# Appendix
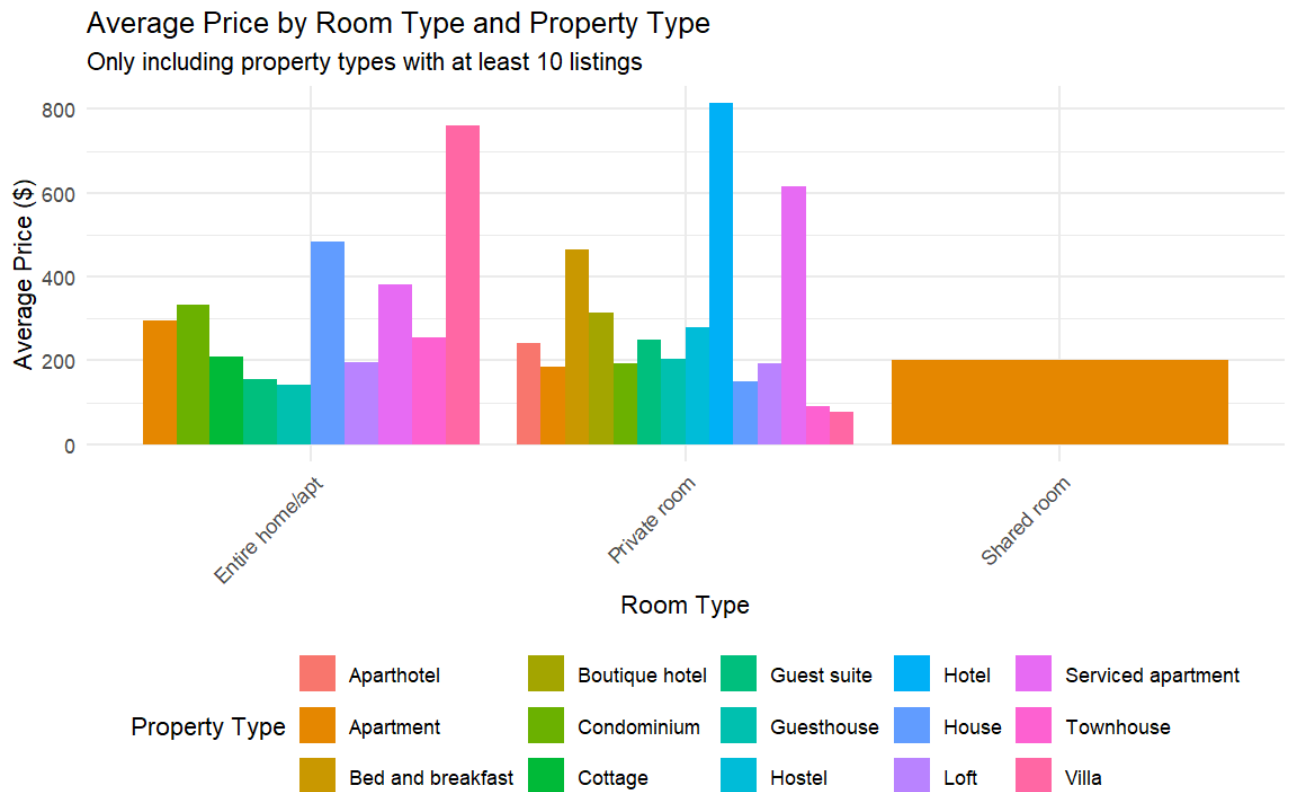
Appendix 1: Description of the "Average Price by Room Type and Property Type" visualization
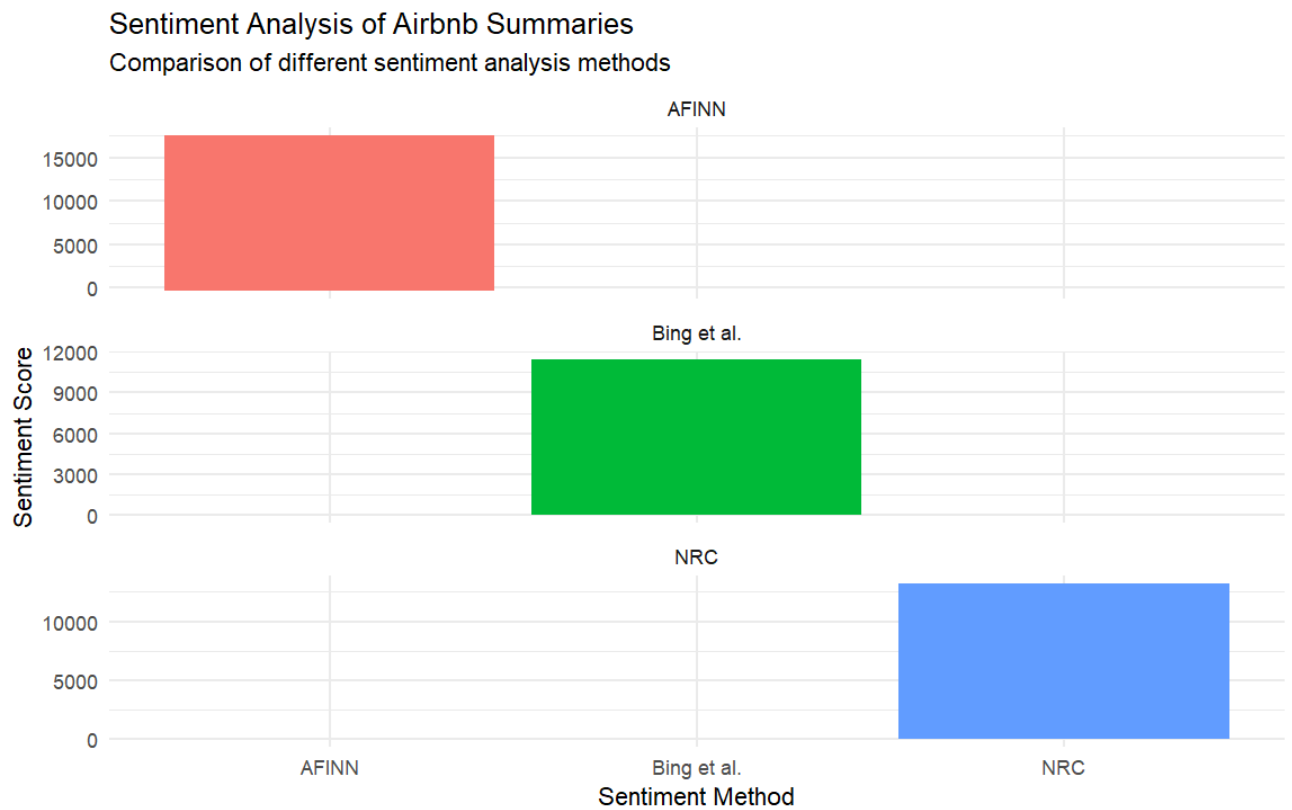


