



HUMBOLDT-UNIVERSITY OF BERLIN  
SCHOOL OF BUSINESS AND ECONOMICS  
LADISLAUS VON BORTKIEWICZ CHAIR OF STATISTICS

STATISTICAL PROGRAMMING LANGUAGES  
SEMINAR PAPER

# Exploratory Data Analysis of Airbnb properties in Berlin

*Silvia Ventoruzzo*  
(592252)

submitted to

Alla PETUKHINA

March 15, 2019

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Data Preparation</b>	<b>2</b>
2.1	Berlin neighbourhoods and districts . . . . .	2
2.2	Berlin VBB Zones . . . . .	4
2.3	Airbnb listings' attributes . . . . .	7
<b>3</b>	<b>Exploratory Data Analysis</b>	<b>9</b>
3.1	Descriptive statistics . . . . .	10
3.2	Distribution plots . . . . .	12
<b>4</b>	<b>Price analysis</b>	<b>14</b>
4.1	Correlation with price . . . . .	14
4.2	Linear regression on price . . . . .	15
<b>5</b>	<b>Clustering</b>	<b>17</b>
5.1	Choosing the number of clusters . . . . .	17
5.2	k-means clustering . . . . .	20
<b>6</b>	<b>ShinyApp</b>	<b>21</b>
<b>7</b>	<b>Results and Conclusions</b>	<b>22</b>
<b>References</b>		<b>23</b>
<b>Appendix</b>		<b>25</b>
<b>A</b>	<b>Code</b>	<b>25</b>
A.1	Data and analysis . . . . .	25
A.2	Functions . . . . .	49
A.3	ShinyApp . . . . .	66

# 1 Introduction

Airbnb (airbnb.com) is now a famous website allowing private people and commercial entities to rent out part of their spaces. It all started in 2007 when two 27-year-olds decided to sublet their living room in their San Francisco apartment during a conference to help pay their rent (Salter, 2012). Now, they are a multibillion-dollar company where properties from all over the world are on rent (Bort, Julie (Business Insider), 2018).

Many studies have already been conducted on price determinants for the hotel industry, like the ones named in Wang and Nicolau (2017), and for Airbnb properties Wang and Nicolau (2017) itself, but none so far specific for the city of Berlin. We therefore attempted an exploratory analysis of Airbnb listings in Berlin and, in particular, of their price. For this purpose linear regression on price was also run and properties were clustered.

As in Wang and Nicolau (2017) summarized, in the case of hotels, hotel prices are negatively influenced by the hotel's location, since shorter distance from the city center, the main attractions and/or important transportation points leads to higher prices. On the other side positive influences on price are hotels' star rating and online customer rating, the services and amenities provided, and by the "presence of car parks and fitness centers". But price determinants could be different in case of different types of accommodations, such as the ones offered on Airbnb. That's why Wang and Nicolau (2017) also derived the drivers of price for Airbnb listings from 33 cities.

What is observed is that the dimension as well as the type of the accommodation positively influences the price. However, the location has almost no impact on the price, except for the case of being inside the circular with respect to the opposite. Clustering on the other hand did not deliver usable results, because of the absence of easily recognizable patterns. Further analysis of the price determinants and the clusters may shed more light into the topic.

The paper is divided into 7 sections. Section 2 will present the data used for this study and explained how it was prepared. Consecutively, 3 will run exploratory data analysis, both in form of tables and plots, on this data. Because of our interest in explaining property price, 4 will focus on this feature and show correlation and regression with respect to it. Furthermore, an attempt at clustering the Airbnb properties will be done in 5. After that it will be explained in 6 why a ShinyApp was developed to present this research. Finally, 7 we will present the results and draw some conclusions.

## 2 Data Preparation

For the analysis explained in this paper data was downloaded for a website independent from airbnb itself. Insideairbnb (Insideairbnb, 2019) scrapes airbnb to get its data, posts it online for the public to use on own analysis, while also providing some analysis of its own.

The data is divided according to cities and for each there is general information about the city's properties and their availability for the next year. For this analysis not all variables are being kept, the focus is indeed on the ones that might have the most affect on price. Moreover, feature engineering will also be performed to extract even more useful information. More details in subsection 2.1.

Since we are dealing with data with coordinates, spatial data is also needed to produce a map of the location. For this purpose further data has been downloaded from the Statistics Office of Berlin-Brandenburg (Amt für Statistik Berlin-Brandenburg, 2019) and Geofabrik (Geofabrik, 2019) and further processed, as explained in subsections 2.1 and 2.2.

To handle data will be make use of functions from the different packages in the `tidyverse`, which allows for easy manipulation of the data and flexible plotting, see Ross et al. (2017). For spatial data the package `sf` has been chosen, since it works well with the `tidyverse` packages and for all the other reasons listed in Pebesma (2018).

### 2.1 Berlin neighbourhoods and districts

Berlin consists of 96 neighbourhoods (Ortsteile), which are grouped into 12 districts (Bezirke). The data downloaded from Amt für Statistik Berlin-Brandenburg (2019) contain one polygon for each neighbourhood and additional information about them, out of which we only kept the district.

The polygons have been extracted from the relative shapefile with the function `st_read` from the `sf` package. For the neighbourhoods we simply rename the variables and keep the ones of interest.

**Listing 1: |berlin\_districts\_neighbourhoods.R|**

```
1
2 # Load shapefiles
3 berlin = sf:::st_read(file.path(getwd(), "Data",
4                         "Berlin-Ortsteile-polygon.shp", fsep="/"))
5
6 # Object with the neighbourhoods (and respective district)
7 berlin_neighbourhood_sf = berlin %>%
```

```

8     dplyr::rename(id      = Name,
9                      group = BEZNAME) %>%

```

However, one of the neighbourhoods, Buckow, is composed of two separate polygons, which we therefore need to unite according to the neighbourhoods' names. Thanks to the flexibility of the `sf` objects, this is done simply with a `summarize` from the package `dplyr`.

**Listing 2:** |berlin\_districts\_neighbourhoods.R|

```

11
12 # Buckow is composed of two separate parts, so we need to join them
13 berlin_neighbourhood_singlebuckow_sf = berlin_neighbourhood_sf %>%

```

For the districts we perform the same procedure as above, but this time we unite the polygons only by their district, which are represented here by the group variable.

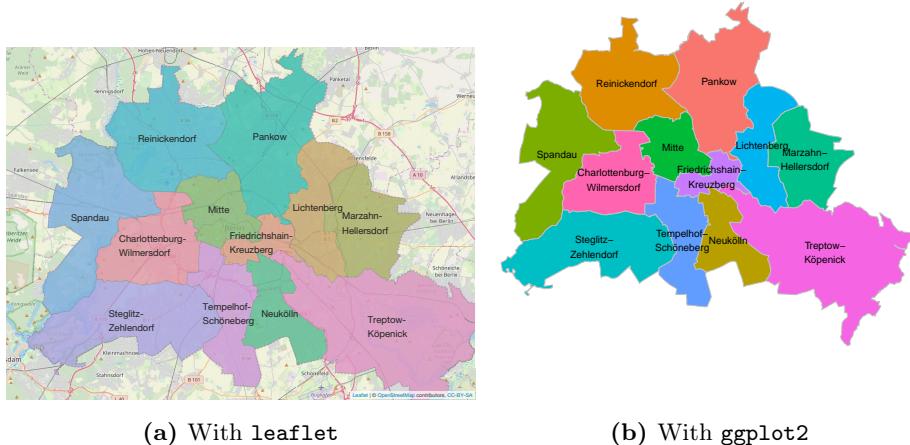
**Listing 3:** |berlin\_districts\_neighbourhoods.R|

```

19
20 # Object with the districts
21 berlin_district_sf = berlin_neighbourhood_sf %>%
22   dplyr::group_by(group) %>%

```

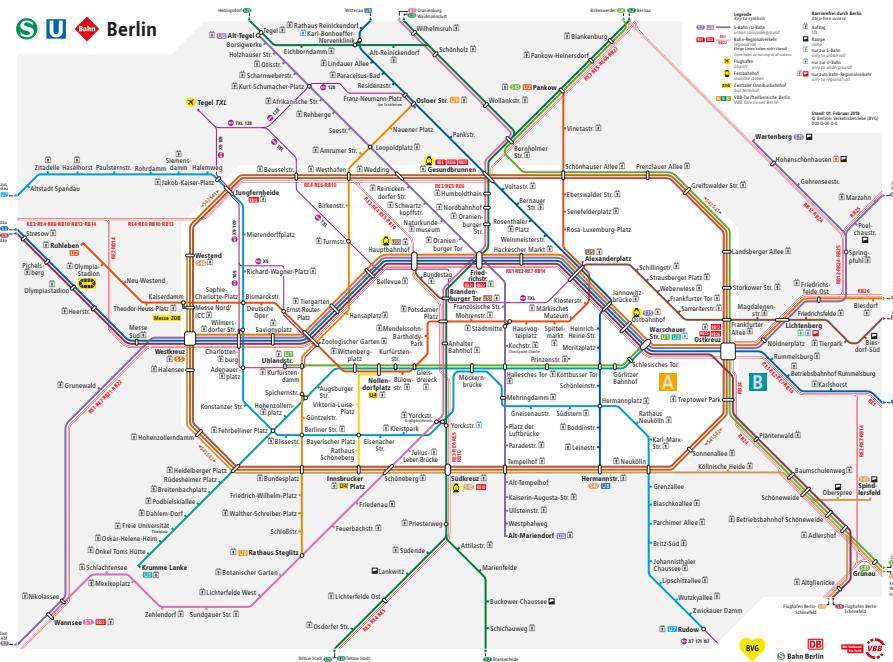
The produced polygons can be used to map Berlin using `ggplot2` or `leaflet`, as shown in figure 1. `leaflet` is more flexible and more suitable for plotting spatial data. Therefore, this package will be used for now on for mapping polygons and coordinate points.



**Figure 1:** Maps of the Berlin Districts  berlin\_districts\_neighbourhoods\_maps.R

## 2.2 Berlin VBB Zones

The VBB (Verkehrsverbund Berlin-Brandenburg) is "the public transport authority covering the federal states of Berlin and Brandenburg" (Verkehrsverbund Berlin-Brandenburg, 2019). The city of Berlin, in particular, is divided in two fare areas: A, covering the center of Berlin up to the circular line (Ringbahn), and B, from the Ringbahn to the border with Brandenburg. After that there is also the area C, which however will not be covered here since we only consider the city of Berlin.



**Figure 2:** Network Map of Berlin Areas A and B (Source: Verkehrsverbund Berlin-Brandenburg (2019))

We tried to replicate these areas by also making use of the spatial points of the Berlin stations.

First of all, we need to create the polygon for area A, which is done by joining the points in the circular line and transforming this into a polygon.

Unfortunately the shapefile with the stations did not contain information about the line to which these stations correspond. Therefore we needed to filter the stations manually with the ones belonging to the circular line. This vector also contains the information in which order the points will be connected. You can notice that the first and last station are the same since the polygon needs to close.

**Listing 4: |berlin\_vbb\_zones.R|**

```

1
2 # Create dataframe with names of stations on the Ringbahn (delimits Area A)
3 ringbahn_names_df = base::data.frame(
4   id = c("Südkreuz", "Schöneberg", "Innsbrucker Platz", "Bundesplatz",
5         "Heidelberger Platz", "Hohenzollerndamm", "Halensee", "Westkreuz",
6         "Messe Nord/ICC", "Westend", "Jungfernheide", "Beusselstraße",
7         "Westhafen", "Wedding", "Gesundbrunnen", "Schönhauser Allee",
8         "Prenzlauer Allee", "Greifswalder Straße", "Landsberger Allee",
9         "Storkower Straße", "Frankfurter Allee", "Ostkreuz",
10        "Treptower Park", "Sonnenallee", "Neukölln", "Hermannstraße",
11        "Tempelhof", "Südkreuz"),
12   stringsAsFactors = FALSE) %>%

```

Since some stations appear multiple times, we firstly filter railway stations, which include both subway and lightrail, and then we calculate the middle point for each station among the ones having the same name. We then join this with the dataframe containing the names of the stations in the circular line, thus filtering the stations to the ones we are interested in. After performing some preparation steps, the function `st_polygon` from the `sf` package was used to create a polygon out of a list of points.

**Listing 5: |berlin\_vbb\_zones.R|**

```

13
14 # Create sf object of Area A
15 berlin_vbb_A_sf = stations %>%
16   dplyr::filter(fclass %like% "railway") %>%
17   dplyr::rename(id = name) %>%
18   dplyr::mutate(id = gsub("Berlin ", "", id),
19                 id = gsub("Berlin-", "", id),
20                 id = gsub(" *\\(.*?\\) *", "", id),
21                 id = gsub("S ", "", id),
22                 id = gsub("U ", "", id)) %>%
23   points_midpoint() %>%
24   dplyr::right_join(ringbahn_names_df, by = "id") %>%
25   dplyr::arrange(order) %>%
26   dplyr::select(long, lat) %>%
27   base::as.matrix() %>%
28   base::list() %>%
29   sf::st_polygon() %>%

```

Second af all, in order to create the polygon for zone B, we need to produce the polygon for entire Berlin. This is done with a procedure already explained in subsection 2.1. In this case we don't perform any grouping, thus uniting all neighbourhoods.

**Listing 6: |berlin\_vbb\_zones.R|**

```
31
32 # Create sf object of entire Berlin
```

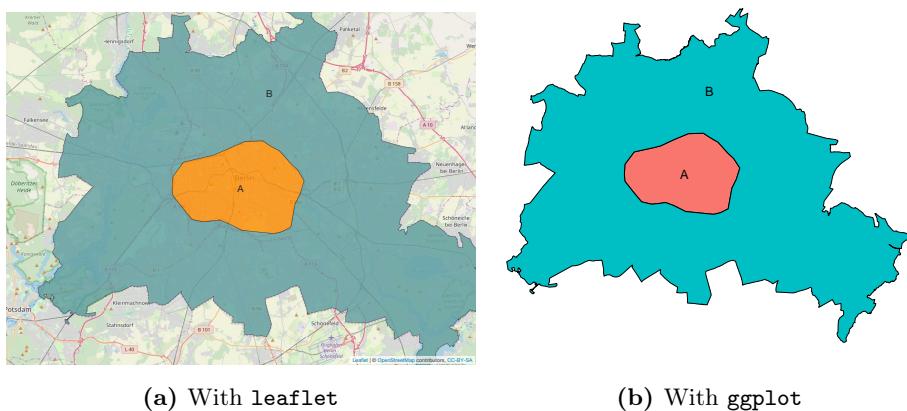
Finally, we bind the two objects by row and calculate their intersections thanks to the function `st_intersection` from the package `sf`. We then define the area names according to how many times the two previous polygons intersect:

- Area A: where polygons intersect ( $n.\text{overlaps} > 1$ )
- Area B: where polygons do not intersect ( $n.\text{overlaps} \leq 1$ )

**Listing 7: |berlin\_vbb\_zones.R|**

```
34
35 # Bind and intersect to create sf object with both VBB Zones (A and B)
36 berlin_vbb_AB_sf = berlin_vbb_A_sf %>%
37   base::rbind(berlin_sf) %>%
38   sf::st_intersection() %>%
39   dplyr::mutate(id = ifelse(n.overlaps > 1, "A", "B")) %>%
```

As in subsection 2.1 the `sf` object can be used to map the VBB zones using `leaflet` or `ggplot2`.



**Figure 3:** Maps of the Berlin VBB Zones  `berlin_vbb_zones_maps.R`

## 2.3 Airbnb listings' attributes

The first part of cleaning the airbnb data consists in joining the two datasets containing general information according to their common variables and correcting some string values. Secondly, we proceed in checking for missing values and deriving that information from other correlated variables. Thirdly, we move on to feature engineering.

We firstly derive the areas where the properties are located thanks to the spatial polygons created before and the function `point_in_polygons`. This function loops through the polygons in the `sf` object and check which points are contained in which polygons. In the end it writes the id associated to the polygon, in our case the area id, in the summary column, which can be named as preferred.

**Listing 8: |point\_in\_polygons.R|**

```
1 point_in_polygons <- function(points_df, polys_sf,
2                                 var_name, join_var = "id") {
3   # Create empty dataframe
4   is_in <- data.frame(matrix(ncol = nrow(polys_sf), nrow = nrow(points_df)))
5   is_in[,var_name] <- NA
6   # Extract coordinates of the polygons
7   coordinates <- as.data.frame(st_coordinates(polys_sf))
8   # Extract names of the polygons
9   name <- as.character(polys_sf$id)
10  # For all polygons check if the points are inside of them
11  for (k in 1:nrow(polys_sf)) {
12    is_in[,k] <- sp::point.in.polygon(point.x = points_df$long,
13                                         point.y = points_df$lat,
14                                         pol.x = coordinates$x,
15                                         [coordinates$L2 == k],
16                                         pol.y = coordinates$y,
17                                         [coordinates$L2 == k])
18    # Get the names of the polygons where the points are in one column
19    is_in[,var_name][is_in[,k] == 1] <- name[k]
20  }
21  # Keep only summary column and add points' names
22  is_in <- is_in %>%
23    dplyr::select(var_name) %>%
24    dplyr::mutate(id = points_df$id)
25  # Add the summary column to the points dataframe
26  points_df <- dplyr::full_join(points_df, is_in, by = join_var)
27  return(points_df)
```

Secondly, for railway stations and tourist attractions we calculate the amount inside a range and the distance to the nearest point using the function `distance_count`. In particular, the following parameters will be used:

- Railway stations: distance = 1000 (1 km)
- (Top 10) attractions: distance = 2000 (2 km)

This function firstly calculates the distance between all properties and all reference points using the Haversine Formula, which "gives minimum distance between any two points on spherical body by using latitude and longitude" (Ingole and Nichat, 2013) according to the following formula:

$$d = 2r \arcsin \left( \sqrt{\sin^2 \left( \frac{\phi_2 - \phi_1}{2} \right) + \cos(\phi_2) \cos(\phi_1) \sin^2 \left( \frac{\psi_2 - \psi_1}{2} \right)} \right) \quad (1)$$

where  $\phi$  corresponds to the latitude and  $\psi$  to the longitude of the two points.

Then it calculates how many reference points are within the set distance and how much is the distance to the nearest reference point. Finally this information is joined into the main dataframe.

**Listing 9: |distance\_count.R|**

```

1 distance_count = function(main, reference, var_name, distance) {
2   # Create variable names
3   var_name_count = paste(var_name, "count", sep = "_")
4   var_name_dist = paste(var_name, "dist", sep = "_")
5   # Calculate distance for each listing to each station
6   point_distance = geosphere::distm(x = main %>%
7     dplyr::select(long, lat),
8     y = reference %>%
9     dplyr::select(long, lat),
10    fun = distHaversine) %>%
11    as.data.frame() %>%
12    data.table::setnames(as.character(reference$id))
13   # Calculate how many "reference" are within "distance"
14   point_distance[,var_name_count] = rowSums(point_distance <= distance)
15   %>%
16   as.factor()

```

```

17 # Calculate the distance to the nearest "reference"
18 point_distance[,var_name_dist] = apply(point_distance[,-ncol(point_
19   distance)], 
20     MARGIN = 1,
21     FUN = min) %>%
22     round(0)
23 
24 # Insert this information into the main DF
25 main = point_distance %>%
26   dplyr::mutate(id = main$id) %>%
27   dplyr::select(id, var_name_count, var_name_dist) %>%
28   dplyr::right_join(main, by = "id")
29
30 return(main)
31 }
```

In the end we use the calendar dataframe to calculate availability of each property in the different seasons.

**Listing 10: |data\_preparation.R|**

```

1 month      = lubridate::month(date, label = TRUE),
2 day        = lubridate::day(date),
3 season     = ifelse((month == "Mar" & day >= 21) |
4   (month == "Apr") |
5   (month == "May") |
6   (month == "Jun" & day < 21), "Spring",
7   ifelse((month == "Jun" & day >= 21) |
8     (month == "Jul") |
9     (month == "Aug") |
10    (month == "Sep" & day < 21), "Summer",
11    ifelse((month == "Sep" & day >= 21) |
12      (month == "Oct") |
13      (month == "Nov") |
14      (month == "Dec" & day < 21), "Fall",
```

### 3 Exploratory Data Analysis

The exploratory data analysis will be divided in two parts: subsection 3.1 will show the descriptive statistics calculated either numeric or categorical variables; subsection 3.2 will then display the distribution plots. Correlation will also be calculated, but only with respect to price. This subject will be dealt with in subsection 4.1.

### 3.1 Descriptive statistics

Descriptive statistics make the interpretation of the data easier by giving grouping it and thus providing a shorter representation of it (Ibe, 2014).

