

# ASSIGNMENT COVER SHEET

SUBJECT NUMBER & NAME	36103 Statistical Thinking in Data Science	
NAME OF STUDENT	Brackenrig, Mark, Campbell, MaryAlice, Cheluget, Linda	
(PRINT CLEARLY - SURNAME, FIRST NAME)	Lehmann, Martin, Xu, Yusheng	
STUDENT ID NUMBER	12964298, 13298029, 13009028, 12543738, 11860870	
	mark.brackenrig@student.uts.edu.au	
	maryalice.campbell@student.uts.edu.au	
STUDENT EMAIL	linda.cheluget@student.uts.edu.au	
	martin.lehmann@student.uts.edu.au	
	yusheng.xu@student.uts.edu.au	
	MB: +61 450 548 776, MC: +61 414 624 429, LC: +61 450 652 779,	
STUDENT CONTACT NUMBER	ML: +61 414 638 601, YX: +61 452 256 666	
DUE DATE	30-Apr-2018 by 23:59	
ASSESSMENT ITEM NUMBER/TITLE	AT2 – Data Analysis Project, Part A: Proposal	

- I confirm that the work submitted conforms with the university's guidelines on academic integrity.
   Refer to the UTS policy on 'Advice to Students on Good Academic Practice':
   <a href="http://www.gsu.uts.edu.au/policies/academicpractice.html">http://www.gsu.uts.edu.au/policies/academicpractice.html</a>
- I am aware of the penalties for plagiarism. This assignment is my own work and I have not handed in this assignment (either part or completely) for assessment in another subject.
- If this assignment is submitted after the due date I understand that it will incur a penalty for lateness unless I have previously had an extension of time approved and have attached the written confirmation of this extension.

Please provide details of extensions granted here if applicable		
Signature of Student:	<b>Date</b> : <u>30 / Apr / 2018</u>	
If submitted electronically tick here to indicate you agree with the above		



## Research Proposal

# AirBnB uptake by Business Sector - Drivers for policy change -

Version: 1.0

**Date:** 30-April-2018

Authors: 3MDL – Data Science Team

### Contents

1 Introduction	3
2 Document Purpose	3
3 Rationale and Research Questions	4
Objectives	4
4 Research Method	4
Data selection	5
Assumptions	5
Techniques	5
5 Challenges	6
6 Timeline	7
Appendix A – Data Acquisition & Merge (code samples)	8
Appendix B – Data Merging code samples	15
Appendix D – Proposed Google Directions Locations	17
Appendix E – Complete list of datasets	18
Selected	18
To Be Confirmed	19
Rejected	19

### **Document Control**

Version	Purpose	Authors	Date
0.1	Draft Template	3MDL – DS team	23-04-2018
0.2	Initial Draft – with stakeholder feedback	3MDL – DS team	24-04-2018
0.3	Incorporating stakeholder feedback	3MDL – DS team	28-04-2018
0.4	Addition of code samples	3MDL – DS team	28-04-2018
1.0	Final version	3MDL – DS team	30-4-2018

### 1 Introduction

The changing nature of tourism accommodation has been enabled by new technology in online marketing, bookings and trust-based systems. Responding to consumer demand, the traditional hotel sector has diversified by offering a range of alternative options such as boutique hotels, private accommodation and last-minute bookings. The 3DML Data Science Team is proposing a research project to investigate the private accommodation market, the demographics of its consumers, and its impact on local communities (see Figure 1). The findings from this research project will be used to make recommendations for changes to existing planning and zoning policy in alignment with 3MDL's strategy vision statement "To hear the community voice".

# Tourism Accommodation Alternative Accommodation Private Accommodation AirBnB Local Community Local Council

Figure 1 - Research Area and Community Participants

### 2 Document Purpose

This document presents an overview of the project rational and articulates the research questions to be investigated. Also presented are the data selection process, proposed modelling techniques and anticipated challenges identified during the preliminary exploratory data analysis (EDA).

### 3 Rationale and Research Questions

This project aims to deliver a foundation for identifying the value to property owners and local-residents in regulating or incentivising AirBnB listings for the business sector. Such policy may encourage listings in areas with under occupied residences and stabilise them in areas with high investment to residential property ratios. This will be achieved by examining the factors driving increased use of AirBnB for business travel in Sydney and its impact on the local communities of AirBnB properties.

### Research questions:

- What are the factors that differentiate AirBnB Business Ready (BR) listings from those who are not Business Ready?
- Is there an association between price and occupancy rates and whether the property has BR status?
- Can we predict the increase in BR listings based on projected increase in visitors traveling on business?
- Identify whether travellers utilise travel information rather than 'as the crow flies' distances to major points of interest?
- Do Australian BR listings adhere to AirBnB rules?
- Do BR listings draw more revenue?
- What features appeal to the BR demographic?
- How does the BR demographic affect social sentiment of the property?
- What kind of reviews appeal to the BR demographic?