**Average Price by Room Type and Property Type**
Only including property types with at least 10 listings

Property Type legend:
- Aparthotel
- Apartment
- Bed and breakfast
- Boutique hotel
- Condominium
- Cottage
- Guest suite
- Guesthouse
- Hostel
- Hotel
- House
- Loft
- Serviced apartment
- Townhouse
- Villa

**Axes:**

- **X-axis**: Represents Room Type. The three room types are "Entire home/apt," "Private room," and "Shared room."

- **Y-axis**: Represents Average Price ($). The scale ranges from 0 to 800.

**Data Representation**:

- Each bar represents a specific combination of room type and property type.

- The height of each bar corresponds to the average price for that combination.

- Bars are color-coded to distinguish between different property types, with a legend provided beneath the chart.

Appendix 2: Description of the sentiment analysis visualization for the Airbnb dataset:

## Sentiment Analysis of Airbnb Summaries
### Comparison of different sentiment analysis methods



The y-axis represents the **Sentiment Score**, and the x-axis represents the **Sentiment Method**.

Appendix 3: Description of the "Correlation Heatmap of Numeric Variables" visualization for the

Airbnb dataset:



Correlation Heatmap of Numeric Variables

**Cells**: Each cell at the intersection of a row and a column represents the correlation coefficient between the two corresponding variables.

**Color Coding**: The cells are colored according to the strength and direction of the correlation:

- **Red** indicates a positive correlation. The intensity of the red color represents the strength of the positive correlation (darker red = stronger positive correlation).

- **White**: Indicates a correlation close to zero (no or very weak correlation).

- **Blue** Indicates a negative correlation. The intensity of the blue color represents the strength of the negative correlation (darker blue = stronger negative correlation).

**Color Scale**: A color scale (legend) is provided to the right of the heatmap, mapping the color intensities to correlation values ranging from -1.0 to 1.0.

Appendix 4: "Superhost vs Non-Superhost Performance: Comparison across key performance metrics", comparing the performance of Airbnb listings based on whether the host is a Superhost (TRUE) or not (FALSE)
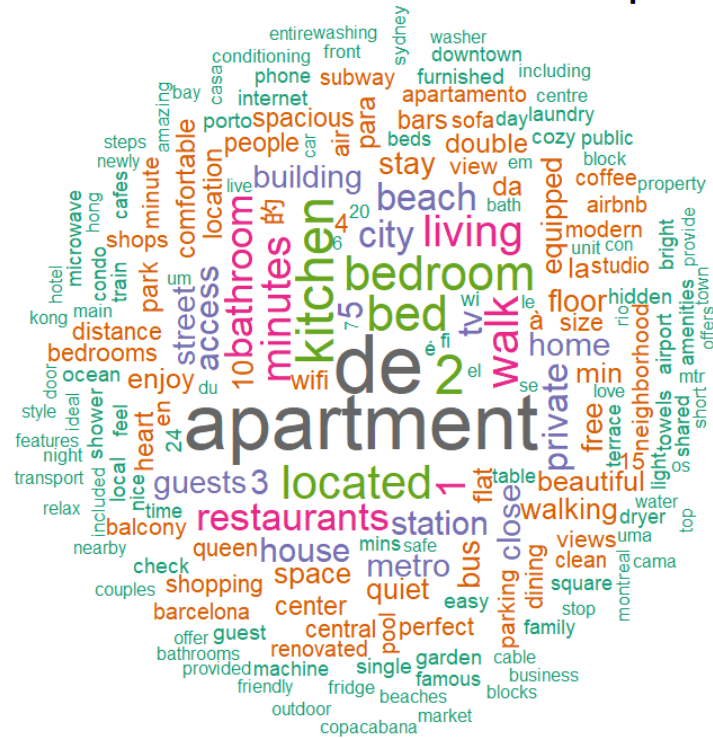
Superhost vs Non-Superhost Performance
Comparison across key performance metrics

| Average Number of Reviews | Average Price ($) | Average Review Score |

**X-axis**: Superhost Status (FALSE/TRUE).

*Average Number of Reviews:* **Y-axis**: *Average Number of Reviews.* The turquoise bar is significantly higher than the coral bar, indicating that listings managed by Superhosts tend to have a much higher average number of reviews than those managed by non-Superhosts.

*Average Price:* **Y-axis**: *Average Price.* On average, the listings managed by non-Superhosts have a higher price than those managed by Superhosts, possibly to break even owing to fewer customers.

