

Is there a correlation between housing prices and living standards in different parts of Victoria?

Authors:

Un Leng Kam 1178863

Ron (Lang) Chen 1181506

Sean (Aoxiang) Xiao 1174270

This project investigated *whether house prices in Victoria are correlated to living standards*. Specifically, we looked at the relationship between the Median House Price (MHP) and living standard using data from each of the 79 LGAs of Victoria. As 'living standard' is a difficult index to measure, with its definition and method of calculation subjective and varying significantly between different sources, one challenge in this project was deriving a function for calculating the living standard that is as objective as possible given the data we could source. We defined living standards by a liveability score: a weighted sum of five different attributes found in each area; its derivation is discussed and justified later. The outcome of our analysis was a scatter plot of MHP to liveability, with linear regression performed as we assumed for such an important investment one would theoretically only pay more for a linear return for their money. Analysis of the regression returned an r^2 value of 0.232, which suggests the existence of a positive correlation between MHP and liveability.

Our research question is related to the themes of liveability and sustainability of communities in Victoria, where communities are represented by LGA areas. Through investigating the correlation between housing prices and living standards, we hypothesised that our results would give insights on the liveability in areas with certain MHPs, and by including open spaces in our living standard calculation, sustainability was also explored - open space area is a typically accepted indicator of how sustainable an area is.

Six datasets in the form of csv and .xlsx were used in this project, sourced from multiple databases (Aurin, Discover Data Victoria¹), of which we extracted the MHP, Crime Rate per 100,000 (CR), Tally of number of Facilities (FT), total Open Area (OA), average ambulance Response Times (RT) and Tally of number of Schools (ST) in each LGA. These five attributes were chosen because we believed they are among the fundamental aspects that homebuyers consider before purchasing a property, and hence the more/better they are in a certain LGA the more demand there would be for a local property and thus the higher its MHP. In our selection process data such as tally of number of open spaces were discarded because they would likely be highly correlated with OA and thus double weigh the impact of sustainability when liveability was calculated.

Major python modules used in the data wrangling and analysis process includes 'Pandas', 're', 'matplotlib', 'Statistics', 'sklearn' and 'NumPy' as well as native datatypes and functions. Excel was used to inspect raw data.

The first wrangling step was creating a csv to map all suburbs to their LGA, and another with an alphabetical order of LGAs so all tallies initiated by it would maintain a unified order that simplified aggregation later. In Victoria, one LGA may consist multiple suburbs, and as most raw datasets came with only LGA or suburbs, the first csv was created to allow them to be tallied by LGA. It is because of this one-to-many relationship that made suburbs an undesirable choice for partitioning Victoria compared to LGAs. Both were created by wrangling and aggregating the 'School.csv' dataset, which for each school provided its suburb, postcode, LGA and LGA code. 're' was used to clean up names (in 'School.csv' all LGA names had a '(' after the name which denoted its type) and for a special case the hyphen in 'Otway-Colac' was removed because some data presented 'Otway Colac', and deleting from the dictionary was easier than adding '-' when wrangling the other

¹ See bibliography for full list

data sets.

All other datasets then followed a similar process of importing into a DataFrame and creating a frequency tally of the attribute by LGAs. In some cases, tally inclusion of a row (object) was also based on values in other columns of the raw dataset, because not every single row represented objects that were sought (i.e. in 'VPA_Open_Space.csv' some rows were schools and thus discarded because it would double-weight schools later). In some files the LGA name had noise and were cleaned up by 're'. When wrangling MHP, we found that our csv which mapped suburbs to LGAs did not include all the Victorian suburbs, likely because in 'School.csv' not every suburb had a school. These were filled in manually in Excel from internet searches as no other 'dictionaries' that had the full set of suburbs-to-LGA could be found. In cases where data of multiple suburbs were merged into one LGA, the median of each merging suburb's values were taken rather than the mean of values because some LGAs had up to 30 suburbs so there was a high chance that the mean would be distorted by high outlier values. The MHP and 5 attributes were then aggregated into one DataFrame.

