

CIVL 4210 – Advanced Construction with AI and Robotics

Guidebook: Housing Price Regression

Prof. Yu
HUANG Xuhong

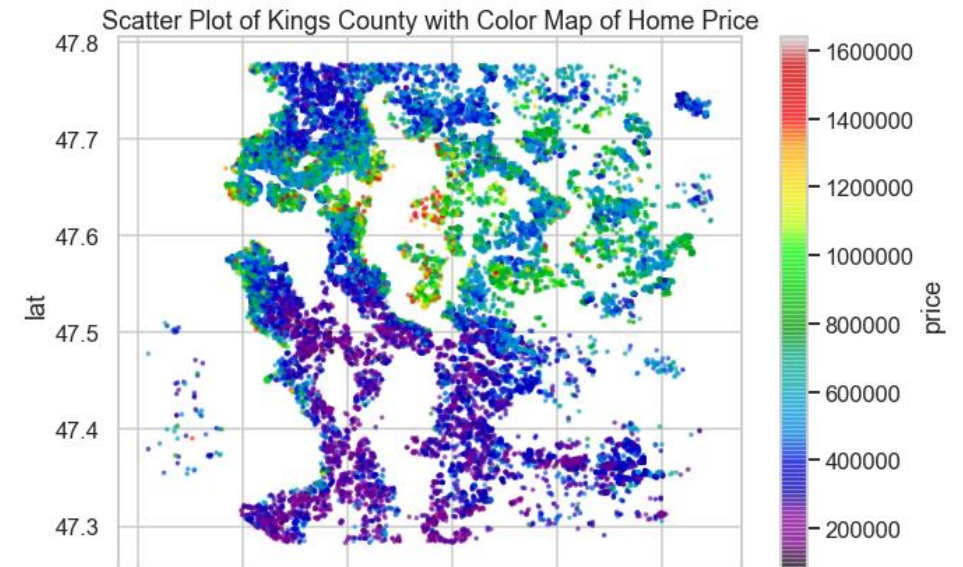
Introduction

The objective is to create a **linear regression model** for a given dataset(House Sales in King County, USA).

The overall idea of regression is to examine two things:

- (1) Does a set of predictor variables do a good job of predicting an outcome (dependent) variable?
- (2) Which variables are significant predictors of the outcome variable, and how do they—indicated by the magnitude and sign of the beta estimates—impact the outcome variable?

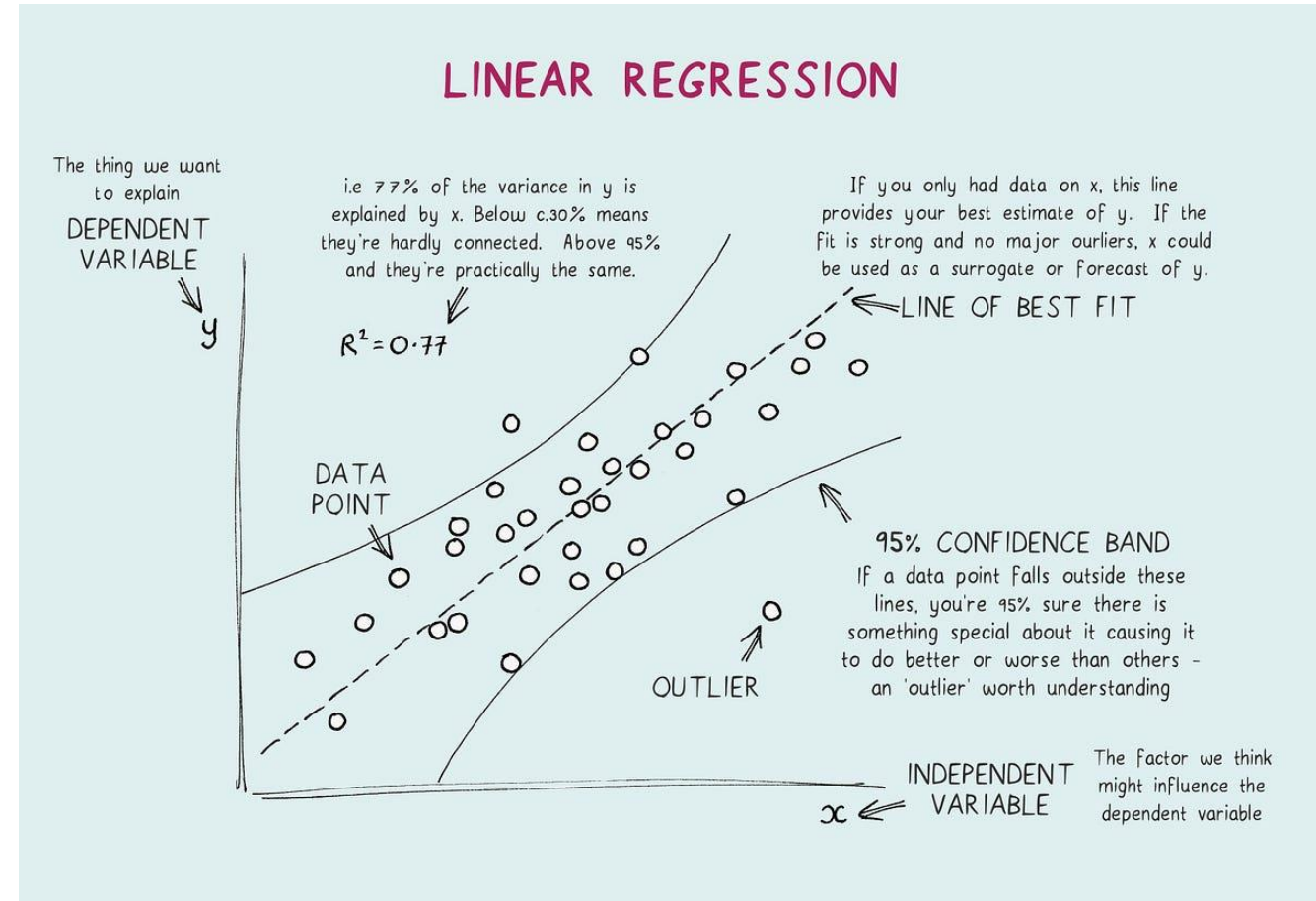
These regression estimates are used to explain the relationship between one dependent variable and one or more independent variables.