*Average Review Score:* **Y-axis**: *Average Review Score.* The bar graph suggests that listings managed by Superhosts generally have higher average review scores.
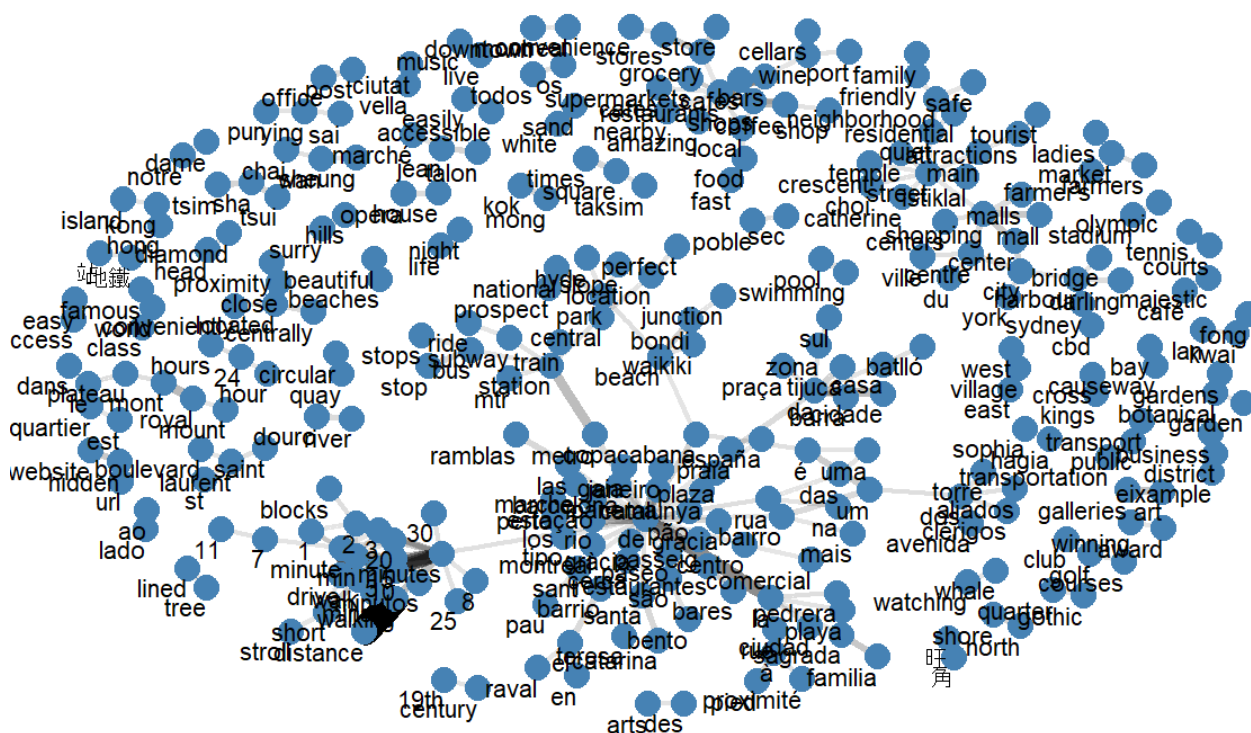
Appendix 5: **Word Cloud** titled "Most Common Words in Airbnb Descriptions."

**Most Common Words in Airbnb Descriptions**



**Dominant Words**: Due to its size, the most prominent word is "apartment." Other significant words include "minutes," "located," "restaurants," "kitchen," and "bedroom." These suggest that descriptions often focus on property type, location, amenities, and the number of bedrooms.

**Descriptive Adjectives**: Adjectives like "beautiful," "cozy," "comfortable," "modern," and "quiet" give insight into how hosts are describing their properties to attract customers.

Appendix 6: **Bigram Network of Neighborhood Descriptions**

Bigram Network of Neighborhood Descriptions
Showing bigrams appearing at least 15 times



**Nodes:** Each word or phrase (bigram) is represented as a node (a blue circle). The label of each node is the bigram itself.

**Edges:** Lines connect the nodes, representing the co-occurrence of these words or phrases within the neighborhood descriptions. The thickness of the line suggests the frequency of their co-occurrence, with thicker lines indicating more frequent pairings.

**Layout:** The nodes are arranged in a force-directed layout where connected nodes are drawn closer together and disconnected nodes are pushed further apart. This allows clusters of related terms to emerge.

Appendix 7: Horizontal bar chart titled "Top 20 Bigrams in Neighborhood Descriptions".



**Y-axis:** The y-axis lists the 20 most frequent bigrams. Each bigram is a pair of words often appearing together in the neighborhood descriptions. Examples include "walking distance," "minute walk," "restaurants bars," and location-specific phrases like "hong kong," "rio de," and "de janeiro." The bigrams are ordered from top to bottom according to frequency (highest to lowest).

