

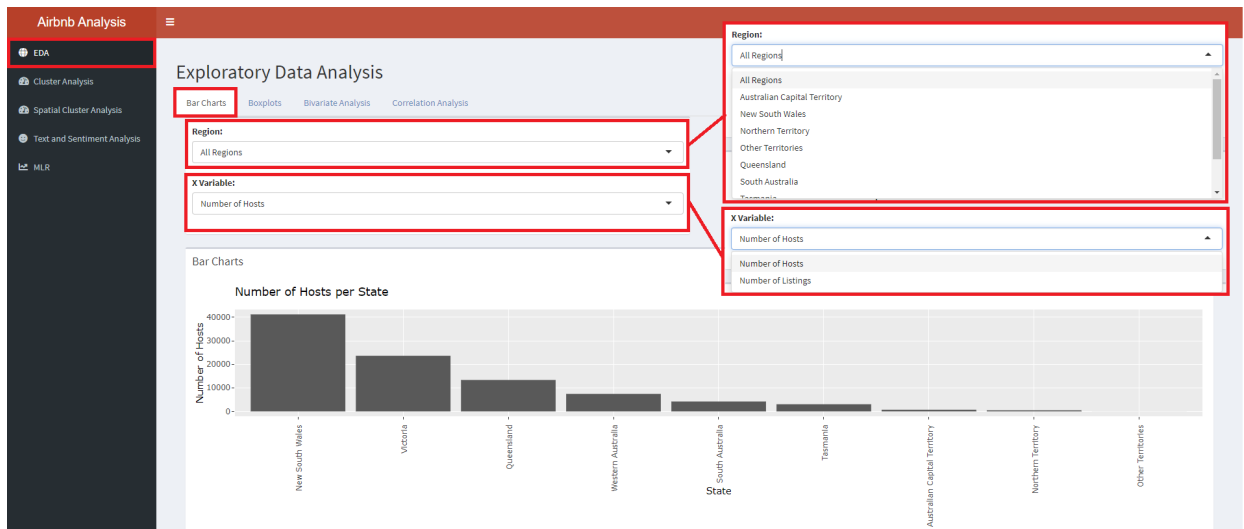
User Guide - Understanding Airbnb listings in Australia

true true true

Exploratory Data Analysis (EDA)

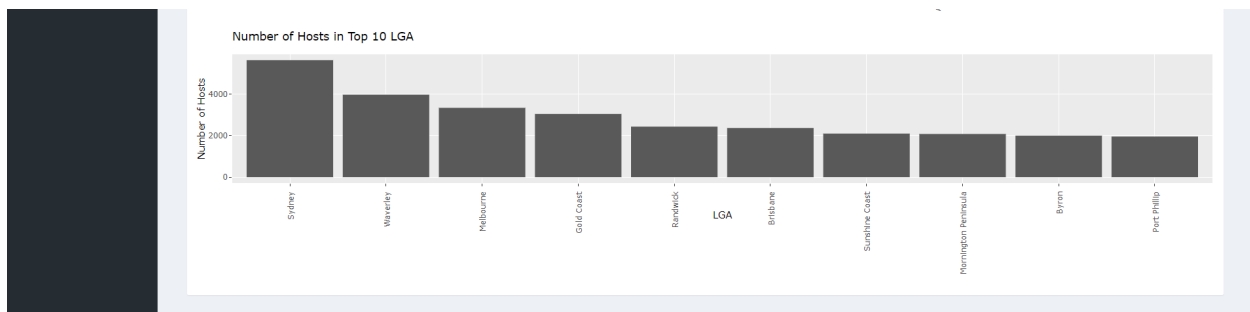
Exploratory Data Analysis allows you to perform initial investigation on the data, so that you may be able to discover patterns, and explore the different variables in the data set. It will allow us to formulate hypothesis and explore different statistics models that could be developed after.

Bar Charts



- Click on “Region:” to explore the different states in Australia
- Click on “X Variable:” to explore “Number of Hosts” and “Number of Listings” in various different States.

The top barchart showcases the “X Variable” per State.

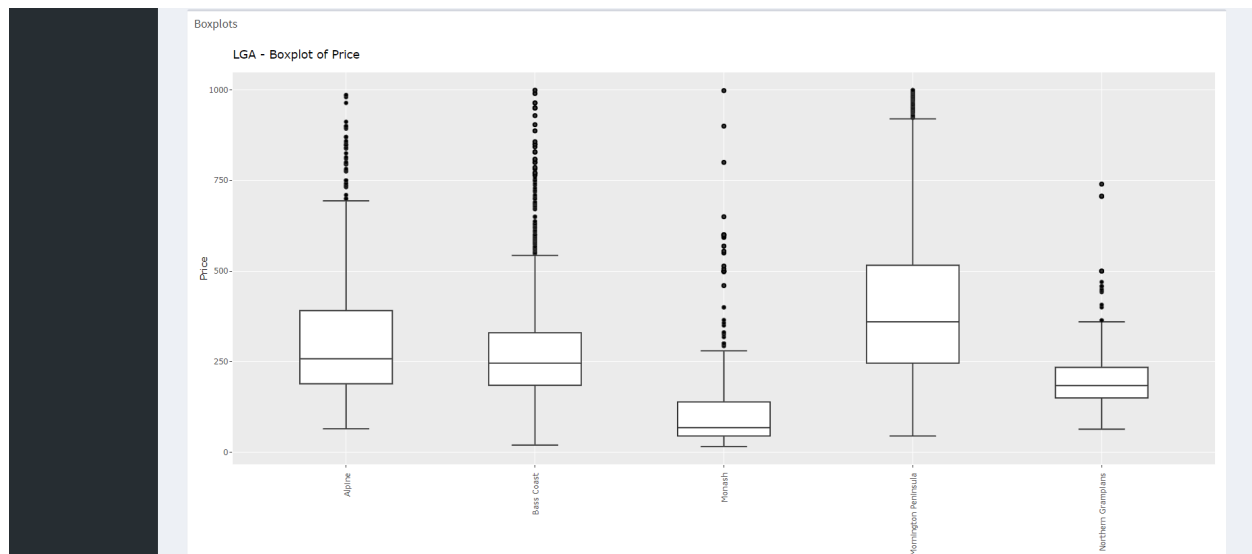


The bottom barchart showcases the “X Variable” per top 10 Local Government Area per State.