In seeking a method to derive liveability from the five wrangled attributes, we decided to use a linear function to calculate a liveability score:

$$\text{liveability score} = a(x_1) + b(x_2) + c(x_3) + d(x_4) + e(x_5)$$

with the x_i values being the normalised scores of each attribute from each LGA and the coefficients being the weighting constants, which were derived by calculating the Normalised Mutual Information (NMI) between the MHP and each attribute. The module 'sklearn.normalized_mutual_info_score' was employed to give NMI scores:

Attribute	NMI scores
CR	0.9919028123572338
FT	0.969466990039968
GA	0.6447755861424472
RT	0.9857433425019687
ST	0.9072602417922371

hence the liveability equation is:

$$\text{liveability score} = 0.9919028123572338(x_1) + 0.969466990039968(x_2) + 0.6447755861424472(x_3) + 0.9857433425019687(x_4) + 0.9072602417922371(x_5),$$

where x_1, x_2, x_3, x_4, x_5 are the normalised values of CR, FT, GA, RT and ST

This was applied to every LGA to get their liveability score², where attributes were each normalised with respect to their entire column. The FT, GA and ST columns were normalised by $(x - \min)/(\max - \min)$ while CR and RT columns were normalised by $(\max - x)/(\max - \min)$. The latter was done differently because for these attributes lower values represented better performances, and it was essential that post-normalisation all attributes had values where higher numbers represented better performance in order for the liveability equation to be effective. Normalisation was critical because some attributes had values that were overall much higher than others and thus would have undermined the weighting system and hence produced inaccurate liveability scores (note performing NMI on MHP and attribute values gives same results as NMI on MHP and normalised attribute values).

Whilst using Pearson Correlation or r^2 of linear regression on the MHP for each of the five attributes to give the weighting coefficients for the liveability equation was considered, NMI was ultimately chosen because upon inspecting the scatter plots (produced by 'matplotlib.pyplot') of MHP to each attribute, none of the five

² See Appendix A for Table of LGA, MHP and Liveability Score

showed a linear relationship, and thus the precondition for using either Pearson or r^2 were not satisfied. Principle Component Analysis was also considered but we believed the data was too complexly related to be reduced to a one-dimensional liveability score.