**X-axis:** The x-axis represents the frequency or count of each bigram in the dataset. The scale extends from 0 to 400.

**Bars:** Each bigram has a horizontal bar extending from the y-axis to a point corresponding to its frequency on the x-axis. The bars are all colored in blue.
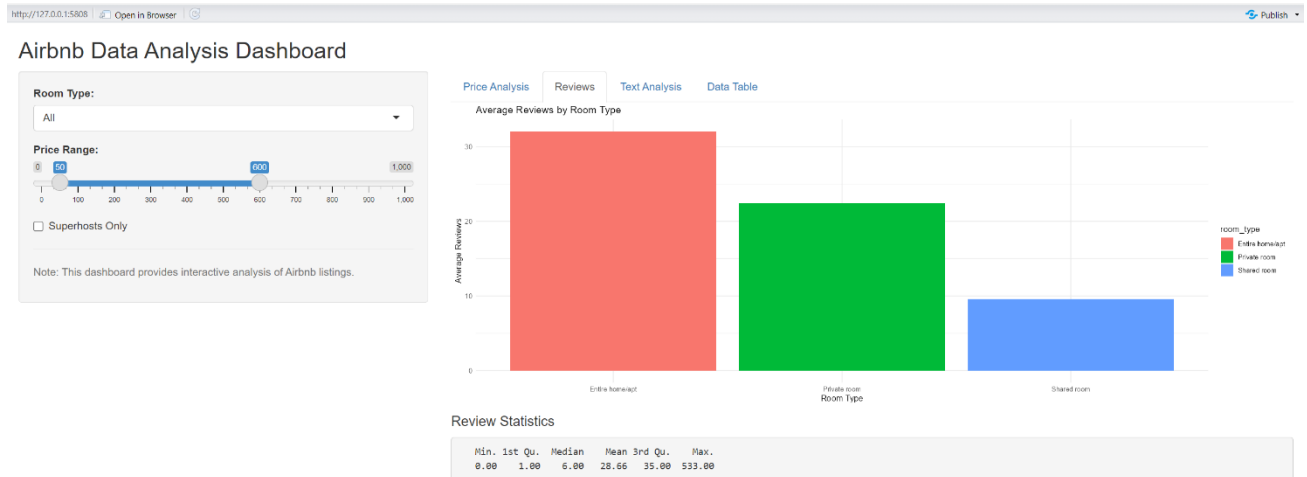
Appendix 8.A: Description of the Airbnb Data Analysis Dashboard
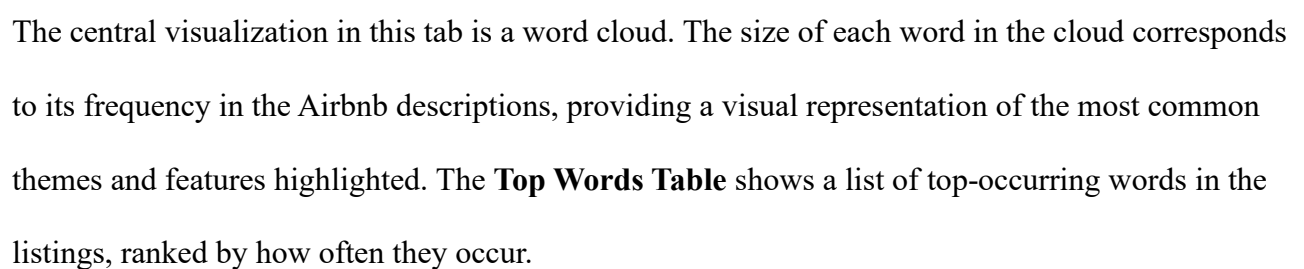
*Price Analysis Tab:*



This tab features a bar chart comparing average prices by room and property types. The chart allows users to identify which combinations command the highest prices quickly. Below the chart are price statistics, such as the minimum, first quartile, median, mean, third quartile, and maximum listing prices.
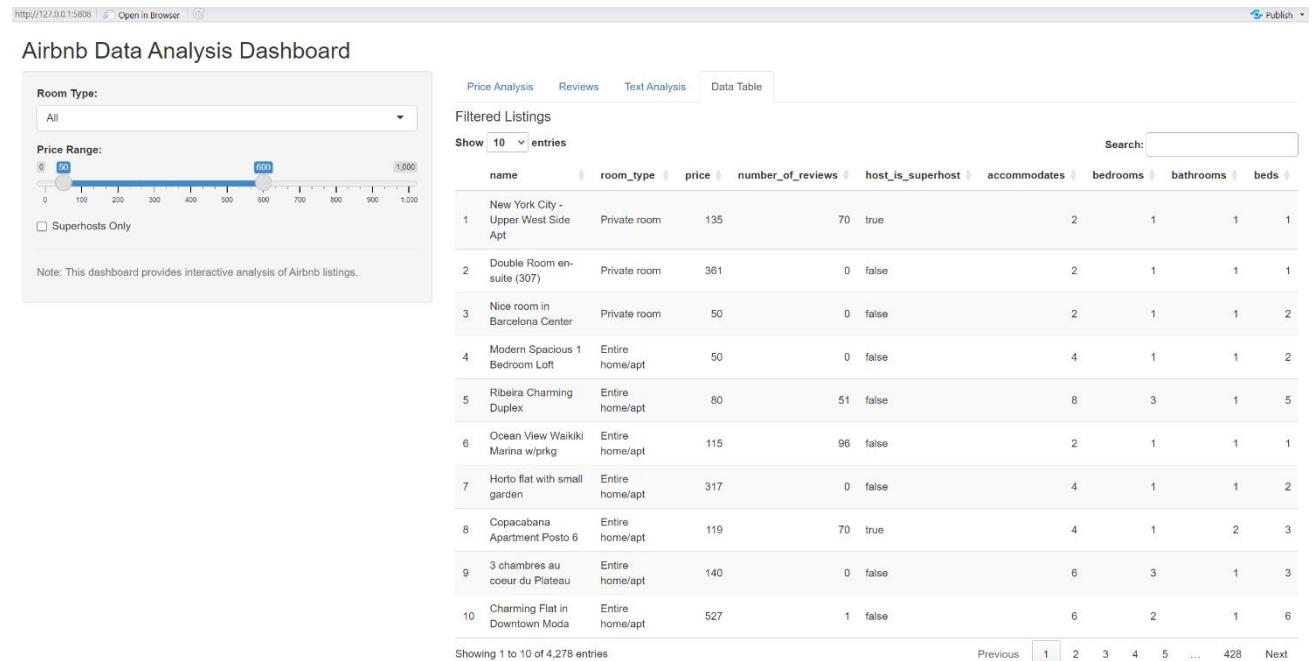
*Reviews Tab:*



The section displays a bar chart showing the average number of reviews across all room types. The