As in Ibe (2014) explained there are three general types of descriptive statistics:

1. Measures of central tendency
2. Measures of spread
3. Graphical displays

This subsection will focus on the first two, while subsection 3.2 will deal with the third.

Most of these statistics are usually applied to continuous data, sometimes even to numerical discrete data, but not to categorical variables. Therefore the function to calculate descriptive statistics has been split in two.

The first part calculate both measures of central tendency and spread of all numerical variables thanks to the function `apply`. The function `apply(X, MARGIN, FUN, ...)` calculates the function in `FUN` for all rows (`MARGIN = 1`) or columns (`MARGIN = 2`) for the data in `X`. One can add additional arguments, like the quantile probabilities in our case.

**Listing 11: |descriptive\_statistics.R|**

```
1 descriptive_statistics = function(df) {  
2  
3   # Descriptive statistics for numeric variables  
4   if (unique(apply(df, 2, function(x) is.numeric(x)))) {  
5  
6     summary = data.frame(  
7       variable = names(df),  
8       min      = apply(df, 2, min),  
9       '1Q'    = apply(df, 2, quantile, probs = 0.25),  
10      median  = apply(df, 2, median),  
11      '3Q'    = apply(df, 2, quantile, probs = 0.75),  
12      max     = apply(df, 2, max),  
13      iqr     = apply(df, 2, IQR),  
14      mean    = apply(df, 2, mean),  
15      sd      = apply(df, 2, sd),  
16      check.names = FALSE) %>%  
17      dplyr::mutate_if(is.numeric, function(x) round(x, 4))
```

Part of the results can be seen in table 1.

variable	min	1Q	median	3Q	max	iqr	mean	sd
price	0.00	30.00	47.00	70.00	9000.00	40.00	67.69	210.53
review_scores_rating	0.00	84.00	95.00	100.00	100.00	16.00	77.26	37.09

**Table 1:** Sample of descriptive table for numeric variables [descriptive\\_statistics.R](#)

For categorical variables the above used statistics do not work, therefore frequencies and proportions of each factor were calculated. The mode is then simply the factor with the highest frequency.

**Listing 12:** [|descriptive\\_statistics.R|](#)

```

19 } else if (unique(apply(df, 2, function(x) is.character(x) | is.factor(x)
20 | is.logical(x)))) {
21
22 # Frequency
23 frequency_list = apply(df, 2, table)
24
25 frequency_df = data.frame(frequency = unlist(frequency_list)) %>%
26   tibble::rownames_to_column(var = "var_fact")
27
28 freq_var = colsplit(string = frequency_df$var_fact, pattern = "\\\\".,
29   names = c("variable", "factor"))
30
31 frequency = freq_var %>%
32   cbind(frequency_df) %>%
33   dplyr::select(-var_fact)
34
35 # Proportion
36 proportion_list = apply(df, 2, function(x) prop.table(table(x)))
37
38 proportion_df = data.frame(proportion = unlist(proportion_list)) %>%
39   tibble::rownames_to_column(var = "var_fact")
40
41 prop_var = colsplit(string = proportion_df$var_fact, pattern = "\\\\".,
42   names = c("variable", "factor"))
43
44 proportion = prop_var %>%
45   cbind(proportion_df) %>%
46   dplyr::select(-var_fact) %>%
47   dplyr::mutate(proportion = round(proportion, 4)*100,

```

```

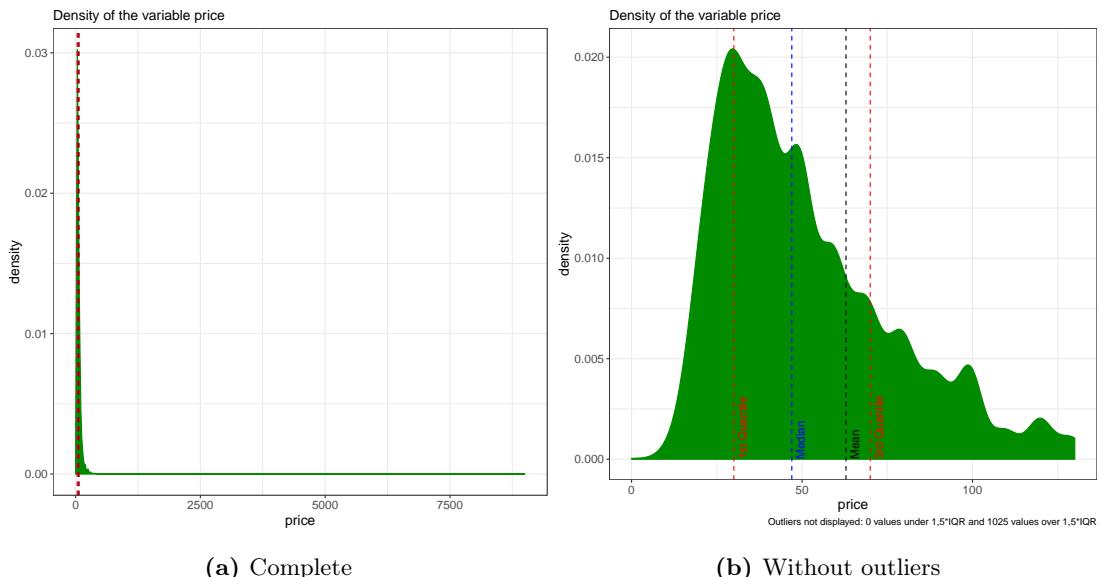
45     proportion = as.character(proportion) %>% paste("%")
46
47     summary = frequency %>%
48         dplyr::inner_join(proportion, by = c("variable", "factor")) %>%
49         dplyr::arrange(variable, desc(frequency))

```

### 3.2 Distribution plots

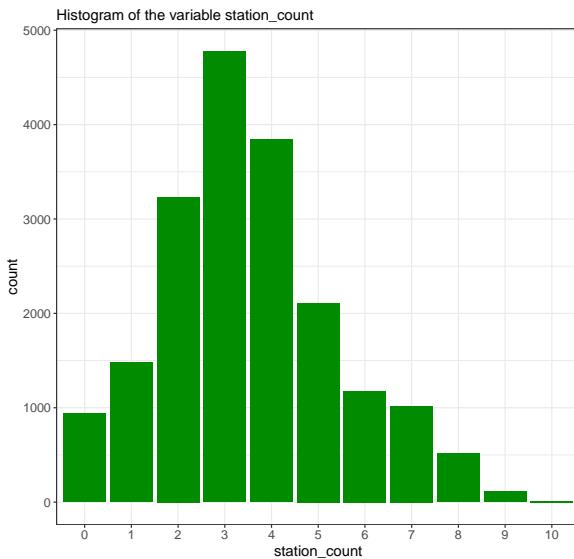
Also for the distribution plots we distinguish between numerical and categorical variables. In the former case a density plot has been chose, while in the latter a bar plot.

Unfortunately, many numeric variables present outliers, which, in the case of very skewed data, have been excluded for better visualization. An example can be seen in figure 4.



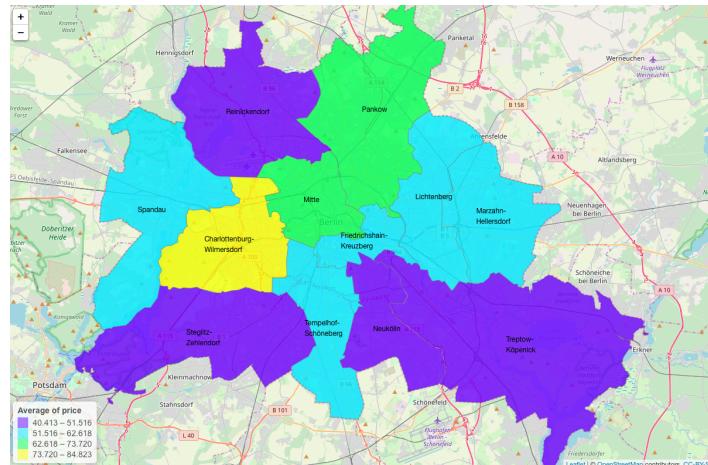
**Figure 4:** Distribution of the variable price distribution\_plot.R

For the categorical variables a bar plot is more appropriate, since variable can only assume a limited amount of values, and usually only a few, like in the case displayed in figure 5.



**Figure 5:** Distribution of the variable station\_count distribution\_plot.R

Having successfully created spatial polygons for Berlin, we can also map the distribution of the variables across the city. For example, in figure 6 one can see how the average of price in each neighbourhood is distributed across neighbourhoods.



**Figure 6:** Distribution of the average of price across Berlin's districts var\_avg\_map.R

## 4 Price analysis

One of the core factors in the choice of the property to book is its price, being one of main drivers of customers' behaviors (Liang et al., 2018). Therefore one may want to try and see what properties' attributes influence its value.

We start in subsection 4.1 by calculating the correlation of all the other variables with respect to price. Then we try in subsection 4.2 to run a linear regression on price to look what variables are statistically relevant and how much they affect the price.

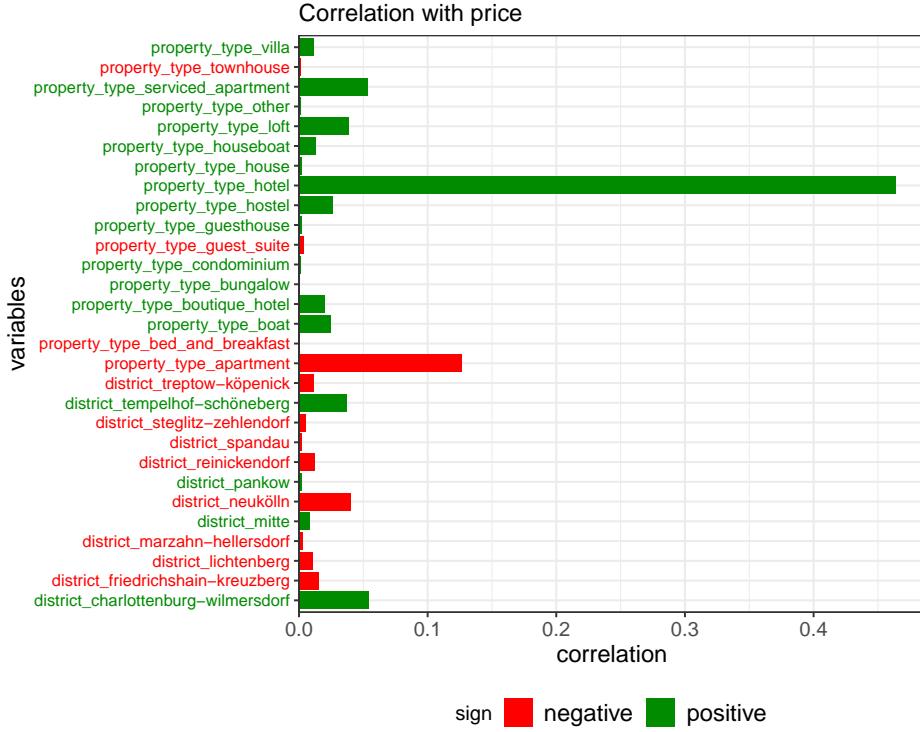
### 4.1 Correlation with price

Since we are only interested in the correlation with price, a classic correlation plot like the one produced by `corrplot` from the package `corrplot` may not be the most easily readable in this case, especially because of the large number of variables.

In fact, categorical variables first need to be transformed into many dummy variables in order to calculate the correlation. The function `from_row_to_col` will be used to achieve this result. It creates dummy variables for all factor levels by creating a column of 1s, spreading it across so many columns as factor levels and filling the empty rows in the columns with 0s.

**Listing 13: |from\_row\_to\_col.R|**

```
1 from_row_to_col = function(df, variable) {
2   df = df %>%
3     dplyr::mutate(yesno = 1) %>%
4     dplyr::distinct() %>%
5     tidyr::spread(variable, yesno, fill = 0)
6   return(df)
7 }
```



**Figure 7:** Sample of plot of correlation with price correlation\_plot.R

## 4.2 Linear regression on price

As pointed out in Wang and Nicolau (2017) linear regression is used to describe possible linear relationships between a dependent variable and one or multiple independent variables.

In this case we will look for the relationship that the properties' attributes have on their price. For this purpose we used again the data with the categorical variables transformed to dummies. The linear regression was then run using all other variables, except *id*, *long*, *lat*, *neighbourhood* and *listing\_url*, as regressors.

On table 2 you can see a sample of the results. A positive coefficient means that the variable has a positive impact on the dependent variable, the higher this regressor is, the higher *price* will be. On the other side, a negative coefficient suggests that a higher value of the independent variable will result in a lower value of *price*. Moreover, a variable is considered significant if the p-value is lower than the significance level, set at 5% in this case, since it rejects the null hypothesis for that regressor coefficient to be equal to zero Moyé (2006).

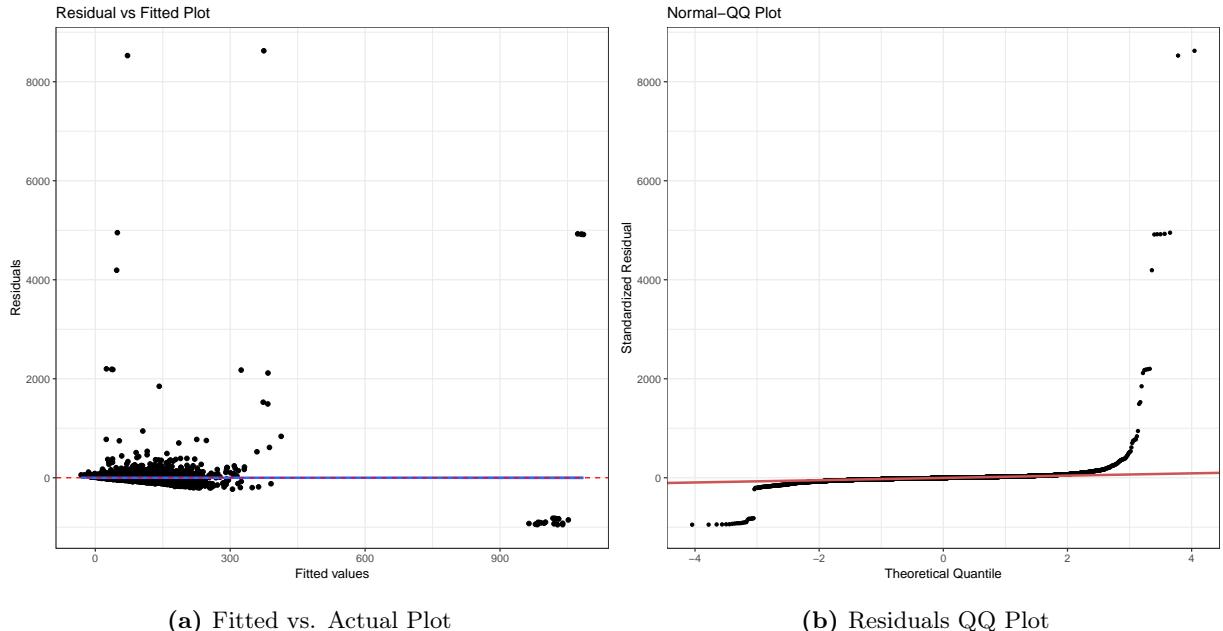
variable	coefficient	significant
(Intercept)	12.69	FALSE
host_listings_count	-0.73	TRUE
accommodates	16.31	TRUE

**Table 2:** Sample of linear regression results

Furthermore, to check the validity of the model one should also look at the distribution of the residuals, which should be normal. One can see from table 3 and picture 8 that this is here not the case.

min	1Q	median	3Q	max	iqr	mean	sd
-947.97	-18.72	-3.31	12.32	8625.31	31.04	0.00	138.02

**Table 3:** Descriptive statistics of regression residuals



**Figure 8:** Plotting the clusters

Finally, the  $R^2$  for this model, that is the proportion of the variance of *price* captured by the regressors, is only 0.1298, meaning that there is much room for improvement. One can test linear regression with other variables using the ShinyApp described in 6.

## 5 Clustering

A further step in data analysis is clustering, that is, trying to find groups in the data where they are similar to each other but different than data in other groups (Kaufman and Rousseeuw, 2009).

In this study we tried to see if Airbnb properties in Berlin can be split into clusters sharing common features. For this purpose we will use k-means clustering, one of the most famous and used partitioning methods.

The basic ideas of this method were developed at the beginning of the second half of the 20<sup>th</sup> century (Bock, 2007). This approach partitions the data to its closest centroid in one of the user-specified  $k$  number of clusters. The process works according to the following steps (Tan, 2018):

1.  $k$  points are randomly selected to be centroids.
2. The rest of the data points are assigned to the respectively closest centroid, i.e. the one which minimizes the sum of the squared error (SSE)

$$SSE = \sum_{i=1}^K \sum_{\mathbf{x} \in C_i} dist(\mathbf{c}_i, \mathbf{x})^2 \quad (2)$$

where  $dist(\mathbf{c}_i, \mathbf{x})$  is the Euclidean distance between the centroid of the  $i^{th}$  cluster and a point  $\mathbf{x}$ .

3. New centroids are selected, being the mean of all points in the respective cluster.
4. Repeat points 2. and 3. until the centroids stop changing or the maximum number of iterations has been reached.

Subsection 5.1 will deal with the process of choosing the appropriate number of clusters to be used in the clustering algorithm, while subsection 5.2 will handle the clustering itself.

### 5.1 Choosing the number of clusters

According to procedure explained above, the first thing to do is to choose the number of clusters  $k$ . As illustrated in Kodinariya and Makwana (2013), there are multiple approaches to selecting the appropriate value of  $k$ . In this study the elbow method was used, since it can be calculated using the k-means algorithm and it does not require too much computational power.

This method is a visual rule of thumb implemented in the function `number_of_clusters` and is Madhulatha (2012). It firstly requires the user to cluster the data multiple times with increasing values of  $k$ , starting at  $k = 2$ .

**Listing 14:** |number\_of\_clusters.R|

```

1 number_of_clusters = function(scaled_df, min = 2, max,
2                               iter_max = 10, plot_breaks) {
3   # Set seed for reproducibility
4   set.seed(900114)
5   # Values to test
6   k_values = min:max
7   # Empty dataframe for the total variance explained
8   tve = data.frame(clusters = k_values,
9                     tve = rep(NA, length(k_values)))
10  # Empty list for the kmeans objects
11  clk = list()
12  # Loop through the possible values of k
13  for (k in k_values) {
14    # Calculate k-means for each k
15    clk[[k-1]] = kmeans(scaled_df, centers = k, iter.max = iter_max)

```

Secondly one calculates for each run the percentage of the total variance explained.

**Listing 15:** |number\_of\_clusters.R|

```

17  # Save the number of clusters
18  names(clk)[k-1] = paste(k, "clusters", sep = " ")
19  # Calculate percentage of total variance explained
20  tve$tve[k-1] = 1 - clk[[k-1]]$tot.withinss / clk[[k-1]]$totss
21  # Print process
22  print(paste("k-means with", k, "clusters done", sep = " "))
23 }

```

One then plots the total variance explained against the number of clusters and chooses the value of  $k$  where adding another cluster does not add sufficient information (Madhulatha, 2012).

**Listing 16:** |number\_of\_clusters.R|

```

25
26  # Plot tve against k values
27  plot = ggplot(data = tve, aes(x = clusters, y = tve)) +
28    geom_line(color = "grey") +

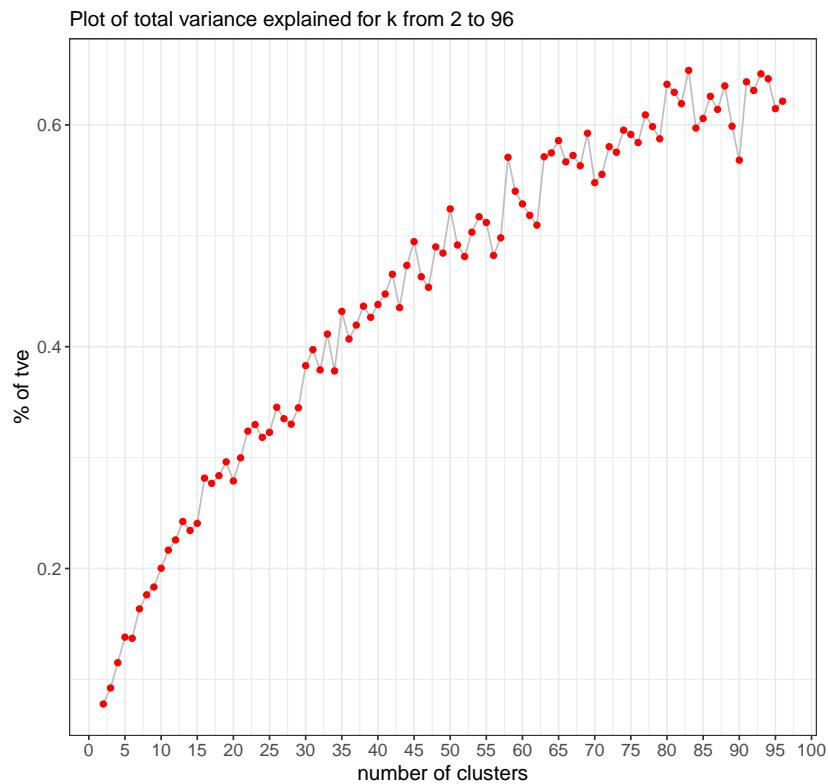
```

```

29     geom_point(color = "red") +
30
31     scale_x_continuous(breaks = plot_breaks) +
32
33     labs(x      = "number of clusters",
34          y      = "% of tve",
35          title = paste("Plot of total variance explained for k from", min,
36                         "to", max, sep = " ")) +
37
38     theme_bw() +
39
40     theme(axis.text.x = element_text(size = rel(1.2)),
41           axis.text.y = element_text(size = rel(1.2)),
42           axis.title.x = element_text(size = rel(1.2)),
43           axis.title.y = element_text(size = rel(1.2)))
44
45
46   return(plot)
47 }

```

The resulting plot can be seen in 9. In this case there is no general best, since the value of *tve* is not continually increasing.