Introduction

Linear Regression Analysis consists of more than just fitting a linear line through a cloud of data points. It consists of 3 stages.

- (1) Analyzing the correlation and directionality of the data
- (2) Estimating the model, i.e., fitting the line
- (3) Evaluating the validity and usefulness of the model



Data

In this dataset, we need to predict the **sales price of houses in King County, Seattle**. It includes homes sold between May 2014 and May 2015.

Before doing anything, we should first know about the dataset, what it contains, what its features are, and what the structure of the data is.

The dataset contains 20 house features plus the price, along with 21613 observations.

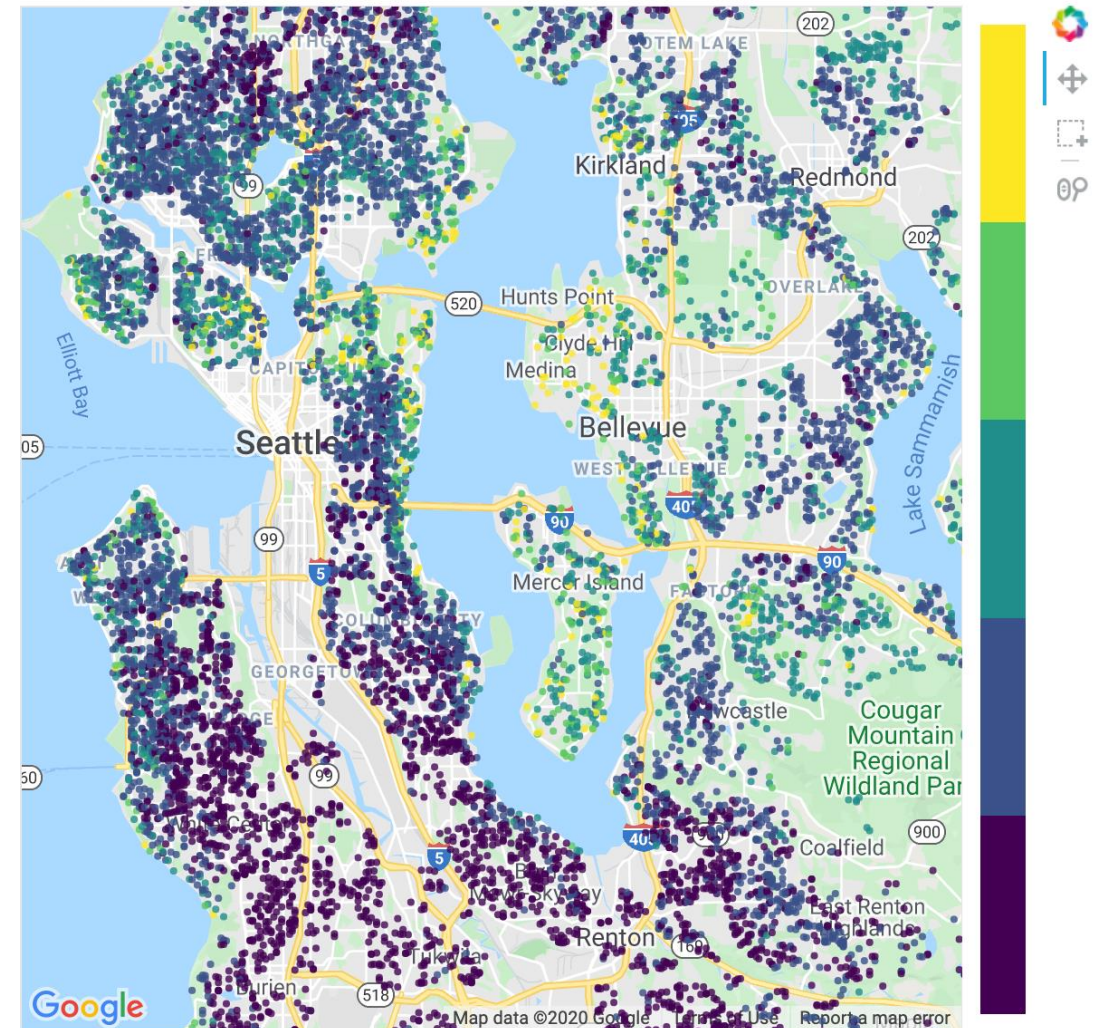


THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY



SCHOOL OF
ENGINEERING

Home Sale Prices in Kings County from May 2014 - May 2015



The description for the 20 features is given below:

- (1) **id** :- It is the unique numeric number assigned to each house being sold.
- (2) **date** :- It is the date on which the house was sold out.
- (3) **price**:- It is the price of house which we have to predict so this is our target variable and apart from it are our features.
- (4) **bedrooms** :- It determines number of bedrooms in a house.
- (5) **bathrooms** :- It determines number of bathrooms in a bedroom of a house.
- (6) **sqft_living** :- It is the measurement variable which determines the measurement of house in square foot.
- (7) **sqft_lot** : It is also the measurement variable which determines square foot of the lot.
- (8) **floors** : It determines total floors means levels of house.
- (9) **waterfront** : This feature determines whether a house has a view to waterfront 0 means no 1 means yes.
- (10) **view** : This feature determines whether a house has been viewed or not 0 means no 1 means yes.

The description for the 20 features is given below:

- (11) **condition** : It determines the overall condition of a house on a scale of 1 to 5.
- (12) **grade** : It determines the overall grade given to the housing unit, based on King County grading system on a scale of 1 to 11
- (13) **sqft_above** : It determines square footage of house apart from basement.
- (14) **sqft_basement** : It determines square footage of the basement of the house.
- (15) **yr_built** : It determines the date of building of the house.
- (16) **yr_renovated** : It determines year of renovation of house.
- (17) **zipcode** : It determines the zipcode of the location of the house.
- (18) **lat** : It determines the latitude of the location of the house.
- (19) **long** : It determines the longitude of the location of the house.
- (20) **sqft_living15** : Living room area in 2015(implies-- some renovations)
- (21) **sqft_lot15** : lotSize area in 2015(implies-- some renovations)

By observing the data, we can know that the price is dependent on various features like **bedrooms**(which is the most dependent feature), **bathrooms**, **sqft_living**(second most important feature), **sqft_lot**, **floors**, etc. The price is also dependent on the location of the house where it is present. The other features, like the waterfront view, are less dependent on the price. Of all the records, there are no missing values, which helps us create a better model.

First, we import the required libraries like pandas, numpy, seaborn, and matplotlib. Now, import the CSV file. Now, we should get to know how the data is and what data type uses the info function. We observe that the date is in 'object' format. To see the no of rows and columns, we use the shape function. Describe the data frame to know the mean, minimum, maximum, standard deviation, and percentiles.

Then, we plot graphs for visualization, and then we do simple regression using 'bedrooms,' multiple regression, and polynomial regression.

Step 1: Check the data in Colab

Step 1: Import and initialization Step 2: Check the data

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

```
csvData =
'https://raw.githubusercontent.com/Yokhong/CIVL4210/main/H3
_housing_price_regression/kc_h
ouse_data.csv'
```

```
df = pd.read_csv(csvData)
```

df.info()

```
>>> <class 'pandas.core.frame.DataFrame'>
RangeIndex: 21613 entries, 0 to 21612
Data columns (total 21 columns):
 #   Column              Non-Null Count  Dtype  
---  --
 0   id                  21613 non-null  int64  
 1   date                21613 non-null  object  
 2   price               21613 non-null  float64 
 3   bedrooms            21613 non-null  int64  
 4   bathrooms           21613 non-null  float64 
 5   sqft_living         21613 non-null  int64  
 6   sqft_lot            21613 non-null  int64  
 7   floors              21613 non-null  float64 
 8   waterfront          21613 non-null  int64  
 9   view                21613 non-null  int64  
10   condition           21613 non-null  int64  
11   grade               21613 non-null  int64  
12   sqft_above          21613 non-null  int64  
13   sqft_basement       21613 non-null  int64  
14   yr_built            21613 non-null  int64  
15   yr_renovated        21613 non-null  int64  
16   zipcode             21613 non-null  int64  
17   lat                 21613 non-null  float64 
18   long                21613 non-null  float64 
19   sqft_living15       21613 non-null  int64  
20   sqft_lot15          21613 non-null  int64  
dtypes: float64(5), int64(15), object(1)
memory usage: 3.5+ MB
```

df.head()

```
>>>
   id            date  price  bedrooms  bathrooms  sqft_living  sqft_lot  floors  waterfront  view  ...
0  7129300520  20141013T000000  221900.0         3         1.00        1180      5650      1.0         0         0  ...
1  6414100192  20141209T000000  538000.0         3         2.25        2570      7242      2.0         0         0  ...
2  5631500400  20150225T000000  180000.0         2         1.00         770     10000      1.0         0         0  ...
3  2487200875  20141209T000000  604000.0         4         3.00        1960      5000      1.0         0         0  ...
4  1954400510  20150218T000000  510000.0         3         2.00        1680      8080      1.0         0         0  ...
5 rows x 21 columns
```

