

# Australian Coronavirus (COVID-19) Analyst

## Introduction

This report is to analyze Australian COVID19 data for general public

Findings of the report is as below

- Most of cases are from NSW and VIC state
- Population, and distance to top affected areas are the major cause contributing to cases
- Lockdown and vaccination is helpful to stop Delta

## Step1: Download data from Elephant DB

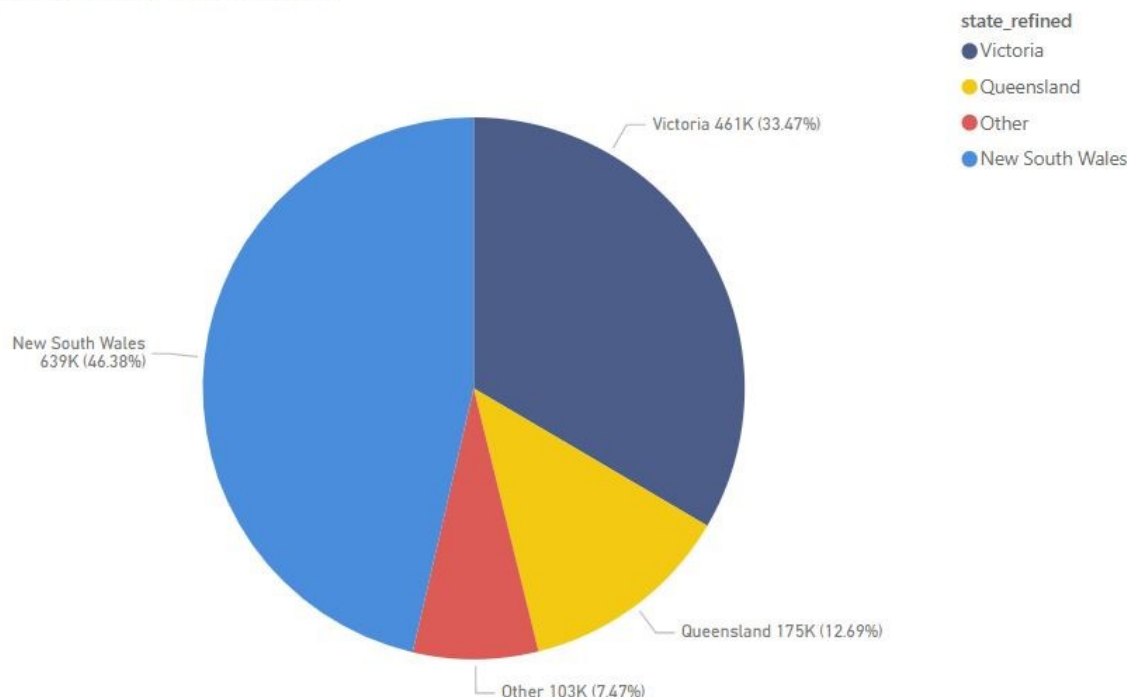
```
postgres://mheodhlf:m6FmQMj_66D6CO7BPmOZAfcUG2-La9Tv@rosie.db.elephantsql.com/mheodhlf
```