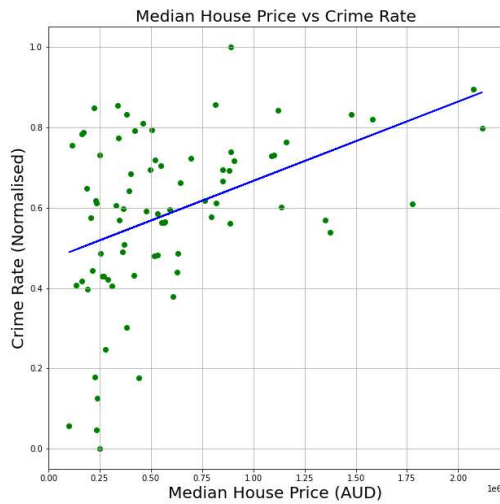


Figure 1 MHP to Normalised CR – looks logarithmic



Figure 2 MHP to Normalised FT – looks logarithmic

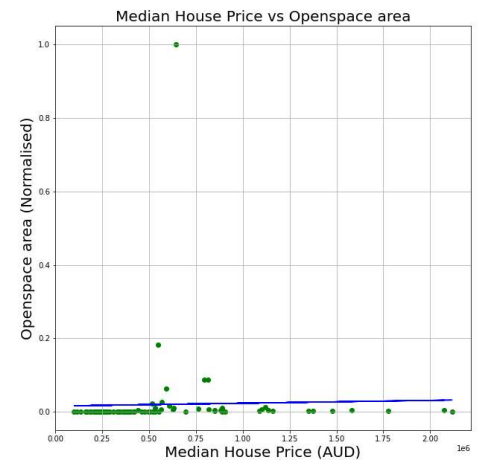


Figure 3 MHP to Normalised OA



Figure 4 MHP to Normalised RT – looks logarithmic

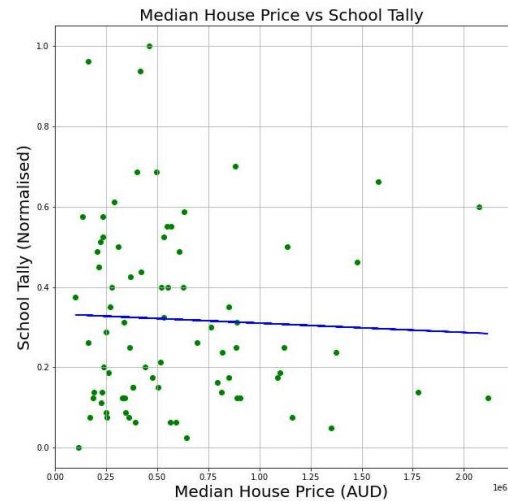


Figure 5 MHP to Normalised ST – looks like hyperbole

Linear regression was performed on the MHP to liveability scatter plot using 'sklearn' and 'NumPy' to give the regression function $y = 0.0000005x + 1.6341$, with r^2 value 0.232.



Figure 6 MHP to Liveability Score (with Linear Regression)



Figure 7 MHP to Liveability Score (With Linear Regression and highlighted outliers)

While 0.232 is typically considered a weak correlation, without domain knowledge (e.g. comparisons to results from this analysis performed on other Australian states) at minimum a positive correlation between MHP and liveability can be concluded; this means generally as house prices goes up the living standard also increases.

The majority of points in Figure 6 are clustered in the lower region of MHPs, with their liveability scores varying between 0.5-2.8 whereas in the higher values (>1,250,000) the points are sparser, varying between 1.7-2.8. The clustering of points in the low MHP region is caused because there are few LGAs in Victoria with high MHPs – most are between 250,000-1,250,000; this is one primary insight from the result. Boroondara is an LGA that fits the trend line well, with MHP 2,075,000 and Liveability Score 2.837; it is observed that it has normalised values over 0.5 for all attributes except 0.0042 for OA³, which is actually rather high for this attribute because OA has a few very high outliers so the majority of normalised scores are below 0.005. This shows that communities with higher MHP have better liveability and sustainability – the latter as the presence of open spaces in the area suggests that over-development is being guarded against by keeping some land ‘construction-less’. Ballarat, Greater Geelong and Yarra Ranges were outliers according to the upper/lower fence outlier analysis performed on residuals. The reason for these outliers is because they have at least two high weighting attributes with very high normalised scores -Ballarat had RT and ST scores at 0.93 and 0.94 which were weighted 0.99 and 0.90 respectively. Ballarat’s high score combined with its comparatively low MHP thus made it an outlier. Stonnington has the highest MHP, but a liveability score of 2.1 which is around the median of all scores. This indicates that at higher MHPs more money does not necessarily return proportionate liveability and sustainability: Stonnington’s normalised value for OS is quite low at 0.0014; it also has limited facilities (0.29) and schools (0.125). This may also suggest that the MHP-liveability relationship may not be linear for all MHP ranges.

Possible beneficiaries of the result and insights are the government and real estate investors. The latter could use it to assist their analysis of whether a property’s value has potential for sharp growth because of the planned improvements related to our five liveability attributes, while the government could use it to understand what the people value having near their homes and improve attributes in areas that are currently under-performing to increase the overall welfare of society. The low OA weighting compared to other attributes may also suggest to the government that sustainability is not a priority to buyers and by inference the public; this may influence their future policymaking. For general homebuyers, the result may help them recognise that the return in liveability for highly priced properties may not be linear.

One limitation of this analysis was we couldn’t find data on other representative attributes such as tally of public transport - including these may have given a better liveability equation as e.g. here buyers would likely favour areas with good access to public transport. Another was that the raw data of OA had ‘missing at random’: many country LGAs had no records (0 on tally)⁴. Mean imputation was not performed as it would have imputed metropolitan data onto country, and the many 0’s likely lowered OA’s weighting. The final model could perhaps be analysed by logarithmic regression, because most points are clustered in the lower MHPs and 16 of them have high liveability scores of >2.25, providing the sharp increase at low values of log functions. Trialling this in Excel returned a higher r^2 at 0.3076. This suggests that at higher MHPs, the current liveability equation either hasn’t captured all attributes that have an impact on property value (i.e. distance from Melbourne CBD) or that the relationship was never linear and our initial belief of linear return for money is incorrect. Nonetheless these limitations do not challenge our primary result of the existence of positive correlation.