■
■
■

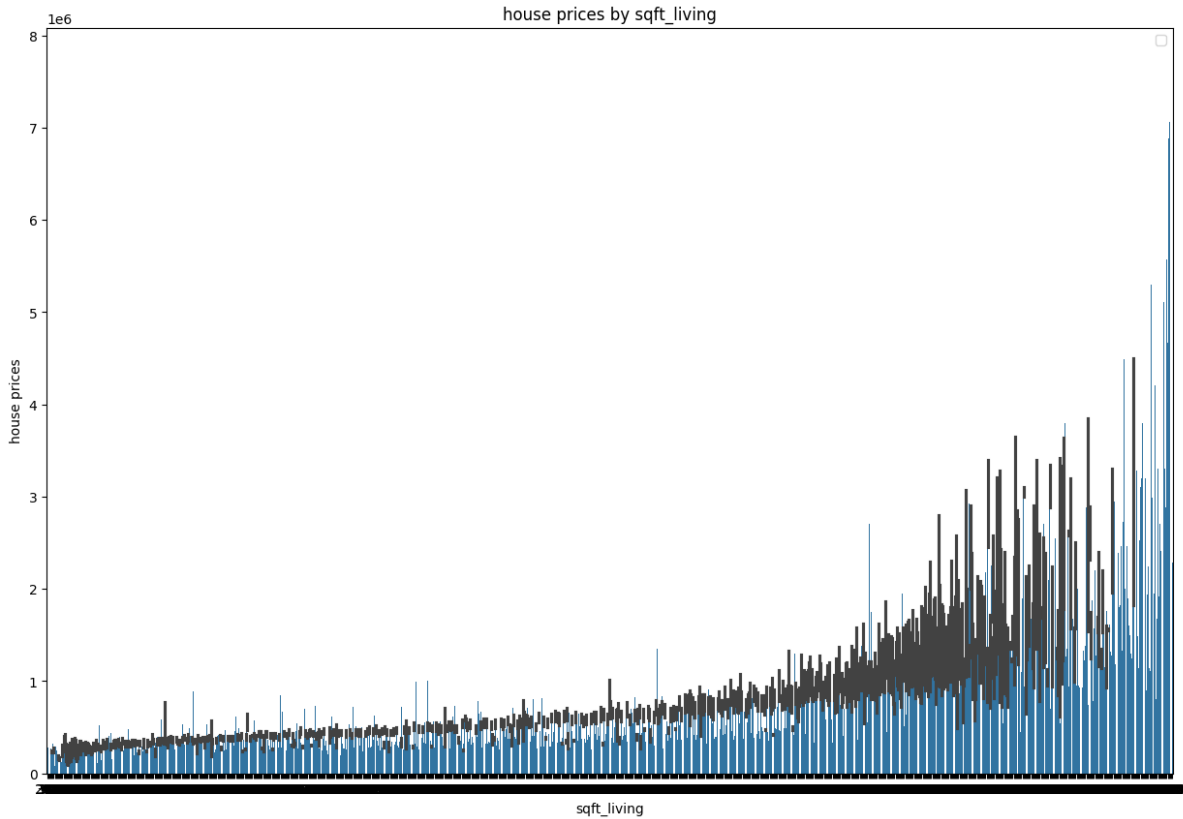
Step 2: Statistics and graphing



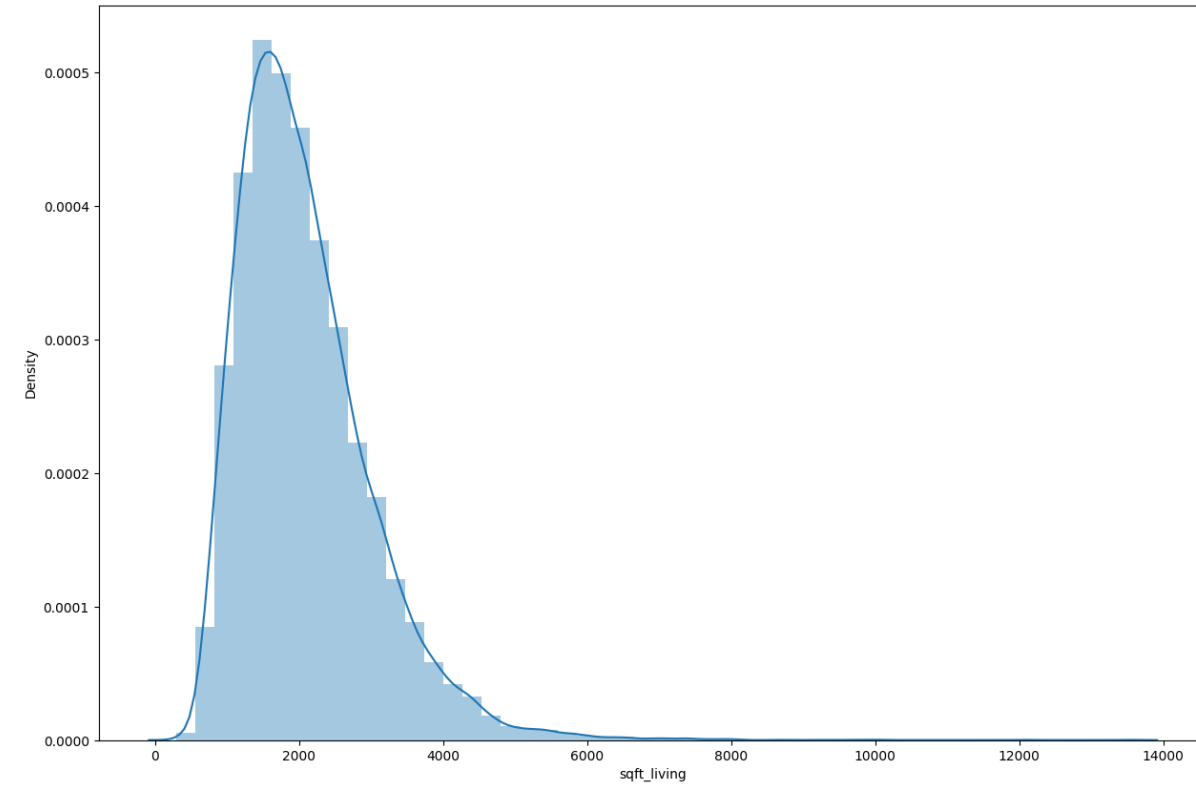
THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY



SCHOOL OF
ENGINEERING



House prices by sqft_living



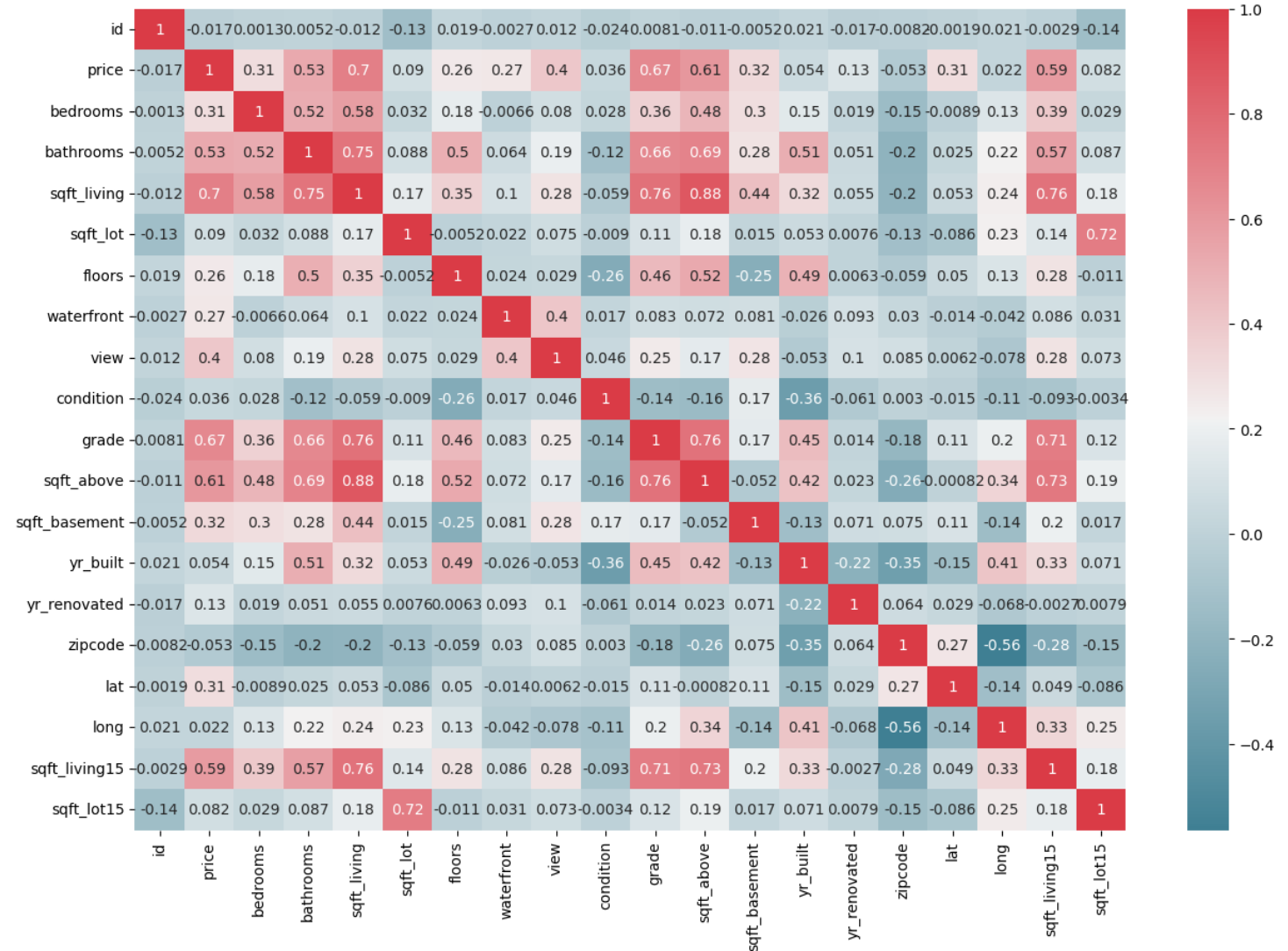
Density by sqft_living

Step 2: Statistics and graphing



What is this?