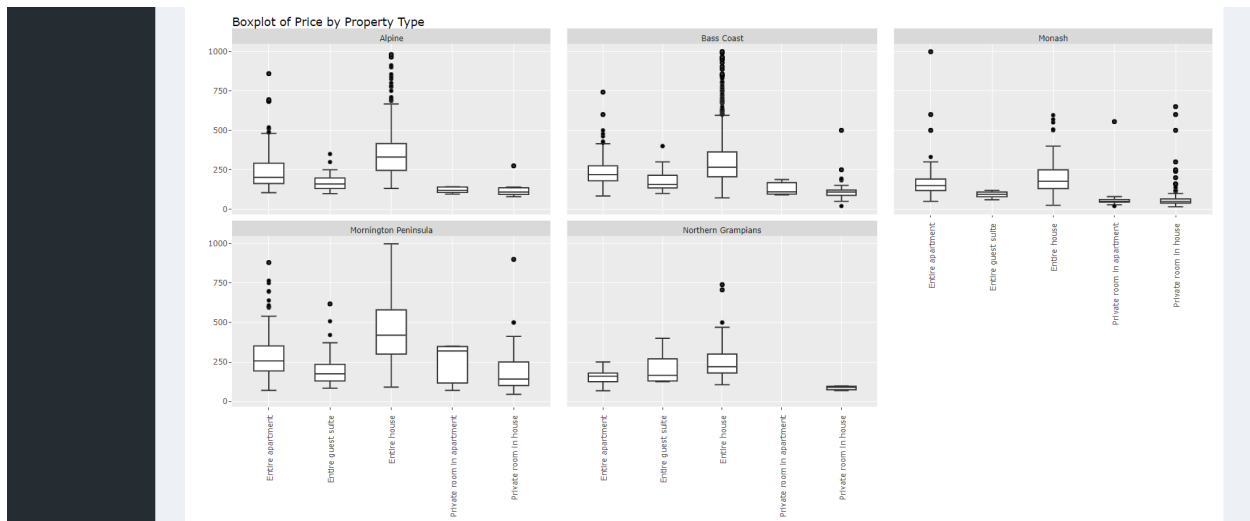
Boxplots

The screenshot shows the Airbnb Analysis dashboard. The 'Exploratory Data Analysis' section is active, with the 'Boxplots' tab selected. The 'Region' dropdown is set to 'Victoria'. The 'Local Government Area' dropdown is set to 'Alpine', 'Bass Coast', 'Monash', 'Mornington Peninsula', and 'Northern Grampians'. The 'Y Variable' dropdown is set to 'Price'. A red box highlights the 'Boxplots' tab and the 'Region', 'Local Government Area', and 'Y Variable' dropdowns. A red arrow points from the 'Boxplots' tab to the resulting boxplot.

- Click on “*Region:*” to explore the different states in Australia.
- Click on “*Local Government Area:*” to explore the different cities in the State. You are able to choose multiple cities to view at a single time through clicking the chosen city in the drop down bar.
- Click on “*Y Variable:*” to view the boxplots of different variables such as *Price*, *Review Scores Ratings*, *Bedrooms*, *Beds* etc.



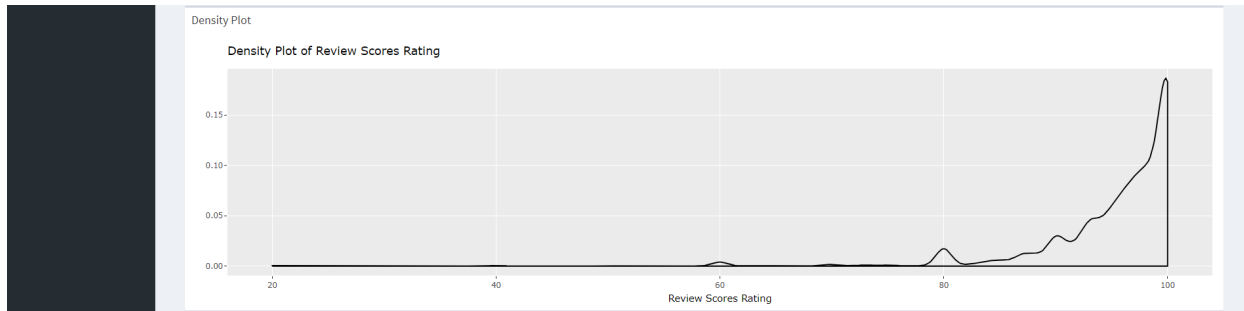
The top boxplot showcases Boxplots of “*Y Variable*” against the Local Government Area.