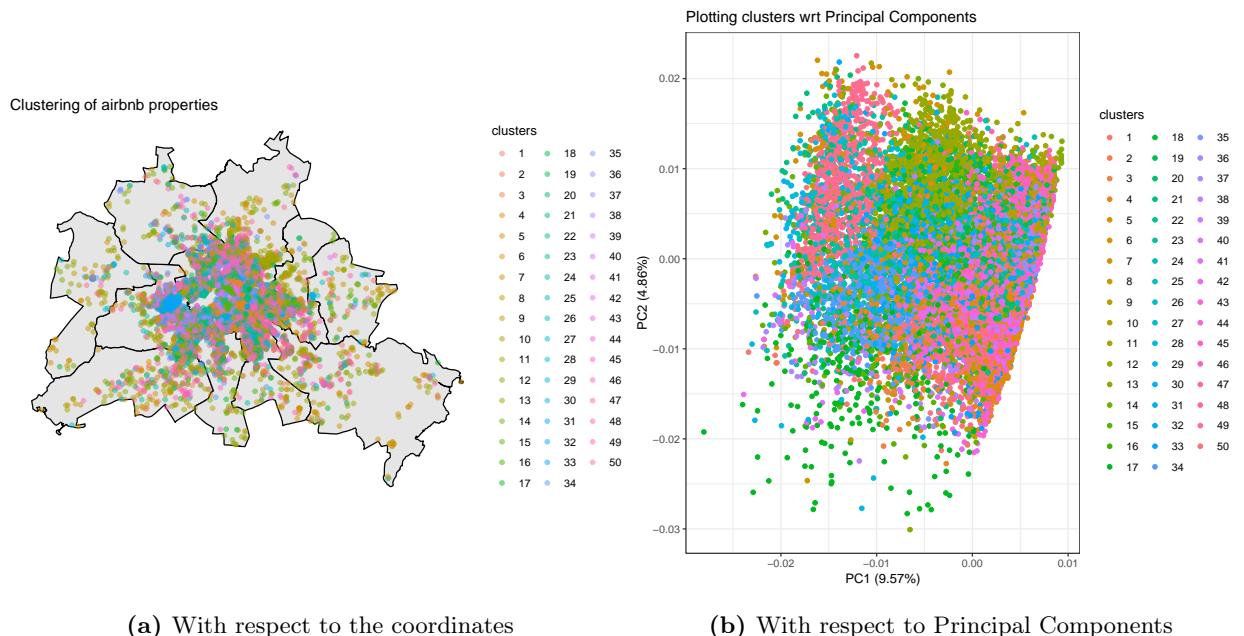


**Figure 9:** Plot of the total variance explained against different values of  $k$  number\_of\_clusters.R

## 5.2 k-means clustering

After selecting the number of clusters to use, one can apply it inside the function `kmeans` and look at the cluster assignments of the different points.

Having multiple variables, plotting the clustered points with respect to all variables is not feasible, one need therefore to find alternatives. Since we have coordinate points, one possibility is to plot the points with respect to their longitude and latitude, like in figure 10a. Another one can be to plot them with respect to the first two principal components, as in figure 10b.



**Figure 10:** Plotting the clusters

Because of the large amount of variables, points and clusters, it is not easy to visually display the groups. This task can be left for further studies.

## 6 ShinyApp

`shiny` is an R package that allows the user to write interactive web applications. These are especially helpful when delivering information to people with no coding experience in a very user-friendly way.

As described in the official site of `shiny`, in its most basic version an App is contained in a single script called `app.R`. As the Apps become more and more complicated one can write the code in two scripts, `ui.R` and `server.R`, or even further split these into thematic scripts.

The App structure is divided into two main elements:

- `ui`: The user interface object determines the appearance of the App itself. Here the possible inputs and the outputs will be defined.
- `server`: The server function uses the input values chosen in the App to produce the outputs indicated in the user interface.

Finally the App will be called using the function `shinyApp(ui, server)`.

Because of its interactiveness and flexibility a ShinyApp was built for this project in order to enable the final user a chance to a first hand analysis of the data. The code for the App can be found in the Appendix at subsection A.3. It is divided into four tabs: the first one is just introductory, presenting the project and the App itself; the second display Berlin maps and descriptive statistics according to section 3; the third relates to price analysis explained in section 4; the fourth and final one enables the user to perform clustering, like illustrated in section 5. The ShinyApp is reachable at this link: <https://silviav.shinyapps.io/airbnb-berlin/>.

## 7 Results and Conclusions

- Give a short summary of what has been done and what has been found.

Airbnb is a relevant subject in this day and age where sharing economy has affected most business fields. Therefore there have been many studies about Airbnb, in relation to hotels (Wang and Nicolau, 2017), and to the rental market (Coles et al., 2017).

In this study we focused on the properties based in Berlin, the capital of Germany. We looked at how the properties are distributed across the cities, noting that most of them are in the central part of town, the one inside the circular line. We further explored the distribution of different attributes, identifying differences across the multiple the districts and neighbourhoods.

As to be expected, most of the properties have at least one station within 1km, being fast to the center an important factor when choosing accomodation.

- Expose results concisely.
- Draw conclusions about the problem studied. What are the implications of your findings?
- Point out some limitations of study (assist reader in judging validity of findings).
- Suggest issues for future research.

## References

- AMT FÜR STATISTIK BERLIN-BRANDENBURG (2019): “Statistik Berlin-Brandenburg BerlinOpenData Geometrien,” <https://www.statistik-berlin-brandenburg.de/produkte/opendata/geometrienOD.asp?Kat=6301>, online; last accessed on 12. November 2018.
- BOCK, H.-H. (2007): “Clustering methods: a history of k-means algorithms,” in *Selected contributions in data analysis and classification*, Springer, 161–172.
- BORT, JULIE (BUSINESS INSIDER) (2018): “Insideairbnb Get the Data,” <https://www.businessinsider.de/airbnb-profit-revenue-2018-2?r=US&IR=T>, online; last accessed on 22. February 2019.
- COLES, P. A., M. EGESDAL, I. G. ELLEN, X. LI, AND A. SUNDARARAJAN (2017): “Airbnb usage across new York City neighborhoods: Geographic patterns and regulatory implications,” *Forthcoming, Cambridge Handbook on the Law of the Sharing Economy*.
- GEOFABRIK (2019): “Geofabrik Downloads: Europe, Germany, Berlin,” <http://download.geofabrik.de/europe/germany/berlin.html>, online; last accessed on 20. November 2018.
- IBE, O. (2014): *Fundamentals of applied probability and random processes*. 2nd ed, Academic Press, chap. 8.
- INGOLE, P. AND M. M. K. NICHAT (2013): “Landmark based shortest path detection by using Dijkstra Algorithm and Haversine Formula,” *International Journal of Engineering Research and Applications*, 3, 162–165.
- INSIDEAIRBNB (2019): “Insideairbnb Get the Data,” <http://insideairbnb.com/get-the-data.html>, online; last accessed on 04. March 2019.
- KAUFMAN, L. AND P. J. ROUSSEEUW (2009): *Finding groups in data: an introduction to cluster analysis*, vol. 344, John Wiley & Sons.
- KODINARIYA, T. M. AND P. R. MAKWANA (2013): “Review on determining number of Cluster in K-Means Clustering,” *International Journal*, 1, 90–95.
- LIANG, L. J., H. C. CHOI, AND M. JOPPE (2018): “Understanding repurchase intention of Airbnb consumers: perceived authenticity, electronic word-of-mouth, and price sensitivity,” *Journal of Travel & Tourism Marketing*, 35, 73–89.

MADHULATHA, T. S. (2012): “An overview on clustering methods,” *arXiv preprint arXiv:1205.1117*.

MOYÉ, L. A. (2006): *Statistical reasoning in medicine: the intuitive P-value primer*, Springer Science & Business Media.

PEBESMA, E. (2018): “Simple Features for R: Standardized Support for Spatial Vector Data,” *The R Journal*, 10, 439–446.

ROSS, Z., H. WICKHAM, AND D. ROBINSON (2017): “Declutter your R workflow with tidy tools,” Tech. rep., PeerJ Preprints.

SALTER, J. (2012): “Airbnb: The story behind the /1.3bnroom – lettingwebsite,” .

TAN, P.-N. (2018): *Introduction to data mining*, Pearson Education India, chap. 8.

VERKEHRSVERBUND BERLIN-BRANDENBURG (2019): “The company VBB,” <https://www.vbb.de/en/about-us/the-company-vbb>, online; last accessed on 14. March 2018.

WANG, D. AND J. L. NICOLAU (2017): “Price determinants of sharing economy based accommodation rental: A study of listings from 33 cities on Airbnb. com,” *International Journal of Hospitality Management*, 62, 120–131.

# A Code

## A.1 Data and analysis

Listing 17: |berlin\_districts\_neighbourhoods.R|

```
1 # Load packages
2 needed_packages = c("dplyr", "data.table", "sf", "tibble")
3 for (package in needed_packages) {
4   if (!require(package, character.only=TRUE))
5     {install.packages(package, character.only=TRUE)}
6   library(package, character.only=TRUE)
7 }
8 rm("needed_packages", "package")
9
10 # Set working directory to the one where the file is located
11
12 # This works when run directly
13 # setwd(dirname(rstudioapi::getActiveDocumentContext()$path))
14
15 # This works when sourced
16 setwd(dirname(sys.frame(1)$ofile))
17
18
19 # Load shapefiles
20 berlin = sf::st_read(file.path(getwd(), "Data",
21                       "Berlin-Ortsteile-polygon.shp", fsep="/"))
22
23 # Object with the neighbourhoods (and respective district)
24 berlin_neighbourhood_sf = berlin %>%
25   dplyr::rename(id      = Name,
26                 group = BEZNAME) %>%
27   dplyr::select(id, group, geometry) %>%
28   dplyr::arrange(group)
29
30 # Buckow is composed of two separate parts, so we need to join them
31 berlin_neighbourhood_singlebuckow_sf = berlin_neighbourhood_sf %>%
32   dplyr::group_by(id, group) %>%
33   dplyr::summarize(do_union = TRUE)
34
35 # Object with the districts
36 berlin_district_sf = berlin_neighbourhood_sf %>%
```

```

37   dplyr::group_by(group) %>%
38   dplyr::summarize(do_union = TRUE) %>%
39   dplyr::mutate(id = group)
40
41 # Create dataframes with the names for plotting
42 # Neighbourhoods
43 berlin_neighbourhoods_names = berlin_neighbourhood_sf %>%
44   sf::st_centroid() %>%
45   sf::st_coordinates() %>%
46   base::as.data.frame() %>%
47   dplyr::rename(long = X,
48                 lat = Y) %>%
49   dplyr::mutate(id = berlin_neighbourhood_sf$id,
50                 group = berlin_neighbourhood_sf$group,
51                 name = gsub("-", "-<br>", berlin_neighbourhood_sf$id))
52
53 # Districts
54 berlin_districts_names = berlin_district_sf %>%
55   sf::st_centroid() %>%
56   sf::st_coordinates() %>%
57   base::as.data.frame() %>%
58   dplyr::rename(long = X,
59                 lat = Y) %>%
60   dplyr::mutate(id = berlin_district_sf$id,
61                 group = berlin_district_sf$group,
62                 name = gsub("-", "-<br>", berlin_district_sf$id))
63
64 # Remove not needed data
65 rm("berlin")

```

Listing 18: |berlin\_districts\_neighbourhoods\_maps.R|

```

1 # Load packages
2 needed_packages = c("leaflet", "ggplot2", "htmltools", "mapview")
3 for (package in needed_packages) {
4   if (!require(package, character.only=TRUE))
5     {install.packages(package, character.only=TRUE)}
6   library(package, character.only=TRUE)
7 }
8 rm("needed_packages")
9
10 # Set working directory to the one where the file is located

```

```

11
12 # This works when run directly
13 setwd(dirname(rstudioapi::getActiveDocumentContext()$path))
14
15 # This works when sourced
16 # setwd(dirname(sys.frame(1)$ofile))
17
18 # Load file where the sf object is created
19 source("berlin_districts_neighbourhoods.R", chdir = TRUE)
20
21 # Colors for polygons
22 district_colors = colorspace::rainbow_hcl(n = 12)
23
24 # Leaflet map of the districts
25 berlin_district_leaflet = leaflet() %>%
26   addTiles() %>%
27   addPolygons(data = berlin_district_sf,
28     weight = 1, smoothFactor = 1, fillOpacity = 0.8,
29     fillColor = district_colors, color = "grey") %>%
30   addLabelOnlyMarkers(data = berlin_districts_names,
31     lng = ~long, lat = ~lat,
32     label = ~lapply(name, htmltools::HTML),
33     labelOptions = labelOptions(noHide = TRUE,
34       direction = 'center',
35       textOnly = TRUE,
36       textSize = "20px")) %>%
37   setView(lng = berlin_neighbourhoods_names$long
38         [berlin_neighbourhoods_names$id == "Mitte"],
39         lat = berlin_neighbourhoods_names$lat
40         [berlin_neighbourhoods_names$id == "Mitte"],
41         zoom = 11)
42 berlin_district_leaflet
43
44 # mapview::mapshot(berlin_district_leaflet,
45 #   file = "berlin_district_leaflet.pdf",
46 #   remove_controls = c("zoomControl", "layersControl",
47 #     "homeButton", "scaleBar"))
48
49
50 # ggplot map
51 ggplot() +

```

```

52   geom_sf(data = berlin_district_sf,
53           show.legend = FALSE, color = "grey") +
54   coord_sf(datum = NA) +
55   aes(fill = district_colors) +
56   geom_text(data = berlin_districts_names,
57             aes(x = long, y = lat,
58                  label = gsub("<br>", "\n", name),
59                  hjust = "center")) +
60   theme_classic() +
61   theme(plot.title = element_text(hjust = 0.5),
62         axis.title.x = element_blank(),
63         axis.title.y = element_blank())
64
65 # dev.copy2pdf(file = "berlin_district_ggplot.pdf")
66 # dev.copy2pdf(file = "../SeminarPaper/berlin_district_ggplot.pdf")

```

Listing 19: |berlin\_vbb\_zones.R|

```

1 # Load packages
2 needed_packages = c("dplyr", "data.table", "sf", "tibble")
3 for (package in needed_packages) {
4   if (!require(package, character.only=TRUE))
5     {install.packages(package, character.only=TRUE)}
6   library(package, character.only=TRUE)
7 }
8 rm("needed_packages", "package")
9
10 # Set working directory to the one where the file is located
11
12 # This works when run directly
13 # setwd(dirname(rstudioapi::getActiveDocumentContext()$path))
14
15 # This works when sourced
16 setwd(dirname(sys.frame(1)$ofile))
17
18
19 # Load helpers
20 source(file.path(getwd(), "Helpers", "points_midpoint.R", fsep="/"))
21
22 # Load shapefiles
23 berlin = sf::st_read(file.path(getwd(), "Data",
24                       "Berlin-Ortsteile-polygon.shp", fsep="/"))

```

```

25 stations = sf::st_read(file.path(getwd(), "Data",
26                         "gis_osm_transport_free_1.shp", fsep="/"))
27
28 # Create dataframe with names of stations on the Ringbahn (delimits Area A)
29 ringbahn_names_df = base::data.frame(
30   id = c("Südkreuz", "Schöneberg", "Innsbrucker Platz", "Bundesplatz",
31         "Heidelberger Platz", "Hohenzollerndamm", "Halensee", "Westkreuz",
32         "Messe Nord/ICC", "Westend", "Jungfernheide", "Beusselstraße",
33         "Westhafen", "Wedding", "Gesundbrunnen", "Schönhauser Allee",
34         "Prenzlauer Allee", "Greifswalder Straße", "Landsberger Allee",
35         "Storkower Straße", "Frankfurter Allee", "Ostkreuz",
36         "Treptower Park", "Sonnenallee", "Neukölln", "Hermannstraße",
37         "Tempelhof", "Südkreuz"),
38   stringsAsFactors = FALSE) %>%
39   tibble::rownames_to_column(var = "order") %>%
40   dplyr::mutate(order = as.numeric(order))
41
42 # Create sf object of Area A
43 berlin_vbb_A_sf = stations %>%
44   dplyr::filter(fclass %like% "railway") %>%
45   dplyr::rename(id = name) %>%
46   dplyr::mutate(id = gsub("Berlin ", "", id),
47                 id = gsub("Berlin-", "", id),
48                 id = gsub(" *\\(.*)?\\) *", "", id),
49                 id = gsub("S ", "", id),
50                 id = gsub("U ", "", id)) %>%
51   points_midpoint() %>%
52   dplyr::right_join(ringbahn_names_df, by = "id") %>%
53   dplyr::arrange(order) %>%
54   dplyr::select(long, lat) %>%
55   base::as.matrix() %>%
56   base::list() %>%
57   sf::st_polygon() %>%
58   sf::st_sfc() %>%
59   sf::st_sf(crs = sf::st_crs(berlin))
60
61 # Create sf object of entire Berlin
62 berlin_sf = berlin %>%
63   dplyr::summarize(do_union = TRUE)
64
65 # Bind and intersect to create sf object with both VBB Zones (A and B)

```

```

66 berlin_vbb_AB_sf = berlin_vbb_A_sf %>%
67   base::rbind(berlin_sf) %>%
68   sf::st_intersection() %>%
69   dplyr::mutate(id = ifelse(n.overlaps > 1, "A", "B")) %>%
70   dplyr::select(-n.overlaps, -origins) %>%
71   dplyr::arrange(desc(id))
72
73 # Create dataframes with coordinates where to show zones' names on map
74 berlin_vbb_A_names = berlin_vbb_A_sf %>%
75   sf::st_centroid() %>%
76   sf::st_coordinates() %>%
77   base::as.data.frame() %>%
78   dplyr::rename(long = X,
79                 lat = Y) %>%
80   dplyr::mutate(id = "A")
81 berlin_vbb_B_names = berlin_sf %>%
82   sf::st_bbox() %>%
83   base::as.matrix() %>%
84   base::t() %>%
85   as.data.frame() %>%
86   dplyr::transmute(long = (xmax + xmin)/2,
87                     lat = (3*ymax + ymin)/4) %>%
88   dplyr::mutate(id = "B")
89 berlin_vbb_AB_names = berlin_vbb_A_names %>%
90   rbind(berlin_vbb_B_names) %>%
91   dplyr::mutate(name = gsub("-", "-<br>", id))
92
93 # Remove not needed data
94 rm("berlin", "berlin_sf", "ringbahn_names_df", "stations",
95    "berlin_vbb_A_sf", "berlin_vbb_A_names", "berlin_vbb_B_names")
96 rm(list=lsf.str()) # All functions

```

Listing 20: |berlin\_vbb\_zones\_maps.R|

```

1 # Load packages
2 needed_packages = c("leaflet", "ggplot2", "htmltools", "mapview")
3 for (package in needed_packages) {
4   if (!require(package, character.only=TRUE))
5     {install.packages(package, character.only=TRUE)}
6   library(package, character.only=TRUE)
7 }
8 rm("needed_packages")

