Does the ratio of bedroom to floors have any impact on price? If yes, what is the suggested ratio?

Does the ratio of bedroom to bathrooms have any impact on price? If yes, what is the suggested ratio?

Data Cleaning:

Irrelevant Data

We considered three columns: **ID**, **sqft_living15**, and **sqft_lot15** to be irrelevant initially, but kept them for following reasons.

ID column: Though this column seemed irrelevant at first, but on analysis we saw that few houses were listed multiple times on different dates with different prices. This led us to keep this column as we thought it might prove to be useful in future analysis.

Figure 1

	id	date	price	bedrooms	...	lat	long	sqft_living15	sqft_lot15
17588	795000620	9/24/2014	115000.0	3	...	47.5045	-122.33	1070.0	6250.0
17589	795000620	12/15/2014	124000.0	3	...	47.5045	-122.33	1070.0	6250.0
17590	795000620	3/11/2015	157000.0	3	...	47.5045	-122.33	1070.0	6250.0

[3 rows x 21 columns]

sqft_living15 and sqft_lot15 columns: These two columns looked like sqft_living and sqft_lot, but still we kept them as in many cases they represented different values. Here is an example:

Figure 2

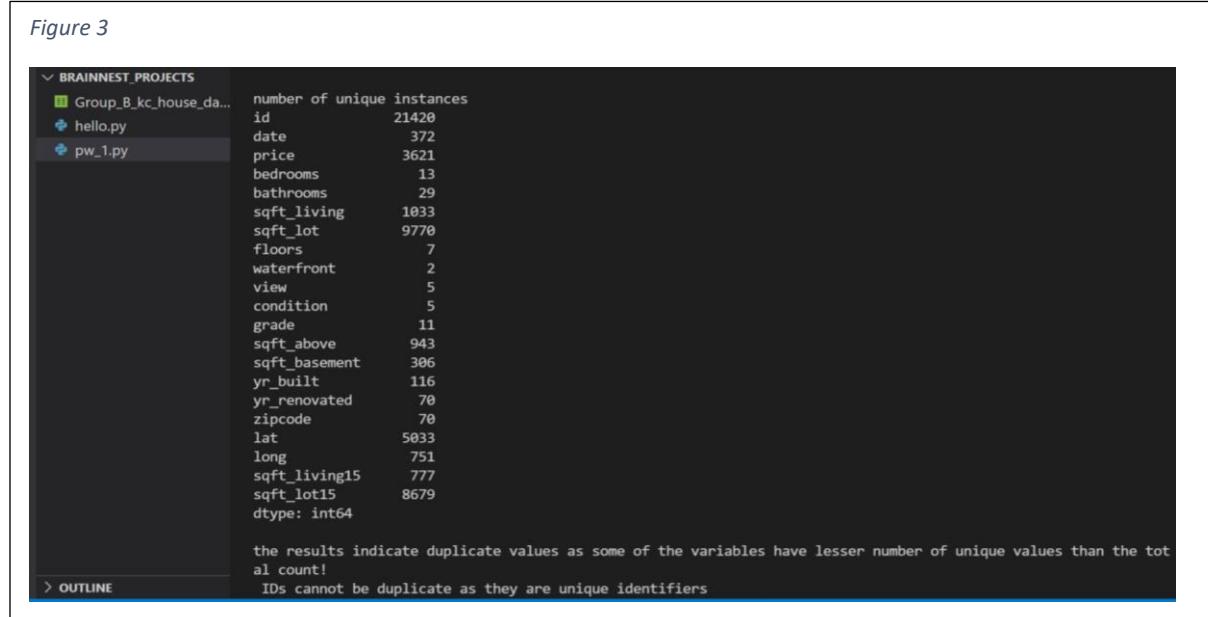
The screenshot shows a Microsoft Excel spreadsheet titled "Group_B_kc_house_data". The ribbon menu is visible at the top, with "Home" selected. The toolbar includes standard functions like Paste, Undo, and Clipboard. The formula bar shows the cell address "7129300520". The main area displays a table with 4 rows and 21 columns. The columns are labeled: id, date, price, bedrooms, bathroom, sqft_living, sqft_lot, floors, waterfront, view, condition, grade, sqft_above, sqft_basement, yr_built, yr_renovated, zipcode, lat, long, sqft_living15, and sqft_lot15. The first four rows of data are:

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
7.129E+09	10/13/2014	221900	3	1	1180	5650	1	0	0	3	7	1180	0	1955	0	98178	47.5112	-122.257	1340	5650
6.414E+09	12/9/2014	538000	3	2.25	2570	7242	2	0	0	3	7	2170	400	1991	98125	47.721	-122.319	1690	7639	
5.632E+09	2/25/2015	180000	2	1	770	10000	1	0	0	3	6	770	0	1933	0	98028	47.7379	-122.233	2720	8062

Duplicate Data

Before getting to completely duplicated rows, we first checked how many unique values each column had.

Below is the image showing that result.



ID column: Figure 3 shows that there are lesser number of IDs than the total rows indicating that there are duplicate IDs. But as indicated already by figure 1, we knew that even though some IDs were duplicated but they still had different values in some other columns.

Latitude and Longitude: In our opinion the combination of latitude and longitude for each house must be unique unless they are apartments in the same building. We saw additional 177 houses with **different IDs** but **duplicated combination of latitude and longitude along** with other similar features. But we still decided to keep them knowing that they might be different houses (as indicated by different IDs) in the same building with similar features.

One such example can be seen below:

Figure 4

```
number of combinations of latitude and longitude with duplicate instances
note True means similar combinations were detected already

False    21420
True     177
dtype: int64
   id      date    price  bedrooms  ...      lat      long  sqft_living15  sqft_lot15
14004  9136103136 3/13/2015 580000.0       2  ...  47.6652 -122.338      1490.0      4013.0
17852  7732418360 8/22/2014 752888.0       3  ...  47.6599 -122.146      2630.0      9000.0
19568  1025049266 9/30/2014 555000.0       2  ...  47.6647 -122.284      1160.0      1327.0
21209  3629700080 1/8/2015 635000.0       3  ...  47.5446 -122.017      2290.0      1407.0
21371  2895800750 4/17/2015 274800.0       3  ...  47.5171 -122.347      1410.0      1899.0
21482  1972200555 7/14/2014      NaN       3  ...  47.6536 -122.354      1570.0      1335.0

[6 rows x 21 columns]
   id      date    price  bedrooms  bathrooms  sqft_living  sqft_lot  floors      lat      long
211  1025049114 7/17/2014 625504.0       3      2.25      1270.0      1566.0      2.0  47.6647 -122.284
13459 1025049115 6/25/2014 594000.0       3      2.25      1270.0      1406.0      2.0  47.6647 -122.284
18297 1025049268 7/21/2014 549900.0       2      1.75      1140.0      936.0      2.0  47.6647 -122.284
19568 1025049266 9/30/2014 555000.0       2      2.25      1160.0      954.0      2.0  47.6647 -122.284
   id      lat      long  waterfront  view  condition  grade  sqft_above  sqft_basement  yr_built
211  47.6647 -122.284       0  0.0      3.0      8.0      1060.0      210.0      2014.0
13459 47.6647 -122.284       0  0.0      3.0      8.0      1060.0      210.0      2014.0
18297 47.6647 -122.284       0  0.0      3.0      8.0      940.0      200.0      2014.0
19568 47.6647 -122.284       0  0.0      3.0      8.0      960.0      200.0      2014.0

PS D:\Brainnest\PW_1\Brainnest_projects> █
```

These analyses still did not rule out the possibility of entire row having exact duplicate, so we checked to see whether that was the case, but we saw no duplication.

Structural Errors

We did not identify any structural errors. However, we found it odd to see Bathrooms and floors to not have discrete values.

Figure 5

```
counts of each value present in column Bathrooms
2.500 5377
1.000 3851
1.750 3047
2.250 2047
2.000 1938
1.500 1445
2.750 1185
3.000 753
3.500 731
3.250 589
3.750 155
4.000 136
4.500 100
4.250 79
0.750 71
4.750 23
5.000 21
5.250 13
5.500 10
1.250 9
6.000 6
0.500 4
5.750 4
6.750 2
8.000 2
6.250 2
6.500 2
```

Figure 6

```
Windows PowerShell
Copyright (C) Microsoft Corporation. All rights reserved.

Try the new cross-platform PowerShell https://aka.ms/pscore6

PS D:\Brainnest\PW_1\Brainnest_projects> & 'C:\Users\Farrukh Nadeem\AppData\Local\Programs\Python\Python37\python.exe' 'c:\Users\Farrukh Nadeem\.vscode\extensions\ms-python.python-2021.8.1159798656\pythonFiles\lib\python\debugpy\launcher' '59686' --- 'd:\Brainnest\PW_1\Brainnest_projects\pw_1_stErr.py'
counts of each value present in column floors
1.000 10672
2.000 8235
1.500 1910
3.000 611
2.500 161
3.500 7
11.000 1
Name: floors, dtype: int64
PS D:\Brainnest\PW_1\Brainnest_projects>
```

We still did not correct them as their high frequency indicates that it cannot be by mistake and might actually mean something important for real estate agents.

Missing Data

We checked to see how many missing values were present in each column, **figure 7** shows the result.