### **Objectives**

- Improve understanding of the value in listing a property on AirBnB with a 'Business Ready' (BR) badge
- Identify the factors that distinguish BR listings and non-BR listings based on attributes of the property
- Identify appropriate 'Business Ready' badge qualifications for Australian business travellers

### 4 Research Method

Of the proposed research questions that can be answered through regression or classification techniques, one of particular interest is:

"What are the factors that differentiate AirBnB Business Ready (BR) listings from those who are not Business Ready?"

A range of potential modelling techniques could be used to determine the potential factors that differentiate business ready properties from those which are not business ready. Logistic regression is an appropriate analysis method as the dependent variable "is\_busines\_ready" is a binary. This method allows an informed statistician/researcher to interpret the results in a way that can be translated for non-technical stakeholders.

Refer to Appendices A-B for data acquisition and data merger code samples.

### Data selection

A range of potential datasets were identified during the preliminary EDA for this project. Each was assessed based on it quality, validity, accessibility and applicability to the research questions. The final stage of selection was to address privacy, ethical and legal considerations.

\*\*Inside AirBnB claims\*\* scrapes of publicly available information are not 'private' data however, the \*\*Australian privacy act\*\* refers to 'personal' rather than 'private'. Consequently, data such as the host's name in the AirBnB dataset should not be included in the analysis. These will be extracted from the listings for cross-referencing for removal during data-cleansing.

Property identification data cannot be removed as they are required for the analysis. For example, exact location and features. These will be retained but not revealed in any published findings.

### Selected datasets:

- Inside AirBnB scrapes of AirBnB listings, reviews and calendars
- TRA visitor and accommodation statistics
- Distances from listed property to points of interest
- Image features content of images in the listings

Refer to Appendix E – Complete list of datasets for list of a complete list datasets considered for this project.

### **Assumptions**

Project assumptions and the basis for forming them, refer to Table 1.

Table 1 – Project assumptions

	Assumption
Proxy for 'vacancy' – number of reviews per month.	A stay results in a review.
Accuracy of location data	Distance data collected is only specified at particular time of day and does not take into account hourly traffic patterns.
AirBnB dataset shows discrepancies	Scrape is an accurate representation of AirBnB listings
Image analysis - multiple suggestions for each detected object	Use item with highest level of confidence
Google distances	Validation of locations against free-text query Google Maps Locations is sufficient.
Focus on Sydney	This focus forms a representative basis for the findings to be extrapolated for Australia.

### **Techniques**

### **Data Acquisition**

- Google distance matrix (API) via "gmapsdistance" R package
- Azure Cognitive Services API (JSON over REST) web service call from R

- AWS Rekognition and S3 API
- Excel API selecting multi-sheet, targeted ranges from R

### **Data merging**

- Using GDrive cache to provide centralised datasets for externally sourced data
- R various munging techniques melt()

### 5 Challenges

Based on issues encountered during our preliminary exploratory data analysis, the following challenges have been identified:

### **Data volumes**

- Free tier services of three cloud platforms (AWS, Google Cloud and Microsoft Azure) were utilised by 3MDL during EDA, however Free Tier Limits placed tight restrictions on the data and images that could be processed. Modelling phase requires appropriate resources and funding.
- TRA datasets are published as large Excel workbooks with multiple tables and tabs. Data acquisition requires targeted selections within and across the tabs. Manual methods are not feasible.
- Sentiment analysis tools providing cost limitations on feature generation for review data

### **Acquisition / Merging datasets**

- Minimal commonality with public datasets hence may require broad assumptions
- Manual downloads required for TRA download site owing to file formats and internal structures

### Privacy/Ethics/Legality

 Data requiring removal from the AirBnB dataset occurs in structured and unstructured content however, owing to the nature of the research, some of this data is required for the analysis. Rationale for retention requires review and sign-off.

### 6 Timeline

The proposed timeline and next steps are outlined in Figure 2 below.

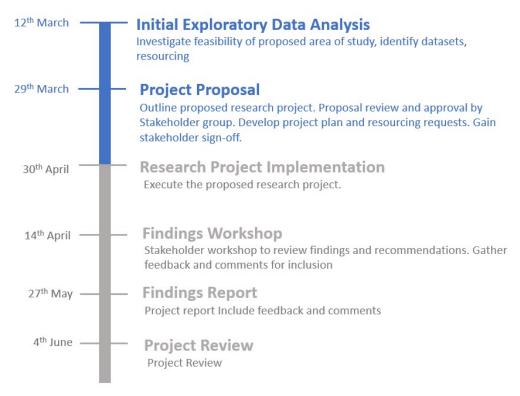


Figure 2 - Research project timeline

### Appendix A – Data Acquisition & Merge (code samples)

### 1A: AirBnB Listing, Review and Calendar Dataset

This code is used for fetching and caching locally the AirBnB data.