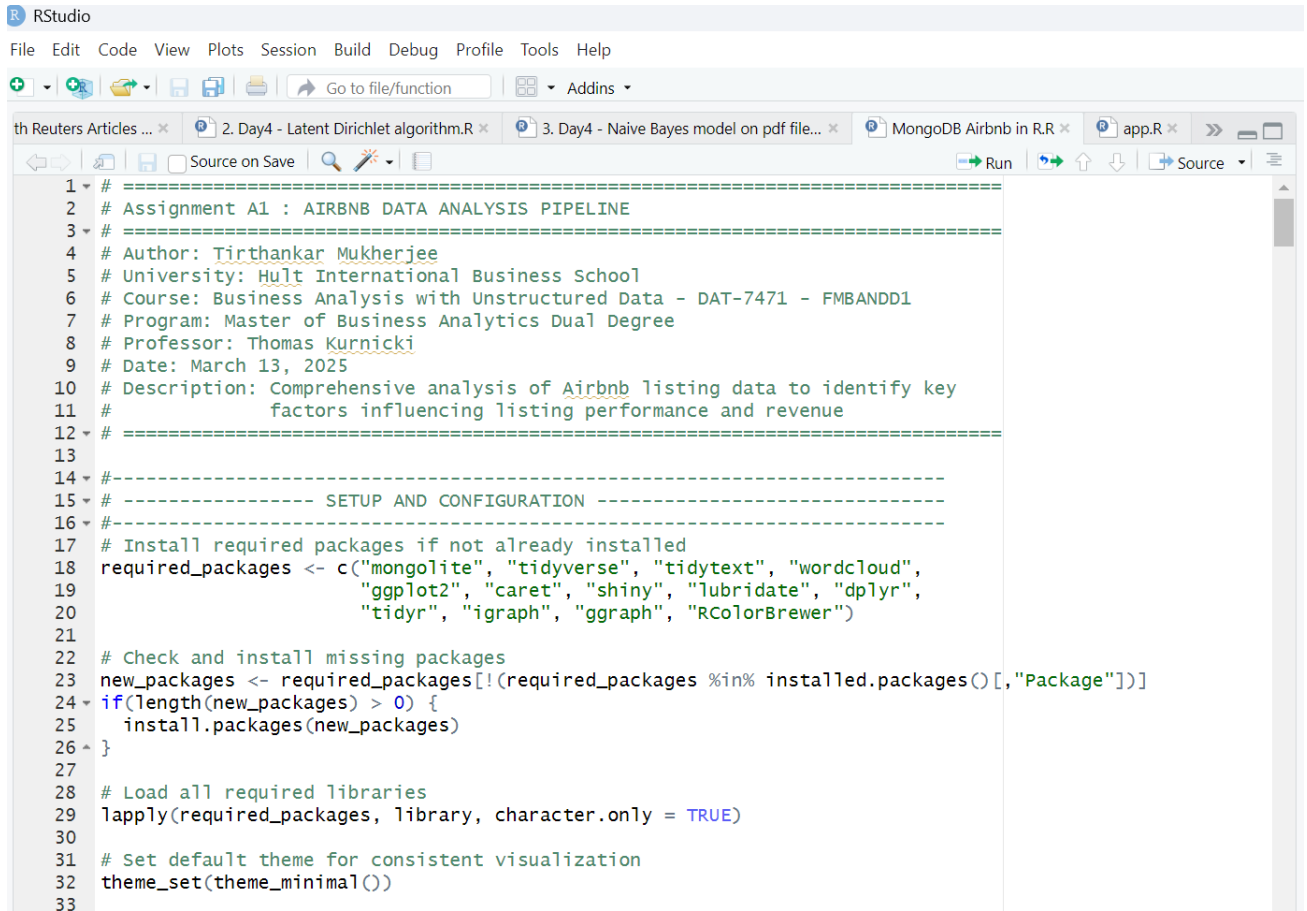Review statistics present descriptive statistics about the reviews.

Appendix 8.B: Description of the Airbnb Data Analysis Dashboard

*Text Analysis Tab:*



The central visualization in this tab is a word cloud. The size of each word in the cloud corresponds to its frequency in the Airbnb descriptions, providing a visual representation of the most common themes and features highlighted. The **Top Words Table** shows a list of top-occurring words in the listings, ranked by how often they occur.

*Data Table Tab:*



Displays raw data in a tabular form, allowing users to browse individual Airbnb listings. The table includes key listing attributes like name, room type, price, number of reviews, and host status.

**Controls** for the data table include filtering and search capabilities, as well as pagination to navigate through multiple pages of listings.

Appendix 9: R-Script

Appendix 9.A: Setup and Configuration



In this code snippet, the required packages are installed and loaded, which are necessary for the

analysis.

Appendix 9.B: Data Connection and Loading

```r
#-----------------------------------------------------------------------
# ---------------- DATA CONNECTION AND LOADING ------------------------
#-----------------------------------------------------------------------
# MongoDB connection details
# Note: In production, consider storing credentials in environment variables
connection_string <- 'mongodb+srv://mukherjeetirthankar:tirthA25@a1textmining.9gdz0.mongodb.net/?retryWr

# Initialize MongoDB connection
airbnb_collection <- mongo(
  collection = "listingsAndReviews",
  db = "sample_airbnb",
  url = connection_string
)

# Function to clean and preprocess Airbnb data
clean_airbnb_data <- function(data) {
  # Data cleaning and type conversion
  data %>%
    mutate(
      # Convert currency string to numeric
      price = as.numeric(gsub("[\\$,]", "", price)),

      # Convert date strings to Date objects
      last_review = as.Date(last_review),
      first_review = as.Date(first_review),

      # Calculate duration between first and last review in days
      review_gap = as.numeric(difftime(last_review, first_review, units = "days")),

      # Extract and clean host response rate
      host_response_rate = as.numeric(gsub("%", "", host$host_response_rate)) / 100,

      # Convert listing details to numeric
      accommodates = as.numeric(accommodates),
      bedrooms = as.numeric(bedrooms),
      bathrooms = as.numeric(bathrooms),
      beds = as.numeric(beds),

      # Extract additional host information
      host_is_superhost = host$host_is_superhost,
      host_identity_verified = host$host_identity_verified,
      host_listings_count = as.numeric(host$host_listings_count)
    )
}

# Fetch and preprocess all data
# Warning: This operation may require significant memory
airbnb <- airbnb_collection$find() %>% clean_airbnb_data()
```

In this code snippet, a MongoDB connection is initiated to load the dataset into R.