Figure 7

```
File Edit Selection View Go Run Terminal Help pw_1.py - Brainnest_projects - Visual Studio Code  
EXPLORER PROBLEMS OUTPUT TERMINAL DEBUG CONSOLE  
BRAINNEST_PROJECTS  
cleaned_house_data...  
Group_B_kc_house_da...  
hello.py  
pw_1_stErr.py  
pw_1.py  
missing values  
id 0  
date 0  
price 9  
bedrooms 0  
bathrooms 1  
sqft_living 10  
sqft_lot 16  
floors 0  
waterfront 0  
view 3  
condition 3  
grade 5  
sqft_above 12  
sqft_basement 2  
yr_built 6  
yr_renovated 0  
zipcode 7  
lat 1  
long 5  
sqft_living15 7  
sqft_lot15 8  
dtype: int64  
PS D:\Brainnest\pw_1\Brainnest_projects> & 'C:\Users\Farrukh Nadeem\AppData\Local\Programs\Python\Python37\python.exe' 'c:\Users\farrukh nadeem\.vscode\extensions\ms-python.python-2021.8.1159798656\pythonFiles\lib\python\debugpy\launcher' '59548' '--' 'd:\Brainnest\pw_1\Brainnest_projects\pw_1.py'  
> OUTLINE  
Python 3.7.7 64-bit 0 △ 0 ⌂  
Ln 13, Col 38 Spaces: 4 UTF-8 CRLF Python 🔍  
37°C Sunny 🔍 ENG 4:05 PM 🔍
```

The percentage of missing values for any column was very low. And except for the columns: zipcode, lat, and long, all other columns were missing completely at random (MCAR) and hence, were imputed with median.

For columns: zipcode, lat, and long, they were missing at random as the missing data could be explained/derived from the other column. Hence, for zipcodes we entered the values of lat, and long on the following website to find their respective zip code.

<https://www.freemaptools.com/convert-us-zip-code-to-lat-lng.htm>

And since missing lat, and long were very few for us to manually check them and replace them we imputed missing values in these columns to the best of our knowledge.

Data Outliers

During analyses we saw outliers in bedrooms and floors columns.

Figure 8

```
PS D:\Brainnest\PW_1\Brainnest_projects> & 'C:\Users\Farrukh Nadeem\AppData\Local\Programs\Python\Python37\python.exe' 'c:\Users\Farrukh Nadeem\.vscode\extensions\ms-python.python-2021.8.1159798656\pythonFiles\lib\python\debugpy\launcher' '59696' '' '' 'd:\Brainnest\PW_1\Brainnest_projects\pw_1_stErr.py'
counts of each value present in column Bedrooms
3      9824
4      6881
6      2760
7      1601
8      272
9      196
10     38
11     13
12     6
13     3
14     1
33     1
40     1
Name: bedrooms, dtype: int64
PS D:\Brainnest\PW_1\Brainnest_projects>
```

Figure 9

```
PS D:\Brainnest\PW_1\Brainnest_projects> & 'C:\Users\Farrukh Nadeem\AppData\Local\Programs\Python\Python37\python.exe' 'c:\Users\Farrukh Nadeem\.vscode\extensions\ms-python.python-2021.8.1159798656\pythonFiles\lib\python\debugpy\launcher' '53745' '' '' 'd:\Brainnest\PW_1\Brainnest_projects\pw_1_stErr.py'
bedrooms greater than 10?
           price  bedrooms  bathrooms  sqft_living  sqft_lot  floors
8748  5200000.000       11       3.000    3000.000   4960.000   2.000
15856 6400000.000       33      1.750   1620.000   6000.000   1.000
16158 5980000.000       40      2.500   3130.000  409180.000   2.000
floors greater than 3?
           price  bedrooms  bathrooms  sqft_living  sqft_lot  floors
10066  4350000.000       3       3.000    1440.000   1350.000   3.500
11582  5440000.000       3       2.500    1760.000   1755.000   3.500
14871  5250000.000       3       3.000    1730.000   1074.000   3.500
15410  4790000.000       2       2.500    1730.000   1037.000   3.500
16146  3650000.000       2       1.000    1250.000   8100.000  11.000
18462  3300000.000       8       4.000    7710.000   11750.000  3.500
20292  5250000.000       2       2.750    1310.000   1268.000  3.500
20756  5635000.000       3       2.500    1400.000   1312.000  3.500
PS D:\Brainnest\PW_1\Brainnest_projects>
```

We believe having 33 or 40 bedrooms in 2 floors with such low number of bathrooms and small sqft_living must be an error. We also found that odd to have only 2 bedrooms in 11 floors! Hence, we decided to replace these values that were more reasonable.

Descriptive Statistics

Figure 10

debugpy\launcher' '53639' '--' 'd:\Brainnest\PW_1\Brainnest_projects\pw_1_stErr.py'									
count	21588.000	21597.000	21596.000	21587.000	21581.000	21597.000			
mean	540382.754	3.375	2.116	2080.185	15113.467	1.495			
std	367405.169	0.959	0.769	917.968	41511.996	0.544			
min	78000.000	1.000	0.500	370.000	520.000	1.000			
25%	322000.000	3.000	1.750	1430.000	5040.000	1.000			
50%	450000.000	3.000	2.250	1910.000	7617.000	1.500			
75%	645000.000	4.000	2.500	2550.000	10677.000	2.000			
max	7700000.000	40.000	8.000	13540.000	1651359.000	11.000			
		waterfront	view	condition	grade	sqft_above	sqft_basement	yr_built	yr_renovated
count	21597.000	21594.000	21594.000	21592.000	21585.000	21595.000	21591.000	21597.000	
mean	0.234	3.410	7.658	1789.258	291.710	1970.998			84.465
std	0.087	0.766	0.651	1.173	834.846	442.664	29.377		401.821
min	0.000	0.000	1.000	3.000	370.000	0.000	1900.000		0.000
25%	0.000	0.000	3.000	7.000	1190.000	0.000	1951.000		0.000
50%	0.000	0.000	3.000	7.000	1560.000	0.000	1975.000		0.000
75%	0.000	0.000	4.000	8.000	2210.000	560.000	1997.000		0.000
max	1.000	4.000	5.000	13.000	18000.000	4820.000	2615.000		2015.000
		zipcode	lat	long	sqft_living15	sqft_lot15			
count	21590.000	21596.000	21592.000	21590.000	21589.000				
mean	98077.950	47.560	-122.214	1986.586	12761.042				
std	53.513	0.139	0.141	685.296	27279.087				
min	98001.000	47.156	-122.519	399.000	651.000				
25%	98033.000	47.471	-122.328	1490.000	5100.000				
50%	98065.000	47.572	-122.231	1840.000	7620.000				
75%	98118.000	47.678	-122.125	2360.000	10083.000				
max	98199.000	47.778	-121.315	6210.000	871200.000				

Feature Engineering

Based on the Project Questions/Objectives we felt a need to create some columns that could best serve our needs.

Figure 11



From figure 11, on the left most group we have two columns, by diving Variables listed in Column 1 by the variables in Column 2 we obtain results that are listed in the right most group's first column i.e., Calculated Field Name.

We believe that comparing price per square foot for different lots is a more relevant measure than the price itself. Thus, we created for ourselves a column named **pricelot**.

We also created the other three columns: **lotbedroom**, **bedroomfloor**, and **bedroomsbathrooms** to facilitate ourselves in trying to answer project questions/objectives 2, 3 ,and 4 respectively.

Normality Tests

We checked the normality of each variable (Column) using following methods:

Frequency Histogram
comparing mean, median, and mode
Skewness and kurtosis
Z-scores (Skewness and kurtosis)
Q-Q Plots
P-P Plots
Box Plots

Here is an example of the tests applied on one of these columns.

Figure 12

Price

Visual and statistical tests like Histogram, Box Plots, Q-Q plots and P-P plots, showed that the price is not normally distributed with skewness of -4.023(SE= 0.017) and kurtosis of 34.536(SE= 0.033).

Test	Normally Distributed
Histogram	No
QQ Plot	No
PP Plot	No
Box Plot	No
Skewness	No
Z score (Skewness)	No (236.64)
Kurtosis	No
Z score (Kurtosis)	No (1046)
Final conclusion	No (8/8)

For the sake of keeping report neat, the following table shows the results of Normality Tests for only the columns that are required to answer our project questions.

Table 1: Normality Test

Variable	Normality Test	Value
Price	Not Normal	100%
Bedrooms	Not Normal	75%
Bathrooms	Not Normal	75%
Sqft_lot	Not Normal	100%
Floors	Approximately Normal	50%
Zipcode	Categorical Variable. Hence, did not test.	

Conclusion on Normality Tests

For Floors we are 50% sure it is Normally distributed

For other variables we are at least 75% sure they are not normally distributed.

So, we suggest going for non-parametric tests in our future analysis.

Hypothesis Testing

As can be seen in Table 2, for every Project Objective (PO.) we created its hypothesis H₀ (null) and H₁ (alternate). Since our Normality Test results suggested us to use non-parametric tests, and for all of our tests we had one dependent variable and one independent variable with more than two groups Kruskal-Wallis test was our best option.

Here is one example of the application of Kruskal-Wallis Test for one of our Hypothesis:

Figure 13

Nonparametric Tests

Hypothesis Test Summary			
	Null Hypothesis	Test	Sig. ^{a,b}
1	The distribution of pricelot is the same across categories of zipcode.	Independent-Samples Kruskal-Wallis Test	<.001
a. The significance level is .050. b. Asymptotic significance is displayed.			

Independent-Samples Kruskal-Wallis Test

pricelot across zipcode

Independent-Samples Kruskal-Wallis Test Summary

Total N	21597
Test Statistic	12075.477 ^a
Degree Of Freedom	69
Asymptotic Sig.(2-sided test)	<.001

a. The test statistic is adjusted for ties.

Independent-Samples Kruskal-Wallis Test

Like this test, for all our hypothesis the p-value was less than 0.05, therefore, we rejected the null hypothesis in favor of alternate hypothesis.