### **Snippet 1A:**

```
1. # Load gzip listings from web and cache locally to team drive
2. listings_url <- "http://data.insideairbnb.com/australia/nsw/sydney/2018-01-
    13/data/listings.csv.gz'
3. airbnb_gdrive_base_path <- "G:/Team Drives/STDS - Assignment 2 - 3MDL/Dataset/AirBnB/"</pre>
4. listings_filename <- "listings.csv"
listings_full_path <- paste(airbnb_gdrive_base_path, listings_filename, sep="")</li>
6.
7. listings <- data.frame()</pre>
8.
9. # If file already cached to g drive then load from there
10. if (file.exists(listings full path)) {
11. listings <- read_csv(listings_full_path)</pre>
12. } else {
13. # Otherwise get from URL and save to drive
14. url = getURL(listings url, encoding="gzip")
15.
      listings_table <- read.table(url)</pre>
16. write_csv(listings_table, listings_full_path)
17.
      listings <- listings_table</pre>
18. }
19.
20. # listings now contains a properly loaded
```

### 2A: GeoSpatial Distance from Tourist Points of Interest (POI) and Business Travel Destinations

The pull from GDM\_API file utilizes parallel processing to call the google distance matrix API to provide the distance (in meters) and time taken (in seconds) to travel between two points by a specified mode of transport (in the sample code chunk, this is from the listing to the Sydney Opera house by public transport). Parallel processing was necessary as there is a requirement to call the API one record at a time.

Midday Saturday (21 May) as used as an arbitrary date to take the observations. This may be extended to include other dates and times to sample from, however it must be a date in the future

### **Snippet 2A:**

```
1. ### Read Google Data###
2. #install.packages("gmapsdistance")
3. #devtools::install_github("rodazuero/gmapsdistance@058009e8d77ca51d8c7dbc6b0e3b622fb7f489a2")
4.
5. library(gmapsdistance)
6.
7. #The gmaps function is as follows
8. #gmapsdistance(origin, destination, combinations, mode, key,
9. #shape, avoid, departure, dep_date, dep_time,
10. #traffic_model, arrival, arr_date, arr_time)
11.
12. #Register do parallel so it is faster
13. library(doParallel)
14.
15. registerDoParallel(cores = detectCores()-1)
16.
```

```
17. #get the unique listing (note they are all unique)
18. Timedata1 <- unique(subset(listings, select = c("latitude", "longitude", "id")))</pre>
19.
20. #Test
21. #gmapsdistance(origin =paste0("-
   33.8","+,+","151"),destination = 'Sydney+Opera+House', departure = as.numeric(as.POSIXct("2018-04-
   21 12:00:00")), combinations = "all", mode = "transit", shape = "wide")
22.
23. ####Public Transport to Opera House####
24. #Opera Public Data <- foreach(i=1:nrow(Timedata)) %dopar% {
25. # library(gmapsdistance)
26.
27. # OperaTimePublic <- gmapsdistance(origin =paste0(Timedata$latitude[i],"+",Timedata$longitude[i]),d
   estination = 'Sydney+Opera+House', departure = as.numeric(as.POSIXct("2018-04-
   21 12:00:00")), combinations = "all", mode = "transit", shape = "wide")
28. #}
29.
31. #Opera_Public <- as.data.frame(matrix(unlist(Opera_Public_Data), ncol = 3,byrow = TRUE))
32. #Opera_Public <- cbind(Opera_Public, listings$id)</pre>
34. #Rename for Cleaning
35. #colnames(Opera_Public) <- c("Time", "Distance", "Status", "id")
37. #Change columns to correct format
38. #Time in Seconds
39. #Opera Public$Time <- as.numeric(as.character(Opera Public$Time))
41. #Distance in Meters
42. #Opera Public$Distance <- as.numeric(as.character(Opera Public$Distance))
44. #write_csv(Opera_Public, "G:/Team Drives/STDS - Assignment 2 - 3MDL/Dataset/AirBnB/Opera House by Pu
   blic Transport.csv")
45.
46. ####Define Grid####
47. # important places to go in sydeny are as follows:
48. # Opera House and Harbour Bridge
49. # Luna Park
50. # Bondi Beach
51. # Manly Beach
52. # Blue mountains (three sisters)
53. # Pokolbin (hunter valley)
55. #We will make the assumption that participants will either want to travel by car or by public transp
   ort, not by walking or cycling
57. #The date and time will be determined on the weekend
58. #Business Travel
59. # Convention centre
60. # Macquarie Park
61. # Bella Vista
62. # Paramatta
63. # Alexandria
64. # UNSW
65. # UTS
67. grid <- expand.grid(destination = c("Sydney+Opera+House", "Bondi+Beach", "Manly+Beach", "Three+Sisters
    ", "Pokolbin+NSW"), mode = c( "transit"))
68.
69. grid <- subset(grid, grid$destination!="Sydney+Opera+House"|grid$mode!="transit")
70.
71. Timedata <- merge(grid, Timedata1) #take every combination of the grid and Timedata
72.
73. ####Public Transport + Drive to all areas####
74. Total Data<- foreach(i=1:nrow(Timedata)) %dopar% {
75.
     library(gmapsdistance)
76.
     set.api.key('')
77.
```

```
TimeTotal <- gmapsdistance(origin =paste0(Timedata$latitude[i],"+",Timedata$longitude[i]),destina
   tion = Timedata$destination[i], departure = as.numeric(as.POSIXct("2018-05-
   21 12:00:00")),combinations = "all",mode = Timedata$mode[i], shape = "wide")
79. }
80. #Write to DF
81. Total Data DF <- as.data.frame(matrix(unlist(Total Data), ncol = 3,byrow = TRUE)) #turn to tidyform
82. Total_Data_DF_Final <- cbind(Total_Data_DF, Timedata) #Match with the other stuff
84. #Rename for Cleaning
85. #colnames(Opera Public) <- c("Time", "Distance", "Status", "id")</pre>
87. #Change columns to correct format
88. #Time in Seconds
89. #Opera Public$Time <- as.numeric(as.character(Opera Public$Time))
91. #Distance in Meters
92. #Opera_Public$Distance <- as.numeric(as.character(Opera_Public$Distance))
93.
94. #write_csv(Opera_Public, "G:/Team Drives/STDS - Assignment 2 - 3MDL/Dataset/AirBnB/Opera House by Pu
   blic Transport.csv")
```