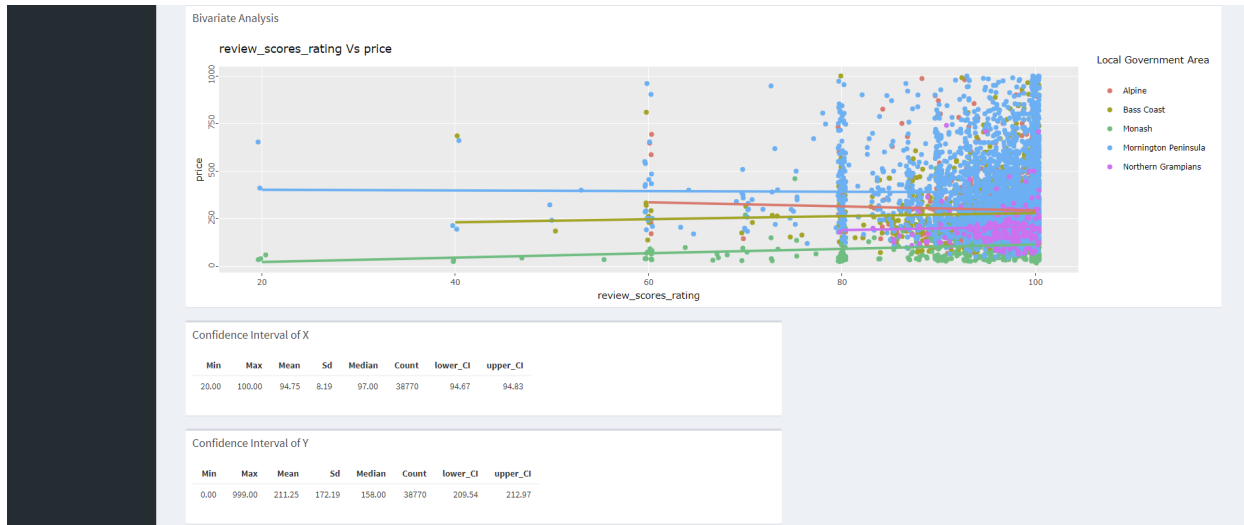
The bottom boxplot showcases “Price” vs “Property Type” with a Facet Wrap of Local Government Area.

Bivariate Analysis

- Click on “Region:” to explore the different states in Australia.
- Click on “Local Government Area” to explore the different cities in the State. You are able to choose multiple cities to view at a single time through clicking the chosen city in the drop down bar.
- Click on “X Variable:” and “Y Variable:” to view the bivariate plot with different variables such as *Price*, *Review Scores Ratings*, *Bedrooms*, *Beds* etc. This will allow us to explore how 2 variables interact with each other



The top plot is a density plot to showcase the density of the “*X Variable*”.



The bottom plot is a bivariate plot to explore the relationship between 2 variables.

The two tables below the bivariate plots showcase the confidence interval of the “*X Variable*” and “*Y Variable*”.

Correlation Analysis

The screenshot shows the 'Airbnb Analysis' interface with the 'Exploratory Data Analysis' section selected. The 'Correlation Analysis' tab is active. Two red boxes highlight the 'Visualisation Method:' and 'Reorder Correlation Matrix:' dropdown menus. Arrows point from these boxes to their respective expanded lists on the right.

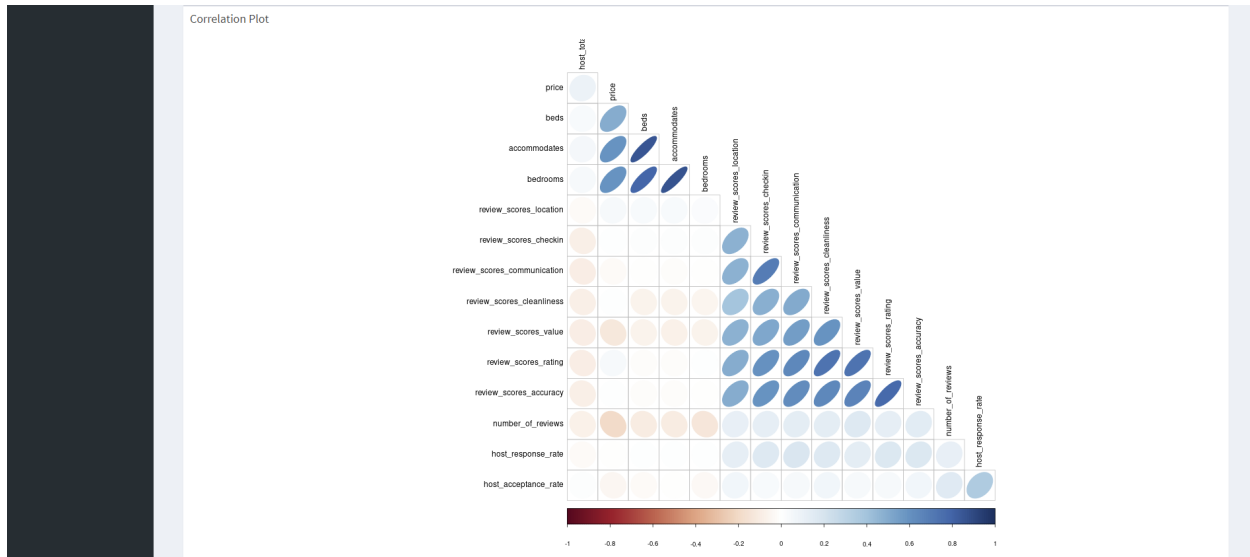
Visualisation Method:

- ellipse
- circle
- square
- ellipse
- number
- shade
- color
- pie

Reorder Correlation Matrix:

- hclust
- AOE
- FPC
- hclust
- alphabet

- Click on “*Visualisation Method:*” to view the correlation plots in different methods such as “*circle, ellipse and number*”.
- Click on “*Reorder Correlation Matrix:*” to view the correlation plots in different orders such as “*hclust, alphabet and AOE*”.