Following is the summary of our results on non-parametric tests:

Table 2: Hypothesis Test

PO.	H0	H1	Applied Test	Result
1	The cost of land per unit area is same across all neighborhoods (i.e., zip codes).	The cost of land per unit area is affected by the neighborhood.	Kruskal-Wallis Test	Accept H1
2	The mean price is not affected by different ratios of lot area to bedroom.	The mean price is affected by lot area to bedroom ratio.		Accept H1
3	The mean price is not affected by different ratios of bedroom to floors.	The mean price is affected by bedroom to floor ratio.		Accept H1
4	The mean price is not affected by different ratios of bedroom to bathrooms.	The mean price is affected by bedroom to bathroom ratio.		Accept H1

Accepting H1 for all tests indicate that our questions were indeed valid and there is a need for in-depth analysis to detect any meaningful patterns.

Pattern Recognition

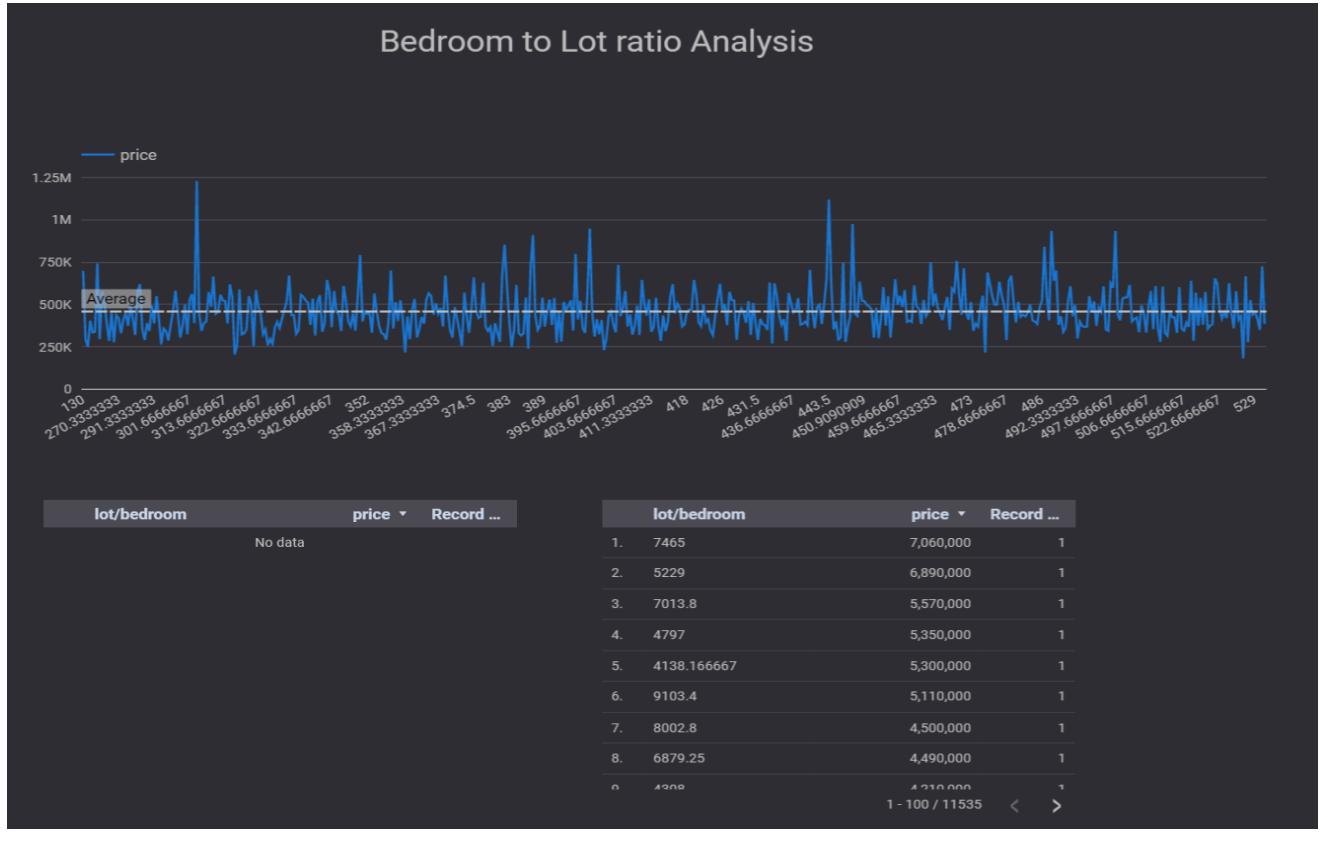
To find patterns between certain variables we first set for ourselves a threshold of at least 2.5% representation for any group in a factor.

PO2, PO3, and PO4

When trying to perform in-depth analysis for PO2, PO3, and PO4 most of the groups failed to meet the threshold. Thus, we decided not to investigate any further as deducing any rules from such small proportion of data would be meaningless and misleading in our opinion.

Below is the example of PO2 where square feet lot to bedroom ratio was to be investigated:

Figure 14



As can be seen from Figure 14, on the bottom left there is an empty table (this is due to the application of filter that restricts groups falling behind the threshold). This can be explained by the line chart on top and table on bottom-right which highlight the fact that there are too many groups within this factor thus rendering it impossible for any group to occupy a larger share of the sample.

Similar scenarios were observed for PO3 and PO4 where bedroom to floors, and bedroom to bathroom ratios were to be investigated respectively. Hence, we decided to end our quest to find any meaningful patterns for these objectives/questions.

PO1: Impact of Location on prices

Unlike the rest of the questions, the number of groups in PO1 were few. This meant that the groups had enough representation from their overall samples, hence, deducing any rules from this analysis was possible.

We conducted in-depth studies using different graphs for price/lot against zipcodes and then selected 7 locations (zipcodes) of interest. We then divided these 7 locations into three groups (G1, G2, and G3) based on their average price per square foot of land, with G1 requiring minimum investment and G3 the highest for a lot of same size.

Figure 15

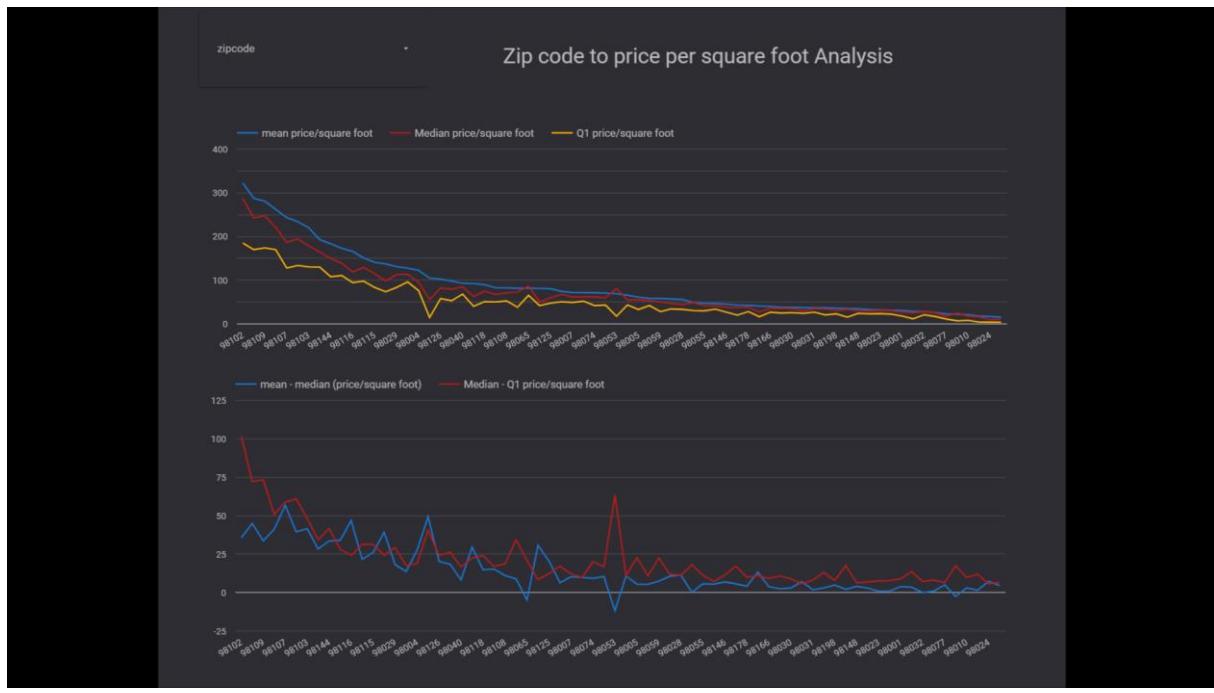


Figure 16

zipcode: 98027, 98053, 98022, 98...		(7)	Zip code to price per square foot Analysis			
zipcode	price/lot	mean - median (price/square foot)	Median - Q1 price/square foot	profit% (price/square feet) (se...	Record Count	
1. 98102	322.56	35.7	101.74	35.47	104	
2. 98119	287.43	45.02	72.33	29.84	184	
3. 98109	281.08	33.69	73.31	29.63	109	
4. 98027	104.82	49.43	40.53	73.17	412	
5. 98053	69.19	-11.78	63.36	78.25	403	
6. 98022	21.99	-2.46	17.53	71.71	233	
7. 98014	17.77	1.6	11.97	74.02	124	

By studying carefully charts from Figures 15, and 16 we decided that our focus should be to try to buy property that is priced under the lower quartile of their respective groups, and then flip them (sell them) at a price equal to or greater than the median price of its respective group.

Though finding property at a price below lower quartile will be hard but it will be a little easy to sell them for median prices as they are mostly closer to average prices and yet cheaper than 30-40% of the houses in same location.