³ See Appendix B for all normalised attribute values of each LGA

⁴ From inspection and domain knowledge of being a Victorian resident

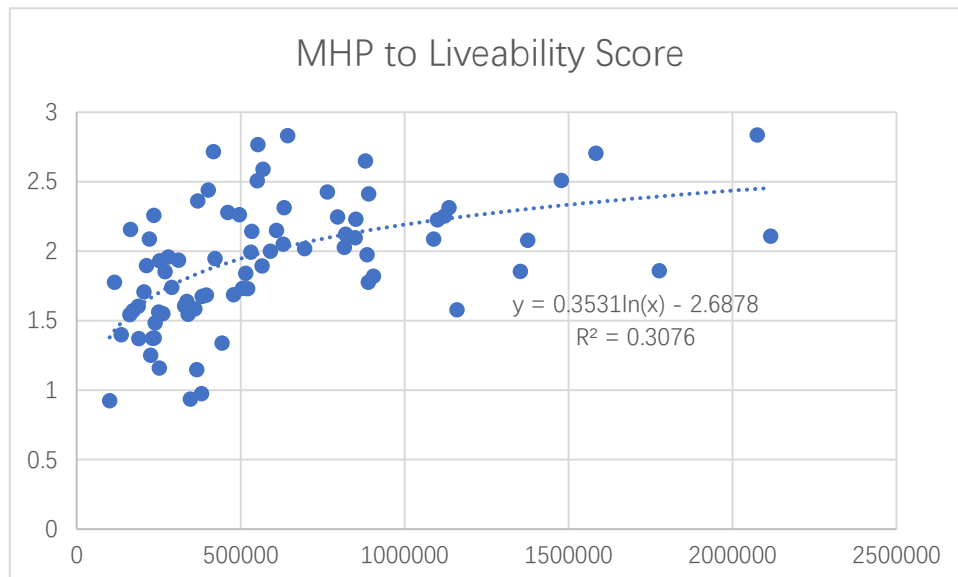


Figure 8 MHP to Liveability Score (Logarithmic Regression) - produced by Excel

An extension/modification to this project is to treat liveability scores as a model and use a method similar to k-fold cross validation but on different combinations of attributes to build and test for the best model of liveability. However, this is more prediction-related and while better serves the purposes of the aforementioned beneficiaries does not answer our research question.

Bibliography

Raw Datasets (italicised filenames as found in the .zip on LMS and on github – may not be same as when first downloaded)

1. Median House Price: discover.data.vic.gov.au/dataset/victoria-property-sales-report-median-by-suburb-time-series1
"Suburb_House_final.csv"
2. Crime Rates: [https://www.crimestatistics.vic.gov.au/family-violence-data-portal/download-data-tables - Table 3](https://www.crimestatistics.vic.gov.au/family-violence-data-portal/download-data-tables-Table-3)
"Ambulance Victoria Data Tables - 2019-20.xlsx"
3. Sporting and Recreational Facilities: Government of Victoria - Department of Environment Land Water and Planning, 2019, VIC DELWP - Vicmap Features of Interest - Sporting Facilities (Polygons), Retrieved from <https://portal.aurin.org.au> on 2021-04-05
"Sport_and_recreation.csv"
4. Open Space: <https://discover.data.vic.gov.au/dataset/open-space>
"VPA_Open_Space.csv"
5. Ambulance Response Time: <https://discover.data.vic.gov.au/dataset/ambulance-victoria-lga-response-time-performance>
"LGA-Response-Time-Performance-FY-2019.csv"
6. Schools: Government of Victoria - Department of Education and Training, 2018, VIC DET - School Locations for Victoria 01/02/2018, Retrieved from <https://portal.aurin.org.au> on 2021-04-05
"School.csv"

Libraries/modules

For the use of libraries and modules please see the README file submitted on LMS or github