The plot represents the correlation plot of all the variables in the Airbnb data set.

Cluster Analysis

Airbnb Analysis

Cluster Analysis

K Means Clustering

Region: Victoria

Cluster Variable: Bedrooms Beds Price Review Scores Rating

Distance Function: euclid

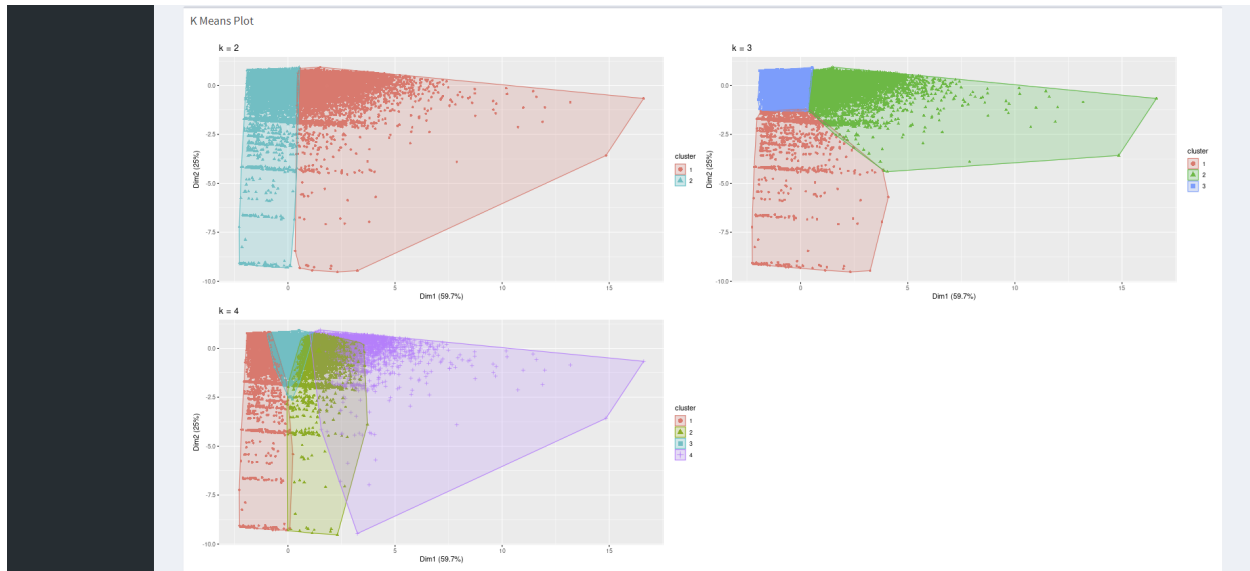
Cluster Size: 4

Distance Function: euclid, euclid, cosine

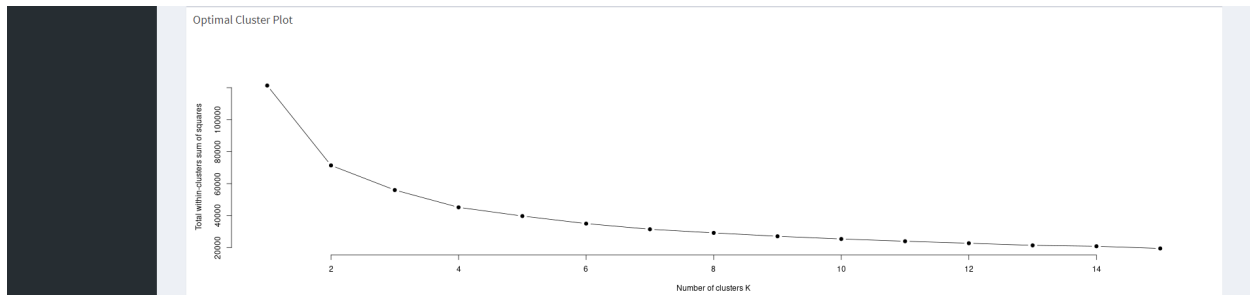
Region: Victoria, New South Wales, Northern Territory, Other Territories, Queensland, South Australia, Tasmania, Victoria

Cluster Variable: Bedrooms Beds Price Review Scores Rating, Accommodates, Host Acceptance Rate, Host Response Rate, Review Scores Accuracy, Review Scores Checkin, Review Scores Cleanliness, Review Scores Communication