Findings

We came to following conclusions/suggestions:

- For G1 (98022, 98014), if investors can buy property under **\$6.9/sqft** and **\$4.2/sqft** respectively and flip them for at least **\$24.4/sqft** and **\$16.2/sqft**, they can expect profit margin between **71-74%**.
- For G2 (98027, 98053), if investors can buy property under **\$14.8/sqft** and **\$17.6/sqft** respectively and flip them for at least **\$55.4/sqft** and **\$81/sqft**, they can expect profit margin between **73-78%**.
- For G3 (98102, 98109), if investors can buy property under **\$185/sqft** and **\$174/sqft** respectively and flip them for at least **\$287/sqft** and **\$247/sqft**, they can expect profit margin between **29-35%**.

Appendix A: Project Presentations

Project Week 1

1. Data Overview

Columns and rows in our data were 21 and 21597 respectively

Nominal data: ID, Waterfront, and Zip code.

Ordinal data: Date, View, Condition, Grade, Year built, and Year renovated.

Ratio: Price, Bedrooms, Bathrooms, sqft_living, sqft_lot, Floors, sqft_above, sqft_basement, sqft_living15, and sqft_lot15.

Interval: Lat, Long

2. Data Cleaning

2.1 Irrelevant Data

- Considered ID column
- Multiple entries with different values in other variables
- Left them untouched

The screenshot shows a Jupyter Notebook interface. On the left is a sidebar with icons for file operations, a search bar, and an 'OUTLINE' section. The main area displays a Pandas DataFrame:

```
id      date    price  bedrooms ...      lat    long  sqft_living15  sqft_lot15
17588  795000620 9/24/2014 115000.0 ... 47.5045 -122.33 1070.0 6250.0
17589  795000620 12/15/2014 124000.0 ... 47.5045 -122.33 1070.0 6250.0
17590  795000620 3/11/2015 157000.0 ... 47.5045 -122.33 1070.0 6250.0
```

[3 rows x 21 columns]

the results above show that though same house has been listed thrice, they were all listed on different dates with different prices!

Hence, I conclude that there are no duplicates in my opinion!

PS D:\Brainnest\PW_1\Brainnest_projects>

At the bottom, the status bar shows: Python 3.7.7 64-bit, Ln 18, Col 38, Spaces: 4, UTF-8, CRLF, Python.

2.1 Irrelevant Data

- Considered sqft_living15 and sqft_lot15 columns
- Different values for same records
- Hence, kept them.

The screenshot shows a Microsoft Excel spreadsheet with a large dataset. The first few rows of the table are:

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	
2	7.129E+09	10/13/2014	221900	3	1	1180	5650	1	0	0	3	7	1180	0	1955	0	98178	47.5112	-122.257	1340	5650
3	6.414E+09	12/9/2014	538000	3	2.25	2570	7242	2	0	0	3	7	2170	400	1991	0	98125	47.721	-122.319	1690	7639
4	5.632E+09	2/25/2015	180000	2	1	770	10000	1	0	0	3	6	770	0	1933	0	98028	47.7379	-122.233	2720	8062
5	4.487E+09	12/9/2014	604000	4	3	1960	5000	1	0	0	5	7	1050	910	1965	0	98136	47.5208	-122.393	1360	5000
6	1.954E+09	2/18/2015	510000	3	2	1680	8080	1	0	0	3	8	1680	0	1987	0	98074	47.6168	-122.045	7503	
7	7.238E+09	5/12/2014	1.23E+06	4	4.5	101930	1	0	0	3	11	3890	1530		0	98053	47.6561	-122.005	4760	101930	
8	1.321E+09	6/27/2014	257500	3	2.25	1715	6819	2	0	0	3	7	1715	0	1995	0	98198	47.4095	-122.315	1650	9711
9	2.008E+09	1/15/2015	291850	3	1.5	1060	9711	1	0	0	3	7	1060	0	1963	0	98007	47.6007		8925	
10	2.415E+09	4/15/2015	229500	3	1	1780	7470	1	0	0	3	7	1050	730	1960	0	98146	47.5123	-122.337	1780	8113
11	3.794E+09	3/12/2015	323000	3	2.5	1890	6560	2	0	0	3	7	1050	0	2003	0	98038	47.3684	-122.031	7570	
12	1.737E+09	4/3/2015	662500	3	2.5	3560	1	0	0	3	8	1860	1700	1965	0	98007	47.6007		8925		
13	9.213E+09	5/27/2014	468000	2	1	1160	6000	1	0	0	4	7	860	300	1942	0	98115	47.69	-122.292	1330	6000
14	114101516	5/28/2014	310000	3	1	1430	19901	1.5	0	0	4	7	1430	0	1927	0	98028	47.755	-122.229	1780	12697
15	6.055E+09	10/7/2014	400000	3	1.75	1370	9680	1	0	0	4	7	1370	0	1977	0	98074	47.6127	-122.045	1370	10208
16	1.175E+09	3/12/2015	530000	5	2	1810	4850	1.5	0	0	3	7	1810	0	1900	0	98107	47.67	-122.394	1360	4850
17	9.297E+09	1/24/2015	650000	4	3	2950	5000	2	0	3	3	9	1980	970		0	98126	47.5114	-122.375	2140	4000
18	1.876E+09	7/31/2014	395000	3	2	14040	2	0	0	3	7	1890	0	1994	0	98019	47.7277	-121.962	1890	14018	
19	6.865E+09	5/29/2014	485000	4	1	1600	4300	1.5	0	0	4	7	1600	0	1916	0	98103	47.6648	-122.343	1610	4300
20	16000397	12/5/2014	189000	2	1	1200	9850	1	0	0	4	7	1200	0	1921	0	98002	47.3089	-122.21	1060	5095
21	7.983E+09	4/24/2015	230000	3	1	1250	1	0	0	4	7	1250	0	1969	0	98003	47.3343	-122.306	1280	8850	
22	6.301E+09	5/14/2014	385000	4	1.75	1620	4980	1	0	0	4	7	860	760	1947	0	98133	47.7025	-122.341	1400	4980

2.2 Duplicate Data

First checked unique values
in each column

```
number of unique instances
id           21420
date         372
price        3621
bedrooms     13
bathrooms    29
sqft_living  1033
sqft_lot     9770
floors        7
waterfront    2
view          5
condition     5
grade         11
sqft_above    943
sqft_basement 306
yr_built      116
yr_renovated  70
zipcode       70
lat           5033
long          751
sqft_living15 777
sqft_lot15    8679
dtype: int64

the results indicate duplicate values as some of the variables
all count!
* cannot be duplicate as they are unique identifier
```

2.2 Duplicate Data

- Considered ID column
- Multiple entries with different values in other variables
- Therefore, did not remove them

```
id      date   price bedrooms ...      lat      long sqft_living15 sqft_lot15
17588 795000620 9/24/2014 115000.0      3 ... 47.5045 -122.33      1070.0      6250.0
17589 795000620 12/15/2014 124000.0      3 ... 47.5045 -122.33      1070.0      6250.0
17590 795000620 3/11/2015 157000.0      3 ... 47.5045 -122.33      1070.0      6250.0
[3 rows x 21 columns]

the results above show that though same house has been listed thrice, they were all listed on different dates with different prices!
Hence, I conclude that there are no duplicates in my opinion!
PS D:\Brainnest\pw_1\Brainnest_projects>
```

Python 3.7.7 64-bit > OUTLINE Ln 18, Col 38 Spaces: 4 UTF-8 CRLF Python

2.2 Duplicate Data

- Considered combination of latitude and longitude
- Houses could be apartments in same building
- Hence kept them

```
      id      date    price bedrooms ...      lat      long sqft_living15 sqft_lot15
14004 9136103136 3/13/2015 580000.0      2 ... 47.6652 -122.338      1490.0      4013.0
17852 7732410360 8/22/2014 752888.0      3 ... 47.6599 -122.146      2630.0      9000.0
19568 1025049266 9/30/2014 555000.0      2 ... 47.6647 -122.284      1160.0      1327.0
21209 3629700080 1/8/2015 635000.0      3 ... 47.5446 -122.017      2290.0      1407.0
21371 2895800750 4/17/2015 274800.0      3 ... 47.5171 -122.347      1410.0      1899.0
21482 1972200555 7/14/2014      NaN      3 ... 47.6536 -122.354      1570.0      1335.0

[6 rows x 21 columns]

      id      date    price bedrooms bathrooms sqft_living sqft_lot floors      lat      long
211 1025049114 7/17/2014 625504.0      3      2.25     1270.0    1566.0    2.0 47.6647 -122.284
13459 1025049115 6/25/2014 594000.0      3      2.25     1270.0    1406.0    2.0 47.6647 -122.284
18297 1025049268 7/21/2014 549900.0      2      1.75     1140.0    936.0    2.0 47.6647 -122.284
19568 1025049266 9/30/2014 555000.0      2      2.25     1160.0    954.0    2.0 47.6647 -122.284
      id      lat      long waterfront view condition grade sqft_above sqft_basement yr_built
211 1025049114 47.6647 -122.284      0 0.0      3.0     8.0    1060.0      210.0 2014.0
13459 1025049115 47.6647 -122.284      0 0.0      3.0     8.0    1060.0      210.0 2014.0
18297 1025049268 47.6647 -122.284      0 0.0      3.0     8.0     940.0      200.0 2014.0
19568 1025049266 47.6647 -122.284      0 0.0      3.0     8.0     960.0      200.0 2014.0

PS D:\Brainnest\PW_1\Brainnest_projects>
```

