

Data Collection

1. Include population
 - 1.1 Request the API from coordinates
 - 1.2 Request related figures as json
 - 1.3 Unpack json and append it as csv format
2. Include deprivation index
 - 2.1 Download csv file from Otago
 - 2.2 Using 'SA1' column as key to acquire corresponding value from the download csv file

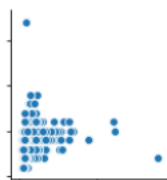
Data Processing

1. Handling NaN issues

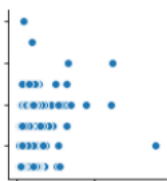
NaN data existed in "Bathrooms" and "Suburbs" column. As "Suburbs" column will not be used later in the processing, only NaN existed in "Bathrooms" column need to be handled. And I choose to remove them as there's 2 data points, and it will not be reasonable to give any estimate on this crucial data.

2. Redundancy data / Noise

Using the histogram, the noise is quite obvious.

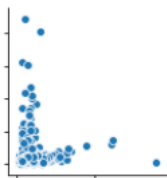


The diagram on the right is graph for the CV-Bedrooms. It is clear that the left top data is a noise. And should be eliminated before moving into analysis.



Similarly, the CV-Bathrooms graph demonstrated the unlikelihood that having more than 6 bathrooms.

The CV-Land area figure also indicates 15000 m² is probably not a regular range that we should estimate.



The similarity of the three graph above is the three data points having CV price more than 10M. This should also be taken out from our analysis as it will create great influence because its excessive numerical value.

3. Correlation Analysis

Bedrooms	1	0.71	0.084	0.22	-0.0071	0.0092	0.032	0.015	-0.051	-0.062	0.029	0.049	-0.027	-0.022	-0.21
Bathrooms	0.71	1	0.05	0.38	0.076	0.045	-0.046	-0.068	-0.061	-0.053	0.025	0.084	0.037	-0.027	-0.33
Land area	0.084	0.05	1	0.11	0.14	-0.029	-0.11	-0.059	-0.11	-0.14	-0.011	0.12	0.05	-0.054	-0.12
CV	0.22	0.38	0.11	1	0.12	0.018	-0.11	-0.16	-0.18	-0.21	-0.045	0.13	0.083	-0.13	-0.38
Latitude	-0.0071	0.076	0.14	0.12	1	-0.31	-0.87	-0.27	-0.16	-0.17	-0.066	-0.014	0.081	-0.16	-0.21
Longitude	0.0092	0.045	-0.029	0.018	-0.31	1	0.53	0.09	-0.041	-0.018	0.023	0.12	0.095	0.063	0.018
SA1	0.032	-0.046	-0.11	-0.11	-0.87	0.53	1	0.3	0.15	0.15	0.065	0.063	-0.0057	0.19	0.21
0-19 years	0.015	-0.068	-0.059	-0.16	-0.27	0.09	0.3	1	0.43	0.6	0.7	0.45	0.075	0.81	0.17
20-29 years	-0.051	-0.061	-0.11	-0.18	-0.16	-0.041	0.15	0.43	1	0.76	0.29	0.16	-0.049	0.68	0.25
30-39 years	-0.062	-0.053	-0.14	-0.21	-0.17	-0.018	0.15	0.6	0.76	1	0.54	0.24	0.051	0.81	0.13
40-49 years	0.029	0.025	-0.011	-0.045	-0.066	0.023	0.065	0.7	0.29	0.54	1	0.53	0.21	0.75	-0.22
50-59 years	0.049	0.084	0.12	0.13	-0.014	0.12	0.063	0.45	0.16	0.24	0.53	1	0.36	0.58	-0.28
60+ years	-0.027	0.037	0.05	0.083	0.081	0.095	-0.0057	0.075	-0.049	0.051	0.21	0.36	1	0.41	-0.19
C18_CURPop	-0.022	-0.027	-0.054	-0.13	-0.16	0.063	0.19	0.81	0.68	0.81	0.75	0.58	0.41	1	0.038
DeplIndex	-0.21	-0.33	-0.12	-0.38	-0.21	0.018	0.21	0.17	0.25	0.13	-0.22	-0.28	-0.19	0.038	1
	Bedrooms	Bathrooms	Land area	CV	Latitude	Longitude	SA1	0-19 years	20-29 years	30-39 years	40-49 years	50-59 years	60+ years	C18_CURPop	DeplIndex

By analysing the correlation between CV and other parameter, the three most related figures are number of bedroom, bathrooms and deprivation index.