```

```

9
10 # Set working directory to the one where the file is located
11
12 # This works when run directly
13 setwd(dirname(rstudioapi::getActiveDocumentContext()$path))
14
15 # This works when sourced
16 # setwd(dirname(sys.frame(1)$ofile))
17
18
19 # Load file where the sf object is created
20 source("berlin_vbb_zones.R", chdir = TRUE)
21
22 # Leaflet map
23 berlin_vbb_zones_leaflet = leaflet() %>%
24   addTiles() %>%
25   addPolygons(data = berlin_vbb_AB_sf,
26               weight = 1, smoothFactor = 1, fillOpacity = 0.8,
27               fillColor = c("cadetblue", "darkorange"),
28               color = "black") %>%
29   addLabelOnlyMarkers(data = berlin_vbb_AB_names,
30                       lng = ~long, lat = ~lat,
31                       label = ~lapply(name, htmltools::HTML),
32                       labelOptions = labelOptions(noHide = TRUE,
33                                         direction = 'center',
34                                         textOnly = TRUE,
35                                         textSize = "20px")) %>%
36   setView(lng = berlin_vbb_AB_names$long
37           [berlin_vbb_AB_names$id == "A"],
38           lat = berlin_vbb_AB_names$lat
39           [berlin_vbb_AB_names$id == "A"],
40           zoom = 11)
41 berlin_vbb_zones_leaflet
42
43 # mapview::mapshot(berlin_vbb_zones_leaflet,
44 #                     file = "berlin_vbb_zones_leaflet.pdf",
45 #                     remove_controls = c("zoomControl", "layersControl",
46 #                                         "homeButton", "scaleBar"))
47 # mapview::mapshot(berlin_vbb_zones_leaflet,
48 #                     file = "../SeminarPaper/berlin_vbb_zones_leaflet.pdf",
49 #                     remove_controls = c("zoomControl", "layersControl",

```

```

50 #                                         "homeButton", "scaleBar"))
51
52
53 # ggplot map
54 ggplot() +
55   geom_sf(data = berlin_vbb_AB_sf, show.legend = FALSE, color = "black") +
56   coord_sf(datum = NA) +
57   aes(fill = c("darkorange", "cadetblue")) +
58   geom_text(data = berlin_vbb_AB_names,
59             aes(x = long, y = lat, label = id, hjust = "center"),
60             size = 5) +
61   theme_classic() +
62   theme(plot.title = element_text(hjust = 0.5),
63         axis.title.x = element_blank(),
64         axis.title.y = element_blank())
65
66 # dev.copy2pdf(file = "berlin_vbb_zones_ggplot.pdf")
67 # dev.copy2pdf(file = "../SeminarPaper/berlin_vbb_zones_ggplot.pdf")

```

**Listing 21:** |data\_preparation.R|

```

1 # Install and load needed packages
2 needed_packages = c("tidyverse",
3                      "geosphere",
4                      "readr",
5                      "rstudioapi",
6                      "Jmisc",
7                      "sp",
8                      "sf")
9 for (package in needed_packages) {
10   if (!require(package, character.only=TRUE)) {
11     install.packages(package, character.only=TRUE)}
12   library(package, character.only=TRUE)
13 }
14 rm("needed_packages", "package")
15
16 # Set working directory to the one where the file is located
17
18 # This works when run directly
19 # setwd(dirname(rstudioapi::getActiveDocumentContext()$path))
20
21 # This works when sourced

```

```

22 setwd(dirname(sys.frame(1)$ofile))
23
24 # Load helper functions and other scripts
25 source(file.path(getwd(), "Berlin_VBB_Zones",
26                   "berlin_vbb_zones.R", fsep="/"), chdir = TRUE)
27 source(file.path(getwd(), "Berlin_Districts_Neighbourhoods",
28                   "berlin_districts_neighbourhoods.R", fsep="/"), chdir =
29                   TRUE)
30 Jmisc::sourceAll(file.path(getwd(), "Helpers", fsep="/"))
31
32 # Load shapefile
33 stations = sf::st_read(file.path(getwd(), "Data", "spatial_data",
34                         "gis_osm_transport_free_1.shp", fsep="/"))
35
36 # Print code development
37 print("Spatial data loaded.")
38
39 # Load dataframes
40 listings_summarized = read_csv(file.path(getwd(), "Data", "airbnb_data",
41                                 "listings.csv", fsep="/"),
42                                 na = c("NA", ""),
43                                 locale = locale(encoding = "UTF-8"))
44 listings_detailed = read_csv(file.path(getwd(), "Data", "airbnb_data",
45                               "listings.csv.gz", fsep="/"),
46                               na = c("NA", ""),
47                               locale = locale(encoding = "UTF-8"))
48 listings_calendar = read_csv(file.path(getwd(), "Data", "airbnb_data",
49                               "calendar.csv.gz", fsep="/"),
50                               na = c("NA", ""),
51                               locale = locale(encoding = "UTF-8"))
52
53 # Print code development
54 print("Airbnb datasets uploaded.")
55
56 # Join dataframes, first clean and keep only variables of interest
57 listings = df_join_clean(df1 = listings_detailed, df2 = listings_summarized)
58
59 # Check if there are missing values
60 apply(listings, 2, function(x) any(is.na(x)))
61

```

```

62 # host_is_superhost, # host_has_profile_pic, host_identity_verified
63 listings = listings %>%
64   dplyr::mutate(host_is_superhost      = ifelse(is.na(host_is_superhost),
65                           FALSE, host_is_superhost),
66   host_has_profile_pic    = ifelse(is.na(host_has_profile_pic
67                           ),
68                           FALSE, host_has_profile_
69                           pic),
70   host_identity_verified = ifelse(is.na(host_identity_
71                           verified),
72                           FALSE, host_identity_
73                           verified))

74 # review_scores_rating
75 # We insert the average of the review_scores_rating
76 listings %>%
77   dplyr::summarize(mean = mean(review_scores_rating, na.rm = TRUE), # 77
78                     median = median(review_scores_rating, na.rm = TRUE), # 95
79                     mode = getmode(review_scores_rating)) %>% # 100
80   round(0)
81 # mean might be influenced by outliers, being it so different from median
82 # and mode, and mode is too high, so I will substitute the missing values
83 # with the median
84 listings = listings %>%
85   dplyr::mutate(review_scores_rating = ifelse(is.na(review_scores_rating),
86                           median(review_scores_rating,
87                           na.rm = TRUE),
88                           review_scores_rating))

89 # bedrooms and beds
90 # Being the amount of missing values in these features relative small, we
91 # can derive their value from the other variables:
92 listings = listings %>%
93   # If beds is NA, but bedrooms has a valid value, we set beds with the
94   # number of bedrooms
95   dplyr::mutate(beds      = ifelse(is.na(beds) & !is.na(bedrooms),
96                           bedrooms, beds),
97   # If bedrooms is NA, but beds has a valid value, we set
98   # bedrooms with the number of beds
99   bedrooms = ifelse(is.na(bedrooms) & !is.na(beds), beds,
100                          bedrooms),

```

```

91      # If both beds and bedrooms are NA: 1
92      # Since in this case the property can accomodate only one
93      # person, we assume that it has 1 bed and 1
94      # bedroom (even if it might not
95      # be separate)
96      beds      = ifelse(is.na(beds) & is.na(bedrooms), 1, beds),
97      bedrooms = ifelse(beds == 1 & is.na(bedrooms), 1, bedrooms))

98
99  # Check if there are missing values
100 apply(listings, 2, function(x) any(is.na(x)))

101 # Listings selection
102 # Keep only listings with at least one customer review, as per Wang and
103 # Nicolau (2017)
104 listings %>%
105   dplyr::filter(number_of_reviews == 0) %>%
106   dplyr::summarize(count = n())
107
108 # Print code development
109 print("Missing values cleaned.")

110 ## New variables
111
112 # Areas: Not doing this with berlin_sf because the polys_sf need to have
113 # geometry sfc_POLYGON
114 # district
115 listings = point_in_polygons(points_df = listings,
116                               polys_sf   = berlin_district_sf,
117                               var_name   = "district")
118 # vbb_zone
119 listings = point_in_polygons(points_df = listings,
120                               polys_sf   = berlin_vbb_AB_sf,
121                               var_name   = "vbb_zone")
122 # neighbourhood
123 listings = point_in_polygons(points_df = listings,
124                               polys_sf   = berlin_neighbourhood_sf,
125                               var_name   = "neighbourhood")
126 # For the sake of consistency, we will delete these listings
127
128 listings = listings %>%

```

```

127 filter(!is.na(district) & !is.na(vbb_zone))

128

129 # Stations

130 railway_stations_df = stations %>%
131   dplyr::filter(fclass %like% "railway") %>%
132   dplyr::rename(id      = name) %>%
133   dplyr::mutate(id      = gsub("Berlin ", "", id),
134                  id      = gsub("Berlin-", "", id),
135                  id      = gsub(" *\\(.*)?\\\" *", "", id),
136                  s_bahn = ifelse(startsWith(id, "S "), TRUE, FALSE),
137                  u_bahn = ifelse(startsWith(id, "U "), TRUE, FALSE),
138                  id      = gsub("S ", "", id),
139                  id      = gsub("U ", "", id)) %>%
140   points_midpoint()

141 railway_stations_df = point_in_polygons(points_df = railway_stations_df,
142                                         polys_sf   = berlin_district_sf,
143                                         var_name   = "district")

144 railway_stations_df = railway_stations_df %>%
145   filter(!is.na(district))

146 listings = distance_count(main = listings, reference = railway_stations_df,
147                           var_name = "station", distance = 1000)

148

149 # (Top 10) attractions

150 attractions_df = data.frame(
151   id    = c("Reichstag", "Brandenburger Tor", "Fernsehturm",
152           "Gendarmenmarkt", "Berliner Dom", "Kurfürstendamm",
153           "Schloss Charlottenburg", "Museuminsel",
154           "Gedenkstätte Berliner Mauer", "Potsdamer Platz"),
155   lat   = c(52.518611, 52.516389, 52.520803, 52.513333, 52.519167,
156           52.500833, 52.521111, 52.520556, 52.535, 52.509444),
157   long  = c(13.376111, 13.377778, 13.40945, 13.393056, 13.401111,
158           13.312778, 13.295833, 13.397222, 13.389722, 13.375833))
159 attractions_df = point_in_polygons(points_df = attractions_df,
160                                     polys_sf   = berlin_district_sf,
161                                     var_name   = "district")

162 listings = distance_count(main = listings, reference = attractions_df,
163                           var_name = "attraction", distance = 2000)

164

165 # season availability

166 listings_calendar = listings_calendar %>%
167   dplyr::rename(id = listing_id) %>%

```

```

168 dplyr::mutate(date = as.Date(date),
169                 year = lubridate::year(date),
170                 month = lubridate::month(date, label = TRUE),
171                 day = lubridate::day(date),
172                 season = ifelse((month == "Mar" & day >= 21) |
173                               (month == "Apr") |
174                               (month == "May") |
175                               (month == "Jun" & day < 21), "Spring",
176                               ifelse((month == "Jun" & day >= 21) |
177                                     (month == "Jul") |
178                                     (month == "Aug") |
179                                     (month == "Sep" & day < 21), "Summer",
180                               ifelse((month == "Sep" & day >= 21) |
181                                     (month == "Oct") |
182                                     (month == "Nov") |
183                                     (month == "Dec" & day < 21), "Fall",
184                               "Winter")))) %>%
185                 factor(levels = c("Spring", "Summer",
186                                   "Fall", "Winter")),
187                 year_season = paste(tolower(season), year, sep = "_")) %>%
188             dplyr::ungroup() %>%
189             dplyr::select(-price)
190
191 # Availability per year_season
192 listings_season_availability = listings_calendar %>%
193   dplyr::group_by(id, year_season) %>%
194   dplyr::summarize(count = sum(available == TRUE),
195                     season_days = n(),
196                     proportion = round(count/season_days*100, 4)) %>%
197   dplyr::arrange(id, year_season) %>%
198   dplyr::ungroup() %>%
199   dplyr::mutate(season_av = paste(year_season, "availability",
200                                     sep = "_")) %>%
201   tolower() %>%
202   factor(levels = paste(year_season,
203                         "availability",
204                         sep = "_")) %>%
205   tolower() %>%
206   unique())) %>%
207   dplyr::select(-count, -year_season, -season_days) %>%
208   tidyr::spread(season_av, proportion)

```

```

209
210 # Print code development
211 print("New variables created.")
212
213 # Join to the listings
214 listings = listings %>%
215   left_join(listings_season_availability, by = "id") %>%
216   replace(is.na(.), 0)
217
218 # Remove unnecessary objects
219 rm("listings_calendar", "listings_season_availability", "stations",
220     "berlin_neighbourhood_sf", "listings_detailed", "listings_summarized")
221 rm(list=lsf.str()) # All functions

```

**Listing 22:** |exploratory\_data\_analysis.R|

```

1 # Install and load needed packages
2 needed_packages = c("tidyverse",
3                      "rstudioapi",
4                      "Jmisc",
5                      "reshape2",
6                      "leaflet",
7                      "mapview",
8                      "xtable")
9
10 for (package in needed_packages) {
11   if (!require(package, character.only=TRUE)) {install.packages(package,
12     character.only=TRUE)}
13   library(package, character.only=TRUE)
14 }
15
16 # Set working directory to the one where the file is located
17
18 # This works when run directly
19 # setwd(dirname(rstudioapi::getActiveDocumentContext()$path))
20
21 # This works when sourced
22 setwd(dirname(sys.frame(1)$ofile))
23
24 # Load helper functions and other scripts
25 source("data_preparation.R", chdir = TRUE)
26 Jmisc::sourceAll(file.path(getwd(), "Helpers", fsep="/"))

```

```

26
27 ## SUMMARY STATISTICS
28
29 # Numeric variables
30 statistics_num = listings %>%
31   dplyr::select(-id, -listing_url, -long, -lat) %>%
32   dplyr::select_if(is.numeric) %>%
33   descriptive_statistics()
34
35 # Categorical variables
36 statistics_cat = listings %>%
37   dplyr::select(-id, -listing_url, -long, -lat) %>%
38   dplyr::select_if(function(x) is.character(x) | is.factor(x)
39               | is.logical(x)) %>%
40   descriptive_statistics()
41
42 # Latex tables
43 # statistics_num %>%
44 #   filter(variable %in% c("price", "review_scores_rating")) %>%
45 #   xtable::xtable() %>%
46 #   print(include.rownames = FALSE)
47 # statistics_cat %>%
48 #   filter(variable == "room_type") %>%
49 #   xtable::xtable() %>%
50 #   print(include.rownames = FALSE)
51
52 ## SUMMARY DATAFRAMES
53
54 # Select numerical and categorical variables
55 num_var = listings %>%
56   dplyr::select(-id, -listing_url, -long, -lat,
57                 -vbb_zone, -district, -neighbourhood) %>%
58   dplyr::select_if(is.numeric) %>%
59   names()
60
61 fact_var = listings %>%
62   dplyr::select(-id, -listing_url, -long, -lat,
63                 -vbb_zone, -district, -neighbourhood) %>%
64   dplyr::select_if(function(x) is.character(x) | is.factor(x)
65                   | is.logical(x)) %>%
66   names()

```

```

67
68 # For all Berlin
69 listings_summary = summarize_df(df = listings,
70                                 vars_mean = num_var,
71                                 vars_count = fact_var)
72
73 # For each VBB Zone
74 listings_vbb_summary = summarize_df(df = listings,
75                                     wrt = "vbb_zone",
76                                     vars_mean = num_var,
77                                     vars_count = fact_var)
78
79 listings_vbb_summary = listings_vbb_summary %>%
80   rename(id = vbb_zone) %>%
81   mutate(view = "VBB Zones",
82         group = id)
83
84 # For each district
85 listings_district_summary = summarize_df(df = listings,
86                                         wrt = "district",
87                                         vars_mean = num_var,
88                                         vars_count = fact_var)
89 listings_district_summary = listings_district_summary %>%
90   rename(id = district) %>%
91   mutate(view = "Districts",
92         group = id)
93
94 # For each neighbourhood
95 listings_neighbourhood_summary = summarize_df(df = listings,
96                                               wrt = "neighbourhood",
97                                               vars_mean = num_var,
98                                               vars_count = fact_var)
99 listings_neighbourhood_summary = berlin_neighbourhoods_names %>%
100   dplyr::select(id, group) %>%
101   dplyr::left_join(listings_neighbourhood_summary,
102                     by = c("id" = "neighbourhood")) %>%
103   dplyr::mutate(view = "Neighbourhoods") %>%
104   replace(is.na(.), 0)
105
106 ## PLOTS
107
```

```

108 # Numeric variables
109
110 # Complete distribution
111 distribution_plot(df = listings, var_name = "price")
112
113 # dev.copy2pdf(file = "./SeminarPaper/price_distribution_complete.pdf")
114 # dev.off()
115
116 # Without outliers
117 distribution_plot(df = listings, var_name = "price", hide_outliers = TRUE)
118
119 # dev.copy2pdf(file = "./SeminarPaper/price_distribution_nooutliers.pdf")
120 # dev.off()
121
122 # Plot factor variables
123 distribution_plot(df = listings, var_name = "station_count")
124
125 # dev.copy2pdf(file = "./SeminarPaper/room_type_distribution.pdf")
126 # dev.off()
127
128 # Maps of variable distribution
129 price_map_distr = var_avg_map(summary_df      = listings_district_summary,
130                                var_name       = "price",
131                                polygon_sf    = berlin_district_sf,
132                                polygon_names_df = berlin_districts_names,
133                                polygon_names_var = "name")
134
135
136 # mapview::mapshot(price_map_distr, file = "./SeminarPaper/price_map_distr.
137 #                         pdf",
138 #                         remove_controls = c("zoomControl", "layersControl", "
139 print("Exploratory Data Analysis done.")
140
141 # Remove unnecessary objects
142 rm("price_map_distr", "fact_var", "num_var")
143 rm(list=lsf.str()) # All functions

```

**Listing 23:** |price\_analysis.R|

```

1 # Install and load needed packages

```

```

2 needed_packages = c("tidyverse",
3
4     "rstudioapi",
5
6     "corrplot",
7
8     "Jmisc",
9
10    "xtable",
11
12    "lindia")
13
14 for (package in needed_packages) {
15
16   if (!require(package, character.only=TRUE)) {
17
18     install.packages(package, character.only=TRUE)}
19
20   library(package, character.only=TRUE)
21
22 }
23
24 rm("needed_packages", "package")
25
26
27 # Set working directory to the one where the file is located
28
29
30 # This works when run directly
31 # setwd(dirname(rstudioapi::getActiveDocumentContext()$path))
32
33
34 # This works when sourced
35 setwd(dirname(sys.frame(1)$ofile))
36
37
38 # Load helper functions and other scripts
39 source("exploratory_data_analysis.R", chdir = TRUE)
40 Jmisc::sourceAll(file.path(getwd(), "Helpers", fsep="/"))
41
42
43 ## Correlation between different variables and price
44
45
46 # Transform categorical variables from rows to column
47 listings_price = listings %>%
48
49   recode_all("room_type") %>% from_row_to_col("room_type") %>%
50
51   recode_all("property_type") %>% from_row_to_col("property_type") %>%
52
53   recode_all("cancellation_policy") %>%
54
55   from_row_to_col("cancellation_policy") %>%
56
57   recode_all("security_deposit") %>% from_row_to_col("security_deposit") %>%
58
59   recode_all("cleaning_fee") %>% from_row_to_col("cleaning_fee") %>%
60
61   recode_all("station_count") %>% from_row_to_col("station_count") %>%
62
63   recode_all("attraction_count") %>% from_row_to_col("attraction_count") %>%
64
65   recode_all("instant_bookable") %>% from_row_to_col("instant_bookable") %>%
66
67   recode_all("host_is_superhost") %>%
68
69   from_row_to_col("host_is_superhost") %>%
70
71   recode_all("host_has_profile_pic") %>%

```

```

43 from_row_to_col("host_has_profile_pic") %>%
44 recode_all("host_identity_verified") %>%
45 from_row_to_col("host_identity_verified") %>%
46 recode_all("district") %>% from_row_to_col("district") %>%
47 recode_all("vbb_zone") %>% from_row_to_col("vbb_zone") %>%
48 dplyr::select(-id, -lat, -long, -listing_url, -neighbourhood)

49

50 # Create dataframe with price correlations
51 listings_price_correlation = cor(listings_price$price, listings_price) %>%
52   as.data.frame() %>%
53   tidyr::gather() %>%
54   dplyr::rename(variable      = key,
55                 correlation = value) %>%
56   dplyr::filter(variable != "price")

57

58 print("Price correlation calculated.")

59

60 # Plot correlation
61 # Example corrplot
62 # listings %>%
63 #   dplyr::select(price, accommodates, bedrooms, beds,
64 #                   minimum_nights, station_dist) %>%
65 #   cor() %>%
66 #   corrplot(method = "circle", type = "upper")
67 # dev.copy2pdf(file = "./SeminarPaper/corrplot.pdf")
68 # dev.off()

69

70 # Correlation with price plot
71 correlation_plot(corr_df      = listings_price_correlation,
72                   variables = c("property_type", "district"))

73

74 # dev.copy2pdf(file = "./SeminarPaper/price_correlation.pdf")
75 # dev.off()

76

77 ## Linear regression

78

79 # Fit linear model
80 listings_price_lm = lm(price ~ ., data = listings_price)

81

82 # Extract coefficients and p-values
83 listings_price_lm_results = data.frame(