2.3 Structural errors

- Odd to see bathrooms and floors not in discrete numbers
- Still left it as it is

```
counts of each value present in column Bathrooms
1. 0.000      227
2. 0.500      3851
3. 1.000      36457
4. 1.500      2047
5. 2.000      33368
6. 2.500      3445
7. 3.000      3385
8. 3.500      723
9. 4.000      735
10. 4.500      589
11. 5.000      125
12. 5.500      136
13. 6.000      1889
14. 6.500      29
15. 7.000      71
16. 7.500      27
17. 8.000      21
18. 8.500      13
19. 9.000      10
20. 9.500      4
21. 10.000      2
22. 10.500      2
23. 11.000      2
24. 11.500      2
25. 12.000      2
26. 12.500      2
27. 13.000      2
28. 13.500      2
29. 14.000      2
30. 14.500      2
31. 15.000      2
32. 15.500      2
33. 16.000      2
34. 16.500      2
35. 17.000      2
36. 17.500      2
37. 18.000      2
38. 18.500      2
39. 19.000      2
40. 19.500      2
41. 20.000      2
42. 20.500      2
43. 21.000      2
44. 21.500      2
45. 22.000      2
46. 22.500      2
47. 23.000      2
48. 23.500      2
49. 24.000      2
50. 24.500      2
51. 25.000      2
52. 25.500      2
53. 26.000      2
54. 26.500      2
55. 27.000      2
56. 27.500      2
57. 28.000      2
58. 28.500      2
59. 29.000      2
60. 29.500      2
61. 30.000      2
62. 30.500      2
63. 31.000      2
64. 31.500      2
65. 32.000      2
66. 32.500      2
67. 33.000      2
68. 33.500      2
69. 34.000      2
70. 34.500      2
71. 35.000      2
72. 35.500      2
73. 36.000      2
74. 36.500      2
75. 37.000      2
76. 37.500      2
77. 38.000      2
78. 38.500      2
79. 39.000      2
80. 39.500      2
81. 40.000      2
82. 40.500      2
83. 41.000      2
84. 41.500      2
85. 42.000      2
86. 42.500      2
87. 43.000      2
88. 43.500      2
89. 44.000      2
90. 44.500      2
91. 45.000      2
92. 45.500      2
93. 46.000      2
94. 46.500      2
95. 47.000      2
96. 47.500      2
97. 48.000      2
98. 48.500      2
99. 49.000      2
100. 49.500      2
101. 50.000      2
102. 50.500      2
103. 51.000      2
104. 51.500      2
105. 52.000      2
106. 52.500      2
107. 53.000      2
108. 53.500      2
109. 54.000      2
110. 54.500      2
111. 55.000      2
112. 55.500      2
113. 56.000      2
114. 56.500      2
115. 57.000      2
116. 57.500      2
117. 58.000      2
118. 58.500      2
119. 59.000      2
120. 59.500      2
121. 60.000      2
122. 60.500      2
123. 61.000      2
124. 61.500      2
125. 62.000      2
126. 62.500      2
127. 63.000      2
128. 63.500      2
129. 64.000      2
130. 64.500      2
131. 65.000      2
132. 65.500      2
133. 66.000      2
134. 66.500      2
135. 67.000      2
136. 67.500      2
137. 68.000      2
138. 68.500      2
139. 69.000      2
140. 69.500      2
141. 70.000      2
142. 70.500      2
143. 71.000      2
144. 71.500      2
145. 72.000      2
146. 72.500      2
147. 73.000      2
148. 73.500      2
149. 74.000      2
150. 74.500      2
151. 75.000      2
152. 75.500      2
153. 76.000      2
154. 76.500      2
155. 77.000      2
156. 77.500      2
157. 78.000      2
158. 78.500      2
159. 79.000      2
160. 79.500      2
161. 80.000      2
162. 80.500      2
163. 81.000      2
164. 81.500      2
165. 82.000      2
166. 82.500      2
167. 83.000      2
168. 83.500      2
169. 84.000      2
170. 84.500      2
171. 85.000      2
172. 85.500      2
173. 86.000      2
174. 86.500      2
175. 87.000      2
176. 87.500      2
177. 88.000      2
178. 88.500      2
179. 89.000      2
180. 89.500      2
181. 90.000      2
182. 90.500      2
183. 91.000      2
184. 91.500      2
185. 92.000      2
186. 92.500      2
187. 93.000      2
188. 93.500      2
189. 94.000      2
190. 94.500      2
191. 95.000      2
192. 95.500      2
193. 96.000      2
194. 96.500      2
195. 97.000      2
196. 97.500      2
197. 98.000      2
198. 98.500      2
199. 99.000      2
200. 99.500      2
201. 100.000      2
202. 100.500      2
203. 101.000      2
204. 101.500      2
205. 102.000      2
206. 102.500      2
207. 103.000      2
208. 103.500      2
209. 104.000      2
210. 104.500      2
211. 105.000      2
212. 105.500      2
213. 106.000      2
214. 106.500      2
215. 107.000      2
216. 107.500      2
217. 108.000      2
218. 108.500      2
219. 109.000      2
220. 109.500      2
221. 110.000      2
222. 110.500      2
223. 111.000      2
224. 111.500      2
225. 112.000      2
226. 112.500      2
227. 113.000      2
228. 113.500      2
229. 114.000      2
230. 114.500      2
231. 115.000      2
232. 115.500      2
233. 116.000      2
234. 116.500      2
235. 117.000      2
236. 117.500      2
237. 118.000      2
238. 118.500      2
239. 119.000      2
240. 119.500      2
241. 120.000      2
242. 120.500      2
243. 121.000      2
244. 121.500      2
245. 122.000      2
246. 122.500      2
247. 123.000      2
248. 123.500      2
249. 124.000      2
250. 124.500      2
251. 125.000      2
252. 125.500      2
253. 126.000      2
254. 126.500      2
255. 127.000      2
256. 127.500      2
257. 128.000      2
258. 128.500      2
259. 129.000      2
260. 129.500      2
261. 130.000      2
262. 130.500      2
263. 131.000      2
264. 131.500      2
265. 132.000      2
266. 132.500      2
267. 133.000      2
268. 133.500      2
269. 134.000      2
270. 134.500      2
271. 135.000      2
272. 135.500      2
273. 136.000      2
274. 136.500      2
275. 137.000      2
276. 137.500      2
277. 138.000      2
278. 138.500      2
279. 139.000      2
280. 139.500      2
281. 140.000      2
282. 140.500      2
283. 141.000      2
284. 141.500      2
285. 142.000      2
286. 142.500      2
287. 143.000      2
288. 143.500      2
289. 144.000      2
290. 144.500      2
291. 145.000      2
292. 145.500      2
293. 146.000      2
294. 146.500      2
295. 147.000      2
296. 147.500      2
297. 148.000      2
298. 148.500      2
299. 149.000      2
300. 149.500      2
301. 150.000      2
302. 150.500      2
303. 151.000      2
304. 151.500      2
305. 152.000      2
306. 152.500      2
307. 153.000      2
308. 153.500      2
309. 154.000      2
310. 154.500      2
311. 155.000      2
312. 155.500      2
313. 156.000      2
314. 156.500      2
315. 157.000      2
316. 157.500      2
317. 158.000      2
318. 158.500      2
319. 159.000      2
320. 159.500      2
321. 160.000      2
322. 160.500      2
323. 161.000      2
324. 161.500      2
325. 162.000      2
326. 162.500      2
327. 163.000      2
328. 163.500      2
329. 164.000      2
330. 164.500      2
331. 165.000      2
332. 165.500      2
333. 166.000      2
334. 166.500      2
335. 167.000      2
336. 167.500      2
337. 168.000      2
338. 168.500      2
339. 169.000      2
340. 169.500      2
341. 170.000      2
342. 170.500      2
343. 171.000      2
344. 171.500      2
345. 172.000      2
346. 172.500      2
347. 173.000      2
348. 173.500      2
349. 174.000      2
350. 174.500      2
351. 175.000      2
352. 175.500      2
353. 176.000      2
354. 176.500      2
355. 177.000      2
356. 177.500      2
357. 178.000      2
358. 178.500      2
359. 179.000      2
360. 179.500      2
361. 180.000      2
362. 180.500      2
363. 181.000      2
364. 181.500      2
365. 182.000      2
366. 182.500      2
367. 183.000      2
368. 183.500      2
369. 184.000      2
370. 184.500      2
371. 185.000      2
372. 185.500      2
373. 186.000      2
374. 186.500      2
375. 187.000      2
376. 187.500      2
377. 188.000      2
378. 188.500      2
379. 189.000      2
380. 189.500      2
381. 190.000      2
382. 190.500      2
383. 191.000      2
384. 191.500      2
385. 192.000      2
386. 192.500      2
387. 193.000      2
388. 193.500      2
389. 194.000      2
390. 194.500      2
391. 195.000      2
392. 195.500      2
393. 196.000      2
394. 196.500      2
395. 197.000      2
396. 197.500      2
397. 198.000      2
398. 198.500      2
399. 199.000      2
400. 199.500      2
401. 200.000      2
402. 200.500      2
403. 201.000      2
404. 201.500      2
405. 202.000      2
406. 202.500      2
407. 203.000      2
408. 203.500      2
409. 204.000      2
410. 204.500      2
411. 205.000      2
412. 205.500      2
413. 206.000      2
414. 206.500      2
415. 207.000      2
416. 207.500      2
417. 208.000      2
418. 208.500      2
419. 209.000      2
420. 209.500      2
421. 210.000      2
422. 210.500      2
423. 211.000      2
424. 211.500      2
425. 212.000      2
426. 212.500      2
427. 213.000      2
428. 213.500      2
429. 214.000      2
430. 214.500      2
431. 215.000      2
432. 215.500      2
433. 216.000      2
434. 216.500      2
435. 217.000      2
436. 217.500      2
437. 218.000      2
438. 218.500      2
439. 219.000      2
440. 219.500      2
441. 220.000      2
442. 220.500      2
443. 221.000      2
444. 221.500      2
445. 222.000      2
446. 222.500      2
447. 223.000      2
448. 223.500      2
449. 224.000      2
450. 224.500      2
451. 225.000      2
452. 225.500      2
453. 226.000      2
454. 226.500      2
455. 227.000      2
456. 227.500      2
457. 228.000      2
458. 228.500      2
459. 229.000      2
460. 229.500      2
461. 230.000      2
462. 230.500      2
463. 231.000      2
464. 231.500      2
465. 232.000      2
466. 232.500      2
467. 233.000      2
468. 233.500      2
469. 234.000      2
470. 234.500      2
471. 235.000      2
472. 235.500      2
473. 236.000      2
474. 236.500      2
475. 237.000      2
476. 237.500      2
477. 238.000      2
478. 238.500      2
479. 239.000      2
480. 239.500      2
481. 240.000      2
482. 240.500      2
483. 241.000      2
484. 241.500      2
485. 242.000      2
486. 242.500      2
487. 243.000      2
488. 243.500      2
489. 244.000      2
490. 244.500      2
491. 245.000      2
492. 245.500      2
493. 246.000      2
494. 246.500      2
495. 247.000      2
496. 247.500      2
497. 248.000      2
498. 248.500      2
499. 249.000      2
500. 249.500      2
501. 250.000      2
502. 250.500      2
503. 251.000      2
504. 251.500      2
505. 252.000      2
506. 252.500      2
507. 253.000      2
508. 253.500      2
509. 254.000      2
510. 254.500      2
511. 255.000      2
512. 255.500      2
513. 256.000      2
514. 256.500      2
515. 257.000      2
516. 257.500      2
517. 258.000      2
518. 258.500      2
519. 259.000      2
520. 259.500      2
521. 260.000      2
522. 260.500      2
523. 261.000      2
524. 261.500      2
525. 262.000      2
526. 262.500      2
527. 263.000      2
528. 263.500      2
529. 264.000      2
530. 264.500      2
531. 265.000      2
532. 265.500      2
533. 266.000      2
534. 266.500      2
535. 267.000      2
536. 267.500      2
537. 268.000      2
538. 268.500      2
539. 269.000      2
540. 269.500      2
541. 270.000      2
542. 270.500      2
543. 271.000      2
544. 271.500      2
545. 272.000      2
546. 272.500      2
547. 273.000      2
548. 273.500      2
549. 274.000      2
550. 274.500      2
551. 275.000      2
552. 275.500      2
553. 276.000      2
554. 276.500      2
555. 277.000      2
556. 277.500      2
557. 278.000      2
558. 278.500      2
559. 279.000      2
560. 279.500      2
561. 280.000      2
562. 280.500      2
563. 281.000      2
564. 281.500      2
565. 282.000      2
566. 282.500      2
567. 283.000      2
568. 283.500      2
569. 284.000      2
570. 284.500      2
571. 285.000      2
572. 285.500      2
573. 286.000      2
574. 286.500      2
575. 287.000      2
576. 287.500      2
577. 288.000      2
578. 288.500      2
579. 289.000      2
580. 289.500      2
581. 290.000      2
582. 290.500      2
583. 291.000      2
584. 291.500      2
585. 292.000      2
586. 292.500      2
587. 293.000      2
588. 293.500      2
589. 294.000      2
590. 294.500      2
591. 295.000      2
592. 295.500      2
593. 296.000      2
594. 296.500      2
595. 297.000      2
596. 297.500      2
597. 298.000      2
598. 298.500      2
599. 299.000      2
600. 299.500      2
601. 300.000      2
602. 300.500      2
603. 301.000      2
604. 301.500      2
605. 302.000      2
606. 302.500      2
607. 303.000      2
608. 303.500      2
609. 304.000      2
610. 304.500      2
611. 305.000      2
612. 305.500      2
613. 306.000      2
614. 306.500      2
615. 307.000      2
616. 307.500      2
617. 308.000      2
618. 308.500      2
619. 309.000      2
620. 309.500      2
621. 310.000      2
622. 310.500      2
623. 311.000      2
624. 311.500      2
625. 312.000      2
626. 312.500      2
627. 313.000      2
628. 313.500      2
629. 314.000      2
630. 314.500      2
631. 315.000      2
632. 315.500      2
633. 316.000      2
634. 316.500      2
635. 317.000      2
636. 317.500      2
637. 318.000      2
638. 318.500      2
639. 319.000      2
640. 319.500      2
641. 320.000      2
642. 320.500      2
```