## Appendix 9.C: Exploratory Data Analysis

```
 87 ▾  #---------------------------------------------------------------------
 88 ▾  #------------------ EXPLORATORY DATA ANALYSIS ------------------------
 89 ▾  #---------------------------------------------------------------------
 90
 91    # Function to summarize dataset
 92 ▾  summarize_dataset <- function(data) {
 93      # Basic dataset summary
 94      cat("Dataset Summary:\n")
 95      cat("Number of listings:", nrow(data), "\n")
 96      cat("Number of variables:", ncol(data), "\n")
 97      cat("Date range:", min(data$first_review, na.rm = TRUE), "to",
 98          max(data$last_review, na.rm = TRUE), "\n")
 99
100      # Missing values summary
101      missing_values <- data %>%
102        summarise(across(everything(), ~sum(is.na(.)))) %>%
103        pivot_longer(everything(), names_to = "variable", values_to = "missing_count") %>%
104        filter(missing_count > 0) %>%
105        arrange(desc(missing_count))
106
107 ▾    if(nrow(missing_values) > 0) {
108        cat("\nVariables with missing values:\n")
109        print(missing_values)
110 ▴    }
111 ▴  }
112
113    # Run initial data summary
114    summarize_dataset(airbnb)
115
```

In this code snippet, an exploratory data analysis of the Airbnb dataset is conducted.

13

Appendix 9.D: Price Analysis

```r
117 ▾ #----------------------------------------------------------------------
118 ▾ # -------------------- PRICE ANALYSIS --------------------------------
119 ▾ #----------------------------------------------------------------------
120   # Analyze how price varies by room type and property type
121   price_analysis <- airbnb %>%
122     # Group data for aggregation
123     group_by(room_type, property_type) %>%
124     # Calculate metrics for each group
125     summarise(
126       avg_price = mean(price, na.rm = TRUE),
127       median_price = median(price, na.rm = TRUE),
128       min_price = min(price, na.rm = TRUE),
129       max_price = max(price, na.rm = TRUE),
130       total_listings = n(),
131       avg_reviews = mean(number_of_reviews, na.rm = TRUE),
132       .groups = "drop"  # Drop grouping after summarization
133     ) %>%
134     # Filter out uncommon property types for clearer visualization
135     filter(total_listings >= 10)
136
137   # Visualization: Average Price by Room Type and Property Type
138 ▾ plot_price_by_room_property <- function(data) {
139     ggplot(data, aes(x = room_type, y = avg_price, fill = property_type)) +
140       geom_bar(stat = "identity", position = "dodge") +
141       labs(
142         title = "Average Price by Room Type and Property Type",
143         subtitle = "Only including property types with at least 10 listings",
144         x = "Room Type",
145         y = "Average Price ($)",
146         fill = "Property Type"
147       ) +
148       theme(
149         axis.text.x = element_text(angle = 45, hjust = 1),
150         legend.position = "bottom",
151         legend.box = "horizontal"
152       )
153 ▴ }
154
155   # Generate and display the price analysis plot
156   plot_price_by_room_property(price_analysis)
157
158
```

Now, the code snippet shows the price analysis and visualizes it with the ggplot function.

Appendix 9.E: Host Performance Analysis

```r
159  #------------------------------------------------------------------------
160  #--------------- HOST PERFORMANCE ANALYSIS -----------------------------
161  #------------------------------------------------------------------------
162  # Analyze how host characteristics affect listing performance
163  host_performance <- airbnb %>%
164    # Filter out entries with missing response rate
165    filter(!is.na(host$host_response_rate)) %>%
166    # Group by superhost status
167    group_by(host_is_superhost) %>%
168    # Calculate metrics for each group
169    summarise(
170      avg_response_rate = mean(host_response_rate, na.rm = TRUE),
171      avg_price = mean(price, na.rm = TRUE),
172      avg_reviews = mean(number_of_reviews, na.rm = TRUE),
173      avg_review_scores = mean(review_scores$review_scores_rating, na.rm = TRUE),
174      total_listings = n(),
175      .groups = "drop"
176    )
177
178  # Visualization: Superhost vs Non-Superhost Performance
179  plot_superhost_performance <- function(data) {
180    # Prepare data for visualization by pivoting to long format
181    plot_data <- data %>%
182      select(host_is_superhost, avg_price, avg_reviews, avg_review_scores) %>%
183      pivot_longer(
184        cols = c(avg_price, avg_reviews, avg_review_scores),
185        names_to = "metric",
186        values_to = "value"
187      ) %>%
188      mutate(
189        # Create readable labels for the metrics
190        metric = case_when(
191          metric == "avg_price" ~ "Average Price ($)",
192          metric == "avg_reviews" ~ "Average Number of Reviews",
193          metric == "avg_review_scores" ~ "Average Review Score",
194          TRUE ~ metric
195        )
196      )
197
198    # Create the plot
199    ggplot(plot_data, aes(x = host_is_superhost, y = value, fill = host_is_superhost)) +
200      geom_bar(stat = "identity") +
201      facet_wrap(~ metric, scales = "free_y") +
202      labs(
203        title = "Superhost vs Non-Superhost Performance",
204        subtitle = "Comparison across key performance metrics",
205        x = "Superhost Status",
206        y = NULL
207      ) +
208      theme(legend.position = "none")
209  }
210
211  # Generate and display the host performance plot
212  plot_superhost_performance(host_performance)
```

The code snippet compares the superhost with the non-superhost across various metrics like Price,

Reviews, and so on. Then, it visualizes it with the ggplot function.