### 3A: Locate required data in Excel Workbooks

This code is used to help identify the location of specific data in an excel workbook. Returns the identify the row and column number for specific strings

Snippet 3A - ## IMPORTANT: snippet only. Has dependencies, uses XLConnect package ##

```
3 # work out where the data is in the source workbook
5 mm <- as.matrix(readWorksheetFromFile("Accommodation/IVS1 YE Dec 2017_UpdatedMar2018.xlsx", sheet=2))
6 class(mm)<-"character" # convert all to character
7
8 # type search string in mm=="" below to work out where the data is in the source Workbook (returns row and column number)
9 rowcol <- which(mm=="VISITOR NIGHTS", arr.ind=T)
10 rowcol
11 # repeat as required
12 ### TO DO: read values from rowcol to feed into other functions
14 # read in blocks of data from the workbook
16 rs <- 8
17 re <- rs + 9
18 rs1 <- re + 3
19 re1 <- rs1 +9
20
21 B1 <- read.table(text=apply(mm[rs:re, 1:12],1,paste, collapse="\t"), sep="\t")
22 B2 <- read.table(text=apply(mm[rs1:re1, 1:12],1,paste, collapse="\t"), sep="\t")
23
24 rs2 <- re1 + 3
25 re2 <- rs2 + 9
26 B3 <- read.table(text=apply(mm[rs2:re2, 1:12],1,paste, collapse="\t"), sep="\t")
28 ### TO DO: loop through all these ... rs2 = re1 + 3 = rs2 + 9 for all data blocks
29 # have a look at what we've got
31 str(B1)
32 View(B1)
33 View(B2)
34 View(B3)
```

### 4A: TRA Visitor and Accommodation Statistics

This extracts specific data blocks across a range of worksheets in an Excel Workbook and loads into R. Resulting data is manipulated into a shape that can be merged with AirBnB dataset, including match on business / non-business.