2.4 Missing Data

```
OUTPUT TERMINAL DEBUG CONSOLE
D:\Brainnest\PW_1\Brainnest_projects> & 'C:\Users\hon.exe' 'c:\Users\Farrukh Nadeem\vscode\extensions\ms-python/debugpy\launcher' '59548' '--' 'd:\Brainnest\PW_1\Brainnest_
missing values
id          0
date        0
price       9
bedrooms    0
bathrooms   1
sqft_living 10
sqft_lot    16
floors      0
waterfront  0
view        3
condition   3
grade       5
sqft_above  12
sqft_basement 2
yr_built     6
yr_renovated 0
zipcode      7
lat         1
long        5
sqft_living15 7
sqft_lot15   8
dtype: int64
PS D:\Brainnest\PW_1\Brainnest_projects>
```

2.5 Data Outliers

33 OR 40 BEDROOMS?

11 FLOORS

```
counts of each value present in column floors
1.000    10672
2.000    8235
1.500    1910
3.000    611
2.500    161
3.500     7
11.000    1
Name: floors, dtype: int64
counts of each value present in column Bedrooms
3       9824
4       6881
2       2760
5       1601
6       272
1       196
7       38
8       13
9       6
10      3
11      1
33      1
40      1
Name: bedrooms, dtype: int64
```

TOO FEW BATHROOMS FOR THAT
MANY BEDROOMS!

11 FLOORS FOR JUST 2 BEDROOMS?

```
      price  bedrooms  bathrooms  sqft_living  sqft_lot  floors
8748  520000.000      11      3.000    3000.000  4900.000  2.000
15856 640000.000      33      1.750    1628.000  6000.000  1.000
16158 598000.000      40      2.500    3130.000  409188.000  2.000
floors greater than 3?

      price  bedrooms  bathrooms  sqft_living  sqft_lot  floors
10066 435000.000      3      3.000    1440.000  1350.000  3.500
11582 544000.000      3      2.500    1760.000  1755.000  3.500
14871 525000.000      3      3.000    1730.000  1874.000  3.500
15410 479000.000      2      2.500    1730.000  1837.000  3.500
16146 365000.000      2      1.000    1250.000  8100.000  11.000
18462 330000.000      8      4.000    7710.000  11750.000  3.500
20292 525000.000      2      2.750    1310.000  1268.000  3.500
20756 563500.000      3      2.500    1400.000  1312.000  3.500
PS D:\Brainnest\PW_1\Brainnest_projects>
```

3. Descriptive Statistics

```

Group_B_kc_house_da...  price  bedrooms  bathrooms  sqft_living  sqft_lot  floors
count    21588.000  21597.000  21596.000  21587.000  21581.000  21597.000
mean     540382.754    3.375    2.116   2080.185  15113.467    1.495
std      367405.169    0.959    0.769   917.968  41511.996    0.544
min      78000.000    1.000    0.500   370.000   520.000    1.000
25%     322000.000    3.000    1.750  1430.000  5040.000    1.000
50%     450000.000    3.000    2.250  1910.000  7617.000    1.500
75%     645000.000    4.000    2.500  2550.000  18677.000    2.000
max     7700000.000   48.000    8.000 13540.000 1651359.000   11.000
waterfront      view  condition  grade  sqft_above  sqft_basement  yr_built  yr_renovated
count    21597.000  21594.000  21594.000  21585.000  21595.000  21591.000  21597.000
mean     0.008    0.234    3.410    7.658   1789.258   291.710  1970.998    84.465
std      0.087    0.766    0.651    1.173   834.846   442.664   29.377   401.821
min      0.000    0.000    1.000    3.000   370.000    0.000  1900.000    0.000
25%     0.000    0.000    3.000    7.000   1190.000    0.000  1951.000    0.000
50%     0.000    0.000    3.000    7.000   1560.000    0.000  1975.000    0.000
75%     0.000    0.000    4.000    8.000   2210.000   560.000  1997.000    0.000
max     1.000    4.000    5.000   13.000  18000.000   4820.000  2015.000   2015.000
 zipcode      lat    long  sqft_living15  sqft_lot15
count    21598.000  21596.000  21592.000  21590.000  21589.000
mean     98077.950    47.560   -122.214   1986.586  12761.042
std      53.513    0.139    0.141   685.296  27279.087
min     98001.000    47.156   -122.519   399.000   651.000
25%     98033.000    47.471   -122.328   1490.000   5100.000
50%     98065.000    47.572   -122.231   1840.000   7620.000
75%     98118.000    47.678   -122.125   2360.000  10083.000
max     98199.000    47.778   -121.315   6210.000  871200.000

```

> OUTLINE

4. Project Questions

What are the ideal locations (Zip codes) to invest in considering profit margins on resale in mind?

Does the ratio of bedroom to lot area have any impact on price? If yes, what is the suggested ratio?

Does the ratio of bedroom to floors have any impact on price? If yes, what is the suggested ratio?

Does the ratio of bedroom to bathrooms have any impact on price? If yes, what is the suggested ratio?