15

Appendix 9.F: Text Mining on Descriptions

```
215  #------------------------------------------------------------------------
216  #---------------- TEXT MINING ON DESCRIPTIONS ---------------------------
217  #------------------------------------------------------------------------
218  # Process and analyze listing descriptions to identify key terms
219
220  # Function to process text and remove stop words
221  process_text <- function(data, text_column, id_column = "name") {
222    data %>%
223      select(!!sym(id_column), !!sym(text_column)) %>%
224      # Filter out missing text values
225      filter(!is.na(!!sym(text_column))) %>%
226      # Tokenize text into individual words
227      unnest_tokens(word, !!sym(text_column)) %>%
228      # Remove stop words that don't add meaning
229      anti_join(stop_words) %>%
230      # Count word frequency
231      count(word, sort = TRUE)
232  }
233
234  # Process descriptions for analysis
235  description_tokens <- process_text(airbnb, "description")
236
237  # Create word cloud of most common words in descriptions
238  create_wordcloud <- function(token_data, max_words = 200, title = "Word Cloud") {
239    # Set up plotting device
240    par(mar = c(0, 0, 2, 0))
241
242    # Generate word cloud
243    wordcloud(
244      words = token_data$word,
245      freq = token_data$n,
246      max.words = max_words,
247      colors = brewer.pal(8, "Dark2"),
248      random.order = FALSE,
249      rot.per = 0.35,
250      scale = c(3, 0.5)
251    )
252
253    # Add title
254    title(main = title)
255  }
256
257  # Generate and display the description word cloud
258  create_wordcloud(
259    description_tokens,
260    max_words = 200,
261    title = "Most Common Words in Airbnb Descriptions"
262  )
263
```

In this code snippet, descriptions of Airbnb properties are processed and then visualized through a word cloud.

Appendix 9.G: Sentiment Analysis

```r
265  #---------------- SENTIMENT ANALYSIS ON LISTING SUMMARIES ----------------
266  #------------------------------------------------------------------------
267  # Analyze the sentiment of listing summaries to identify emotional tone
268
269  # Function to perform sentiment analysis using multiple lexicons
270  analyze_sentiment <- function(token_data) {
271    # AFINN sentiment analysis (numerical scores)
272    afinn_sentiment <- token_data %>%
273      inner_join(get_sentiments("afinn")) %>%
274      group_by(name) %>%
275      summarise(sentiment = sum(value), .groups = "drop") %>%
276      mutate(method = "AFINN")
277
278    # Bing and NRC sentiment analysis (categorical)
279    bing_and_nrc_sentiment <- bind_rows(
280      # Bing lexicon (positive/negative)
281      token_data %>%
282        inner_join(get_sentiments("bing")) %>%
283        mutate(method = "Bing et al."),
284
285      # NRC lexicon (positive/negative only)
286      token_data %>%
287        inner_join(get_sentiments("nrc") %>%
288                   filter(sentiment %in% c("positive", "negative"))) %>%
289        mutate(method = "NRC")
290    ) %>%
291      # Count positive and negative words
292      count(method, sentiment) %>%
293      # Reshape data for visualization
294      pivot_wider(names_from = sentiment, values_from = n, values_fill = 0) %>%
295      # Calculate net sentiment
296      mutate(sentiment = positive - negative)
297
298    # Combine results from different methods
299    bind_rows(afinn_sentiment, bing_and_nrc_sentiment)
300  }
301
302  # Process and tokenize summaries
303  summary_tokens <- airbnb %>%
304    select(name, summary) %>%
305    filter(!is.na(summary)) %>%
306    unnest_tokens(word, summary) %>%
307    filter(!is.na(word)) %>%
308    anti_join(stop_words)
309
310  # Perform sentiment analysis
311  combined_sentiments <- analyze_sentiment(summary_tokens)
312
313  # Visualization: Sentiment Analysis Results
314  plot_sentiment_analysis <- function(sentiment_data) {
315    ggplot(sentiment_data, aes(method, sentiment, fill = method)) +
316      geom_col(show.legend = FALSE) +
317      facet_wrap(~method, ncol = 1, scales = "free_y") +
318      labs(
319
517:1   (Untitled)                                                                    R Script
```

A sentiment analysis is performed on the tokenized description column in the code snippet.

17

Appendix 9.H: Neighborhood Bigram Analysis

```
331 ▾ #--------------------- NEIGHBORHOOD BIGRAM ANALYSIS --------------------------
332 ▾ #------------------------------------------------------------------------------
333   # Analyze neighborhood descriptions using bigrams to identify location features
334
335   # Function to process text into n-grams
336 ▾ process_ngrams <- function(data, text_column, n = 2, id_column = "name", min_count = 1) {
337     # Create ngrams from text
338     ngrams <- data %>%
339       filter(!is.na(!!sym(text_column))) %>%
340       unnest_tokens(
341         ngram,
342         !!sym(text_column),
343         token = "ngrams",
344         n = n
345       )
346
347     # For bigrams, separate into component words and filter stop words
348 ▾   if (n == 2) {
349       ngrams <- ngrams %>%
350         separate(ngram, c("word1", "word2"), sep = " ") %>%
351         filter(!word1 %in% stop_words$word) %>%
352         filter(!word2 %in% stop_words$word) %>%
353         filter(!is.na(word1) & !is.na(word2)) %>%
354         count(word1, word2, sort = TRUE) %>%
355         filter(n >= min_count)
356 ▾   } else {
357       # For other n-grams, just count occurrences
358       ngrams <- ngrams %>%
359         count(ngram, sort = TRUE) %>%
360         filter(n >= min_count)
361 ▴   }
362
363     return(ngrams)
364 ▴ }
365
366   # Process neighborhood descriptions into bigrams
367   neighborhood_bigrams <- process_ngrams(
368     airbnb,
369     "neighborhood_overview",
370     n = 2,
371     min_count = 10
372   )
373
374   # Visualization: Top Bigrams Bar Chart
375 ▾ plot_top_bigrams <- function(bigram_data, top_n = 20) {
376     bigram_data %>%
377       head(top_n) %>%
378       unite(bigram, word1, word2, sep = " ") %>%
379       ggplot(aes(x = reorder(bigram, n), y = n)) +
380       geom_col(fill = "steelblue") +
381       coord_flip() +
382       labs(
383         title = paste("Top", top_n, "Bigrams in Neighborhood Descriptions"),
384         x = "Bigram",
```