Appendix A

<u>LGA Name</u>	<u>Median House Price</u>	<u>Liveability Score</u>
ALPINE	460000	2.279017856451208
ARARAT	163250	2.156227507844175
BALLARAT	416000	2.715337431545911
BANYULE	880500	2.650137192237664
BASS COAST	495500	2.2624904920603903
BAW BAW	401250	2.4392453058358505
BAYSIDE	1582500	2.7040439085699037
BENALLA	290000	1.739107387612497
BOROONDARA	2075000	2.8368436256478136
BRIMBANK	632500	2.3126065593780973
BULOKE	135500	1.399155772145344
CAMPASPE	234500	2.2574981381775667
CARDINIA	550000	2.507555786556805
CASEY	567500	2.5903105680780767

CENTRAL GOLDFIELDS	235500	1.3763201587999694
COLAC OTWAY	530250	1.993244858979771
CORANGAMITE	221500	2.0881095462474115
DAREBIN	1135000	2.312900073883726
EAST GIPPSLAND	310000	1.9351826044607967
FRANKSTON	608750	2.149084342082875
GANNAWARRA	205000	1.7065373023940116
GLEN EIRA	1477750	2.508541533252486
GLENELG	212750	1.8950754089430983
GOLDEN PLAINS	421500	1.9468128633004147
GREATER BENDIGO	368750	2.362089981610652
GREATER DANDENONG	629000	2.050355371724189
GREATER GEELONG	552500	2.766546728484037
GREATER SHEPPARTON	279250	1.9589742822693
HEPBURN	521250	1.731252297553474
HINDMARSH	100500	0.9232552798950614
HOBSONS BAY	850750	2.229605965690888
HORSHAM	269000	1.8523082802029143
HUME	533500	2.1422621159216613
INDIGO	335500	1.6371074508014265
KINGSTON	890250	2.411518603017613
KNOX	763500	2.424973980723895
LATROBE	250000	1.5630712337383574
LODDON	161750	1.5432433202421507
MACEDON RANGES	695000	2.018901742491153
MANNINGHAM	1121000	2.2519375032496507
MANSFIELD	365750	1.1485055580864514
MARIBYRNONG	885000	1.9751810758665256
MAROONDAH	819250	2.121887851495286
MELBOURNE	1375000	2.0784089943949478
MELTON	515000	1.8407566888554507
MILDURA	239000	1.4848009921631191
MITCHELL	442500	1.3381691157435764
MOIRA	262500	1.5513108842527197
MONASH	1100000	2.2247534602091124
MOONEE VALLEY	1088500	2.0871231817153233
MOORABOOL	477500	1.687042963304442
MORELAND	849250	2.097552124204076
MORNINGTON PENINSULA	795000	2.2465036008755965
MOUNT ALEXANDER	505000	1.7321940038389902
MOYNE	382000	1.675348363714706
MURRINDINDI	380500	0.97557993491256
NILLUMBIK	815500	2.0263814069655566
NORTHERN GRAMPIANS	189000	1.3712036614723058
PORT PHILLIP	1776250	1.8599249784535532
PYRENEES	231000	1.3739751945208492
QUEENSCLIFFE	888750	1.7762092395683893
SOUTH GIPPSLAND	328750	1.6068044229687402
SOUTHERN GRAMPIANS	186750	1.6030831375529997

STONNINGTON	2116250	2.108344659706451
STRATHBOGIE	339500	1.5453839892393293
SURF COAST	905000	1.8195476444503096
SWAN HILL	225000	1.2500399183520183
TOWONG	251500	1.1583313255658538
WANGARATTA	345750	0.9361798788937388
WARRNAMBOOL	360000	1.5858065426312844
WELLINGTON	252500	1.9316711048384558
WEST WIMMERA	171250	1.570497748249457
WHITEHORSE	1159500	1.5783227537831162
WHITTLESEA	590500	2.0008594549663754
WODONGA	395000	1.6839158235564338
WYNDHAM	565000	1.89392751574602
YARRA	1352500	1.855080142348744
YARRA RANGES	642500	2.831952643734072
YARRIAMBIACK	114500	1.7755305119836042