## Step2: Visualize the data, and conduct data analysis

Data in states: it finds that NSW and VIC has most of cases in Australia

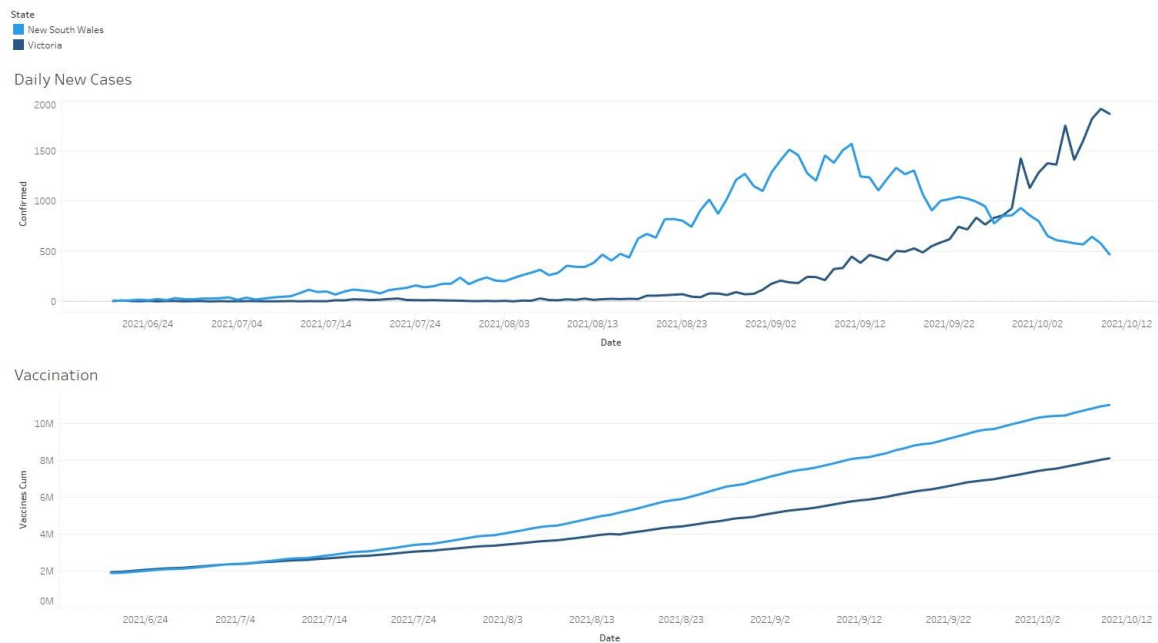
- NSW: 46.3%, VIC: 33.4%
- QLD and the other states are about 20% in total

Total confirmed cases over time



Hyphothesis 1: vaccination is related with covid cases during lockdown

- Check data for NSW and VIC during the Lockdown for Delta (2021-06-20 to 2021-10-10)



OLS Analysis for NSW vaccines and cases

R-squared is 0.6, this indicates strong correlation between vaccinations and cases in NSW

#### OLS Regression Results

<b>Dep. Variable:</b>	confirmed	<b>R-squared:</b>	0.609
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.606
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	173.2
<b>Date:</b>	Sun, 23 Jan 2022	<b>Prob (F-statistic):</b>	2.09e-24
<b>Time:</b>	23:41:15	<b>Log-Likelihood:</b>	-805.29
<b>No. Observations:</b>	113	<b>AIC:</b>	1615.
<b>Df Residuals:</b>	111	<b>BIC:</b>	1620.
<b>Df Model:</b>	1		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>Intercept</b>	-188.7203	63.539	-2.970	0.004	-314.627	-62.814
<b>vaccines_cum</b>	0.0001	1e-05	13.162	0.000	0.000	0.000

<b>Omnibus:</b>	3.208	<b>Durbin-Watson:</b>	0.077
<b>Prob(Omnibus):</b>	0.201	<b>Jarque-Bera (JB):</b>	2.664
<b>Skew:</b>	0.269	<b>Prob(JB):</b>	0.264
<b>Kurtosis:</b>	3.526	<b>Cond. No.</b>	1.41e+07

Notes:

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

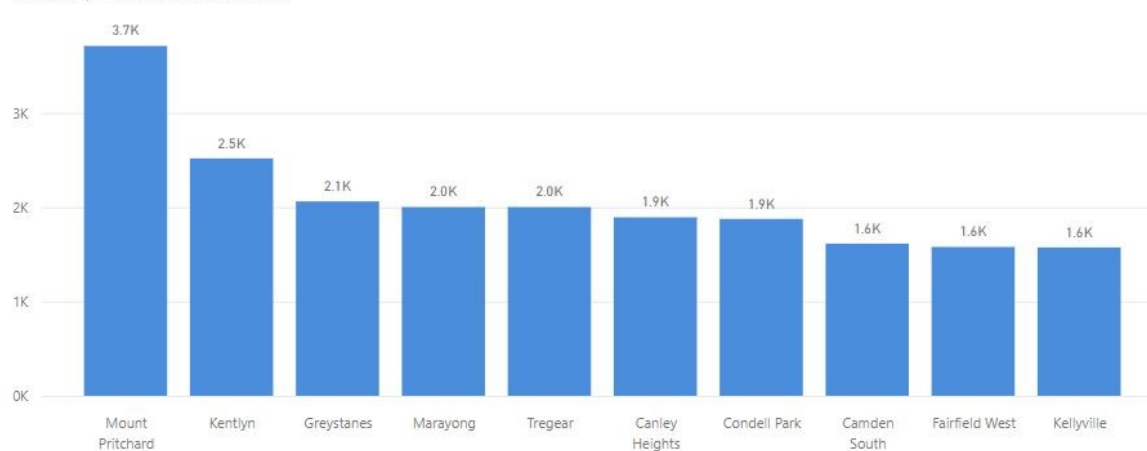
[2] The condition number is large, 1.41e+07. This might indicate that there are

strong multicollinearity or other numerical problems.

Hypothesis 2: distance to CBD or top1 areas is related with covid cases

- Check data for NSW and VIC, which areas were most of cases

NSW top10 areas in this week

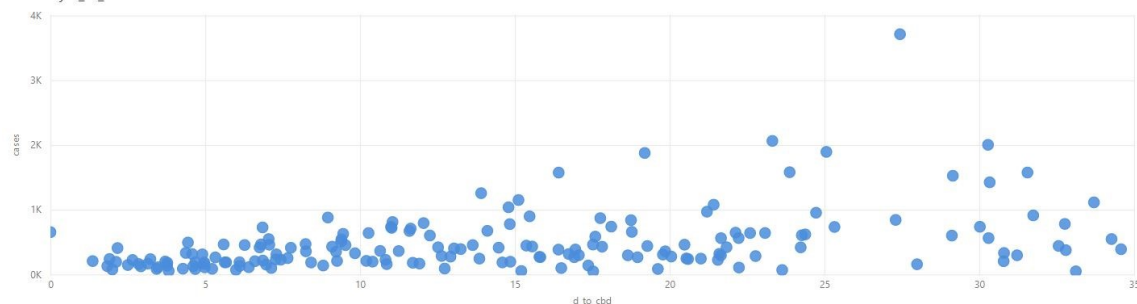


VIC top10 areas in this week

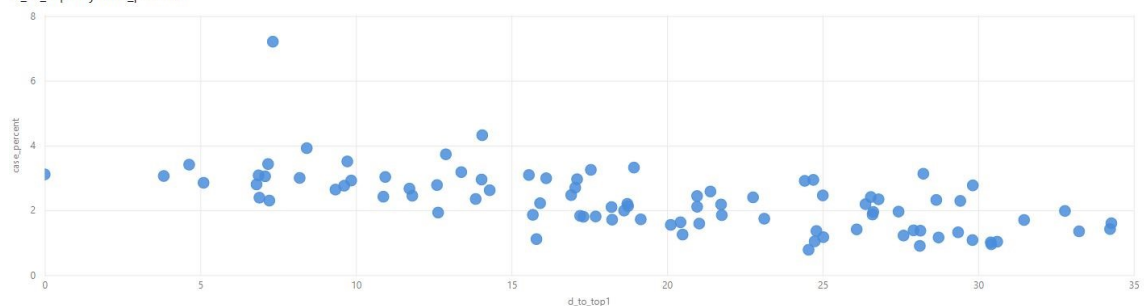


- Calculate distance of each suburb to top\_one suburb, and to CBD -- NSW CBD: 2000, VIC CBD: 3000 -- NSW top1: 2170, VIC top1: 3029