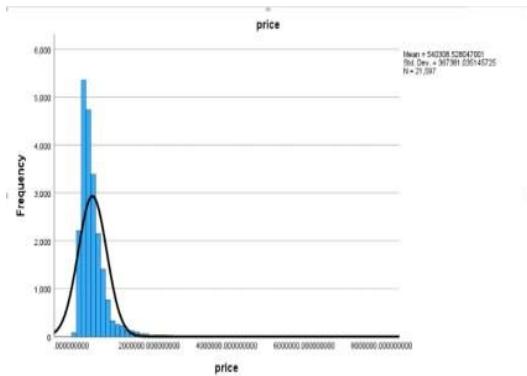
Project week 2

Skewness and Kurtosis

variable name	Skewness			Kurtosis		
	Statistic	Std. Error	Stat/std. error	Statistic	Std. Error	Stat/std. error
id	0.243	0.017	14.29411765	-1.261	0.033	-38.21212121
price	4.023	0.017	236.6470588	34.536	0.033	1046.545455
bedrooms	0.549	0.017	32.29411765	1.797	0.033	54.45454545
bathrooms	0.52	0.017	30.58823529	1.279	0.033	38.75757576
sqft_living	1.471	0.017	86.52941176	5.253	0.033	159.1818182
sqft_lot	13.075	0.017	769.1176471	285.571	0.033	8653.666667
floors	0.615	0.017	36.17647059	-0.491	0.033	-14.87878788
waterfront	11.275	0.017	663.2352941	125.139	0.033	3792.090909
view	3.395	0.017	199.7058824	10.89	0.033	330
condition	1.028	0.017	60.47058824	0.536	0.033	16.24242424
grade	0.76	0.017	44.70588235	1.245	0.033	37.72727273
sqft_above	1.745	0.017	102.6470588	9.77	0.033	296.0606061
sqft_basement	1.577	0.017	92.76470588	2.713	0.033	82.21212121
yr_built	-20.444	0.017	-1202.588235	1062.416	0.033	32194.42424
yr_renovated	4.548	0.017	267.5294118	18.684	0.033	566.1818182
zipcode	-84.549	0.017	-4973.470588	7162.595	0.033	217048.3333
lat	-77.329	0.017	-4548.764706	6358.629	0.033	192685.7273
long	83.655	0.017	4920.882353	7062.384	0.033	214011.6364
sqft_living15	1.103	0.017	64.88235294	1.595	0.033	48.33333333
sqft_lot15	9.525	0.017	560.2941176	151.406	0.033	4588.060606

Price

Histogram

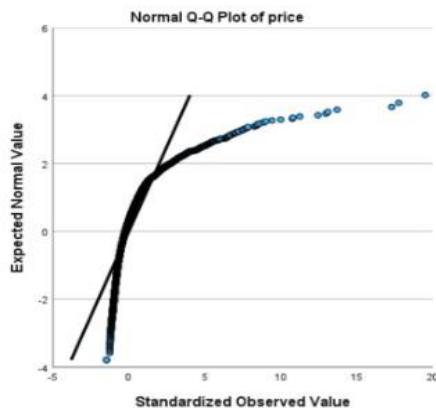


- We can see that the tail is skewing more to the right
- The distribution isn't symmetric
- The median is going also to the right.
- If Median=M, Mean=x, and Mode=m, then $x \neq M \neq m$

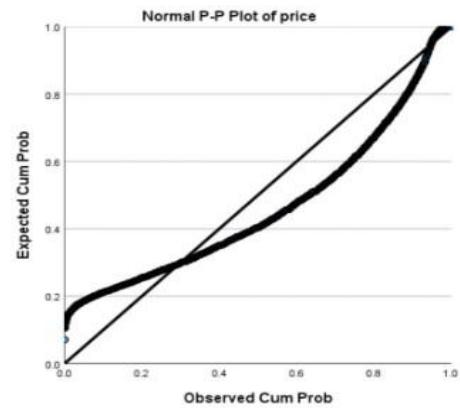
Price

We can see that the distribution is not normal, many of the values are far away from the straight line.

Q-Q Plot



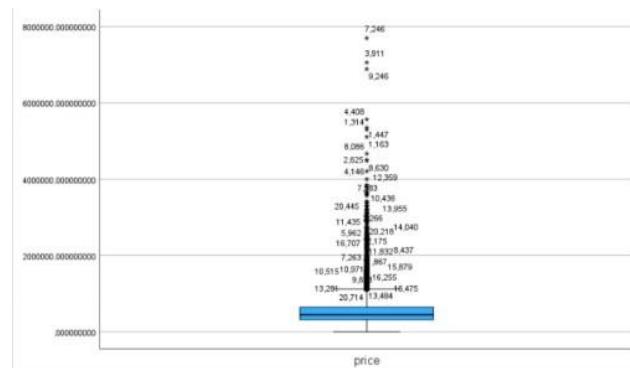
P-P Plot



Price

Box Plot

- We can see that the whiskers of the box plot are almost equals but the boxes aren't.
- The distribution is skewed up (positive)
- Non symmetric distribution
- There are outliers. There are two different types. The stars are the Extreme outliers, and the circles are the mild outliers



Price

Visual and statistical tests like Histogram, Box Plots, Q-Q plots and P-P plots, showed that the price is not normally distributed with skewness of -4.023(SE= 0.017) and kurtosis of 34.536(SE= 0.033).

Test	Normally Distributed
Histogram	No
QQ Plot	No
PP Plot	No
Box Plot	No
Skewness	No
Z score (Skewness)	No (236.64)
Kurtosis	No
Z score (Kurtosis)	No (1046)
Final conclusion	No (8/8)

Bathrooms

Visual and statistical tests like Histogram, Box Plots, Q-Q plots and P-P plots, showed that the Bathrooms is not normally distributed with skewness of 0.52(SE=0.017) and kurtosis of 1.279 (SE=0.033).

Test	Normally Distributed
Histogram	No
QQ Plot	No
PP Plot	Yes
Box Plot	No
Skewness	Close to 0
Z score (Skewness)	30.58
Kurtosis	Not that close to 0
Z score (Kurtosis)	38.75
Final conclusion	Not normal(6/8)

Bedroom

Visual and statistical tests like Histogram, Box Plots, Q-Q plots and P-P plots, showed that the Bedroom is not normally distributed with skewness of 0.549(SE=0.017) and kurtosis of 1.797(SE=0.033).

Test	Normally Distributed
Histogram	No
QQ Plot	No
PP Plot	Yes
Box Plot	No
Skewness	Is close to 0
Z score (Skewness)	(32.29) Big
Kurtosis	1.797
Z score (Kurtosis)	(54.45)
Final conclusion	Non normal (6/8)

Square Feet Living

Visual and statistical tests like Histogram, Box Plots, Q-Q plots and P-P plots, showed that the Square Feet Living is not normally distributed with skewness of 1.471(SE=0.017) and kurtosis of 5.253 (SE=0.033).

Test	Normally Distributed
Histogram	No
QQ Plot	No
PP Plot	No
Box Plot	No
Skewness	No
Z score (Skewness)	No(86.52)
Kurtosis	No
Z score (Kurtosis)	No(159.18)
Final conclusion	No

Square Feet Lot

Visual and statistical tests like Histogram, Box Plots, Q-Q plots and P-P plots, showed that the Square Feet Lot is not normally distributed with skewness of 13.075(SE=0.017) and kurtosis of 285.571 (SE=0.033).

Test	Normally Distributed
Histogram	No
QQ Plot	No
PP Plot	No
Box Plot	No
Skewness	No
Z score (Skewness)	No
Kurtosis	No
Z score (Kurtosis)	No
Final conclusion	No

Floors

Visual and statistical tests like Histogram, Box Plots, Q-Q plots and P-P plots, showed that the Floors is approximately normally distributed with skewness of 0.615 (SE=0.017) and kurtosis of -0.491(SE=0.033).

Test	Normally Distributed
Histogram	No
QQ Plot	Yes
PP Plot	Yes
Box Plot	No
Skewness	Close to zero
Z score (Skewness)	36.1764
Kurtosis	Close to zero
Z score (Kurtosis)	-14.8787
Conclusion	Approximately normally distributed (4/8)

Square Feet Above

Visual and statistical tests like Histogram, Box Plots, Q-Q plots and P-P plots, showed that the Square Feet Above variable is not normally distributed with skewness of -1.745 (SE=0.017) and kurtosis of 9.77 (SE=0.033).

Test	Normally Distributed
Histogram	No
QQ Plot	No
PP Plot	No
Box Plot	No
Skewness	Not close to 0
Z score (Skewness)	No (102.64)
Kurtosis	Not close to 0
Z score (Kurtosis)	No (296.06)
Final conclusion	No (8/8)

Square Feet Basement

Visual and statistical tests like Histogram, Box Plots, Q-Q plots and P-P plots, showed that the Square Feet Basement is not normally distributed with skewness of 1.577(SE=0.017) and kurtosis of 2.713(SE=0.033).

Test	Normally Distributed
Histogram	No
QQ Plot	No
PP Plot	No
Box Plot	No
Skewness	Not close to 0
Z score (Skewness)	No (92.76)
Kurtosis	Not close to 0
Z score (Kurtosis)	No (82.21)
Final conclusion	No (8/8)

Year Built

Visual and statistical tests like Histogram, Box Plots, Q-Q plots and P-P plots, showed that the year built is not normally distributed with skewness of -20.44 (SE=0.017) and kurtosis of 1062 (SE=0.033).

Test	Normally Distributed
Histogram	No
QQ Plot	Yes
PP Plot	Yes
Box Plot	No
Skewness	No
Z score (Skewness)	No
Kurtosis	No
Z score (Kurtosis)	No
Final conclusion	NO

Year Renovated

Visual and statistical tests like Histogram, Box Plots, Q-Q plots and P-P plots, showed that the year renovated is or not normally distributed with skewness of 4.548 (SE=0.017) and kurtosis of 18.684 (SE=0.033).

Test	Normally Distributed
Histogram	No
QQ Plot	No
PP Plot	No
Box Plot	No
Skewness	No
Z score (Skewness)	No (4.548)
Kurtosis	No
Z score (Kurtosis)	No (18.684)
Final conclusion	No 8/8

Sqft_living15

Visual and statistical tests like Histogram, Box Plots, Q-Q plots and P-P plots, showed that the {Column Name} is {approximately normally distributed or not normally distributed} with skewness of 1.103 (SE = 0.017) and kurtosis of 1.595 (SE=0.033).