### Code 4A:

```
1. # This file takes the tourism spreadsheet data and turns into a tidy dataset
2.
3. ## Note: these require JAVA to be installed. Using v8+ in this code
4.
5. # check java version ... returns error if JAVA not installed so go and get it!
6. system("java -version")
7.
8. # set environment variable for java - this is not required if already been added to o/s system variables
9. Sys.setenv(JAVA_HOME = "C:/Program Files/Java/jdk/")
10.
11. # url <- "https://www.tra.gov.au/ArticleDocuments/233/IVS1%20YE%20Dec%202017.xlsx.aspx"
12. # Need to get data manually and place in drive, since website does support direct access (instead stream
  s data through aspx page)
13.
14. # Install required packes if not already installed:
15. # XLConnect, XLConnectJars - called by XLConnect, reshape2, stringr, saRifx
16.
17. install.packages("reshape2")
                                       #has recast()
18. if(!"XLConnect" %in% rownames(installed.packages())) {
     install.packages("XLConnect")
19.
20.}
21.
22. if(!"reshape2" %in% rownames(installed.packages())) {
23.
     install.packages("reshape2")
24. }
25.
26. # Use XLConnect to deal with multiple worksheets
27. library(tidyverse) # mutate
28. library(XLConnect) # excel
29. library(reshape2)
                       # melt
30. library(stringr)
                       # regex
31. library(taRifx)
                        # destring
32.
33. tourism base path <- "G:/Team Drives/STDS - Assignment 2 - 3MDL/Dataset/Accommodation/"
34. tourism filename <- "IVS1 YE Dec 2017 UpdatedMar2018.xlsx"
35. tourism_dest_filename <- "sydney_tourism_201803.csv"
36. tourism_fullpath <- paste(tourism_base_path, tourism_filename, sep="")</pre>
37. tourism_dest_fullpath <- paste(tourism_base_path, tourism_dest_filename, sep="")</pre>
38.
39. # Load tourism workbook
40. tourism_workbook <- loadWorkbook(tourism_fullpath)
41.
42. # Get sheet names
43. sheet_names <- getSheets(tourism_workbook)</pre>
44.
45. # We are only interested in country sheets
46. sheet_names <- sheet_names[!sheet_names %in% c("Contents", "Total")]
47.
48. # Spreadsheet configuration details
49.
50. # There are 7 visit dimensions located in Cell A8 and increasing by 12 rows with last dimensions in Cell
    Δ80
51. sheet block range = 1:7
52. sheet_year_range = 2007:2017
53.
54. visitor dimension startRow = 8
55. visitor dimension col = 1
56. visitor dimension nextRowStep = 12
57.
```

```
58. # Each visit dimension has its own associated block of data
59. data startRow = 8
60. data startCol = 1
61. data endRow = 17
62. data endCol = 12
63. data nextRowStep = 12
64.
65. # Helper functions
66.
67. # Visit dimension name cleaned for friendly format i.e. friendly feature name
68. # features are currently of format UPPER.WORDS.AND.SO..ON.(000)
69. cleanVisitDimensionName <- function(s) {</pre>
70. s %>%
       str_replace_all("\\.", " ") %>%
str_replace_all(" ", " ") %>%
str_extract("[A-Z ]+") %>%
71.
                                            # Replace all . with spaces
72.
                                        # Collapse double spaces
73.
                                            # Get words only
74.
        trimws() %>%
                                            # Remove leading and trailing whitespace
        str_replace_all(" ", "_") %>%
75.
                                            # Replace spaces with underscore
76.
        tolower()
                                            # Lower case
77.}
78.
79. # Get Visit Dimension Name from sheet
80. getVisitDimensionName <- function(workbook, sheet, block_index) {</pre>
     visit_dimension_name = readWorksheet(workbook, sheet,
81.
82.
                                             startRow = visitor_dimension_startRow + ((block_index-

    * visitor dimension nextRowStep),

83.
                                             startCol = visitor_dimension_col,
                                             endRow = visitor_dimension_startRow + ((block_index-
84.

    * visitor dimension nextRowStep),

                                             endCol = visitor dimension col)
85.
86.
      cleanVisitDimensionName(names(visit dimension name))
87.}
88.
89. # Get Visit Dimension Block Data from sheet
90. getVisitDimensionBlock <- function(workbook, sheet, block_index) {</pre>
91.
     readWorksheet(workbook, sheet,
92.
                    startRow = data startRow + ((block index-
   1) * data nextRowStep), startCol = data startCol,
                    endRow = data_endRow + ((block_index-1) * data_nextRowStep), endCol = data endCol)
93.
94. }
95.
96. # Final Tourism data frame
97. tourism.data = data.frame()
98.
99. # Build dataset across all sheets
100. for (sheet in sheet names){
101.
102.
         # Context data for current sheet
103.
         sheet_data <- data.frame()</pre>
104.
105.
         # Loop through each block in the sheet as they are iterated over each visit dimension
106.
         for (block_index in sheet_block_range) {
107.
            # get sheet data for visit dimension Metric
108.
109.
            sheet_block_data <- getVisitDimensionBlock(tourism_workbook, sheet, block_index)</pre>
110.
111.
            # Change names of features to "purpose" and from 2007, to 2017
112.
            names(sheet_block_data) <- c("purpose", sheet_year_range)</pre>
113.
            # visit dimension name, want in format lower case words seperated by underscore
114.
115.
            visit_dimension_name = getVisitDimensionName(tourism_workbook, sheet, block_index)
116.
117.
            # Melt the years features into a single year feature
118.
            sheet block data %>%
119.
              melt(id.vars=c("purpose"),
                   variable.name="year",
120.
121.
                   value.name=visit_dimension_name) -> sheet_block_data
122.
```

```
# Each block is a feature, inner join on purpose and year
124.
           if (nrow(sheet_data) == 0) {
125.
              sheet data <- sheet block data # first time
126.
            } else {
127.
              sheet data <- merge(sheet data, sheet block data, by=c("purpose","year"))</pre>
128.
129.
         }
130.
131.
         # add country name, this is the current sheet name
132.
         sheet data %>%
133.
           mutate(country = as.factor(sheet)) -> sheet data
134.
135.
         # need to append now this dataset
136.
         tourism.data <- rbind(tourism.data, sheet data)</pre>
137.
       }
138.
139.
       # Clean Tourism.data
140.
141.
       # Strip leading and trailing whitespace
142.
       tourism.data$purpose <- trimws(tourism.data$purpose)</pre>
143.
144.
145.
       tourism.data <- tourism.data[!(tourism.data$purpose == "Total"),]</pre>
146.
147.
       # Strip backpacker data - this data duplicates the by reason data
148.
       ### Consider commenting these lines and split this into separate data frame
       tourism.data <- tourism.data[!(tourism.data$purpose == "Backpackers"),]</pre>
149.
       tourism.data <- tourism.data[!(tourism.data$purpose == "Non backpackers"),]
150.
151.
152.
       # Convert to factor
153.
       tourism.data$purpose <- as.factor(tourism.data$purpose)</pre>
154.
155.
       # Clean metric data, and put as proper type
156.
       dim index <- which(names(tourism.data) %in% c("purpose", "year", "country"))</pre>
       metric columns <- names(tourism.data[-dim index])</pre>
157.
158.
       for (col in metric columns) {
159.
160.
         tourism.data[,col] <- as.numeric(destring(tourism.data[,col])) # convert to number, and set np to
   NΑ
161.
       }
162.
163.
       # Check structure and content are as expected
164.
       str(tourism.data)
165.
       head(tourism.data)
166.
167.
       # the sort the rows so values are in a predicable order
168.
       tourism.data <- select(tourism.data, purpose, country, year, everything())</pre>
169.
       View(tourism.data)
170.
       # Write data to shared drive
171.
172.
       write.csv(tourism.data, file=tourism_dest_filename)
173.
174.
       # Summarize totals for business and non business
175.
176.
       # create is_business flag which is set if visit purpose was Business Travel or Employment
177.
       tourism.data %>%
178.
         mutate(is_business = ifelse(purpose %in% c("Employment","Business Travel"), 1, 0)) -
   > tourism.data
179.
180.
       # Sum across all metric columns by Business, Year and Country.
181.
       aggregate(tourism.data[metric columns],
182.
                  bv=list(
183.
                    year=tourism.data$year,
                    is business=tourism.data$is business,
184.
185.
                    country=tourism.data$country), FUN=sum, na.rm=TRUE)
```