cases by d\_to\_cbd



d\_to\_top1 by case\_percent



OLS Analysis for NSW distance and cases

R-squared is 0.45, this indicates a correlation between distance and cases

#### OLS Regression Results

<b>Dep. Variable:</b>	case_percent	<b>R-squared:</b>	0.451
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.448
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	142.9
<b>Date:</b>	Sun, 23 Jan 2022	<b>Prob (F-statistic):</b>	1.96e-24
<b>Time:</b>	23:41:16	<b>Log-Likelihood:</b>	-160.30
<b>No. Observations:</b>	176	<b>AIC:</b>	324.6
<b>Df Residuals:</b>	174	<b>BIC:</b>	330.9
<b>Df Model:</b>	1		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>Intercept</b>	3.4339	0.130	26.468	0.000	3.178	3.690
<b>d_to_top1</b>	-0.0691	0.006	-11.956	0.000	-0.081	-0.058

<b>Omnibus:</b>	122.269	<b>Durbin-Watson:</b>	1.450
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	1446.254
<b>Skew:</b>	2.391	<b>Prob(JB):</b>	0.00
<b>Kurtosis:</b>	16.204	<b>Cond. No.</b>	64.0

Notes:

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

Hyphothesis 3: Population is related with covid cases

- Check data for NSW and VIC, population and cases of each suburb



OLS Analysis for NSW distance and cases

R-squared is 0.82, this indicates strong correlation between population and cases

#### OLS Regression Results

<b>Dep. Variable:</b>	cases	<b>R-squared:</b>	0.829
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.828
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	843.7
<b>Date:</b>	Sun, 23 Jan 2022	<b>Prob (F-statistic):</b>	1.22e-68
<b>Time:</b>	23:41:16	<b>Log-Likelihood:</b>	-1173.7
<b>No. Observations:</b>	176	<b>AIC:</b>	2351.
<b>Df Residuals:</b>	174	<b>BIC:</b>	2358.
<b>Df Model:</b>	1		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>Intercept</b>	-140.3022	25.697	-5.460	0.000	-191.021	-89.584
<b>population</b>	0.0275	0.001	29.047	0.000	0.026	0.029

<b>Omnibus:</b>	16.303	<b>Durbin-Watson:</b>	0.995
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	32.526
<b>Skew:</b>	0.416	<b>Prob(JB):</b>	8.65e-08
<b>Kurtosis:</b>	4.935	<b>Cond. No.</b>	4.82e+04

Notes:

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

[2] The condition number is large, 4.82e+04. This might indicate that there are

strong multicollinearity or other numerical problems.

Hypothesis 4: Lockdown is correlated with cases in NSW

AS Covid-19 is becoming a 'New Normal' of our life, this visualisation aims to demonstrate how the currently active cases are located in NSW. Do people who live close to CBD with more people living around have a higher risk of catching covid-19 today? To answer this question, we will compare the overall performance of each suburb during and after the lockdown. The answer is no during the lockdown and yes when all the restrictions have been eased.

Load data from elephant db to local.

Get the data for the covid pandemic in NSW during the lockdown in 2021 from June 24 to December 15. Then rank the covid cases per 100 population in each postcode from 0 to 1 by summing up each of their ranks. We can estimate their overall ranking during the lockdown.

	postcode	lat	lng	suburb	population	total_rank
0	2338	-31.72628	150.79265	Ardglen	1428	0.0

Calculate the distance from each suburb to Sydney CBD.