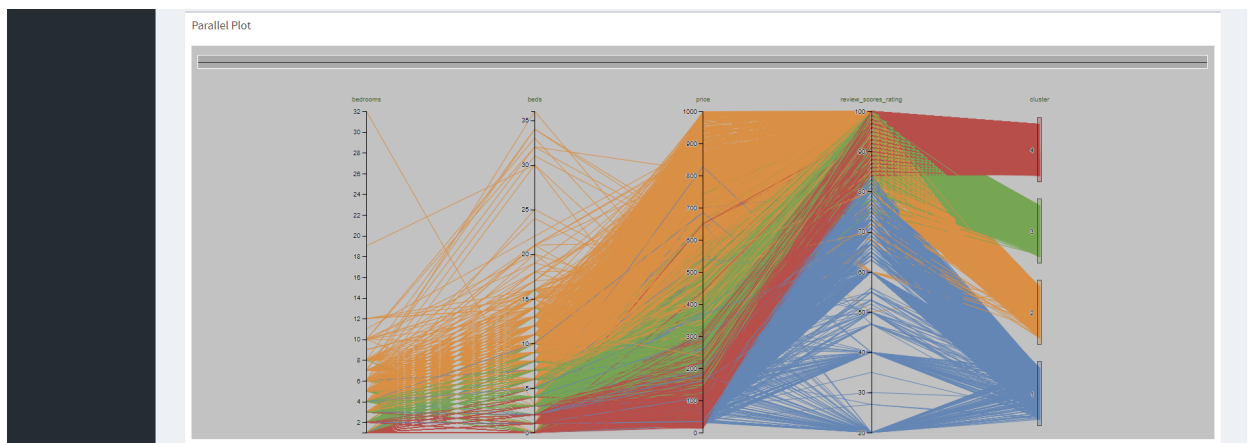
- Click on “*Region:*” to explore the different states in Australia.
- Click on “*Cluster Variable:*” to explore the different variables you would like to cluster. Types of variables are “*Price, Bedrooms, Beds, Review Score Ratings*” etc. You are also able to click to choose more variables that would be used for clustering.
- Click on “*Distance Function:*” to choose between “*euclid*” and “*cosine*” distance for the clusters.
- Toggle “*Cluster Size*” between the sizes 2 and 15, to see different cluster sizes in the data set.



The first plot is the kmeans plot. it showcases the different cluster sizes.



The second plot is an Optimal Cluster Plot, that gives you an indication the number of clusters that will be best for the data set used.



The third plot is a parallel plot that showcases the relationship of the variables and the clusters. It gives an indication of the characteristics of the cluster that was formed.

Exploratory Spatial Data Visualization and Analysis

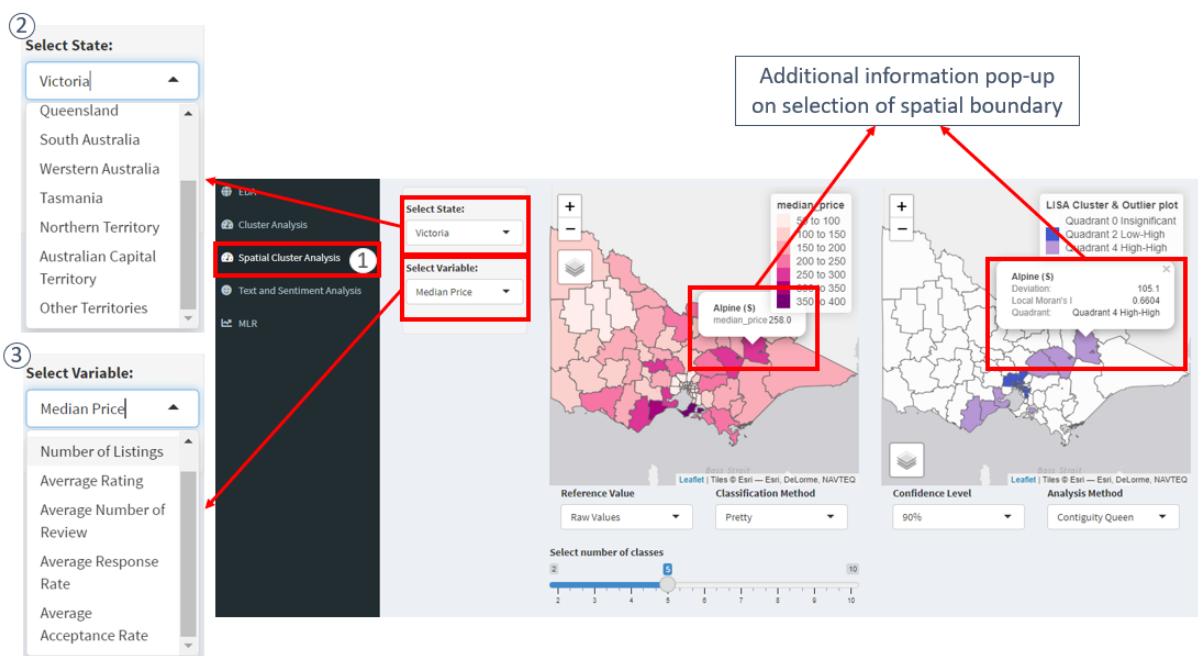
This tab can be used for exploratory spatial data visualization and analysis on the Local Government Area (LGA) level for each of the States of Australia.

Step 1) by clicking on the Spatial Cluster Analysis main left-hand panel, we enter this tab, with two maps coming into view (Note: it might take awhile for the viz to load). The left-hand viz provides an overview of the distribution of either the raw values (default) or clustering statistics for the LGAs within a selected area (State) of Australia. The right-hand viz displays the spatial cluster analysis results - showing areas of high/low clusters, and outliers

Step 2) clicking on the drop down list on the “Select State:” tab allows you to select the state of interest from a list of nine different states in Australia (*Note: while the visualization allows for user to select the “All” state option, it is highly discouraged as some of the spatial analyses do not have meaningful interpretations for areas with considerable missing data - for which at the whole-of-Australia level, data are more sparsed. When analyzing the data, users are also advised to conduct due diligence in determining whether the selected data range of interest has sufficiently satisfied analysis requirements.*)

Step 3) allows you to select the variable of interest from a list of six..

The resultant maps also allows you to click on the area within each contiguous boundary to have more information, such as the name of the LGA, the status (e.g. (C) is a city, (A) an area, (T) towns, (M) municipalities, and (S) shire), and the value of the variable of interest on the first geographical visualization, and the clustering analysis statistics on the second geo-viz on the right.