Appendix B

<u>LGA Name</u>	<u>CR</u>	<u>FT</u>	<u>OA</u>	<u>RT</u>	<u>ST</u>
ALPINE	0.809079854	0.156976744	0.0	0.423076923	1.0
ARARAT	0.416700446	0.174418605	0.0	0.710702341	0.9625
BALLARAT	0.431698419	0.540697674	0.0	0.925585284000 0001	0.9375
BANYULE	0.693149574000 0001	0.404069767	0.0056071830 00000001	0.945652174	0.7
BASS COAST	0.693960276	0.22965116300 000002	0.0	0.738294314	0.6875
BAW BAW	0.684637211	0.42151162799 999997	0.0	0.738294314	0.6875
BAYSIDE	0.820429672	0.395348837	0.005305291	0.915551839	0.6625
BENALLA	0.421970004	0.09011627900 000001	0.0	0.68729097	0.6125
BOROONDARA	0.895419538	0.502906977	0.00423198	0.927257525	0.6
BRIMBANK	0.485610053	0.38662790700 000005	0.0102045809 99999999	0.929765886	0.5875
BULOKE	0.40656668	0.203488372	0.0	0.280936455	0.575
CAMPASPE	0.612079449	0.398255814	0.0	0.753344482	0.575
CARDINIA	0.704499392	0.39244186	0.183467522	0.822742475	0.55
CASEY	0.565058776	0.630813953	0.0260270570 00000003	0.915551839	0.55
CENTRAL GOLDFIELDS	0.047426024000 000004	0.16569767400 000002	0.0	0.702341137	0.525
COLAC OTWAY	0.585731657999 9999	0.26744186	0.0	0.686454849000 0001	0.525
CORANGAMITE	0.848804215999 9999	0.276162791	0.0	0.52090301	0.5125
DAREBIN	0.601134982	0.328488372	0.005129389	0.954849498	0.5
EAST GIPPSLAND	0.404945277999 99996	0.424418605	0.0	0.678093645	0.5

FRANKSTON	0.378597487	0.39244186	0.01739442	0.953177258	0.4875
GANNAWARRA	0.575597892	0.209302326	0.0	0.497491639	0.4875
GLEN EIRA	0.832590191	0.33430232600 000004	0.0029072359 999999997	0.950668895999 9999	0.4625
GLENELG	0.443453587000 00004	0.27906976699 999997	0.0	0.787625418	0.45
GOLDEN PLAINS	0.792055128	0.27906976699 999997	0.0	0.50083612	0.4375
GREATER BENDIGO	0.508309688	0.65406976700 00001	0.0	0.850334448	0.425
GREATER DANDENONG	0.43899473	0.313953488	0.0061021040 000000006	0.95735786	0.4
GREATER GEELONG	0.563437373	1.0	0.0	0.887959866	0.4
GREATER SHEPPARTON	0.247263883000 00002	0.470930233	0.0	0.907190635	0.4
HEPBURN	0.719092015	0.151162791	0.0	0.515886287999 9999	0.4
HINDMARSH	0.056749087999 999996	0.104651163	0.0	0.431438127	0.375
HOBSONS BAY	0.665585732	0.36046511600 000003	0.004974073	0.912207358	0.35
HORSHAM	0.429671666	0.215116279	0.0	0.913043478	0.35
HUME	0.482367247999 99997	0.470930233	0.0115048739 99999998	0.918060201	0.325
INDIGO	0.854073774000 0001	0.127906977	0.0	0.387959866000 00004	0.3125
KINGSTON	0.738548845	0.470930233	0.010405871	0.945652174	0.3125
KNOX	0.618565059	0.60755814	0.008838097	0.95819398	0.3
LATROBE	0.0	0.450581395	0.0	0.877926421	0.2875
LODDON	0.783948115	0.18023255800 000001	0.0	0.357859532000 00004	0.2625
MACEDON RANGES	0.723145521	0.313953488	0.0	0.77006689	0.2625
MANNINGHAM	0.843129307	0.34011627899 999997	0.01192645	0.863712375	0.25
MANSFIELD	0.597486826	0.063953488	0.0	0.27090301	0.25
MARIBYRNONG	0.561815971	0.261627907	0.001811823	0.949832775999 9999	0.25
MAROONDAH	0.612079449	0.331395349	0.005907211	0.988294314	0.2375
MELBOURNE	0.538710985	0.354651163	0.003592891	0.996655518	0.2375
MELTON	0.479935144	0.28488372100 000003	0.0227359379 99999997	0.893812709000 0001	0.2125
MILDURA	0.125253344000 00002	0.302325581	0.0	0.898829430999 9999	0.2
MITCHELL	0.176732873999 99998	0.244186047	0.004505155	0.752508361	0.2
MOIRA	0.429671666	0.316860465	0.0	0.657190635	0.1875
MONASH	0.731657884000 0001	0.447674419	0.006114921	0.903846154	0.1875