Summing up the monthly ranking during the lockdown, each suburb received a mark to rank their overall performance during the lockdown, and lower is better. Below are the 20 worst performed suburbs.

	postcode	suburb	population	total_rank	distance_cbd
612	2168	Busby	43449	6.158283	30.755024
611	2191	Belfield	6322	5.987127	12.120299
610	2190	Mount Lewis	25568	5.984774	15.672347
609	2192	Belmore	12718	5.922951	12.651825
608	2174	Edmondson Park	2271	5.855576	34.803631

5 best performed suburbs within 50km from CBD.

	postcode	suburb	population	total_rank	distance_cbd
40	2555	Badgerys Creek	225	0.000000	43.033619
143	2083	Bar Point	1524	0.358779	41.236992
149	2071	Killara	13552	0.454194	12.005354
167	2082	Berowra Waters	5402	0.617326	30.533327
171	2072	Gordon	7668	0.630901	13.270306

R-squared represents how good distance to CBD and population can explain the performance of epidemic prevention during the lockdown. In this case, the correlation is not strong enough. This means, during NSW 2021 lockdown, the distance from a suburb to Sydney CBD and its population doesn't strongly affect the level of risk of catching covid-19.

#### OLS Regression Results

<b>Dep. Variable:</b>	total_rank	<b>R-squared:</b>	0.365
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.363
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	175.3
<b>Date:</b>	Sun, 23 Jan 2022	<b>Prob (F-statistic):</b>	6.97e-61
<b>Time:</b>	23:41:23	<b>Log-Likelihood:</b>	-992.04
<b>No. Observations:</b>	613	<b>AIC:</b>	1990.

<b>Df Residuals:</b>	610			<b>BIC:</b>	2003.	
<b>Df Model:</b>	2					
<b>Covariance Type:</b>	nonrobust					
	<b>coef</b>	<b>std err</b>	<b>t</b>	<b>P&gt; t </b>	<b>[0.025</b>	<b>0.975]</b>
<b>Intercept</b>	2.0383	0.097	20.914	0.000	1.847	2.230
<b>distance_cbd</b>	-0.0032	0.000	-12.696	0.000	-0.004	-0.003
<b>population</b>	3.052e-05	3.95e-06	7.733	0.000	2.28e-05	3.83e-05
<b>Omnibus:</b>	50.711	<b>Durbin-Watson:</b>		0.647		
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>		61.482		
<b>Skew:</b>	0.769	<b>Prob(JB):</b>		4.46e-14		
<b>Kurtosis:</b>	3.209	<b>Cond. No.</b>		3.60e+04		

Notes:

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

[2] The condition number is large, 3.6e+04. This might indicate that there are

... ..

Get the data for the covid pandemic in NSW after eased restrictions in 2021 from DEC 15 to current. We can estimate their overall ranking during the lockdown by summing up their relative position each day.

	postcode	lat	lng	suburb	population	total_rank
0	2611	-35.66563	148.70878	Cooleman	52	0.0

Add distance to CBD into the table

Like the calculation above, we summed up the daily ranking for each suburb after the lockdown. Below are the 20 worst performed suburbs in NSW.

	postcode	suburb	population	total_rank	distance_cbd
608	2026	North Bondi	32488	32.383562	6.507199
609	2762	Schofields	4983	32.433877	35.942849
610	2020	Mascot	14772	32.692533	9.076371
611	2481	Broken Head	11772	33.570675	614.497119
612	2174	Edmondson Park	2271	33.754071	34.803631

Below are the 5 best performed suburbs in NSW after restriction eased.

	postcode	suburb	population	total_rank	distance_cbd
0	2611	Cooleman	52	0.000000	303.551512
1	2649	Nurenmermong	31	0.000000	342.085900
2	2668	Barmedman	459	0.000000	349.789948
3	2356	Gwabegar	162	0.004098	421.904178
4	2735	Koraleigh	451	0.007590	725.843478

5 best performed suburbs in NSW within 50km from Sydney CBD.

	postcode	suburb	population	total_rank	distance_cbd
223	2563	Menangle Park	257	5.675193	49.109234
228	2083	Bar Point	1524	5.971365	41.236992
242	2082	Berowra Waters	5402	6.651336	30.533327

263	2080	Mount Kuring-Gai	1708	7.804422	25.830430
267	2105	Leppington	1854	7.822197	26.821494