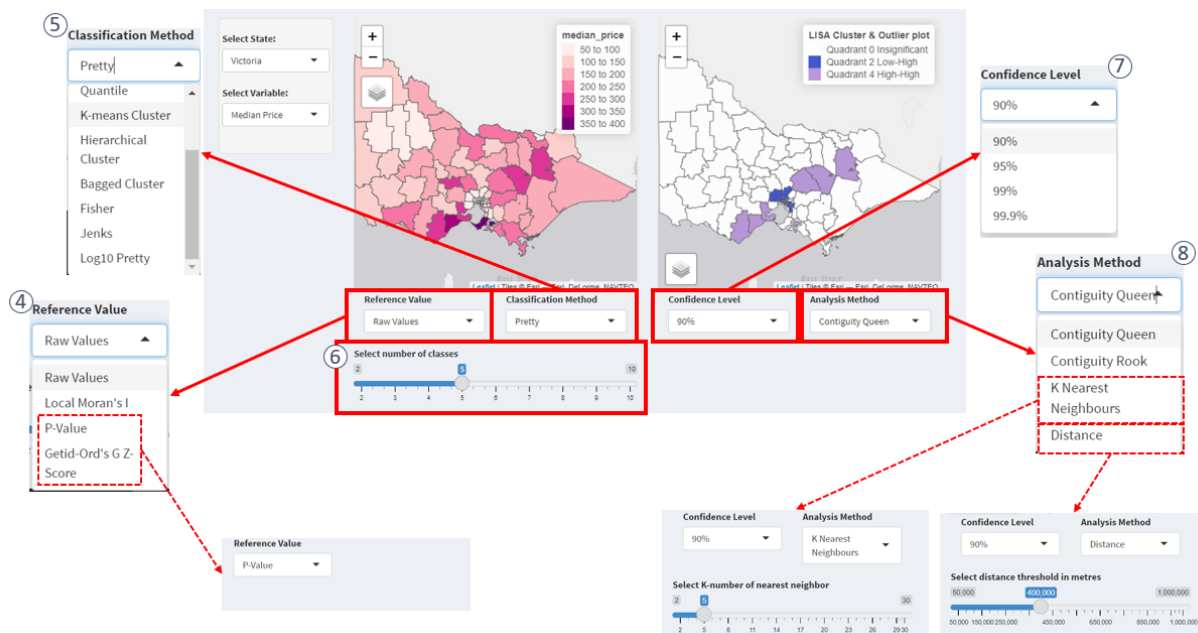
Step 4) we go on to the options for individual geoviz. For the left-hand viz, user can first select whether you want the distribution of raw values displayed, or that of the local indicators of spatial association's statistics displayed. Note that if the Moran's I's p-value or the Getis-Ord's G z-scores are selected, selection for other options (see step 4 and 5 below) will be disabled as these statistical results are to be interpreted on pre-determined intervals.

Step 5) allows you to select the binning methods - in other words, how should the values be split up, be it by equal intervals, by natural breaks (Jenks' and Fisher's), standard deviation bins, etc. A total of 10 options are available under the dropdown.

Step 6) allows user to select the number of intervals to have - i.e. how many different classes should the data be split into. User can shift the slider manually to select the desired class. Note that the selection is very critical (compulsory) for classification methods such as the K-Mean selection; practical for methods such as Equal, Pretty, Jenks and Fisher; but not as useful for methods such as log10 Pretty or even standard deviation.

Step 7) on the right-hand viz, user can select variables to alter the clustering parameters. Firstly, the level of confidence - i.e. how much confidence do we need to reject the null hypothesis that there are no spatial clustering/outlier at the areas identified to have such cluster/outlier.

Step 8) in geospatial analysis, there are also different methods to select the definition of 'neighbor' - this can be selected by changing the "Analysis methods:" panel options. Note that if "K Nearest Neighbors" is selected, an additional slider requesting user to input the number of nearest neighbors to be considered will appear. Same thing, if "Distance" is selected, a slider to select the distance radius to be considered neighbor will appear as well.



Sentiment Analysis

The sentiment analysis module of this Shiny app includes two analysis models: a word cloud and a topic model.