```

```

84 variable      = names(listings_price_lm$coefficients
85                         [!is.na(listings_price_lm$coefficients)]),
86 coefficient    = listings_price_lm$coefficients
87                         [!is.na(listings_price_lm$coefficients)],
88 significant     = ifelse(summary(listings_price_lm)
89                         $coefficients [,4] < 0.5, TRUE, FALSE),
90 model_r_squared = summary(listings_price_lm)$r.squared) %>%
91 dplyr::mutate_if(is.numeric, function(x) round(x, 4))

92
93 # listings_price_lm_results %>%
94 #   dplyr::filter(variable %in% c("(Intercept)", "accommodates", "host_
95 #   listings_count")) %>%
96 #   dplyr::select(-model_r_squared) %>%
97 #   xtable::xtable() %>%
98 #   print(include.rownames = FALSE)

99 # Descriptive statistics of residuals
100 listings_price_lm_residuals = data.frame(
101   residuals = listings_price_lm$residuals
102 )
103
104 # descriptive_statistics(listings_price_lm_residuals) %>%
105 #   dplyr::select(-variable) %>%
106 #   xtable::xtable() %>%
107 #   print(include.rownames = FALSE)

108
109 # Fitted vs actual plot
110 ggplot(listings_price_lm, aes(x = .fitted, y = .resid)) +
111   geom_point() +
112   geom_smooth(method = lm) +
113   geom_hline(yintercept=0, col="red", linetype="dashed") +
114   xlab("Fitted values") +
115   ylab("Residuals") +
116   ggtitle("Residual vs Fitted Plot") +
117   theme_bw()
118 # dev.copy2pdf(file = "./SeminarPaper/fittedactual.pdf")
119 # dev.off()

120
121 # Residuals QQ Plot
122 lindia::gg_qqplot(listings_price_lm) +
123   theme_bw()

```

```

124 # dev.copy2pdf(file = "./SeminarPaper/qqplot.pdf")
125 # dev.off()
126
127 print("Linear regression on price calculated.")
128
129 # Remove objects not further needed
130 rm("listings_price_lm", "listings_price_lm_results")
131 rm(list=lsf.str()) # All functions

```

**Listing 24: |clustering.R|**

```

1 # Install and load needed packages
2 needed_packages = c("tidyverse",
3                     "rstudioapi",
4                     "Jmisc",
5                     "ggfortify")
6 for (package in needed_packages) {
7   if (!require(package, character.only=TRUE)) {
8     install.packages(package, character.only=TRUE)}
9   library(package, character.only=TRUE)
10 }
11 rm("needed_packages", "package")
12
13 # Set working directory to the one where the file is located
14
15 # This works when run directly
16 # setwd(dirname(rstudioapi::getActiveDocumentContext()$path))
17
18 # This works when sourced
19 setwd(dirname(sys.frame(1)$ofile))
20
21 # Load helper functions and other scripts
22 source("price_analysis.R", chdir = TRUE)
23 Jmisc::sourceAll(file.path(getwd(), "Helpers", fsep="/"))
24
25 ## Calculate number of clusters
26
27 # Scale dataframe
28 listings_scaled = listings %>%
29   recode_all("room_type") %>% from_row_to_col("room_type") %>%
30   recode_all("property_type") %>% from_row_to_col("property_type") %>%
31   recode_all("cancellation_policy") %>%

```

```

32 from_row_to_col("cancellation_policy") %>%
33 recode_all("security_deposit") %>% from_row_to_col("security_deposit") %>%
34 recode_all("cleaning_fee") %>% from_row_to_col("cleaning_fee") %>%
35 recode_all("station_count") %>% from_row_to_col("station_count") %>%
36 recode_all("attraction_count") %>% from_row_to_col("attraction_count") %>%
37 recode_all("instant_bookable") %>% from_row_to_col("instant_bookable") %>%
38 recode_all("host_is_superhost") %>%
39 from_row_to_col("host_is_superhost") %>%
40 recode_all("host_has_profile_pic") %>%
41 from_row_to_col("host_has_profile_pic") %>%
42 recode_all("host_identity_verified") %>%
43 from_row_to_col("host_identity_verified") %>%
44 dplyr::select(-id, -lat, -long, -listing_url,
45               -district, -vbb_zone, -neighbourhood) %>%
46 scale()

47

48 # Total variance explained
49 number_of_clusters(scaled_df      = listings_scaled,
50                     max          = 96,
51                     iter_max     = 20,
52                     plot_breaks = seq(0, 100, by = 5))

53

54 # dev.copy2pdf(file = "./SeminarPaper/numclusters.pdf")
55 # dev.off()

56

57 ## Clustering
58 n_clusters = 50
59 listings_kmeans = kmeans(listings_scaled,
60                          centers = n_clusters, iter.max = 20)

61

62 listings_clusters = listings_kmeans$cluster %>% as.factor()

63

64 # Plot clusters in map
65 ggplot() +
66   geom_sf(data = berlin_district_sf,
67            show.legend = FALSE, color = "black") +
68   coord_sf(datum = NA) +
69   geom_point(data = listings,
70              aes(x = long, y = lat, color = listings_clusters),
71              alpha = 0.5) +
72   theme(plot.title = element_text(hjust = 0.5)) +

```

```

73     ggtile("Clustering of airbnb properties") +
74     labs(color = "clusters") +
75     theme_void()
76
77 # dev.copy2pdf(file = "./SeminarPaper/cluster_map.pdf")
78 # dev.off()
79
80 # Principal components
81 listings_pc = prcomp(x = listings_scaled, center = FALSE, scale = FALSE)
82
83 ggplot2::autoplot(listings_pc, data = listings_kmeans, colour = "cluster") +
84   ggtile("Plotting clusters wrt Principal Components") +
85   labs(color = "clusters") +
86   theme_bw()
87
88 # dev.copy2pdf(file = "./SeminarPaper/cluster_plot.pdf")
89 # dev.off()
90
91 # Remove objects not further needed
92 rm("n_clusters")
93 rm(list=lsf.str()) # All functions

```

**Listing 25:** |shiny\_preparation.R|

```

1 # Install and load needed packages
2 needed_packages = c("tidyverse",
3                      "rstudioapi",
4                      "sf")
5 for (package in needed_packages) {
6   if (!require(package, character.only=TRUE)) {
7     install.packages(package, character.only=TRUE)}
8   library(package, character.only=TRUE)
9 }
10 rm("needed_packages", "package")
11
12 # Set working directory to the one where the file is located
13
14 # This works when run directly
15 setwd(dirname(rstudioapi:::getActiveDocumentContext()$path))
16
17 # This works when sourced
18 # setwd(dirname(sys.frame(1)$ofile))

```

```

19
20 # Load scripts
21 source("clustering.R", chdir = TRUE)
22
23 # Calculate the center of Berlin polygon
24 berlin_center = berlin_neighbourhood_singlebuckow_sf %>%
25   dplyr::group_by() %>%
26   dplyr::summarize(do_union = TRUE) %>%
27   sf::st_centroid() %>%
28   sf::st_coordinates() %>%
29   base::as.data.frame() %>%
30   dplyr::rename(long = X,
31                 lat = Y) %>%
32   dplyr::select(lat, long)
33
34 # Join area summaries
35 listings_area_summary = rbind(listings_district_summary, listings_vbb_
36   summary, listings_neighbourhood_summary)
37
38 # Join sf objects
39 berlin_district_sf = berlin_district_sf %>%
40   dplyr::mutate(view = "Districts")
41 berlin_neighbourhood_singlebuckow_sf =
42   berlin_neighbourhood_singlebuckow_sf %>%
43   dplyr::mutate(view = "Neighbourhoods")
44 berlin_vbb_AB_sf = berlin_vbb_AB_sf %>%
45   dplyr::mutate(view = "VBB Zones",
46                 group = id)
47 berlin_sf = rbind(berlin_district_sf,
48                   berlin_neighbourhood_singlebuckow_sf,
49                   berlin_vbb_AB_sf)
50
51 # Join names dataframes
52 berlin_districts_names = berlin_districts_names %>%
53   dplyr::mutate(view = "Districts")
54 berlin_neighbourhoods_names = berlin_neighbourhoods_names %>%
55   dplyr::mutate(view = "Neighbourhoods")
56 berlin_vbb_AB_names = berlin_vbb_AB_names %>%
57   dplyr::mutate(view = "VBB Zones",
58                 group = id)
59 berlin_names = rbind(berlin_districts_names,

```

```

59         berlin_neighbourhoods_names,
60         berlin_vbb_AB_names)
61
62 # Recode and transform from row to columns factor variables
63 listings_recoded = listings %>%
64   recode_all("room_type") %>% from_row_to_col("room_type") %>%
65   recode_all("property_type") %>% from_row_to_col("property_type") %>%
66   recode_all("cancellation_policy") %>%
67   from_row_to_col("cancellation_policy")
68
69 # Remove objects not further needed
70 rm("listings_district_summary", "listings_vbb_summary",
71   "listings_neighbourhood_summary", "berlin_district_sf",
72   "berlin_neighbourhood_singlebuckow_sf", "berlin_vbb_AB_sf",
73   "berlin_districts_names", "berlin_neighbourhoods_names",
74   "berlin_vbb_AB_names")
75 rm(list=lsf.str()) # All functions
76
77 # Save workspace
78 # save.image(file = file.path(getwd(), "ShinyApp", "workspace_for_shinyapp.RData", fsep="/"))

```

## A.2 Functions

**Listing 26:** |capwords.R|

```

1 # Capitalize the first letter of every word
2 # from the help for toupper()
3 capwords <- function(s, strict = FALSE) {
4   cap <- function(s) paste(toupper(substring(s, 1, 1)),
5                           {s <- substring(s, 2);
6                            if(strict) tolower(s) else s},
7                           sep = "", collapse = " "))
8   sapply(strsplit(s, split = " "), cap, USE.NAMES = !is.null(names(s)))
9 }

```

**Listing 27:** |correlation\_plot.R|

```

1 correlation_plot = function(corr_df, variables = corr_df$variable,
2                             variables_original = TRUE) {
3
4   correlation = corr_df %>%

```

```

5   dplyr::mutate(sign      = ifelse(correlation > 0, "positive",
6                   ifelse(correlation < 0, "negative",
7                           "null")),
8                   correlation = abs(correlation)) %>%
9   dplyr::arrange(variable)
10
11  if (all(correlation$variable == corr_df$variable)) {
12
13    # If variables are listed with their names from the original dataframe
14    # (before dummies)
15
16    if (variables_original) {
17
18      correlation = correlation %>%
19          dplyr::filter(grepl(paste(variables, collapse = " | "),
20                          corr_df$variable))
21
22    # If variables are listed with their names from corr_df
23    } else {
24
25      correlation = correlation %>%
26          dplyr::filter(variable %in% variables)
27
28    }
29
30  corr_plot = correlation %>%
31    ggplot(aes(x = variable, y = correlation, fill = sign)) +
32    geom_bar(stat="identity") +
33    scale_fill_manual("sign",
34                      values = c("positive" = "green4",
35                                "negative" = "red1",
36                                "null" = "blue3")) +
37    coord_flip() +
38    labs(x = "variables", y = "correlation",
39         title = "Correlation with price") +
40    theme_bw() +
41    theme(axis.text.y      = element_text(
42                                      color = ifelse(
43                                          correlation$sign == "positive",
44                                          "green4",

```

```
45
46         ifelse(
47             correlation$sign == "negative",
48             "red1",
49             "blue3"))),
50
51     axis.text.x      = element_text(size = rel(1.2)),
52
53     axis.title.x    = element_text(size = rel(1.2)),
54
55     axis.title.y    = element_text(size = rel(1.2)),
56
57     legend.text     = element_text(size = rel(1.2)),
58
59     legend.position = "bottom") +
60
61     scale_y_continuous(expand = expand_scale(mult = c(0, 0.05)))
62
63
64
65
66     return(corr_plot)
67
68 }
```

**Listing 28:** |descriptive statistics.R|

```

1 descriptive_statistics = function(df) {
2
3   # Descriptive statistics for numeric variables
4   if (unique(apply(df, 2, function(x) is.numeric(x)))) {
5
6     summary = data.frame(
7       variable = names(df),
8       min = apply(df, 2, min),
9       '1Q' = apply(df, 2, quantile, probs = 0.25),
10      median = apply(df, 2, median),
11      '3Q' = apply(df, 2, quantile, probs = 0.75),
12      max = apply(df, 2, max),
13      iqr = apply(df, 2, IQR),
14      mean = apply(df, 2, mean),
15      sd = apply(df, 2, sd),
16      check.names = FALSE) %>%
17      dplyr::mutate_if(is.numeric, function(x) round(x, 4))
18
19   # Descriptive statistics for categorical variables
20 } else if (unique(apply(df, 2, function(x) is.character(x) | is.factor(x)
21 | is.logical(x)))) {
22
23   # Frequency
24   frequency_list = apply(df, 2, table)

```

```

25 frequency_df = data.frame(frequency = unlist(frequency_list)) %>%
26   tibble::rownames_to_column(var = "var_fact")
27
28 freq_var = colsplit(string = frequency_df$var_fact, pattern = "\\\\".,
29   names = c("variable", "factor"))
30
31 frequency = freq_var %>%
32   cbind(frequency_df) %>%
33   dplyr::select(-var_fact)
34
35 # Proportion
36 proportion_list = apply(df, 2, function(x) prop.table(table(x)))
37
38 proportion_df = data.frame(proportion = unlist(proportion_list)) %>%
39   tibble::rownames_to_column(var = "var_fact")
40
41 prop_var = colsplit(string = proportion_df$var_fact, pattern = "\\\\".,
42   names = c("variable", "factor"))
43
44 proportion = prop_var %>%
45   cbind(proportion_df) %>%
46   dplyr::select(-var_fact) %>%
47   dplyr::mutate(proportion = round(proportion, 4)*100,
48                 proportion = as.character(proportion) %>% paste("%"))
49
50 summary = frequency %>%
51   dplyr::inner_join(proportion, by = c("variable", "factor")) %>%
52   dplyr::arrange(variable, desc(frequency))
53
54 } else {
55
56   stop("Check that all variables are either numeric or categorical.")
57
58 }
59
60 return(summary)

```

**Listing 29:** |df\_join\_clean.R|

```

1 # Join dataframes and cleaning

```

```

2 df_join_clean = function(df1, df2) {
3
4   df1 = df1 %>%
5     # price and fee columns
6   dplyr::mutate_if(grepl("price|cleaning|deposit|people" ,
7                         names(.)), funs(gsub("\\\\$", "", .))) %>%
8   dplyr::mutate_if(grepl("price|cleaning|deposit|people" ,
9                         names(.)), funs(gsub("\\\\,", "", .))) %>%
10  dplyr::mutate_if(grepl("price|cleaning|deposit|people" ,
11                     names(.)), funs(as.numeric(.)))
12
13  df2 = df2 %>%
14    # price and fee columns
15  dplyr::mutate_if(grepl("price|cleaning|deposit|people" ,
16                        names(.)), funs(gsub("\\\\$", "", .))) %>%
17  dplyr::mutate_if(grepl("price|cleaning|deposit|people" ,
18                        names(.)), funs(gsub("\\\\,", "", .))) %>%
19  dplyr::mutate_if(grepl("price|cleaning|deposit|people" ,
20                      names(.)), funs(as.numeric(.)))
21
22  df_joined = full_join(df1, df2, by = c("id", "host_id", "room_type", "price",
23                           "minimum_nights",
24                           "number_of_reviews",
25                           "latitude", "longitude",
26                           "availability_365", "calculated_host_listings_count")) %>%
27
#First selection
28  dplyr::select(id, longitude, latitude, price, property_type, room_type,
29                 instant_bookable, listing_url, host_is_superhost,
30                 security_deposit, cleaning_fee,
31                 calculated_host_listings_count, host_has_profile_pic, host
32                 _identity_verified,
33                 accommodates, bedrooms, beds, minimum_nights,
34                 review_scores_rating, number_of_reviews, cancellation_
35                 policy,
36                 availability_30, availability_60, availability_90,
37                 availability_365) %>%
38
39  dplyr::rename(lat           = latitude,
40                long          = longitude,
41                host_listings_count = calculated_host_listings_count) %>%

```

```

38 # property_type
39 dplyr::group_by(property_type) %>%
40 dplyr::mutate(property_type_count = n()) %>%
41 dplyr::ungroup() %>%
42 dplyr::mutate(property_type      = as.factor(ifelse(property_type_count
43             <= 10,
44                         "Other", property_type)),
45
46 # security_deposit
47             security_deposit = ifelse(security_deposit == 0 |
48                                         is.na(security_deposit), FALSE,
49                                         TRUE),
50
51 # cleaning_fee
52             cleaning_fee     = ifelse(cleaning_fee == 0 |
53                                         is.na(cleaning_fee), FALSE,
54                                         TRUE),
55
56 # availabilities
57             availability_30   = round(availability_30/30*100, 4),
58             availability_60   = round(availability_60/60*100, 4),
59             availability_90   = round(availability_90/90*100, 4),
60             availability_365  = round(availability_365/365*100, 4))
61             %>%
62
63 dplyr::select(-property_type_count)
64
65
66 return(df_joined)
67
68 }

```

**Listing 30:** |distance\_count.R|

```

1 # Calculate the distance from the nearest reference point (from reference df
2
3 # and the number of reference points within a certain distance
4
5 distance_count = function(main, reference, var_name, distance) {
6
7     # Create variable names
8
9     var_name_count = paste(var_name, "count", sep = "_")
10    var_name_dist = paste(var_name, "dist", sep = "_")
11
12    # Calculate distance for each listing to each station
13
14    point_distance = geosphere::distm(x = main %>%
15
16                                dplyr::select(long, lat),
17
18                                y = reference %>%
19
20                                dplyr::select(long, lat),
21
22                                fun = distHaversine) %>%
23
24                                as.data.frame() %>%

```

```

14     data.table::setnames(as.character(reference$id))
15 
16 # Calculate how many "reference" are within "distance"
17 point_distance[,var_name_count] = rowSums(point_distance <= distance)
18 %>% as.factor()
19 
20 # Calculate the distance to the nearest "reference"
21 point_distance[,var_name_dist] = apply(point_distance[,-ncol(point_
22 distance)], MARGIN = 1,
23 FUN = min) %>%
24 round(0)
25 
26 # Insert this information into the main DF
27 main = point_distance %>%
28   dplyr::mutate(id = main$id) %>%
29   dplyr::select(id, var_name_count, var_name_dist) %>%
30   dplyr::right_join(main, by = "id")
31 
32 return(main)
33 }
```

**Listing 31: |distribution\_plot.R|**

```

1 distribution_plot = function(df, var_name, hide_outliers = FALSE, tilt_text_
2   x = FALSE) {
3 
4   variable = sym(var_name)
5 
6   # Plot for numeric variables
7   if (is.numeric(df[, var_name])) {
8 
9     numvar_plot = df %>%
10       dplyr::transmute(
11         var = !!variable,
12         median = median(var),
13         mean = mean(var),
14         '1Q' = quantile(var, probs = 0.25),
15         '3Q' = quantile(var, probs = 0.75),
16         IQR = IQR(var),
17         upper = (IQR*1.5) + '3Q',
18         lower = '1Q' - (IQR*1.5),
19         outlier_upper = ifelse(var > unique(upper), 1, 0),
20         outlier_lower = ifelse(var < unique(lower), 1, 0))
21 
22     if (!hide_outliers) {
23       numvar_plot = numvar_plot %>%
24         filter(var > lower & var < upper)
25     }
26 
27     if (tilt_text_x) {
28       numvar_plot = numvar_plot %>%
29         mutate(tilt_text_x = ifelse(var == '1Q', 1.5, 0))
30     }
31 
32     numvar_plot %>%
33       ggplot(aes(x = var)) +
34       geom_boxplot() +
35       geom_point(data = main, aes(x = id, y = var_name_dist))
36   }
37 }
```

```

19         outlier_lower = ifelse(var < unique(lower), 1, 0),
20         outlier        = ifelse(outlier_upper == 1 | outlier_lower
21             == 1, 1, 0),
22         count_upper    = sum(outlier_upper),
23         count_lower    = sum(outlier_lower)
24     )
25
26 if (hide_outliers) {
27
28     plot = numvar_plot %>%
29         dplyr::filter(outlier == 0) %>%
30         ggplot() +
31         geom_density(aes(x = var), fill = "green4", color = "green4") +
32         geom_vline(aes(xintercept = unique(median)),
33                     color = "blue", linetype = "dashed") +
34         annotate(geom = "text", x = (unique(numvar_plot$median) + 2.6), y =
35             0,
36                     label = "Median",
37                     colour = "blue",
38                     angle = 90, hjust = 0) +
39         geom_vline(aes(xintercept = unique(mean)),
40                     color = "black", linetype = "dashed") +
41         annotate(geom = "text", x = (unique(numvar_plot$mean) + 2.6), y = 0,
42                     label = "Mean",
43                     colour = "black",
44                     angle = 90, hjust = 0) +
45         geom_vline(aes(xintercept = '1Q'),
46                     color = "red", linetype = "dashed") +
47         annotate(geom = "text", x = (unique(numvar_plot$'1Q') + 2.6), y = 0,
48                     label = "1st Quantile",
49                     colour = "red",
50                     angle = 90, hjust = 0) +
51         geom_vline(aes(xintercept = '3Q'),
52                     color = "red", linetype = "dashed") +
53         annotate(geom = "text", x = (unique(numvar_plot$'3Q') + 2.6), y = 0,
54                     label = "3rd Quantile",
55                     colour = "red",
56                     angle = 90, hjust = 0) +
57     labs(caption = paste("Outliers not displayed:",
58                         numvar_plot$count_lower, "values under 1,5*IQR
59                         and",
60