### **5A: AirBnB Listing Profile Image Analysis**

This function is intended to get sentiment data for the reviews from AirBnb. The purpose is to feature engineer a new feature against each AirBnB listing to understand the overall sentiment against that listing. API Notes: The API is a RESTFul webservice using JSON as data format exchange. There is a max of 5k per month uses. Need to distribute around team to run.

### Snippet 5A - \*\* The following is more of a process instruction rather than an r script\*\*

```
5 # Extract id andpicture url from listings file
6 #id picture url
7 #13467416 https://a0.muscache.com/im/pictures/778d441e-2c94-431a-8307-b32eaec69930.jpg?aki_policy=large
8 #15882089 https://a0.muscache.com/im/pictures/3e9e50cf-2949-4a9e-9782-fecf327ea75a.jpg?aki_policy=large
10 ## Remove '?aki_policy=large' on all picture url to the image direct link
11 #id picture url
12 #13467416 https://a0.muscache.com/im/pictures/778d441e-2c94-431a-8307-b32eaec69930.jpg
13 #15882089 https://a0.muscache.com/im/pictures/3e9e50cf-2949-4a9e-9782-fecf327ea75a.jpg
15
16 #Process the picture to get the context of each photo to get the context confidencelevel
17 #id context confidencelevel
18 #13467416 Building 0.923
19 #13467416 Housing 0.923
20 #13467416 Apartment Building 0.547
21 #13467416 Architecture 0.512
22 #15882089 Building 0.962
23 #15882089 Housing 0.962
32 #....
33
34 ## Analysis the output and find the most appeared context in this case 'Building' and 'Housing', and some interesting features,
high rise 35 #id context confidencelevel Count of context
36 #13467416 Building 0.923 4
37 #15882089 Building 0.962 4
40 #13467416 Housing 0.923 4
41 #15882089 Housing 0.962 4
44 #13467416 Architecture 0.512 2
45 #13467416 Apartment Building 0.547
53 #15882089 High Rise 0.52 2
54 #...
57
58 ## Convert the most appeared words into features for adding back to the lising
59 #id has building has housing has high rise has city
60 #13467416 1 1 1 1
61 #15882089 1 1 1 1
62 #16197117 1 1 0 0
63 #22298997 1 1 0 0
65 ## R code for Merge new feature columns back to the main listings list for regression analysis
66 library(tidyverse)
68 image <- read.csv("image.csv")
69 listings <- read.csv("listings.csv")
71 listings2 <- merge(listings, image)
72
```

### Appendix B – Data Merging code samples

### 1B: Merging AirBnB Listings dataset with Google Distance Matrix

The Google Distance Matrix dataset was constructed using the longitude and latitude from the listings file as the "origin" for the API GET call. Because of this, we were able link the listing ID directly to each API call. Thus a merge between the two files is relatively straightforward.

### **Snippet 1B:**

```
1. # Load Libraries
2. library(tidyverse)
3. library(httr)
                        # For constructing POST request
4. library(jsonlite) # For converting JSON to data frame
5. library(data.table)
6. library(dplyr)
8. # Params for API update
9. # Can only run 5000 Queries per month, so each person needs to update the start and end to distribu
   te.
10. \text{ start} = 0
11. end = 5000
12.
13. # Get Review Data
14. reviews_filename = 'G:/Team Drives/STDS - Assignment 2 - 3MDL/Dataset/AirBnB/reviews.csv'
15. reviews <- read_csv(reviews_filename)</pre>
16.
17. df <- reviews[start:end,c('id','comments')]</pre>
19. # Azure API expects request format params of 'id' and 'text'
20. df <- rename(df, text = comments)</pre>
21. data <-list(documents=df)</pre>
22.
23. # Call Sentiment API
24. api_key = 'xxx'
25. sentiment_api_url = "https://westcentralus.api.cognitive.microsoft.com/text/analytics/v2.0/sentimen
26.
27. # Call APT
28. response <- POST(sentiment_api_url,</pre>
                     add headers(`Ocp-Apim-Subscription-Key`=api key),
29.
30.
                     body=toJSON(data))
31. # Get Response
32. response_content <- content(response, as="text", encoding = "utf-8")
34. # convert JSON to data.frame
35. fromJSON(response_content)$documents %>%
36. mutate(id=as.numeric(id)) ->
37.
     responses
39. # Create table with Review ID and Sentiment Score only
40. df %>% left_join(responses, by = 'id') %>% select(c('id', 'score')) -> dfs
42. # Write out which records of dataset were done so not overlapping
43. write_csv(dfs, sprintf("data\\review_sentiment_start_%s_end_%s.csv",start,end))
```