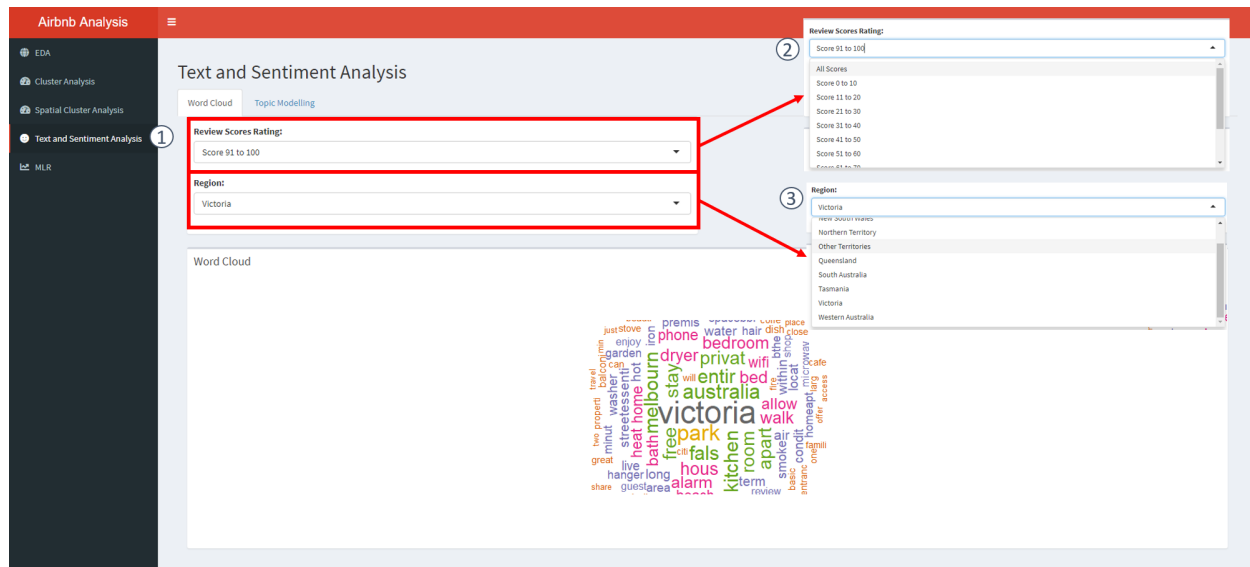
Word Cloud

Step 1) Click on Text and Sentiment Analysis to open the module.

Step 2) Select the range of review score ratings that you want to focus on. By default, all scores will be included.

Step 3) Select the region of focus. In this case, it is Victoria.

The word cloud will be generated. The size of the font represents the relative frequency of the word in the description section of Airbnb listings that fall in the category of the filter settings above (the word in the largest font is the most common word).



Topic Modelling

Topic modeling is an unsupervised machine learning technique that detects word and phrase patterns in documents and clusters them into groups known as topics.

Latent Dirichlet Allocation (LDA) is one common topic modeling technique. The basic assumption of LDA is that similar topics make use of similar words (i.e. distributional hypothesis). The purpose of LDA is to map the corpus to topics covering a significant number of words in the documents in the corpus.

LDA assigns topics to arrangements of words for example, the best word for a topic related to accommodation. This is based on the assumption that documents are written with a certain arrangement of words and that those arrangements will determine the topics. LDA assumes that all words in the document can be assigned a probability of belonging to a topic. As such, the goal of LDA is to determine the mixture of topics that a document contains.

The navigation is similar to the Word Cloud pane.

Step 1) Click on Text and Sentiment Analysis to open the module.

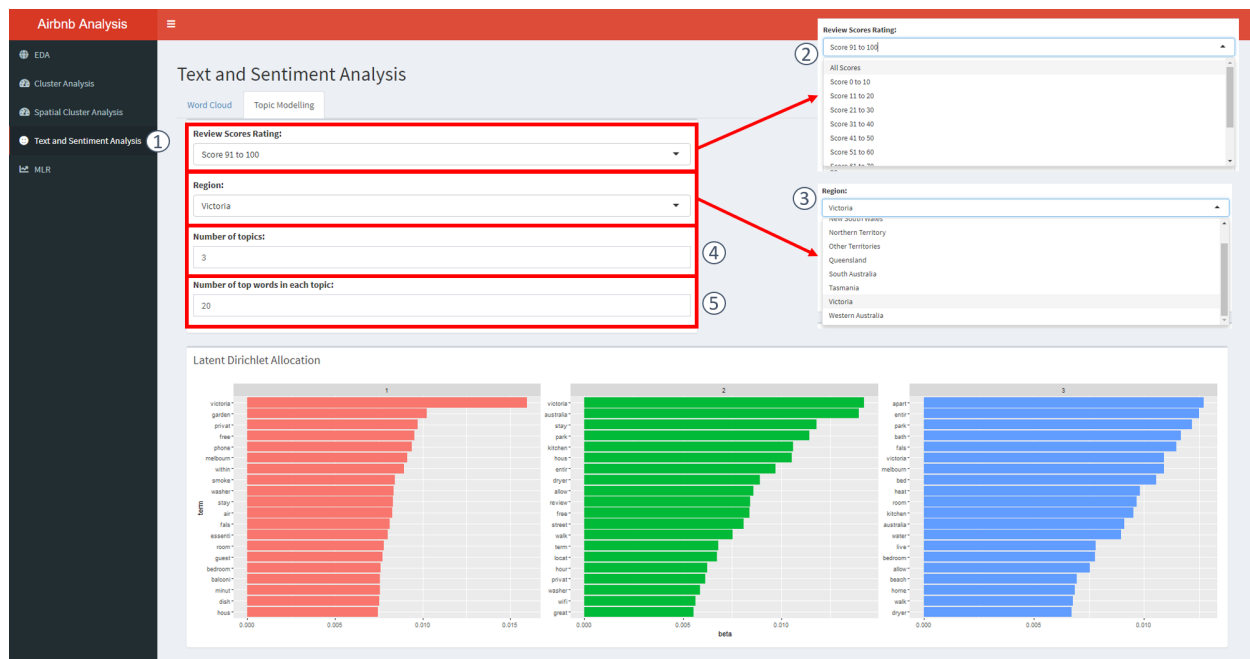
Step 2) Select the range of review score ratings that you want to focus on. By default, all scores will be included.

Step 3) Select the region of focus. In this case, it is Victoria.

Step 4) Enter the number of topics you want identified. The default is three topics.