A house price correlation heatmap shows the **relationships between different factors** that can affect house prices, such as location, size, number of bedrooms, or amenities. It uses colors to indicate the strength and direction of these relationships: for example, **dark colors** might show a **strong positive correlation**, while **lighter colors** indicate **weak or no correlation**. This helps identify which factors are most related to house prices.



Step 2: Statistics and graphing

What is this?

1. Variables Analyzed:

Common variables include square footage, number of bedrooms, bathrooms, age of the home, location, and features like pools or garages.

2. Correlation Coefficient:

The heatmap uses a correlation coefficient (usually between -1 and 1) to show relationships:

1 means a perfect positive correlation (as one increases, the other does too).

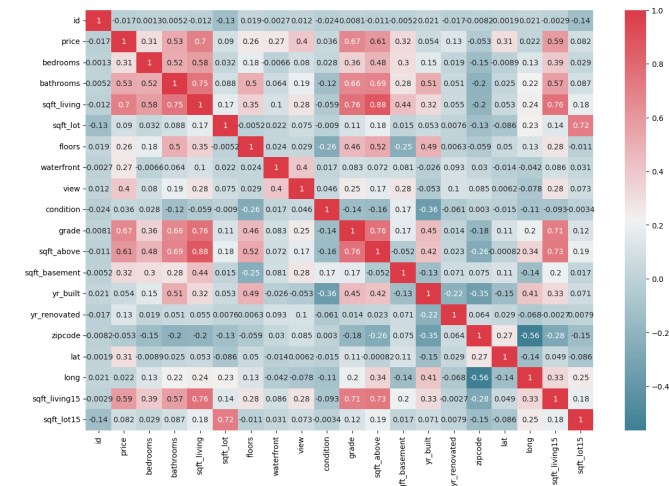
-1 means a perfect negative correlation (as one increases, the other decreases).

0 means no correlation.

3. Color Coding:

Different colors represent different levels of correlation. For example:

- Darker colors** may indicate stronger correlations.
- Lighter colors** may indicate weaker correlations.



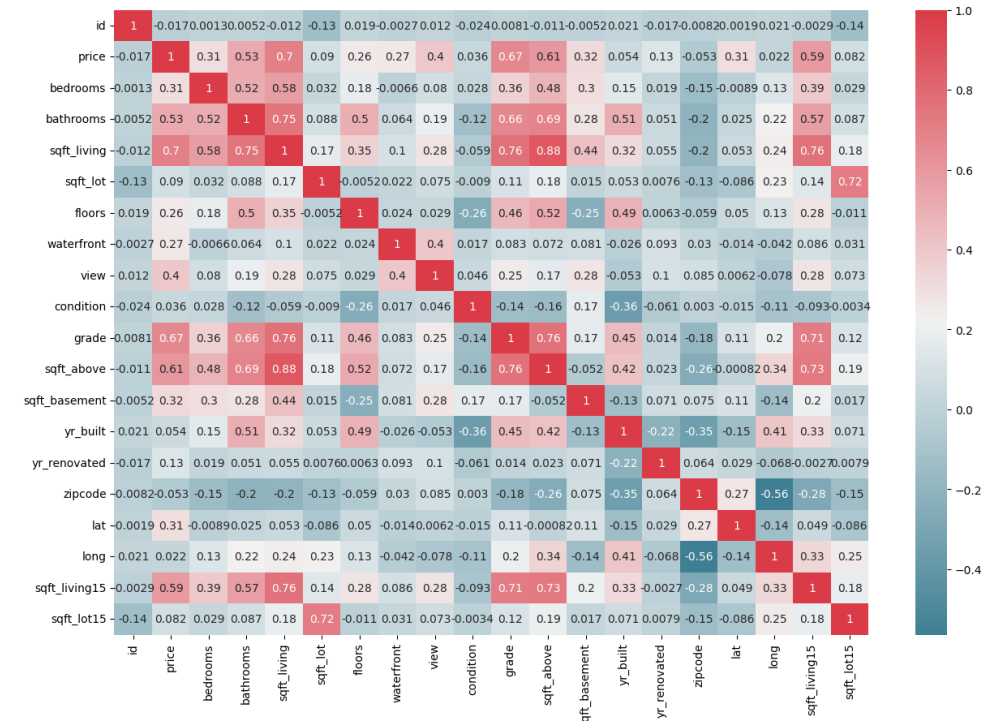
Step 2: Statistics and graphing



What can it do?

It helps real estate professionals and buyers understand which features most influence house prices. Buyers can use this information to make informed decisions about what features to prioritize. Sellers can understand what aspects of their home might increase its value.

For example, if the heatmap shows a strong correlation between square footage and price, it suggests larger homes tend to sell for more.



Step 3: Simple Linear Regression

How can we build this?



1. Collect Data: (done)

Gather data that includes the independent variable and the dependent variable. For example, as shown on the right, you might have living space(X) and house price (Y).

2. Visualize the Data: (done)

Examining data features, counting the density distribution of each feature, plotting statistics between two features to determine if there is a linear relationship, and creating heat maps to determine the relationship between different features.