```

```

57             numvar_plot$count_upper, "values over 1,5*IQR",
58             sep = " "),
59         title = paste("Density of the variable", variable, sep = " "),
60         x = var_name) +
61     theme_bw() +
62     theme(axis.text.x = element_text(size = rel(1.2)),
63           axis.text.y = element_text(size = rel(1.2)),
64           axis.title.x = element_text(size = rel(1.2)),
65           axis.title.y = element_text(size = rel(1.2)))
66 } else {
67
68   plot = numvar_plot %>%
69     ggplot() +
70     geom_density(aes(x = var), fill = "green4", color = "green4") +
71     geom_vline(aes(xintercept = unique(median)),
72                 color = "blue", linetype = "dashed") +
73     geom_vline(aes(xintercept = unique(mean)),
74                 color = "black", linetype = "dashed") +
75     geom_vline(aes(xintercept = '1Q'),
76                 color = "red", linetype = "dashed") +
77     geom_vline(aes(xintercept = '3Q'),
78                 color = "red", linetype = "dashed") +
79     labs(x = var_name,
80           title = paste("Density of the variable", variable, sep = " "))
81     +
82     theme_bw() +
83     theme(axis.text.x = element_text(size = rel(1.2)),
84           axis.text.y = element_text(size = rel(1.2)),
85           axis.title.x = element_text(size = rel(1.2)),
86           axis.title.y = element_text(size = rel(1.2)))
87 }
88
89 if (tilt_text_x) {
90
91   plot = plot +
92     theme(axis.text.x = element_text(angle = 90, hjust = 1, vjust = 0.5)
93       )
94 }
95

```

```

96  # Plot for categorical variables
97 } else if (is.character(df[, var_name]) |
98   is.factor(df[, var_name]) |
99   is.logical(df[, var_name])) {
100
101 factvar_plot = listings %>%
102   dplyr::rename(var = !!variable) %>%
103   dplyr::group_by(var) %>%
104   dplyr::summarize(count = n())
105
106 plot = ggplot(factvar_plot) +
107   geom_bar(aes(x = var, y = count), fill = "green4", stat = "identity")
108   +
109   labs(x      = variable,
110         title = paste("Histogram of the variable", variable, sep = " "))
111   +
112   theme_bw() +
113   theme(axis.text.x = element_text(size = rel(1.2)),
114         axis.text.y = element_text(size = rel(1.2)),
115         axis.title.x = element_text(size = rel(1.2)),
116         axis.title.y = element_text(size = rel(1.2)))
117
118 if (tilt_text_x) {
119
120   plot = plot +
121     theme(axis.text.x = element_text(angle = 90, hjust = 1, vjust = 0.5))
122
123 }
124
125 if (hide_outliers) {
126
127   stop("hide_outliers cannot be TRUE for categorical variables.")
128
129 }
130
131 } else {
132
133   stop("Check what type of variable you area dealing with.")

```

```

134 }
135
136 plot
137
138 }
```

**Listing 32: |from\_row\_to\_col.R|**

```

1 # Spread rows into columns using just one variable
2 from_row_to_col = function(df, variable) {
3   df = df %>%
4     dplyr::mutate(yesno = 1) %>%
5     dplyr::distinct() %>%
6     tidyr::spread(variable, yesno, fill = 0)
7   return(df)
8 }
```

**Listing 33: |getmode.R|**

```

1 # getmode function
2 # from: https://www.tutorialspoint.com/r/r_mean_median_mode.htm
3 getmode <- function(v) {
4   uniqv <- unique(v)
5   uniqv[which.max(tabulate(match(v, uniqv)))]
6 }
```

**Listing 34: |number\_of\_clusters.R|**

```

1 # Calculate plot of tve to choose number of clusters for k-means clustering
2 number_of_clusters = function(scaled_df, min = 2, max,
3                               iter_max = 10, plot_breaks) {
4   # Set seed for reproducibility
5   set.seed(900114)
6   # Values to test
7   k_values = min:max
8   # Empty dataframe for the total variance explained
9   tve = data.frame(clusters = k_values,
10                     tve = rep(NA, length(k_values)))
11  # Empty list for the kmeans objects
12  clk = list()
13  # Loop through the possible values of k
14  for (k in k_values) {
15    # Calculate k-means for each k
```

```

16 clk[[k-1]] = kmeans(scaled_df, centers = k, iter.max = iter_max)
17 # Save the number of clusters
18 names(clk)[k-1] = paste(k, "clusters", sep = " ")
19 # Calculate percentage of total variance explained
20 tve$tve[k-1] = 1 - clk[[k-1]]$tot.withinss / clk[[k-1]]$totss
21 # Print process
22 print(paste("k-means with", k, "clusters done", sep = " "))
23 }

24

25 # Plot tve against k values
26 plot = ggplot(data = tve, aes(x = clusters, y = tve)) +
27   geom_line(color = "grey") +
28   geom_point(color = "red") +
29   scale_x_continuous(breaks = plot_breaks) +
30   labs(x      = "number of clusters",
31        y      = "% of tve",
32        title = paste("Plot of total variance explained for k from", min,
33                      "to", max, sep = " ")) +
34   theme_bw() +
35   theme(axis.text.x = element_text(size = rel(1.2)),
36         axis.text.y = element_text(size = rel(1.2)),
37         axis.title.x = element_text(size = rel(1.2)),
38         axis.title.y = element_text(size = rel(1.2)))
39
40 return(plot)
}

```

Listing 35: |point\_in\_polygons.R|

```

1 # Tells in which polygon a point is
2 point_in_polygons <- function(points_df, polys_sf,
3                                var_name, join_var = "id") {
4
5   # Create empty dataframe
6   is_in <- data.frame(matrix(ncol = nrow(polys_sf), nrow = nrow(points_df)))
7   is_in[,var_name] <- NA
8
9   # Extract coordinates of the polygons
10  coordinates <- as.data.frame(st_coordinates(polys_sf))
11
12  # Extract names of the polygons
13  name <- as.character(polys_sf$id)
14
15  # For all polygons check if the points are inside of them
16  for (k in 1:nrow(polys_sf)) {
17    is_in[,k] <- sp::point.in.polygon(point.x = points_df$long,
18

```

```

14         point.y = points_df$lat,
15         pol.x   = coordinates$X
16                           [coordinates$L2 == k],
17         pol.y   = coordinates$Y
18                           [coordinates$L2 == k])
19
20     # Get the names of the polygons where the points are in one column
21     is_in[,var_name][is_in[,k] == 1] <- name[k]
22   }
23
24   # Keep only summary column and add points' names
25   is_in <- is_in %>%
26     dplyr::select(var_name) %>%
27     dplyr::mutate(id = points_df$id)
28
29   # Add the summary column to the points dataframe
30   points_df <- dplyr::full_join(points_df, is_in, by = join_var)
31
32   return(points_df)
33 }
```

**Listing 36:** |points\_midpoint.R|

```

1 # Calculate midpoint from a sf point object
2 points_midpoint <- function(points_sf, point_name = "id") {
3
4   # Create symbol version of point_name
5   point_name <- sym(point_name)
6
7   # Get coordinates from sf point objects
8   coordinates <- sf::st_coordinates(points_sf) %>%
9     base::as.data.frame() %>%
10    dplyr::rename(lat = Y,
11                  long = X)
12
13   # Bind coordinates with sf point object and calculate mean of lat and long
14   midpoints_df <- points_sf %>%
15     base::as.data.frame() %>%
16     dplyr::select(-geometry) %>%
17     base::cbind(coordinates) %>%
18     dplyr::group_by(!!point_name) %>%
19     dplyr::summarize(lat      = mean(lat),
20                       long     = mean(long),
21                       s_bahn  = any(s_bahn),
22                       u_bahn  = any(u_bahn))
23
24   return(midpoints_df)
25 }
```

**Listing 37:** |recode\_all.R|

```

1 # Recode variable: rename all factor levels of that variable
2 recode_all = function(df, variable) {
3   variable_name = sym(variable)
4   df = df %>%
5     mutate (!!variable_name := as.factor (!!variable_name))
6   levels(df[,variable]) = gsub(" ", "_", levels(df[,variable])) %>%
7     tolower() %>%
8     paste(variable, ., sep = "_")
9   return(df)
10 }

```

**Listing 38:** |summarize\_df.R|

```

1 # Summary information for numeric and factor variables
2 summarize_df <- function(df, wrt = NULL, vars_mean, vars_count) {
3   # Create empty vector
4   summary_df <- c()
5   # For the case in which the summary is general
6   if(is.null(wrt)) {
7     # Variables of which we calculate the mean
8     for (var in vars_mean) {
9       var_name <- sym(var)
10      mean_df <- df %>% dplyr::summarize (!!var_name := mean (!!var_name))
11      summary_df <- c(summary_df, mean_df)
12    }
13    # Variables of which we count the unique values of the factors
14    for (var in vars_count) {
15      var_name <- sym(var)
16      count_df <- df %>%
17        recode_all(var) %>%
18        group_by (!!var_name) %>%
19        dplyr::summarize(count = n()) %>%
20        spread(key = !!var_name, value = count)
21      summary_df <- summary_df %>% cbind(count_df)
22    }
23    # For the case in which the summary is done wrt a variable
24  } else {
25    wrt_name <- sym(wrt)
26    # Variables of which we calculate the mean
27    for (var in vars_mean) {
28      var_name <- sym(var)

```

```

29   mean_df <- df %>%
30     group_by(!!wrt_name) %>%
31       dplyr::summarize (!!var_name := mean(!!var_name))
32   summary_df <- c(summary_df, mean_df)
33   summary_df[duplicated(summary_df)] <- NULL
34 }
35 # Variables of which we count the unique values of the factors
36 for (var in vars_count) {
37   var_name <- sym(var)
38   count_df <- df %>%
39     recode_all(var) %>%
40     group_by(!!var_name, !!wrt_name) %>%
41       dplyr::summarize(count = n()) %>%
42         spread(key = !!var_name, value = count)
43   summary_df <- summary_df %>% as.data.frame() %>%
44     dplyr::mutate(!!wrt_name := as.character(!!wrt_name)) %>%
45     full_join(count_df, by = wrt) %>%
46     replace(is.na(.), 0)
47 }
48 # Percentage of the rows
49 percentage_df <- df %>%
50   group_by(!!wrt_name) %>%
51     dplyr::summarize(percentage = round(n()/nrow(.)*100, 2))
52   summary_df <- summary_df %>% full_join(percentage_df, by = wrt)
53 }
54 return(summary_df)
55 }
```

**Listing 39: |summary\_stat\_fact.R|**

```

1 summary_stat_fact = function(df, var, constraint_var = NULL, constraint_
2   value = NULL) {
3
4   if (is.null(constraint_var) & is.null(constraint_value)) {
5     factor = df[, var]
6   } else if (is.null(constraint_var) | is.null(constraint_value)) {
7     stop("constraint_var or constraint_value is not defined. Review function
8       values.")
9   } else {
10    constraint = sym(constraint_var)
11    df          = df %>%
12      dplyr::filter (!!constraint == constraint_value)
13  }
14}
```

```

11     factor    = df[, var]
12 }
13
14 frequency_df = table(factor) %>%
15   as.data.frame() %>%
16   dplyr::rename(frequency = Freq)
17
18 proportion_df = prop.table(table(factor)) %>%
19   as.data.frame() %>%
20   dplyr::rename(proportion = Freq) %>%
21   dplyr::mutate(proportion = round(proportion, 4)*100,
22                 proportion = as.character(proportion) %>% paste("%"))
23
24 summary = frequency_df %>%
25   dplyr::inner_join(proportion_df, by = "factor") %>%
26   dplyr::arrange(desc(frequency))
27
28 return(summary)
29
30 }
```

**Listing 40: |summary\_stat\_num.R|**

```

1 summary_stat_num = function(df, var, constraint_var = NULL, constraint_value
2                           = NULL) {
3
4   if (is.null(constraint_var) & is.null(constraint_value)) {
5     variable = df[, var]
6   } else if (is.null(constraint_var) | is.null(constraint_value)) {
7     stop("constraint_var or constraint_value is not defined. Review function
8          values.")
9   } else {
10     constraint = sym(constraint_var)
11     df         = df %>%
12       dplyr::filter (!!constraint == constraint_value)
13     variable   = df[, var]
14   }
15
16   summary = data.frame(
17     min      = min(variable),
18     `1Q`    = quantile(variable, probs = 0.25),
19     median  = median(variable),
```

```

18     `3Q`    = quantile(variable, probs = 0.75),
19     max     = max(variable),
20     iqr     = IQR(variable),
21     mean    = mean(variable),
22     sd      = sd(variable),
23     check.names = FALSE) %>%
24   dplyr::mutate_if(is.numeric, function(x) round(x, 4))
25
26 return(summary)
27
28 }

```

**Listing 41:** |var\_avg\_map.R|

```

1 var_avg_map = function(summary_df, var_name,
2                           polygon_sf, polygon_names_df, polygon_names_var,
3                           bins = TRUE, nsplits = 4, pretty = FALSE, zoom_level
4                           = 10,
5                           title = paste("Average of", var_name, sep = " ")) {
6
7   if (bins) {
8
9     colors = colorBin(topo.colors(nsplits), summary_df[, var_name], bins =
10       nsplits, pretty = pretty)
11
12   } else {
13
14     if (pretty) {
15
16       stop("Parameter pretty not used in colorQuantile.")
17
18     }
19
20     colors = colorQuantile(topo.colors(nsplits), summary_df[, var_name],
21       bins = nsplits)
22
23   }
24
25   map = leaflet() %>%
26     addTiles() %>%
27     addPolygons(data = polygon_sf, weight = 1, smoothFactor = 1, fillOpacity

```

```

= 0.8,
      color = "grey",
      fillColor = ~colors(summary_df[, var_name])) %>%
addLabelOnlyMarkers(data = polygon_names_df,
      lng = ~long, lat = ~lat, label = ~lapply(polygon_
names_df[,polygon_names_var], htmltools::HTML),
      labelOptions = labelOptions(noHide = TRUE, direction
= 'center',
      textOnly = TRUE,
      textsize = 20,
      style = list("color" =
black))) %>%
addLegend(position = "bottomleft", pal = colors,
      values = summary_df[, var_name],
      labFormat = labelFormat(transform = function(x) sort(x,
decreasing = FALSE)),
      title = title)

return(map)
}

```

### A.3 ShinyApp

Listing 42: |ui.R|

```

1 ui = dashboardPage(
2   skin = "red",
3   dashboardHeader(title = "Airbnb in Berlin"),
4   dashboardSidebar(
5     sidebarMenu(
6       id = "tabs",
7       menuItem("Intro",
8         tabName = "intro",
9         icon = icon("play-circle")),
10      menuItem("Maps",
11        tabName = "maps",
12        icon = icon("map")),
13      menuItem("Price",
14        tabName = "price",
15        icon = icon("dollar")),

```

```

16     menuItem("Clustering",
17               tabName = "cluster",
18               icon = icon("object-group"))
19 )
20 ),
21 dashboardBody(
22   tabItems(
23
24     tabItem(tabName = "intro",
25       fluidRow(
26         box(
27           title = "Welcome to the ShinyApp on Airbnb properties in
28             Berlin!",
29           width = 12,
30           status = "danger",
31           uiOutput("intro_ui"),
32           br(),
33           column(width = 4,
34             h4("Maps Tab"),
35             textOutput("intro_maps_text"),
36             br(),
37             actionButton("switch_maps", "Go to Maps Tab")),
38           column(width = 4,
39             h4("Price Tab"),
40             textOutput("intro_price_text"),
41             br(),
42             actionButton("switch_price", "Go to Price Tab")),
43           column(width = 4,
44             h4("Clustering Tab"),
45             textOutput("intro_cluster_text"),
46             br(),
47             actionButton("switch_cluster",
48                         "Go to Clustering Tab")))
49 )
50
51   ## TAB2 (Maps)
52   tabItem(tabName = "maps",
53     fluidRow(
54       tabBox(id = "berlin",
55             width = 12,

```

```

56   tabPanel(
57     title = "Berlin map",
58     textOutput("berlin_text"),
59     br(),
60     box(
61       width = 12,
62       solidHeader=TRUE,
63       column(
64         width = 6,
65         selectInput(
66           inputId = "view",
67           label = "Choose view",
68           choices = berlin_sf %>%
69             as.data.frame() %>%
70             select(view) %>%
71             unique() %>%
72             t() %>%
73             c() %>%
74             setdiff("Neighbourhoods"))),
75       column(
76         width = 6,
77         selectInput(
78           inputId = "variable1",
79           label = "Choose variable",
80           choices = c("none", listings_area_summary %>%
81             colnames() %>%
82             setdiff(c("id", "view", "group")) %>%
83             gsub("_", " ", .)))),
84     ),
85     box(
86       width = 12,
87       solidHeader=TRUE,
88       column(
89         width = 4,
90         checkboxInput(inputId = "main_properties",
91                       label = "Display properties"
92                       ,
93                       value = FALSE)),
94       column(
95         width = 4,

```

```

95         checkboxInput(inputId = "main_stations",
96                           label    = "Display train/
97                               subway stations",
98                           value    = FALSE)),
99
100        column(
101            width = 4,
102            checkboxInput(
103                inputId = "main_sights",
104                label    = "Display top 10 attractions",
105                value    = FALSE))
106        ),
107
108        box(
109            width      = 12,
110            align      = "center",
111            solidHeader = TRUE,
112            h4(textOutput("title_main_map")),
113            leafletOutput("main_map") %>% withSpinner()
114        ),
115
116        box(
117            width      = 12,
118            align      = "center",
119            solidHeader = TRUE,
120            column(
121                width = 4,
122                h4(textOutput("title_main_table")),
123                tableOutput("main_table1"),
124                tableOutput("main_table2")
125            ),
126            column(
127                width = 8,
128                h4(textOutput("title_main_plot")),
129                plotOutput("main_plot")
130            )
131        ),
132
133        tabPanel(
134            title = "District map",
135            textOutput("district_text"),
136            br(),
137            box(
138                width = 12,

```

```

135     solidHeader=TRUE,
136
137     column(
138         width = 6,
139         selectInput(inputId = "district",
140                     label = "Choose district",
141                     choices = berlin_names %>%
142                         filter(view == "Districts"
143                               ) %>%
144                         select(id) %>% unique()
145                         %>% t() %>% c(),
146                         selected = "Mitte")),
147
148     column(
149         width = 6,
150         selectInput(inputId = "variable2",
151                     label = "Choose variable",
152                     choices = c("none", listings_area_
153                               summary %>%
154                               colnames() %>%
155                               setdiff(c("id", "view", "group"))
156                               ) %>%
157                               gsub("_", " ", .)
158                               ))))
159
160     ),
161     box(
162         width = 12,
163         solidHeader=TRUE,
164         column(width = 4,
165                 checkboxInput(inputId = "district_
166                               properties",
167                               label = "Display
168                               properties",
169                               value = FALSE)),
170
171         column(
172             width = 4,
173             checkboxInput(inputId = "district_stations"
174                           ",
175                           label = "Display train/
176                           subway stations",
177                           value = FALSE)),
178
179         column(