Test	Normally Distributed
Histogram	No
QQ Plot	No
PP Plot	No
Box Plot	No
Skewness	No
Z score (Skewness)	Yes (1.103)
Kurtosis	No
Z score (Kurtosis)	Yes (1.595)
Final conclusion	No 6/8

Sqft_lot15_1

Visual and statistical tests like Histogram, Box Plots, Q-Q plots and P-P plots, showed that the {Column Name} is {approximately normally distributed or not normally distributed} with skewness of 9.525 (SE=0.017) and kurtosis of 151.406 (SE= 0.033).

Test	Normally Distributed
Histogram	No
QQ Plot	No
PP Plot	No
Box Plot	No
Skewness	No
Z score (Skewness)	No (9.525)
Kurtosis	No
Z score (Kurtosis)	No (151.406)
Final conclusion	No (8/8)

Overall Picture

Variable	Remarks	Value
Id	Unique identifiers no need to test normality	
Date	Considered them to be categorical so no need to test normality	
Price	Not Normal	100%
Bedrooms	Not Normal	75%
Bathrooms	Not Normal	75%
sqft_living	Not Normal	100%
sqft_lot	Not Normal	100%
Floors	Approximately normal	50%
Waterfront	Categorical variable. Hence, no need to test normality.	

Overall Picture Continued...

Variable	Remarks	Value
View	Categorical variable. Hence, no need to test normality.	
Condition	Categorical variable. Hence, no need to test normality.	
Grade	Categorical variable. Hence, no need to test normality.	
sqft_above	Not Normal	100%
sqft_basement	Not Normal	100%
yr_built	Not Normal	100%
yr_renovated	Not Normal	100%
Zipcode	Categorical variable. Hence, no need to test normality.	
Lat	Unique identifiers no need to test normality	

Overall Picture Continued...

Variable	Remarks	Value
Long	Unique identifiers no need to test normality	
sqft_living15	Not Normal	75%
sqft_lot15	Not Normal	100%

- For Floors we are 50% sure it is Normally distributed
- For other variables we are at least 75% sure they are not normally distributed.
- So we suggest going for nonparametric tests in our future analysis.

Final Presentation

Project Objectives

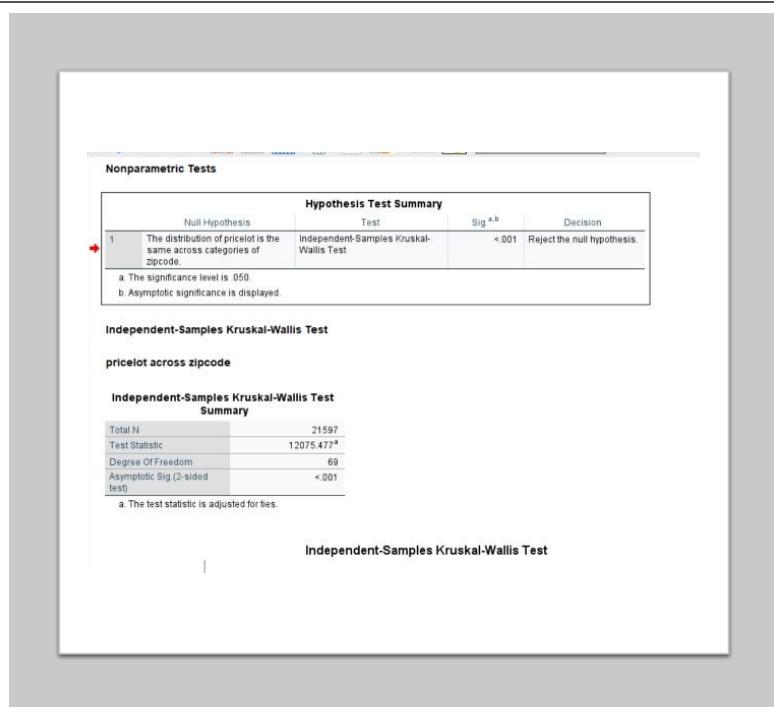
PO.	Objective
1	What are the ideal locations (Zip codes) to invest in considering profit margins on resale in mind?
2	Does the ratio of bedroom to lot area have any impact on price? If yes, what is the suggested ratio?
3	Does the ratio of bedroom to floors have any impact on price? If yes, what is the suggested ratio?
4	Does the ratio of bedroom to bathrooms have any impact on price? If yes, what is the suggested ratio?

Hypothesis Testing

PO.	H0	H1	Applied Test	Result
1	The cost of land per unit area is same across all neighborhoods (i.e., zip codes).	The cost of land per unit area is affected by the neighborhood.		Accept H1
2	The price is not affected by different ratios of lot area to bedroom.	The price is affected by lot area to bedroom ratio.	Kruskal-Wallis Test	Accept H1
3	The price is not affected by different ratios of bedroom to floors.	The price is affected by bedroom to floor ratio.		Accept H1
4	The price is not affected by different ratios of bedroom to bathrooms.	The price is affected by bedroom to bathroom ratio.		Accept H1

Hypothesis Test for PO1

- The p-value is less than 0.05, hence we reject null hypotheses and conclude that:
- Cost of land per square foot is indeed affected by the location (i.e., Zip codes).



lot/bedroom	price	Record ...	lot/bedroom	price	Record ...
No data			1. 7465	7,060,000	1
			2. 5229	6,890,000	1
			3. 7013.8	5,570,000	1
			4. 4797	5,350,000	1
			5. 4138.166667	5,300,000	1
			6. 9103.4	5,110,000	1
			7. 8002.8	4,500,000	1
			8. 6879.25	4,490,000	1
			9. 4209	4,210,000	1
					1 - 100 / 11535 < >

PO2

- Results from hypotheses testing did indicate possible existence of patterns.

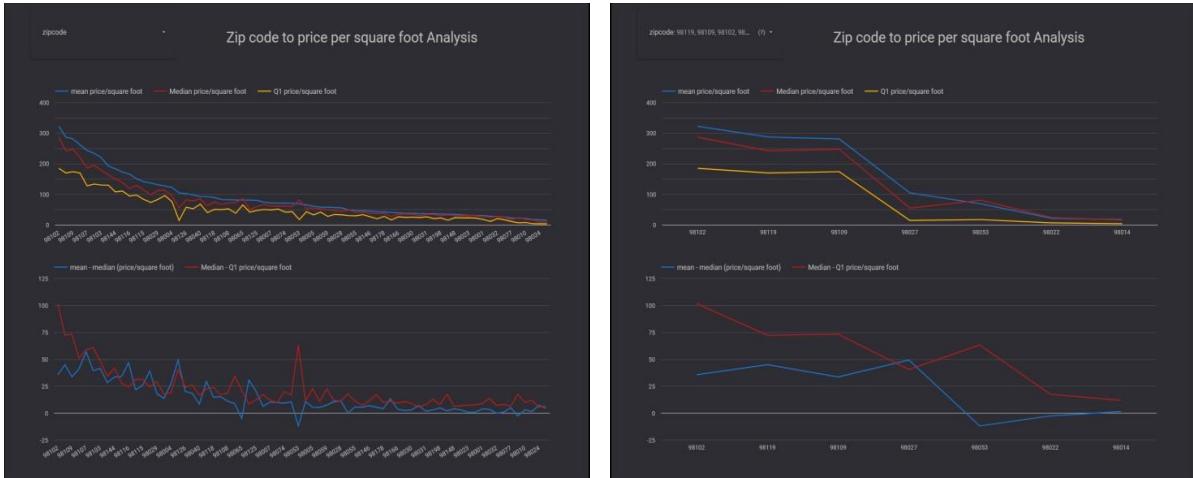
PO3

- The group sizes within respective samples did not meet the minimum threshold of at least 2.5% representation of the total sample size.

PO4

- Thus, detecting any pattern from such small groups is meaningless in our opinion.

PO1: What are the ideal locations (Zip codes) to invest in considering profit margins on resale in mind?



zipcode: 98027, 98053, 98022, 98... (7) ▾ Zip code to price per square foot Analysis

zipcode	mean price/square fo...	Median price/square fo...	Q1 price/square fo...	Median - Q1 price/square fo...	profit% (price/square feet) (selling at me...
1. 98102	322.56	286.86	185.12	101.74	35.47
2. 98119	287.43	242.41	170.08	72.33	29.84
3. 98109	281.08	247.39	174.08	73.31	29.63
4. 98027	104.82	55.39	14.86	40.53	73.17
5. 98053	69.19	80.97	17.61	63.36	78.25
6. 98022	21.99	24.45	6.92	17.53	71.71
7. 98014	17.77	16.17	4.2	11.97	74.02

1 - 7 / 7 < >

- We divided our zones of interest into three groups G1, G2, and G3; with G1 requiring minimum investment and G3 the highest for same property size.
- Buy property at a price per square foot that lies within first quartile of its respective zip-code
- Try to sell them at a price equal to or greater than Median to earn desired profit.

PO1

PO1: Report/Insights

- For G1 (98022, 98014), if investors can buy property under **\$6.9/sqft** and **\$4.2/sqft** respectively and flip them for at least **\$24.4/sqft** and **\$16.2/sqft**, they can expect profit margin between **71-74%**.
- For G2 (98027, 98053), if investors can buy property under **\$14.8/sqft** and **\$17.6/sqft** respectively and flip them for at least **\$55.4/sqft** and **\$81/sqft**, they can expect profit margin between **73-78%**.
- For G3 (98102, 98109), if investors can buy property under **\$185/sqft** and **\$174/sqft** respectively and flip them for at least **\$287/sqft** and **\$247/sqft**, they can expect profit margin between **29-35%**.

Appendix B: Tools Used

Throughout the project we used following tools:

Python (Data cleaning and EDA)

SPSS (Statistical analysis, and some visualizations)

Google Data Studio (Data Visualization)