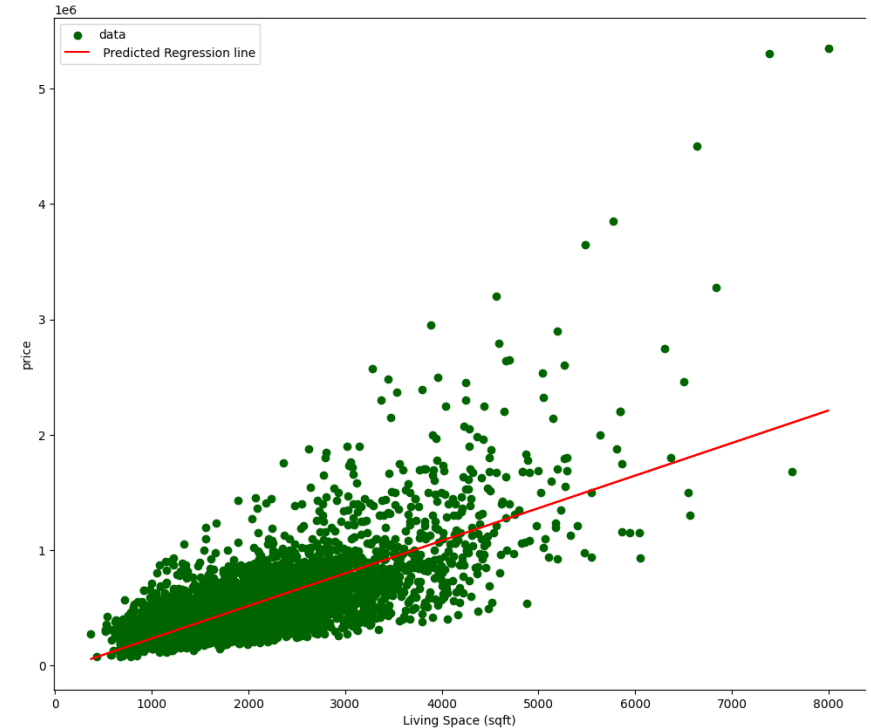
3. Prepare the Data: (done)

Clean the data by handling missing values and outliers if necessary.

4. Split the Data:

Divide your data into training and testing sets (e.g., 80% for training and 20% for testing).

```
train_data,test_data=train_test_split(df,train_size=0.8,random_state=3)
```



Step 3: Simple Linear Regression



How can we build this?

5. Fit the Model:

Use statistical software or programming language (like Python or R) to fit the linear regression model.

The formula for simple linear regression is: $[Y = b_0 + b_1X]$ Where:

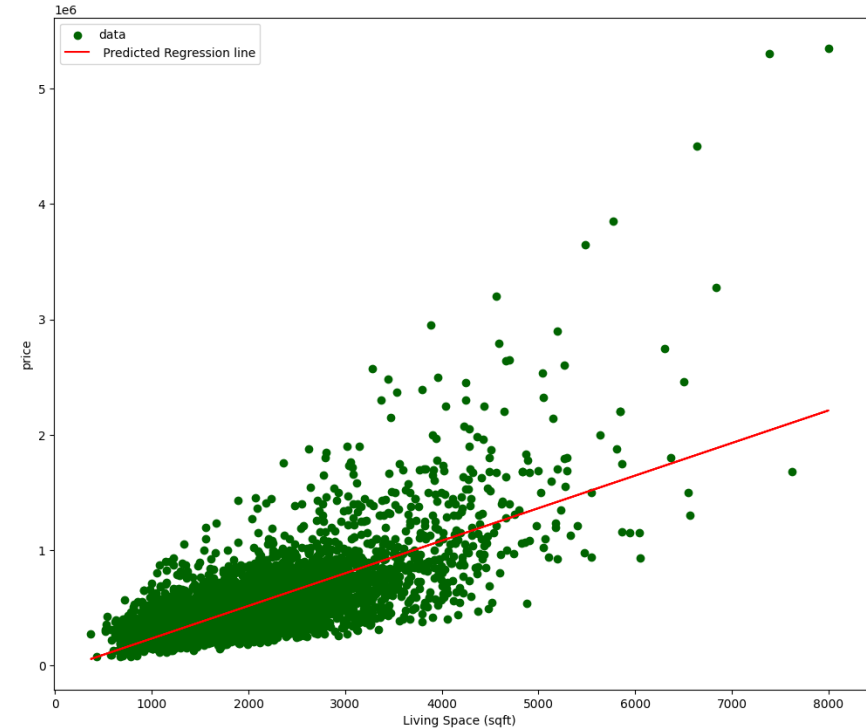
1. (Y) is the predicted value (house price).
2. (b_0) is the y-intercept.
3. (b_1) is the slope of the line (how much Y changes for a unit change in X).
4. (X) is the independent variable (square footage).

6. Evaluate the Model:

Check the model's performance using R-squared, Mean Absolute Error (MAE), or Mean Squared Error (MSE). In our case, we use **R-squared** to evaluate the model.

The R-squared value represents the proportion of the variance in the target variable explained by the model's independent variables.

The R-squared value can range from **0** (can not explain any variance) to **1** (perfectly fitting), with higher values indicating better performance.



7. Make Predictions:

Use the model to make predictions on new data.