5 worst performed suburbs in NSW within 50km from Sydney CBD.

	postcode	suburb	population	total_rank	distance_cbd
607	2179	Leppington	6522	32.001325	39.792150
608	2026	North Bondi	32488	32.383562	6.507199
609	2762	Schofields	4983	32.433877	35.942849
610	2020	Mascot	14772	32.692533	9.076371
612	2174	Edmondson Park	2271	33.754071	34.803631

The correlation between suburb performance and CBD distance and suburb population is stronger after restriction eased. After NSW 2021 lockdown, the distance from a suburb to Sydney CBD and the number of people to an extent affects the likelihood of catching covid-19.

#### OLS Regression Results

<b>Dep. Variable:</b>	total_rank	<b>R-squared:</b>	0.525
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.523
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	337.1
<b>Date:</b>	Sun, 23 Jan 2022	<b>Prob (F-statistic):</b>	2.48e-99
<b>Time:</b>	23:41:24	<b>Log-Likelihood:</b>	-2059.7
<b>No. Observations:</b>	613	<b>AIC:</b>	4125.
<b>Df Residuals:</b>	610	<b>BIC:</b>	4139.
<b>Df Model:</b>	2		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>Intercept</b>	16.1465	0.556	29.029	0.000	15.054	17.239
<b>distance_cbd</b>	-0.0281	0.001	-19.331	0.000	-0.031	-0.025
<b>population</b>	0.0002	2.25e-05	8.436	0.000	0.000	0.000

<b>Omnibus:</b>	44.868	<b>Durbin-Watson:</b>	0.857
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	54.062
<b>Skew:</b>	0.654	<b>Prob(JB):</b>	1.82e-12
<b>Kurtosis:</b>	3.636	<b>Cond. No.</b>	3.60e+04

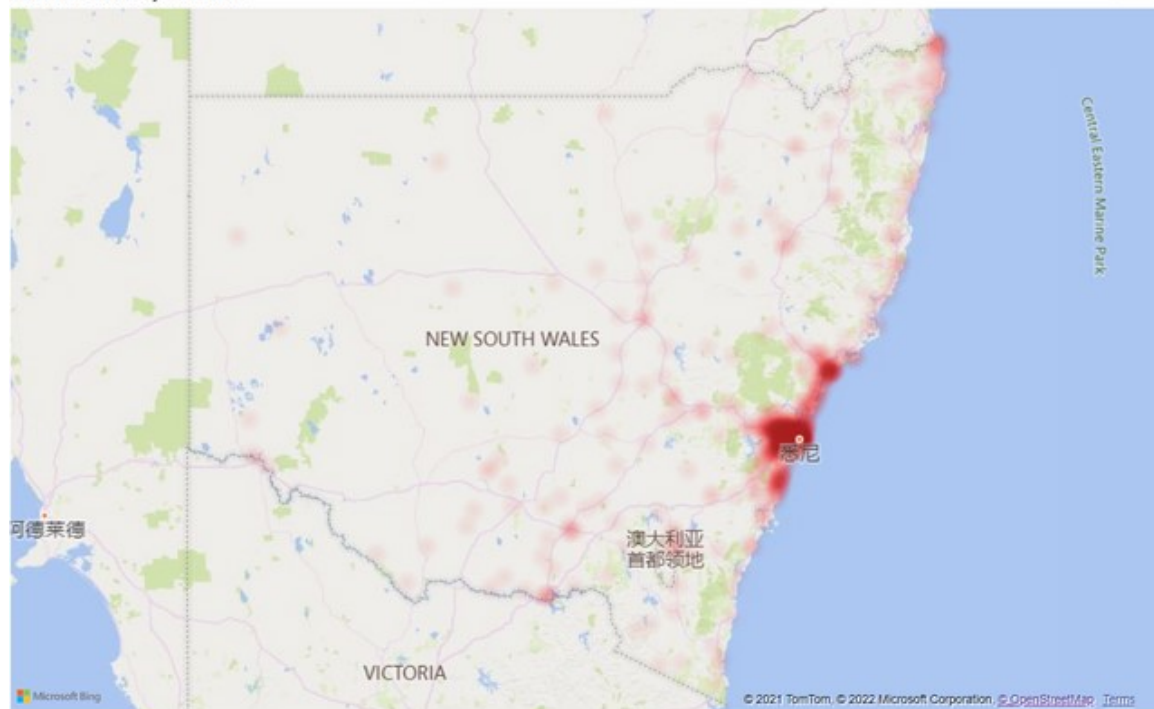
#### Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 3.6e+04. This might indicate that there are strong multicollinearity or other numerical problems.