Data Analysis

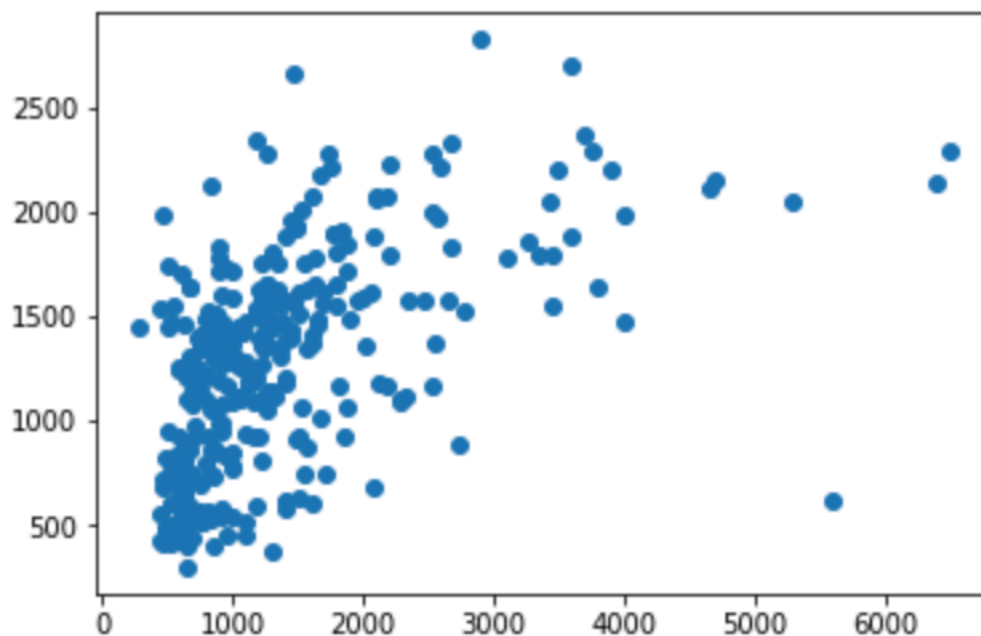
Linear Regression

Linear Regression isn't the best fit model in the situation. As it can hardly justify the impact of area code, such as "SA1", and how longitude and Latitude contributed to the final output. The location has strong impact on people decision to purchase houses but it can hardly to be linearly quantify.

The coefficient of show a different indication as what the heat map indicates.

```
1.18475309e+04, 2.80552411e+05, -7.31442092e+00, -5.96570762e+04,  
3.88543570e+05, -3.43679564e+01, 3.12183490e+03, 3.94489978e+03,  
-6.72674231e+03, -6.20698879e+03, 1.30256811e+04, 2.68116981e+03,  
-2.11761880e+03, -9.44175305e+04 ] )
```

The scatter plot show some prediction efficiency between 1000 – 3000 K, but miss out a lot more point in price higher than 3000K. The model generally underestimate the housing price.

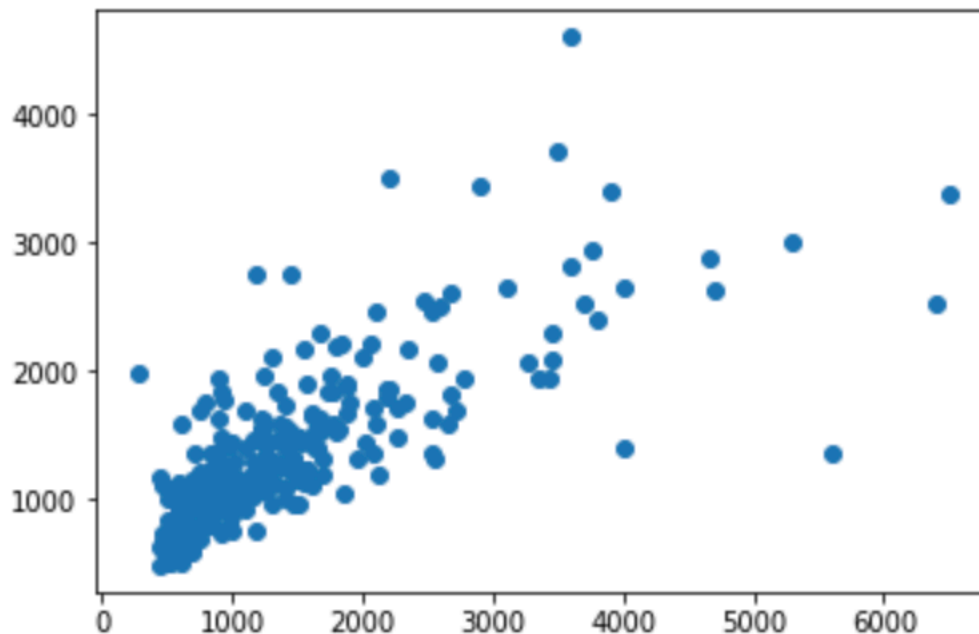


```
linear_model.score(test_x,test_y)
```

```
0.28981868462074956
```

Random Forest Regression

Random Forest Regression will address the problem listing above. Because it has decision trees to help analyse the influence of area code. And the random forest Regression allow a small sample splits, thus, it can better handle extreme training case compare to Linear Regression.



The scatter plot show a more concentration display.

```
RFR_model.score(test_x,test_y)
```

0.5629860202264829

And it give a more accurate estimate of the housing price.