Assignment requirement

Q1 (**1 point**): count the number of occurrences of each unique value in "condition"

Q2 (**2 points**): please draw a **bar plot** of 'house prices by sqft_above' and a **density plot** of sqft_above

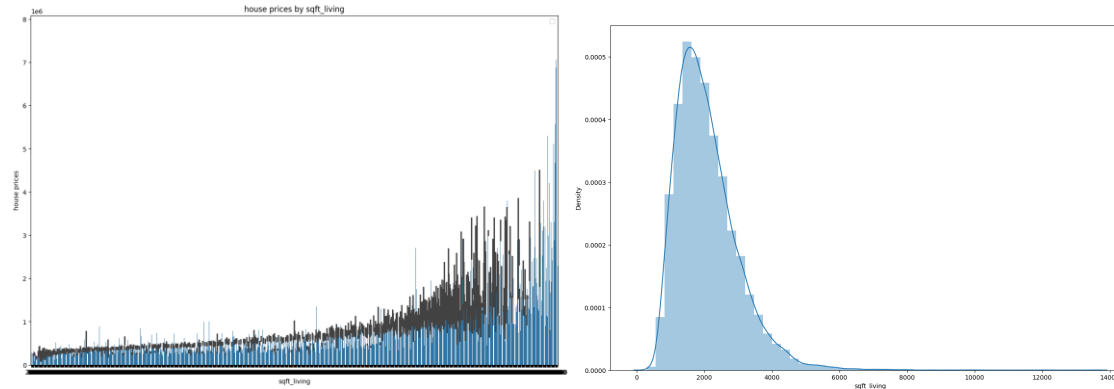
Q3 (**2 points**): please draw a Simple Linear Regression plot of 'house prices by sqft_above' and a Simple Linear Regression plot of 'house prices by bathrooms'

```
Out[8]:
```

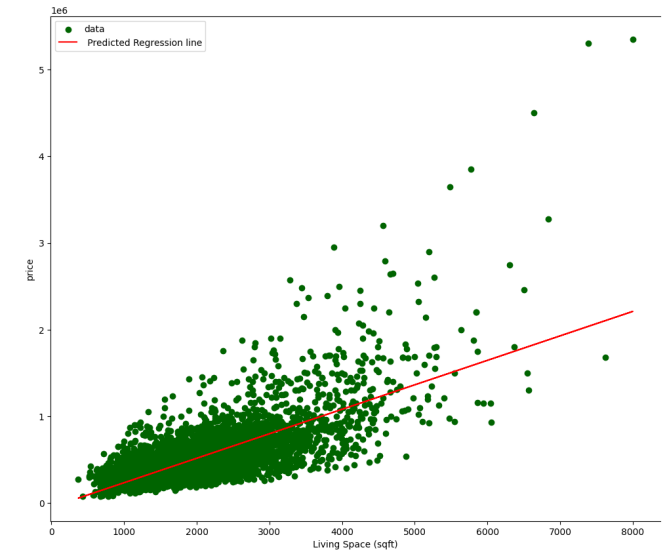
	count
bedrooms	
3	9824
4	6882
2	2760
5	1601
6	272
1	199
7	38
0	13
8	13
9	6
10	3
11	1
33	1

dtype: int64

Q1 sample



Q2 sample



Q3 sample

Assignment requirement



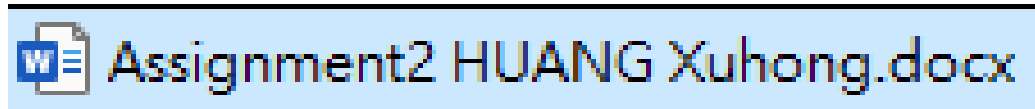
THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY



SCHOOL OF
ENGINEERING

What you need to submit to Canvas is a **PDF file** named **“Assignment 2 + your name”**.

名稱	修改日期	類型	大小
3.ipynb	28/9/2024 13:06	Jupyter 來源檔案	789 KB
3_full.ipynb	28/9/2024 12:59	Jupyter 來源檔案	901 KB
3_old.ipynb	18/9/2023 6:08	Jupyter 來源檔案	883 KB
Guidebook_Housing_Price_Regression...	28/9/2024 13:13	Microsoft Power...	3,626 KB
Housing_price_prediction_update.ipynb	28/9/2024 12:50	Jupyter 來源檔案	651 KB
kc_house_data.csv	18/9/2023 6:08	Microsoft Excel ...	2,457 KB
W6_Assignment_requirement.docx	28/9/2024 13:23	Microsoft Word ...	123 KB
W6_Assignment_temp.docx	28/9/2024 13:30	Microsoft Word ...	17 KB
Assignment2 HUANG Xuhong.docx	28/9/2024 13:30	Microsoft Word ...	17 KB



Assignment 2:

Q1 (1 point): count the number of occurrences of each unique value in "condition."

Your code:

(Copy your core code here)

Your result:

(Screenshot your results here)

Q2 (2 points): please draw a bar plot of 'house prices by sqft_above' and a density plot of sqft_above.

Your code:

(Copy your core code here)

Your result:

(Screenshot your results here)

Q3 (2 points): please draw a Simple Linear Regression plot of 'house prices by sqft_above' and a Simple Linear Regression plot of 'house prices by bathrooms.'

Your code:

(Copy your core code here)

Your result:

(Screenshot your results here)

Assignment 2 template