### 2B: AirBnB Reviews: Feature Construction with Sentiment Analysis

Sentiment Analysis was proposed using Microsoft Azure Cognitive Services, Indico.ai Sentiment API, and potentially others. This was because all API service offerings had a limited amount of API calls per month. This EDA acquisition was abandoned as the AirBnB review dataset was 335k+ records. Acquisition costs would be in the range of \$1500+ AUD. It is recommended to re-attempt with an on-premise sentiment analysis library.

The merge approach was to pass the review ID and associated review text to the API and then merge by Review ID to the original dataset.

### **Snippet 2B:**

```
    # Load Libraries

library(tidyverse)
                        # For constructing POST request
3. library(httr)
4. library(jsonlite) # For converting JSON to data frame
5. library(data.table)
6. library(dplyr)
8. # Params for API update
9. # Can only run 5000 Queries per month, so each person needs to update the start and end to distribu
   te.
10. \text{ start} = 0
11. end = 5000
12.
13. # Get Review Data
14. reviews_filename = 'G:/Team Drives/STDS - Assignment 2 - 3MDL/Dataset/AirBnB/reviews.csv'
15. reviews <- read_csv(reviews_filename)
16.
17. df <- reviews[start:end,c('id','comments')]</pre>
18.
19. # Azure API expects request format params of 'id' and 'text'
20. df <- rename(df, text = comments)</pre>
21. data <-list(documents=df)</pre>
22.
23. # Call Sentiment API
24. api_key = 'xxx'
25. sentiment_api_url = "https://westcentralus.api.cognitive.microsoft.com/text/analytics/v2.0/sentimen
   t"
26.
27. # Call API
28. response <- POST(sentiment_api_url,</pre>
                     add headers(`Ocp-Apim-Subscription-Key`=api key),
29.
30.
                     body=toJSON(data))
31. # Get Response
32. response_content <- content(response, as="text", encoding = "utf-8")</pre>
33.
34. # convert JSON to data.frame
35. fromJSON(response_content)$documents %>%
36. mutate(id=as.numeric(id)) ->
37.
     responses
38.
39. # Create table with Review ID and Sentiment Score only
40. df %>% left_join(responses, by = 'id') %>% select(c('id', 'score')) -> dfs
42. # Write out which records of dataset were done so not overlapping
```

### Appendix D – Proposed Google Directions Locations

The proposed locations of interest to be used in the Derived Distances dataset are listed in Table 2 below. The number of actual locations that can be processed is dependent on access to appropriate resources. Final selection will be determined prior to commencing this research project.