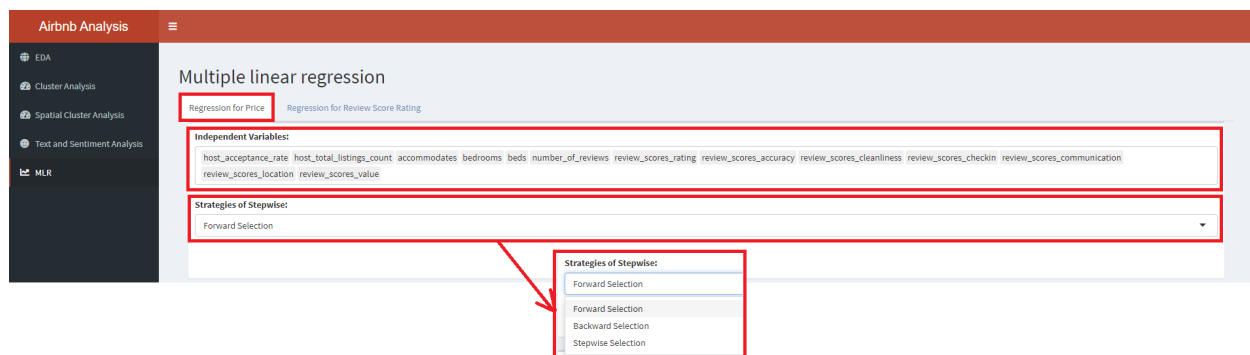
Step 5) Enter the number of words you would like to be displayed for each topic. The default is 10 words.

The LDA topic model will be generated. The bar beside the corresponding word represents the probability (beta) of that word appearing in that particular topic. The longer the bar, the higher the probability.



Multiple Linear Regression

Regression for Price



- Click on “*Independent Variables:*” to choose the multiple variables that you would like to apply in the Multiple Linear Regression (MLR). The choices of variables could be influenced through the correlation plot found in the EDA. It is better to choose variables with low correlation with each other.
- Click on “*Strategies of Stepwise:*” to choose whether the MLR will move in a “*Forward*”, “*Backward*”, or “*Stepwise*” direction.

Regression Summary:

```
Start: AIC=968813.2
price ~ host_acceptance_rate + host_total_listings_count + accommodates +
        bedrooms + beds + number_of_reviews + review_scores_rating +
        review_scores_accuracy + review_scores_cleanliness + review_scores_checkin +
        review_scores_communication + review_scores_location + review_scores_value
```

		Df	Sum of Sq	RSS	AIC
<none>				1642073670	968813
+ review_scores_accuracy		1	75596	1643404570	968816
+ review_scores_checkin		1	506755	1643480425	968842
+ host_acceptance_rate		1	1166489	1644148170	968882
+ review_scores_communication		1	1580770	1644813440	969183
+ review_scores_cleanliness		1	2874678	1645848348	968185
+ review_scores_location		1	7157613	1650171283	968447
+ host_total_listings_count		1	7888776	1650968682	968495
+ beds		1	9561806	1652534676	968598
+ number_of_reviews		1	21871218	1671044888	969788
+ review_scores_rating		1	31235463	1674339133	969889
+ bedrooms		1	66751130	1789728800	971982
+ review_scores_value		1	79538175	1713583545	972101
+ accommodates		1	87318428	1738818298	973163

Call:
lm(formula = price ~ host_acceptance_rate + host_total_listings_count +
accommodates + bedrooms + beds + number_of_reviews + review_scores_rating +
review_scores_accuracy + review_scores_cleanliness + review_scores_checkin +
review_scores_communication + review_scores_location + review_scores_value,
data = Rdata)

Residuals:

	Min	1Q	Median	3Q	Max
	-2056.64	-73.13	-21.18	47.83	916.96

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-83.583406	0.405028	-9.935	< 2e-16 ***
host_acceptance_rate	-15.878689	1.863749	-8.412	< 2e-16 ***
host_total_listings_count	0.108167	0.004959	22.012	< 2e-16 ***
accommodates	24.864277	0.341966	72.706	< 2e-16 ***
bedrooms	42.236676	0.663693	63.639	< 2e-16 ***
beds	-8.831159	0.367989	-24.004	< 2e-16 ***
number_of_reviews	-0.292998	0.007100	-41.268	< 2e-16 ***
review_scores_rating	4.688513	0.107520	43.532	< 2e-16 ***
review_scores_accuracy	1.952084	0.900719	2.146	0.0319 *
review_scores_cleanliness	9.054620	0.685639	13.286	< 2e-16 ***
review_scores_checkin	-5.328460	0.968998	-5.485	2.95e-06 ***
review_scores_communication	-0.864118	0.947713	-0.904	< 2e-16 ***
review_scores_location	19.518758	0.934067	20.897	< 2e-16 ***
review_scores_value	-47.811855	0.738914	-65.414	< 2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 128.4 on 99676 degrees of freedom
Multiple R-squared: 0.4425, Adjusted R-squared: 0.4424
F-statistic: 6885 on 13 and 99676 DF, p-value: < 2.2e-16

The “*Regression Summary*” showcases the MLR calculation and gives us an indication of the significance of each variable in the regression with respect to “*Price*”.

Regression for Review Score Ratings

Airbnb Analysis

Multiple linear regression

Regression for Price Regression for Review Score Rating

Independent Variables:

host_acceptance_rate host_total_listings_count accommodates bedrooms beds number_of_reviews review_scores_rating review_scores_accuracy review_scores_cleanliness review_scores_checkin review_scores_communication review_scores_location review_scores_value

Strategies of Stepwise:

Forward Selection

Strategies of Stepwise:

Forward Selection
Forward Selection
Backward Selection
Stepwise Selection

- Click on “*Independent Variables:*” to choose the multiple variables that you would like to apply in the Multiple Linear Regression (MLR). The choices of variables could be influenced through the correlation plot found in the EDA. It is better to choose variables with low correlation with each other.
- Click on “*Strategies of Stepwise:*” to choose whether the MLR will move in a “*Forward*”, “*Backward*”, or “*Stepwise*” direction.

The “*Regression Summary*” showcases the MLR calculation and gives us an indication of the significance of each variable in the regression with respect to “*Review Score Ratings*”.