MOONEE VALLEY	0.729631131	0.33430232600 000004	0.0030042579 999999997	0.891304348	0.175
MOORABOOL	0.590595865	0.162790698	0.0	0.795986622	0.175
MORELAND	0.694770977	0.34011627899 999997	0.002707699	0.931438127000 0001	0.175
MORNINGTON PENINSULA	0.576813944	0.610465116	0.0876387590 0000001	0.891304348	0.1625
MOUNT ALEXANDER	0.79327118	0.220930233	0.0	0.60367893	0.15
MOYNE	0.832590191	0.177325581	0.0	0.549331104	0.15
MURRINDINDI	0.301175517	0.19476744199 999999	0.0	0.357023411	0.15
NILLUMBIK	0.856505878	0.238372093	0.0883022719 9999999	0.775083612	0.1375
NORTHERN GRAMPIANS	0.396432913999 99997	0.154069767	0.0	0.714046823	0.1375
PORT PHILLIP	0.610052696	0.186046512	0.002896546	0.961538462	0.1375
PYRENEES	0.618565059	0.09883720900 000001	0.0	0.547658863	0.1375
QUEENSCLIFFE	1.0	0.0	0.0	0.680602007000 0001	0.125
SOUTH GIPPSLAND	0.604783137000 0001	0.375	0.0	0.537625418	0.125
SOUTHERN GRAMPIANS	0.647750304	0.220930233	0.0	0.642140468	0.125
STONNINGTON	0.798540737999 9999	0.293604651	0.001376566	0.930602007000 0001	0.125
STRATHBOGIE	0.772598298	0.12209302300 000001	0.0	0.555183946	0.125
SURF COAST	0.71625456	0.34011627899 999997	0.0	0.675585284000 0001	0.125
SWAN HILL	0.177948926	0.191860465	0.0	0.796822742	0.1125
TOWONG	0.732063235	0.069767442	0.0	0.289297659	0.0875
WANGARATTA	0.568301581000 0001	0.302325581	0.0	0.0	0.0875
WARRNAMBOOL	0.49047426	0.244186047	0.0	0.806020067000 0001	0.075
WELLINGTON	0.486420754	0.441860465	0.0	0.966555184000 0001	0.075
WEST WIMMERA	0.788406972	0.104651163	0.0	0.627926421000 0001	0.075
WHITEHORSE	0.764085934	0.40988372100 000003	0.0034940220 000000003	0.357859532000 00004	0.075
WHITTLESEA	0.595460072999 9999	0.36046511600 000003	0.0641965970 0000001	0.976588629000 0001	0.0625
WODONGA	0.641264694	0.14244186	0.0	0.865384615	0.0625
WYNDHAM	0.563032023	0.38662790700 000005	0.007328721	0.912207358	0.0625
YARRA	0.569517633	0.348837209	0.002563835	0.918060201	0.05
YARRA RANGES	0.662748277	0.537790698	1.0	1.0	0.025

YARRIAMBIA	0.755573570999 9999	0.200581395	0.0	0.843645485	0.0
-------------------	------------------------	-------------	-----	-------------	-----

Appendix c

Raw data files, data files created in the intermediate steps and .ipynb used to achieve the results are submitted via the .zip file in the LMS. The order they are to be run to be guaranteed to work and input/output they require/produce are specified in the README file.