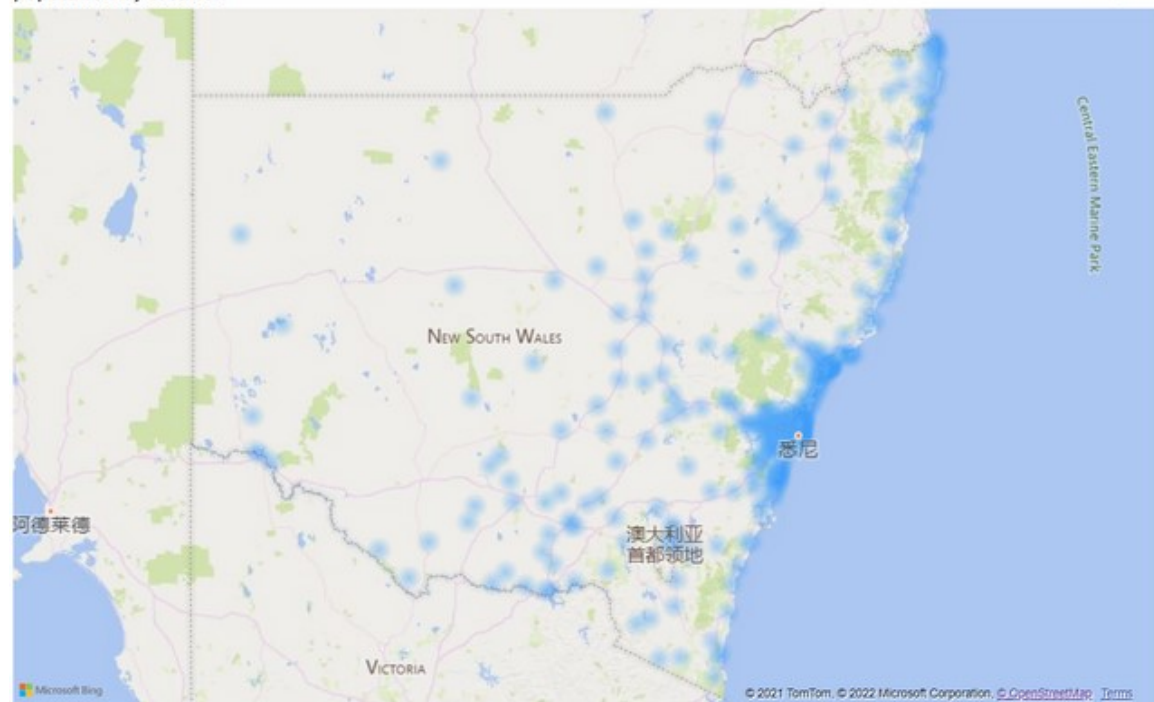
Converting the active case and population into a heat map in PowerBI allows us to visualise this correlation.



Active cases by location



population by location



Conclusion: Total cases shows a strong correlation with population, and little correlation with distance to Mount Partchard, Roxbargh Park or CBD

- If your living area has a population, and close to above areas, better to stay at home to keep safe
- New Cases in NSW and VIC continues and doesnt show decreasing trend.
- After comparing OLS model between Lockdown and Eased, the effectiveness of lockdown is proved.