```

```

166             width = 4,
167             checkboxInput(inputId = "district_sights",
168                           label = "Display top 10
169                           attractions",
170                           value = FALSE))
171
172         ),
173         box(width = 12,
174              align = "center",
175              solidHeader = TRUE,
176              h4(textOutput("title_secondary_map")),
177              leafletOutput("secondary_map") %>%
178              withSpinner()
179         ),
180         box(width = 12,
181              align = "center",
182              solidHeader = TRUE,
183              column(
184                  width = 4,
185                  h4(textOutput("title_secondary_table")),
186                  tableOutput("secondary_table1"),
187                  tableOutput("secondary_table2")
188              ),
189              column(
190                  width = 8,
191                  h4(textOutput("title_secondary_plot")),
192                  plotOutput("secondary_plot")
193              )
194          )
195      )
196  )
197 )
198
199 ### TAB3 (Price)
200 tabItem(tabName = "price",
201         fluidRow(
202             box(
203                 status = "danger",
204                 width = 12,

```

```

205     column(
206       width = 6,
207       textOutput(outputId = "price_text")
208     ),
209     column(
210       width = 6,
211       pickerInput(inputId = "variable3", label = "Choose variable
(s)",
212                     choices = listings_price_correlation %>%
213                     select(variable) %>%
214                     unique() %>%
215                     t() %>%
216                     c() %>%
217                     gsub("_", " ", .),
218                     multiple = TRUE,
219                     selected = listings_price_correlation %>%
220                     select(variable) %>%
221                     unique() %>%
222                     t() %>%
223                     c() %>%
224                     gsub("_", " ", .),
225                     options = list('actions-box' = TRUE))
226       )
227
228     )
229   ),
230   fluidRow(
231     tabBox(id = "price_box",
232       width = 12,
233       tabPanel(
234         title = "Correlation with price",
235         plotOutput(outputId = "corr_plot")
236       ),
237       tabPanel(
238         title = "Linear regression on price",
239         tabBox(id = "lm_model",
240           width = 12,
241           tabPanel(
242             title = "Model",
243             box(

```

```

245         width = 12,
246         align = "center",
247         solidHeader = TRUE,
248         column(
249             width = 3,
250             h4("R-Squared of the model"),
251             textOutput(outputId = "price_lm_r2")
252         ),
253         column(
254             width = 9,
255             h4("Coefficients of the linear regression
256                 on price."),
257             tableOutput(outputId = "price_lm_
258                 coefficients")
259         )
260     ),
261     tabPanel(
262         title = "Residuals",
263         box(width = 12,
264             solidHeader = TRUE,
265             align = "center",
266             h4("Summary statistics of the residuals"),
267             tableOutput(outputId = "price_lm_residuals
268                 ")
269         ),
270         box(width = 12,
271             solidHeader = TRUE,
272             align = "center",
273             h4("Residuals plots"),
274             column(
275                 width = 6,
276                 plotOutput(outputId = "price_lm_resplot"
277                     ) %>% withSpinner()
278             ),
279             column(
280                 width = 6,
281                 plotOutput(outputId = "price_lm_qqplot")
282                     %>% withSpinner()
283             )
284     )

```

```

281
282
283
284
285
286
287 ### TAB4 (Clustering)
288 tabItem(tabName = "cluster",
289     fluidRow(
290         box(status = "danger",
291             width = 12,
292             uiOutput("cluster_ui")),
293         tabBox(id = "cluster_test",
294             width = 12,
295             tabPanel(title = "Find number of clusters",
296                 textOutput("ctest_text"),
297                 br(),
298                 sliderInput(inputId = "ctest",
299                     label = "Choose which number of
300                         clusters to test",
301                         min = 2, max = 96, value = c(2, 40),
302                         step = 1,
303                         width = "100%"),
304                 actionButton("ctest_start", "Calculate"),
305                 progressBar(id = "ctest_progress",
306                     status = "danger",
307                     value = 0, total = 100,
308                     display_pct = TRUE),
309                 plotOutput(outputId = "tve")
310             ),
311
312             tabPanel(title = "Cluster Airbnb properties",
313                 uiOutput("cnum_ui"),
314                 br(),
315                 numericInput(inputId = "cnum",
316                     label = "Choose number of clusters for
317                         clustering",
318                     min = 2, max = 96, value = 12, step =
319                         1,
320                     width = "30%"),
321                 actionButton("cnum_start", "Calculate"),
322             )
323     )
324 )
325 )
326 )
327

```

```

318     box(
319         width      = 12,
320         align      = "center",
321         solidHeader = TRUE,
322         column(
323             width = 6,
324             h4("Clusters on coordinates"),
325             plotOutput(outputId = "clustercoord") %>%
326                 withSpinner()
327         ),
328         column(
329             width = 6,
330             h4("Clusters on Principal Components"),
331             plotOutput(outputId = "clustercpc") %>%
332                 withSpinner()
333         )
334     )
335 )
336
337
338 )
339 )
340 )
341 )

```

**Listing 43:** |server.R|

```

1 server = function(input, output, session) {
2
3     # Tab 1
4     introduction_server(input, output, session)
5
6     # Tab 2
7
8     # Panel 1
9     maps_berlin_server(input, output, session)
10
11    # Panel 2
12    maps_district_server(input, output, session)
13

```

```

14 # Tab 3
15
16 # Panel 1
17 price_corr_server(input, output, session)
18
19 # Panel 2
20 price_lm_server(input, output, session)
21
22 # Tab 4
23 clustering_server(input, output, session)
24
25 }

```

**Listing 44: |introduction\_server.R|**

```

1 ##########
2 # Introduction Server - file #####
3 # Author: Silvia Ventoruzzo #####
4 ##########
5
6 introduction_server = function(input, output, session) {
7
8   output$intro_ui = renderUI({
9
10   text1 = "This ShinyApp was developed by "
11   url1 = a("Silvia Ventoruzzo",
12           href="https://www.linkedin.com/in/silvia-ventoruzzo-1
13           a042110a/")
14   text2 = " for the university course 'Statistical Programming
15   Languages' at the "
16   url2 = a("Humboldt-University of Berlin",
17           href="https://www.hu-berlin.de/")
18   text3 = " in the Winter Term 2018/19. It shows some Exploratory Data
19   Analysis of the Airbnb properties
20   in Berlin updated at the 14. January 2019."
21
22   part1 = paste(text1, url1, text2, url2, text3, sep = "") %>%
23     paste("<br/>", sep = "<br/>")
24   html1 = HTML(part1)
25
26   text4 = "The data for this project has been downloaded from the
27   following sources:"

```

```

24 url4  = a("Inside Airbnb",
25           href="http://insideairbnb.com/")
26 url5  = a("Statistics Office of Berlin-Brandenburg",
27           href="https://www.statistik-berlin-brandenburg.de/produkte/
28             opendata/geometrienOD.asp?Kat=6301")
29 url6  = a("Geofabrik",
30           href="http://download.geofabrik.de/europe/germany/berlin.
31             html")
32
33 part2 = paste(text4, "<ul><li>", sep = "<br/>") %>%
34   paste(url4, sep = "") %>%
35   paste(url5, sep = "</li><li>") %>%
36   paste(url6, sep = "</li><li>") %>%
37   paste("",     sep = "</li></ul><br/>")
38
39 html2 = HTML(part2)
40
41
42
43 text5 = "The complete code for this project can be found on"
44 url7  = a("GitHub",
45           href="https://github.com/silvia-ventoruzzo/SPL-WISE-2018")
46 text6 = "."
47
48
49
50 })
51
52
53 output$intro_maps_text = renderText({
54   "Map of Berlin or one of its district and variable distribution."
55 })
56
57 output$intro_price_text = renderText({
58   "Variable correlation with price and linear regression on price."
59 })
60
61 output$intro_cluster_text = renderText({
62   "Clustering of Airbnb properties."
63 })

```

```

63
64     # Connections to other tabs
65     observeEvent(input$switch_maps, {
66         updateTabItems(session, inputId = "tabs", selected = "maps")
67     })
68
69     observeEvent(input$switch_price, {
70         updateTabItems(session, inputId = "tabs", selected = "price")
71     })
72
73     observeEvent(input$switch_cluster, {
74         updateTabItems(session, inputId = "tabs", selected = "cluster")
75     })
76
77 }

```

**Listing 45:** |maps\_berlin\_server.R|

```

1 ##########
2 # Maps Berlin Server - file #####
3 # Author: Silvia Ventoruzzo #####
4 #####
5
6 maps_berlin_server = function(input, output, session) {
7
8     output$berlin_text <- renderText({
9         "In this tab you can view a map of entire Berlin divided in its
10            different areas. You can choose which view
11            to show and also which variable to show the average in the different
12            areas. Moreover you can also display
13            the properties themselves, the railway stations and the top 10 tourist
14            attractions."
15
16     })
17
18     ### LEAFLET MAP
19
20     ## TITLE
21     title_main_map_change = reactive({
22
23         if (variable1() == "none") {
24
25             paste("Map of Berlin's", input$view, sep = " ")
26
27         }
28
29     })
30
31     # TABS
32     tabItems(
33         tabItem(tabName = "maps", title_main_map_change(),
34                 # Map
35                 leafletOutput("map", width = "100%", height = 500),
36
37                 # Buttons
38                 sidebarMenu(
39                     menuItem("View", tabName = "view", icon = icon("list-ul")),
40                     menuItem("Variable", tabName = "variable", icon = icon("list-ul"))
41                 )
42             )
43         ,
44         tabItem(tabName = "price", title_main_map_change(),
45                 # Map
46                 leafletOutput("map", width = "100%", height = 500),
47
48                 # Buttons
49                 sidebarMenu(
50                     menuItem("View", tabName = "view", icon = icon("list-ul")),
51                     menuItem("Variable", tabName = "variable", icon = icon("list-ul"))
52                 )
53             )
54         ,
55         tabItem(tabName = "cluster", title_main_map_change(),
56                 # Map
57                 leafletOutput("map", width = "100%", height = 500),
58
59                 # Buttons
60                 sidebarMenu(
61                     menuItem("View", tabName = "view", icon = icon("list-ul")),
62                     menuItem("Variable", tabName = "variable", icon = icon("list-ul"))
63                 )
64             )
65     )
66 }

```

```

22
23     } else if (grepl(paste(factor_vars(), collapse = " | "), variable1())) {
24
25     variable = stringr::str_match(variable1(), paste(factor_vars(),
26                                         collapse = " | ")) %>%
27                                         c() %>% gsub("_", " ", .) %>% capwords()
28
29     variable %>%
30         paste("Count of the variable", ., "in Berlin's", input$view, sep =
31             " ")
32
33 } else {
34
35     variable = capwords(input$variable1)
36
37     variable %>%
38         paste("Average of the variable", ., "in Berlin's", input$view, sep
39             = " ")
40
41 }
42
43 observeEvent({c(input$view, input$variable1)},
44               {output$title_main_map = renderText({title_main_map_change
45                                         ()})},
46               ignoreNULL = FALSE)
47
48 ## REACTIVE ELEMENTS
49
50
51 # SF Dataframe to plot polygons
52 polygons = reactive({berlin_sf %>% filter(view == input$view)})
53
54 # Dataframe with coordinates of the areas' names
55 area_names = reactive({berlin_names %>% filter(view == input$view)})
56
57 # Number of colors depend on how many areas to display
58 colors = reactive({
59     colorspace::rainbow_hcl(n = berlin_names %>%
60                             filter(view == input$view) %>%
61                             dplyr::summarize(count = n()) %>%

```

```

59             as.numeric())
60         })
61
62     # Rename variable1
63     variable1 = reactive({
64
65     gsub(" ", "_", input$variable1)
66
67   })
68
69   var_table_plot = reactive({
70
71     if (variable1() == "none") { return() }
72
73     if (grepl(paste(factor_vars(), collapse = "|"), variable1())) {
74
75       variable = stringr::str_match(variable1(), paste(factor_vars(),
76                                         collapse = "|")) %>% c()
77
78     } else {
79
80       variable = variable1()
81
82     }
83
84   })
85
86   # Number of areas in selected view
87   area_num = reactive({
88
89     berlin_names %>%
90       dplyr::filter(view == input$view) %>%
91       dplyr::summarize(count = n()) %>%
92       as.numeric()
93
94   })
95
96   # Colors for average maps
97   colors_var = reactive({
98     leaflet::colorNumeric(

```

```

99      domain  = listings_area_summary %>%
100          dplyr::filter(view == input$view) %>%
101          dplyr::select(variable1()),
102      n       = area_num(),
103      reverse = TRUE)})
```

104

105 *# Text size depend on how many areas to display*

```
106 text_size = reactive({ifelse(input$view == "Districts", "11px", "20px")
107 })
```

108 *## MAP*

109

```
110 output$main_map = renderLeaflet({
111
112     main_icons = awesomeIconList(
113         "main_properties" = makeAwesomeIcon(icon = "home", markerColor = "red"),
114         "main_stations"   = makeAwesomeIcon(icon = "subway", library = "fa",
115             markerColor = "green", iconColor = "#FFFFFF"),
116         "main_sights"     = makeAwesomeIcon(icon = "eye-open", markerColor =
117             "blue"))
118
119
120     main_map = leaflet() %>%
121         addTiles() %>%
122         addLabelOnlyMarkers(data = area_names(),
123             lng = ~long, lat = ~lat, label = ~lapply(name,
124                 htmltools::HTML),
125                 labelOptions = labelOptions(noHide = TRUE,
126                     direction = 'center',
127                         textOnly = TRUE,
128                         textsize = text_
129                             size())) %>%
130
131         addAwesomeMarkers(data = listings,
132             lng = ~long, lat = ~lat,
133             icon = main_icons[["main_properties"]],
134             group = "main_properties",
135             popup = ~lapply(paste("<b><a href=", listing_url,
136                 "> ID: ", id, "</a></b>", sep = "") %>%
137                     paste("Price: ", sep = "<br/>"))
138             %>%
```

```

130          paste(price, "$", sep = "") %>%
131          paste(property_type, sep = "<br/>") %>%
132          paste(room_type, sep = " - ") %>%
133          paste("Accomodates", sep = "<br/>") %>%
134          paste(accommodates, ifelse(
135              accommodates == 1, "person",
136              people), sep = " "),
137          htmltools::HTML),
138      clusterId = "1",
139      clusterOptions = markerClusterOptions(
140          iconCreateFunction=JS("function (cluster) {
141              var childCount = cluster.
142                  getChildCount();
143              if (childCount > 1) {
144                  c = 'rgba(235, 0, 0, 1)';
145              }
146              return new L.DivIcon({
147                  html: '<div style="'
148                      background-color:' + c + '"';
149                  '><span>' + childCount +
150                      '</span></div>',
151                  className: 'marker-
152                  cluster', iconSize: new
153                  L.Point(40, 40) });
154          })",
155          singleMarkerMode = FALSE)) %>%
156          addAwesomeMarkers(data = railway_stations_df,
157              lng = ~long, lat = ~lat,
158              icon = main_icons[["main_stations"]],
159              popup = ~htmltools::htmlEscape(
160                  paste(ifelse(s_bahn == TRUE, "S", ""),
161                      ifelse(s_bahn == TRUE & u_bahn == TRUE, "+",
162                          "")),
163                      ifelse(u_bahn == TRUE, "U", ""),
164                      sep = " ") %>%
165                  paste(id, sep = " ")),
166              group = "main_stations",
167              clusterId = "2",
168              clusterOptions = markerClusterOptions(
169                  iconCreateFunction = JS("function (cluster) {

```

```

157                         var childCount = cluster
158                             .getChildCount();
159                             if (childCount > 1) {
160                               c = 'rgba(7, 187, 31, 1)
161                               ;'}
162                               return new L.DivIcon({
163                                 html: '<div style='
164                                   'background-color:' + c +
165                                   '\"><span>' +
166                                   childCount + '</span><
167                                   /div>', className: '
168                                   marker-cluster',
169                                   iconSize: new L.Point
170                                   (40, 40 )});}),
171
172   singleMarkerMode = FALSE)) %>%
173
174 addAwesomeMarkers(data = attractions_df,
175                     lng = ~long, lat = ~lat,
176                     icon = main_icons[["main_sights"]],
177                     popup = ~htmltools::htmlEscape(id),
178                     group = "main_sights",
179                     clusterId = "3",
180                     clusterOptions = markerClusterOptions(
181                       iconCreateFunction = JS("function (cluster) {
182                         var childCount = cluster
183                             .getChildCount();
184                             if (childCount > 1) {
185                               c = 'rgba(11, 200, 213,
186                               1);'}
187                               return new L.DivIcon({
188                                 html: '<div style='
189                                   'background-color:' + c +
190                                   '\"><span>' +
191                                   childCount + '</span><
192                                   /div>', className: '
193                                   marker-cluster',
194                                   iconSize: new L.Point
195                                   (40, 40 )});}),
196
197   singleMarkerMode = FALSE)) %>%
198
199 hideGroup(group = c("main_properties", "main_stations", "main_sights
200   ")) %>%
201
202 setView(lng = berlin_center$long, lat = berlin_center$lat, zoom =

```

```

10)

177
178 if (variable1() == "none") {
179
180   main_map %>%
181     addPolygons(data = polygons(), weight = 1, smoothFactor = 1,
182                 fillOpacity = 0.8,
183                 color = "grey", fillColor = colors())
184
185 } else {
186
187   colors_vars = colors_var()
188
189   main_map %>%
190     addPolygons(data = polygons(), weight = 1, smoothFactor = 1,
191                 fillOpacity = 0.8,
192                 color = "grey",
193                 fillColor = ~colors_vars(listings_area_summary[
194                   listings_area_summary$view == input$view, variable1
195                   ()]])) %>%
196     addLegend(position = "bottomleft", pal = colors_var(),
197               values = listings_area_summary[listings_area_summary$  

198                 view == input$view, variable1()],
199               labFormat = labelFormat(transform = function(x) sort(x,
200                             decreasing = TRUE)),
201               title = variable1())
202
203   }
204
205   })
206
207
208   # checkboxInputs
209   observeEvent({c(input$main_properties, input$main_stations, input$main_sights)}, {
210
211     proxy = leafletProxy("main_map")
212
213     if (input$main_properties) {proxy %>% showGroup("main_properties")}
214     } else {proxy %>% hideGroup("main_properties")}
215
216
217     if (input$main_stations) {proxy %>% showGroup("main_stations")}
218     } else {proxy %>% hideGroup("main_stations")}

```

```

210
211     if (input$main_sights) {proxy %>% showGroup("main_sights")
212   } else {proxy %>% hideGroup("main_sights")}
213
214 })
215
216
217 ### DESCRIPTIVE TABLE
218
219 ## TITLE
220 title_main_table_change = reactive({
221
222
223   if (variable1() == "none") { return() }
224
225   if (grepl(paste(factor_vars(), collapse = "|"), variable1())) {
226
227     variable = stringr::str_match(variable1(), paste(factor_vars(),
228                                     collapse = "|")) %>%
229     c() %>% gsub("_", " ", .) %>% capwords()
230
231   } else {
232
233     variable = capwords(input$variable1)
234
235   }
236
237   variable %>%
238     paste("Descriptive statistics of the variable", ., sep = " ")
239
240 })
241
242 observeEvent({c(input$view, input$variable1)},
243               {output$title_main_table = renderText({title_main_table_
244                                         change()})},
245               ignoreNULL = FALSE)
246
247 ## REACTIVE ELEMENTS
248
249 # Categorical variables
factor_vars = reactive({

```

```

249
250   listings %>%
251     dplyr::select(-id, -listing_url, -long, -lat, -vbb_zone, -district,
252       -neighbourhood) %>%
253       dplyr::select_if(function(x) is.character(x) | is.factor(x) | is.
254         logical(x)) %>%
255       names()
256
257
258 var_table_plot = reactive({
259
260   if (variable1() == "none") { return() }
261
262   if (grepl(paste(factor_vars(), collapse = "|"), variable1())) {
263
264     variable = stringr::str_match(variable1(), paste(factor_vars(),
265       collapse = "|")) %>% c()
266
267   } else {
268
269     variable = variable1()
270
271   }
272
273 }
274
275 ## TABLES
276
277
278 output$main_table1 = renderTable({
279
280   if (variable1() == "none") {
281
282     return()
283
284   } else if (grepl(paste(factor_vars(), collapse = "|"), variable1())) {
285
286     main_df = statistics_cat %>%
287       dplyr::filter(variable == var_table_plot()) %>%
288       dplyr::select(-variable)

```

```

287
288     } else {
289
290         main_df = statistics_num %>%
291             dplyr::filter(variable == var_table_plot()) %>%
292             dplyr::select(min, '1Q', median, '3Q', max)
293     }
294
295     return(main_df)
296
297 })
298
299 output$main_table2 = renderTable({
300
301     if (variable1() == "none") {
302
303         return()
304
305     } else if (grepl(paste(factor_vars(), collapse = " | "), variable1())) {
306
307         return()
308
309     } else {
310
311         main_df = statistics_num %>%
312             dplyr::filter(variable == var_table_plot()) %>%
313             dplyr::select(mean, sd)
314     }
315
316     return(main_df)
317
318 })
319
320
321 ### DISTRIBUTION PLOT
322
323 ## TITLE
324 title_main_plot_change = reactive({
325
326     if (variable1() == "none") { return() }
327