Table 2 - Points of Interest

Point of Interests -1	Precinct	LGA 👱	Suburb	Postcoc*
Business Park	Airgate Business Park	Sydney	Mascot	2020
Business Park	Auslink Business Park	Northern Beaches	Belrose	2085
Business Park	Australian Technology Park	Sydney	Eveleigh	2015
Business Park	Botany Industrial Park	Sydney	Hillsdale	2036
Business Park	Box Hill Business Park	The Hills	Box Hill	2765
Business Park	Central Hills Business Park	Camden/Liverpool	Narellan	2567
Business Park	City Views Business Park	Ryde	Macquarie Park	2113
Business Park	Forestridge Business Park	Northern Beaches	Frenchs Forest	2086
Business Park	Gateway Business Park	Parramatta/Canada Bay	Silverwater	2148
Business Park	Greystanes Business Park	Cumberland	Greystanes	2145
Business Park	Heritage Business Park	Sydney	Mascot	2020
Business Park	Huntley Business Park	Sydney	Alexandria	2015
Business Park	Kingsgrove Business Park	Canterbury-Bankstown	Kingsgrove	2208
Business Park	Lakes Business Park	Sydney	Botany	2019
Business Park	Lane Cove Business Park	Ryde	Lane Cove West	2066
Business Park	Leppington North Business Park	Camden/Liverpool	Leppington North	2179
Business Park	Link Business Park	Ryde	Macquarie Park	2113
Business Park	Newington Business Park	Parramatta/Canada Bay	Newington	2127
Business Park	Norwest Business Park	The Hills	Baulkham Hills	2153
Business Park	Orion Business Park	Svdnev	Sydney	2066
Business Park	Potts Hill Business Park	Canterbury-Bankstown	Potts Hill	2143
Business Park	Prince Henry Centre	Sydney	Little Bay	2015
Business Park	Riverstone West	Blacktown	Riverstone West	2765
Business Park	Ryde Road, Pymble	Ku-ring-gai	Pymble	2073
Business Park	St Peters Business Park	Sydney	St Peters	2044
Business Park	Sydney Business Park	Blacktown	Marsden Park	2765
Business Park	Sydney Corporate Park	Sydney	Alexandria	2015
Business Park	Sydney Olympic Park/Rhodes	Parramatta/Canada Bay		2127
Business Park	Sydney Park Pavilion	Sydney	Alexandria	2015
Business Park	Sydney Science Park	Penrith	Luddenham	2745
Business Park	Talavera Business Park	Ryde	Macquarie Park	2113
Business Park	Warriewood Business Park	Northern Beaches	Warriewood	2102
Business Park	Waterloo Business Park	Ryde	Macquarie Park	2113
Business Park	West Pennant Hills Business Park	Hornsby	West Pennant Hills	2125
	All Sorts Function Centre	Sydney	Alexandria	2015
	Christie Spaces and Conferencing Centre	Sydney	North Sydney	2060
Conference Centre		Sydney	Sydney	2000
	Figtree Conference Centre	Parramatta/Canada Bay	, ,	2127
	Greater Sydney Conference	Ryde	Epping	2121
	International Convention Centre Sydney	Sydney	Sydney	2000
Conference Centre		Sydney	The Rocks	2000
	Northside Conference Centre	Sydney	St Leonards	2065
	Sydney Conference & Training Centre	-,,		
	Wesley Conference Centre	Northern Beaches	Ingleside Sydney	2101
		Sydney		2135
University	Australian Catholic University	Sydney	Strathfield	
University	Charles Darwin University	Sydney	Haymarket	2043
University	CQUniversity	Sydney	Sydney	2000
University	Curtin University	Sydney	Sydney	2000
University	Macquarie University	Sydney	North Ryde	2109
University	The University of Notre Dame Australia	Sydney	Chippendale	2008
University	University of New South Wales	Sydney	Kensington	2052
University	University of Sydney	Sydney	Camperdown/Darling	2006
University	University of Technology Sydney	Sydney	Ultimo	2007
University	Western Sydney University	Parramatta/Canada Bay	Parramatta	2150

### Appendix E – Complete list of datasets

### Selected

More information about the datasets selected for this project are detailed in Table 3 below.

Table 3 - Datasets (selected and to be confirmed)

Source	Data set	Basis for Selection
Inside AirBnB  http://insideairbnb.com/get- the-data.html	<ul><li>listings</li><li>reviews</li><li>calendar</li></ul> grain: month	Richness of dataset and applicability to the research questions.
Travel Research Australia (TRA) https://www.tra.gov.au/  Accessed under license issued to UTS Libary. http://databases.lib.uts.edu.au/search_databases.py?search=tourism+research+australia+database	International Visitor statistics - trips, nights, average stay and spend by	Specific subsets TRA published datasets based on commonality with available AirBnB data and applicability to the research questions.  Statistics published by Tourism Research Australia (TRA) will be used for industry benchmarking and validate these findings.
Derived distances <a href="https://developers.google.com/maps/documentation/distance-matrix/">https://developers.google.com/maps/documentation/distance-matrix/</a>	<ul> <li>distance</li> <li>time</li> <li>transport mode</li> <li>grain: latitude / longitude</li> </ul>	Points of interest to:

### To Be Confirmed

Source	Data set	Basis for Selection
CoreData (Industry Partners) <a href="https://coredata.com.au/">https://coredata.com.au/</a>	short market research questionnaire about business travel	data set will provide context about what business travellers look for when selecting a place to stay
Australian Bureau of Statistics <a href="http://www.abs.gov.au/AUSSTATS/abs@.nsf/DetailsPage/863">http://www.abs.gov.au/AUSSTATS/abs@.nsf/DetailsPage/863</a> <a href="mailto:5.02015-16?OpenDocument">5.02015-16?OpenDocument</a>	Tourist Accommodation, Australia, 2015-16 Establishments, rooms, bed spaces, occupancy, takings	Specific TRA published datasets were selected based on commonality with available AirBnB data and applicability to the research questions.

### Rejected

The datasets that were investigated but not selected for this research project are detailed in Table 4. Reasons for rejection include insufficient granularity and/or keys for linking with primary datasets, high level of sparsity, prohibitive volumes.

Table 4 - Datasets (rejected)

Tom Slee (scrapes of AirBnB listings)	<ul> <li>cleaner but less rich dataset</li> <li>reduced options for investigation</li> <li>dataset will not be updated in the future - reduced value for future research</li> </ul>
Image analysis	funding dependent
Rental price (benchmark)	insufficient time for acquisition
Obike data (impact on community)	<ul> <li>insufficient time for acquisition</li> <li>tenuously related to research questions</li> </ul>
Weather Data	<ul> <li>insufficient time for acquisition</li> <li>not strongly related to research questions</li> </ul>