This code snippet shows the bigram analysis and visualizes the top 20 bigrams with a Bar Chart.

```
392  # Visualization: Bigram Network Graph
393 ▾ plot_bigram_network <- function(bigram_data, min_count = 10) {
394    # Filter bigrams by minimum count
395    filtered_bigrams <- bigram_data %>%
396      filter(n >= min_count)
397
398    # Create network graph
399    bigram_graph <- filtered_bigrams %>%
400      graph_from_data_frame()
401
402    # Plot the graph
403    ggraph(bigram_graph, layout = "fr") +
404      geom_edge_link(aes(edge_alpha = n, edge_width = n), show.legend = FALSE) +
405      geom_node_point(color = "steelblue", size = 5) +
406      geom_node_text(aes(label = name), vjust = 1.5, hjust = 1) +
407      labs(
408        title = "Bigram Network of Neighborhood Descriptions",
409        subtitle = paste("Showing bigrams appearing at least", min_count, "times")
410      ) +
411      theme_void()
412 ▴ }
413
414  # Generate and display the bigram network
415  plot_bigram_network(neighborhood_bigrams, min_count = 15)
416
```

This follow-up code on Bigram Analysis shows the visualization of a bigram network with the help of the bigram data created with the help of the ggraph function.

19

# Appendix 9.I: TF-IDF Analysis

```
418 ▾ #--------------------------------------------------------------------------
419 ▾ #----------------------- TF-IDF ANALYSIS ------------------------------
420 ▾ #--------------------------------------------------------------------------
421
422   # Identify distinctive terms in neighborhood descriptions using TF-IDF
423
424   # Function to perform TF-IDF analysis on bigrams
425 ▾ analyze_tfidf <- function(data, text_column, grouping_var = "neighborhood_cleansed") {
426     # Process text into bigrams
427     bigrams <- data %>%
428       filter(!is.na(!!sym(text_column)) & !is.na(!!sym(grouping_var))) %>%
429       unnest_tokens(
430         bigram,
431         !!sym(text_column),
432         token = "ngrams",
433         n = 2
434       ) %>%
435       separate(bigram, c("word1", "word2"), sep = " ") %>%
436       filter(!word1 %in% stop_words$word) %>%
437       filter(!word2 %in% stop_words$word) %>%
438       filter(!is.na(word1) & !is.na(word2)) %>%
439       unite(bigram, word1, word2, sep = " ")
440
441     # Calculate TF-IDF
442     bigrams %>%
443       count(!!sym(grouping_var), bigram) %>%
444       bind_tf_idf(bigram, !!sym(grouping_var), n) %>%
445       arrange(desc(tf_idf))
446 ▴ }
447
448   # Analyze TF-IDF if neighborhood_cleansed is available
449 ▾ if ("neighborhood_cleansed" %in% names(airbnb)) {
450     neighborhood_tfidf <- analyze_tfidf(
451       airbnb,
452       "neighborhood_overview",
453       "neighborhood_cleansed"
454     )
455
456     # Display top TF-IDF terms by neighborhood
457     top_tfidf_terms <- neighborhood_tfidf %>%
458       group_by(neighborhood_cleansed) %>%
459       slice_max(order_by = tf_idf, n = 5) %>%
460       ungroup()
461
462     print(top_tfidf_terms)
463 ▴ }
464
465
```

This code snippet shows the UDF of the TF-IDF analysis to identify distinctive terms in neighborhood descriptions.

# Appendix 9.J: Correlation Analysis

```r
466 ▾ #-------------------------------------------------------------------------
467 ▾ #----------------------- CORRELATION ANALYSIS ----------------------------
468 ▾ #-------------------------------------------------------------------------
469   # Analyze correlations between numeric variables
470
471   # Function to calculate correlations between numeric variables
472 ▾ calculate_correlations <- function(data) {
473     # Select numeric columns
474     numeric_cols <- data %>%
475       select(where(is.numeric)) %>%
476       # Remove columns with too many NAs
477       select_if(function(x) mean(!is.na(x)) > 0.5)
478
479     # Calculate correlation matrix
480     cor_matrix <- cor(numeric_cols, use = "pairwise.complete.obs")
481
482     # Convert to long format for visualization
483     cor_data <- as.data.frame(cor_matrix) %>%
484       rownames_to_column("variable1") %>%
485       pivot_longer(-variable1, names_to = "variable2", values_to = "correlation")
486
487     return(cor_data)
488 ▴ }
489
490   # Calculate correlations
491   correlations <- calculate_correlations(airbnb)
492
493   # Visualization: Correlation Heatmap
494 ▾ plot_correlation_heatmap <- function(cor_data) {
495     ggplot(cor_data, aes(x = variable1, y = variable2, fill = correlation)) +
496       geom_tile() +
497       scale_fill_gradient2(
498         low = "blue",
499         mid = "white",
500         high = "red",
501         midpoint = 0,
502         limits = c(-1, 1)
503       ) +
504       theme(
505         axis.text.x = element_text(angle = 45, hjust = 1),
506         axis.title = element_blank()
507       ) +
508       labs(
509         title = "Correlation Heatmap of Numeric Variables",
510         fill = "Correlation"
511       )
512 ▴ }
513
514   # Generate and display correlation heatmap
515   plot_correlation_heatmap(correlations)
```

This code snippet shows the correlation analysis between all the numeric variables and visualizes the relationship between these variables with a heat map.