```

```

328 if (grepl(paste(factor_vars(), collapse = " | "), variable1())) {
329
330   variable = stringr::str_match(variable1(), paste(factor_vars(),
331                                     collapse = " | ")) %>%
332     c() %>% gsub("_", " ", .) %>% capwords()
333
334 } else {
335
336   variable = capwords(variable1())
337
338 }
339
340   variable %>%
341     paste("Distribution of the variable", ., sep = " ")
342 )
343
344 observeEvent({c(input$view, input$variable1),
345               {output$title_main_plot = renderText({title_main_plot_
346                                         change()})},
347               ignoreNULL = FALSE)
348
349
350   ## REACTIVE ELEMENTS
351
352   ## PLOT
353
354   output$main_plot = renderPlot({
355
356     if (variable1() == "none") {
357
358       return()
359
360     } else if (grepl(paste(factor_vars(), collapse = " | "), variable1())) {
361
362       variable = stringr::str_match(variable1(), paste(factor_vars(),
363                                     collapse = " | ")) %>%
364         c() %>% gsub("_", " ", .) %>% capwords()
365
366       plot = distribution_plot(listings, var_table_plot(), tilt_text_x =
367                                 TRUE)
368
369     }
370   })

```

```

365     } else {
366
367         plot = distribution_plot(listings, var_table_plot(), hide_outliers =
368             TRUE)
369     }
370
371     return(plot)
372 }
373
374 }
375 }
```

**Listing 46: |maps\_district\_server.R|**

```

1 ##########
2 # Maps District Server-file #####
3 # Author: Silvia Ventoruzzo #####
4 ##########
5
6 maps_district_server = function(input, output, session) {
7
8     output$district_text <- renderText({
9
10         "This tab is similar to the previous one. You can still choose which
11             variable to show the average
12             in the different areas and display various properties, railway stations
13             and tourist attractions.
14
15         But here you see the city at the district's level, showing therefore the
16             neighbourhoods in the
17             different districts."
18
19     })
20
21     #### LEAFLET MAP
22
23     ## TITLE
24
25
26     title_secondary_map_change = reactive({
27
28         if (variable2() == "none") {
29
30             paste("Map of", input$district, sep = " ")
31             %>%
32             paste("'s", sep = "") %>%
33
34         }
35
36     })
37
38 }
```

```

25     paste("Neighbourhoods", sep = " ")
26
27 } else if (grepl(paste(factor_vars(), collapse = "|"), variable2())) {
28
29     variable = stringr::str_match(variable2(), paste(factor_vars(),
30                                     collapse = "|")) %>%
31                                     c() %>% gsub("_", " ", .) %>% capwords()
32
33     variable %>%
34         paste("Count of the variable", ., "in the district", input$district,
35               sep = " ")
36
37 } else {
38
39     variable = capwords(input$variable2)
40
41     variable %>%
42         paste("Average of the variable", ., "in the district", input$district,
43               sep = " ")
44
45 }
46
47 observeEvent({c(input$district, input$variable2)},
48               {output$title_secondary_map = renderText({title_secondary_map
49                                         _change()})}, ignoreNULL = FALSE)
50
51
52 ## REACTIVE ELEMENTS
53
54 # SF Dataframe to plot polygons
55
56 neighbourhoods = reactive({berlin_sf %>%
57     filter(view == "Neighbourhoods",
58           group == input$district)})
59
60
61 # Dataframe with coordinates of the areas' names
62
63 neighbourhood_names = reactive({berlin_names %>%
64     filter(view == "Neighbourhoods",
65           group == input$district)})
66
67
68 # Number of colors depend on how many areas to display

```

```

62  neighbourhood_colors = reactive({
63    colorspace::rainbow_hcl(n = berlin_names %>%
64      filter(view == "Neighbourhoods",
65             group == input$district) %>%
66      dplyr::summarize(count = n()) %>%
67      as.numeric())
68  })
69
70  # Rename variable2
71  variable2 = reactive({
72
73    gsub(" ", "_", input$variable2)
74
75  })
76
77  # Number of areas in selected view
78  neighbourhood_area_num = reactive({
79
80    berlin_names %>%
81      dplyr::filter(view == "Neighbourhoods",
82                    group == input$district) %>%
83      dplyr::summarize(count = n()) %>%
84      as.numeric()
85
86  })
87
88  # Colors for average maps
89  neighbourhood_colors_var = reactive({
90    leaflet::colorNumeric(
91      palette = topo.colors(area_num()),
92      domain = listings_area_summary %>%
93        dplyr::filter(view == input$view) %>%
94        dplyr::select(variable1()),
95      n       = area_num(),
96      reverse = TRUE)})
97
98
99  # When choosing a variable (input$variable2), the color should represent
100 # the variable's distribution
101 neighbourhood_colors_var = reactive({
102   leaflet::colorNumeric(

```

```

102 palette = topo.colors(neighbourhood_area_num()),
103 domain = listings_area_summary %>%
104   dplyr::filter(view == "Neighbourhoods",
105                 group == input$district) %>%
106   dplyr::select(variable2()),
107 n      = neighbourhood_area_num(),
108 reverse = TRUE)
109
110 })
111
112
113 ## MAP
114
115 output$secondary_map = renderLeaflet({
116
117   district_icons = awesomeIconList(
118     "district_properties" = makeAwesomeIcon(icon = "home", markerColor = "red"),
119     "district_stations" = makeAwesomeIcon(icon = "subway", library = "fa",
120       markerColor = "green", iconColor = "#FFFFFF"),
121     "district_sights" = makeAwesomeIcon(icon = "eye-open", markerColor =
122       "blue")
123   )
124
125   secondary_map = leaflet() %>%
126     addTiles() %>%
127     addLabelOnlyMarkers(data = neighbourhood_names(),
128       lng = ~long, lat = ~lat, label = ~lapply(name,
129         htmltools::HTML),
130       labelOptions = labelOptions(noHide = TRUE,
131         direction = 'center',
132         textOnly = TRUE,
133         textSize = "11px"))
134
135   secondary_map %>%
136     addAwesomeMarkers(data = listings %>% filter(district == input$district),
137       lng = ~long, lat = ~lat,
138       icon = district_icons[["district_properties"]],
139       group = "district_properties",
140       popup = ~lapply(paste("<b><a href=", listing_url, ">ID: ", id, "</a></b>", sep = "") %>%

```

```

134                               paste("Price: ", sep = "<br/>") %>%
135                               paste(price, "$", sep = "") %>%
136                               paste(property_type, sep = "<br/>") %>%
137                               paste(room_type, sep = " - ") %>%
138                               paste("Accomodates", sep = "<br/>") %>%
139                               paste(accommodates, ifelse(
140                                 accommodates == 1, "person",
141                                 people), sep = " "),
142                               htmltools::HTML),
143                               clusterId = "1",
144                               clusterOptions = markerClusterOptions(
145                                 iconCreateFunction=JS("function (cluster) {
146                                   var childCount = cluster.
147                                     getCount();
148                                   if (childCount > 1) {
149                                     c = 'rgba(235, 0, 0, 1);'
150                                     return new L.DivIcon({ html:
151                                       '<div style='background-
152                                         color:' +c+ '><span>' +
153                                         childCount + '</span></div>',
154                                         className: 'marker-
155                                         cluster', iconSize: new L.
156                                         Point(40, 40) });
157                                   }
158                                 singleMarkerMode = FALSE)) %>%
159                               addAwesomeMarkers(data = railway_stations_df %>% filter(district ==
160                                 input$district),
161                                 lng = ~long, lat = ~lat,
162                                 icon = district_icons[["district_stations"]],
163                                 popup = ~htmltools::htmlEscape(
164                                   paste(ifelse(s_bahn == TRUE, "S", ""),
165                                     ifelse(s_bahn == TRUE & u_bahn == TRUE, "+",
166                                       "")),
167                                     ifelse(u_bahn == TRUE, "U", ""),
168                                     sep = "")) %>%
169                               paste(id, sep = " ")),
170                               group = "district_stations",
171                               clusterId = "2",
172                               clusterOptions = markerClusterOptions(

```

```

161         iconCreateFunction = JS("function (cluster) {
162             var childCount = cluster.
163                 getChildCount();
164             if (childCount > 1) {
165                 c = 'rgba(7, 187, 31, 1)
166                 ;'}
167             return new L.DivIcon({
168                 html: '<div style="'
169                     background-color:' + c + '"'
170                     ><span>' + childCount +
171                     '</span></div>',
172                     className: 'marker-
173                         cluster', iconSize: new
174                         L.Point(40, 40 )});}"),
175             singleMarkerMode = FALSE)) %>%
176             addAwesomeMarkers(data = attractions_df %>% filter(district == input$district),
177                               lng = ~long, lat = ~lat,
178                               icon = district_icons[["district_sights"]],
179                               popup = ~htmltools::htmlEscape(id),
180                               group = "district_sights",
181                               clusterId = "3",
182                               clusterOptions = markerClusterOptions(
183                                   iconCreateFunction = JS("function (cluster) {
184                                       var childCount = cluster.
185                                           getChildCount();
186                                           if (childCount > 1) {
187                                               c = 'rgba(11, 200, 213, 1)
188                                               ;'}
189                                           return new L.DivIcon({
190                                               html: '<div style="'
191                                                 background-color:' + c + '"'
192                                                 ><span>' + childCount +
193                                                 '</span></div>',
194                                                 className: 'marker-
195                                         cluster', iconSize: new
196                                         L.Point(40, 40 )});}"),
197                                         singleMarkerMode = FALSE)) %>%
198                                         hideGroup(group = c("district_properties", "district_stations", "district_sights")) %>%
199                                         addMiniMap(centerFixed      = berlin_center %>% t() %>% c(),

```

```

182         zoomLevelFixed = 8,
183         zoomAnimation = FALSE)
184
185 # if no variable is selected
186 if (variable2() == "none") {
187
188     secondary_map %>%
189     addPolygons(data = neighbourhoods(), weight = 1, smoothFactor = 1,
190                 fillOpacity = 0.8,
191                 color = "grey", fillColor = neighbourhood_colors())
192
193 # if variable is selected
194 } else {
195
196     neighbourhood_colors_vars = neighbourhood_colors_var()
197
198     secondary_map %>%
199     addPolygons(data = neighbourhoods(), weight = 1, smoothFactor = 1,
200                 fillOpacity = 0.8,
201                 color = "grey",
202                 fillColor = ~neighbourhood_colors_vars(
203                     listings_area_summary
204                     [listings_area_summary$view == "
205                         Neighbourhoods" &
206                         listings_area_summary$group ==
207                         input$district,
208                         variable2())) %>%
209     addLegend(position = "bottomleft", pal = neighbourhood_colors_var(),
210               values = listings_area_summary[listings_area_summary$view
211 == "Neighbourhoods" &
212                     listings_area_summary$group == input$district
213
214                     ,
215                     variable2()],
216               labFormat = labelFormat(transform = function(x) sort(x,
217                             decreasing = TRUE)),
218               title = variable2())
219
220     }
221 }
222 })

```

```

215 # checkboxInputs
216 observeEvent({c(input$district_properties, input$district_stations, input$district_sights, input$district)}, {
217
218   proxy = leafletProxy("secondary_map")
219
220   if (input$district_properties) {proxy %>% showGroup("district_properties")}
221   } else {proxy %>% hideGroup("district_properties")}
222
223   if (input$district_stations) {proxy %>% showGroup("district_stations")}
224   } else {proxy %>% hideGroup("district_stations")}
225
226   if (input$district_sights) {proxy %>% showGroup("district_sights")}
227   } else {proxy %>% hideGroup("district_sights")}
228
229 })
230
231 #### SUMMARY TABLE
232
233 ## TITLE
234
235 title_secondary_table_change = reactive({
236
237   if (variable2() == "none") { return() }
238
239   if (grepl(paste(factor_vars(), collapse = "|"), variable2())) {
240
241     variable = stringr::str_match(variable2(), paste(factor_vars(),
242       collapse = "|")) %>%
243       c() %>% gsub("_", " ", .) %>% capwords()
244
245   } else {
246
247     variable = capwords(input$variable2)
248
249   }
250
251   variable %>%
252     paste("Summary statistics of the variable", ., "in the district",
253       input$district, sep = " ")

```

```

252 })
253
254 observeEvent({c(input$district, input$variable2)},
255   {output$title_secondary_table = renderText({title_secondary_
256     table_change()})}, ignoreNULL = FALSE)
257
258 ## REACTIVE ELEMENTS
259
260 # Categorical variables
261 factor_vars = reactive({
262
263   listings %>%
264     dplyr::select(-id, -listing_url, -long, -lat, -vbb_zone, -district, -
265       neighbourhood) %>%
266     dplyr::select_if(function(x) is.character(x) | is.factor(x) | is.
267       logical(x)) %>%
268     names()
269
270 })
271
272 var_table_plot2 = reactive({
273
274   if (variable2() == "none") { return() }
275
276   if (grepl(paste(factor_vars(), collapse = "|"), variable2())) {
277
278     variable = stringr::str_match(variable2(), paste(factor_vars(),
279       collapse = "|")) %>% c()
280
281   } else {
282
283     variable = variable2()
284   }
285
286 })
287
288 table_df = reactive({

```

```

289     if (variable2() == "none") { return() }
290
291     df = listings %>%
292       dplyr::filter(district == input$district) %>%
293       dplyr::select(var_table_plot2())
294
295     return(df)
296
297   })
298
299   ## TABLE
300
301   # Descriptive table
302   output$secondary_table1 = renderTable({
303
304     if (variable2() == "none") {
305
306       return()
307
308     } else if (grepl(paste(factor_vars(), collapse = "|"), variable2())) {
309
310       secondary_df = descriptive_statistics(table_df())
311
312     } else {
313
314       secondary_df = descriptive_statistics(table_df()) %>%
315         dplyr::select(min, '1Q', median, '3Q', max)
316
317     }
318
319     return(secondary_df)
320   })
321
322   output$secondary_table2 = renderTable({
323
324     if (variable2() == "none") {
325
326       return()
327
328     } else if (grepl(paste(factor_vars(), collapse = "|"), variable2())) {

```

```

330
331     return()
332
333 } else {
334
335     secondary_df = descriptive_statistics(table_df()) %>%
336         dplyr::select(mean, sd)
337
338 }
339
340 return(secondary_df)
341
342 })
343
344 ### DISTRIBUTION PLOT
345
346 ## TITLE
347 title_secondary_plot_change = reactive({
348
349     # factor_vars() <- c("room_type", "property_type", "room_type")
350
351     if (variable2() == "none") { return() }
352
353     if (grepl(paste(factor_vars(), collapse = "|"), variable2())) {
354
355         variable = stringr::str_match(variable2(), paste(factor_vars(),
356                                         collapse = "|")) %>%
357             c() %>% gsub("_", " ", .) %>% capwords()
358
359     } else {
360
361         variable = capwords(input$variable2)
362
363     }
364
365     variable %>%
366         paste("Distribution of the variable", ., "in the district", input$district,
367               sep = " ")
368 }
369
370 observeEvent({c(input$view, input$variable2)},

```

```

369     {output$title_secondary_plot = renderText({title_secondary_
370         plot_change()}},
371         ignoreNULL = FALSE)
372
373 ## REACTIVE ELEMENTS
374
375 plot_df2 = reactive({
376
377     if (variable2() == "none") {
378
379         return()
380
381     } else if (grepl(paste(factor_vars(), collapse = "|"), variable2())) {
382
383         df = listings %>%
384             filter(district == input$district)
385
386     } else {
387
388         variable = var_table_plot2() %>% sym()
389
390         df = listings %>%
391             filter(district == input$district)
392
393     }
394
395     return(df)
396 })
397
398 ## PLOT
399
400 output$secondary_plot = renderPlot({
401
402     if (variable2() == "none") {
403
404         return()
405
406     } else if (grepl(paste(factor_vars(), collapse = "|"), variable2())) {
407
408         plot2 = distribution_plot(plot_df2(), var_table_plot2(), tilt_text_x =

```

```
    TRUE)

409

410    } else {

411

412        plot2 = distribution_plot(plot_df2(), var_table_plot2(), hide_outliers
413                                = TRUE)

414    }

415

416    return(plot2)

417

418    })
419
420 }
```

**Listing 47:** |price corr server.R|

```

23 validate(
24   need(!is.null(input$variable3), "Please select variable(s)")
25 )
26
27 listings_price_correlation %>%
28   filter(variable %in% gsub(" ", "_", input$variable3))
29 }
30
31 ## PLOT
32
33 output$corr_plot = renderPlot({
34
35   correlation_plot(price_correlation())
36
37 })
38
39 }
```

**Listing 48: |price\_lm\_server.R|**

```

1 ##########
2 # Price Server-file #####
3 # Author: Silvia Ventoruzzo #####
4 ##########
5
6 price_lm_server = function(input, output, session) {
7
8   ## REACTIVE ELEMENTS
9
10  # For plot in tab3-panel1, check that at least one variable (input$variable3) is selected
11  # If not show an error message
12  price_lm_df = reactive({
13    validate(
14      need(!is.null(input$variable3), "Please select variable(s)")
15    )
16
17    listings_price %>%
18      dplyr::select(price, gsub(" ", "_", input$variable3))
19
20  })
```

```

22 price_lm = reactive({
23
24   lm(price ~ ., data = price_lm_df())
25
26 })
27
28 ## TEXT WITH R-SQUARED
29
30 output$price_lm_r2 = renderText({
31
32   summary(price_lm())$r.squared %>% round(4)
33
34 })
35
36 ## TABLE WITH COEFFICIENTS
37
38 output$price_lm_coefficients = renderTable({
39
40   coefficients = data.frame(
41     variable = c("(Coefficient)", gsub(" ", "_", input$variable3)),
42     coefficient = price_lm()$coefficients) %>%
43     tidyrr::drop_na() %>%
44     dplyr::mutate(num = 1:nrow(.)) %>%
45     dplyr::mutate_if(is.numeric, function(x) round(x, 4))
46
47   significance = data.frame(
48     significant = ifelse(summary(price_lm())$coefficients[,4] < 0.5,
49                           TRUE, FALSE)) %>%
50     dplyr::mutate(num = 1:nrow(.))
51
52   lm = coefficients %>%
53     dplyr::inner_join(significance, by = "num") %>%
54     dplyr::select(-num)
55
56   return(lm)
57
58 })
59
60 ## SUMMARY TABLE OF RESIDUALS
61
62 output$price_lm_residuals = renderTable({

```

```

63
64     lm = data.frame(
65         residuals = price_lm()$residuals
66     )
67
68     summary_stat_num(df = lm,
69                        var = "residuals")
70
71 })
72
73 ## FITTED VS ACTUAL PLOT
74
75 output$price_lm_resplot = renderPlot({
76
77     ggplot(price_lm(), aes(x = .fitted, y = .resid)) +
78         geom_point() +
79         geom_smooth(method = lm) +
80         geom_hline(yintercept=0, col="red", linetype="dashed") +
81         xlab("Fitted values") +
82         ylab("Residuals") +
83         ggtitle("Residual vs Fitted Plot") +
84         theme_bw()
85
86 })
87
88 ## QQ PLOT RESIDUALS
89
90 output$price_lm_qqplot = renderPlot({
91
92     gg_qqplot(price_lm()) +
93     theme_bw()
94
95 })
96
97 }

```

**Listing 49: |clustering\_server.R|**

```

1 ######
2 # Clustering Server - file #####
3 # Author: Silvia Ventoruzzo #####
4 #####

```

```

5
6 clustering_server = function(input, output, session) {
7
8   output$cluster_ui <- renderUI({
9
10   text1 = "In this tab you can cluster Airbnb properties in Berlin. The
11     process is the following:"
12   text2 = "Choose the number of the clusters according to the elbow-
13     method"
14   text2a = "Choose the number of clusters to test"
15   text2b = "Calculate for each number the percentage of the total variance
16     explained"
17   text2c = "Plot the total variance explained against the number of
18     clusters"
19   text2d = "Choose the number of clusters where adding another does not
20     add sufficient informatoin"
21   text3 = "Cluster the properties with the selected number of clusters"
22   text4 = "Plot the clusters"
23
24   intern = paste(text2, "<ul><li>", sep = "<br/>") %>%
25     paste(text2a, sep = "") %>%
26     paste(text2b, sep = "</li><li>") %>%
27     paste(text2c, sep = "</li><li>") %>%
28     paste(text2d, sep = "</li><li>") %>%
29     paste("", sep = "</li></ul>")
30
31   extern = paste(text1, "<ul><li>", sep = "<br/>") %>%
32     paste(intern, sep = "") %>%
33     paste(text3, sep = "</li><li>") %>%
34     paste(text4, sep = "</li><li>") %>%
35     paste("", sep = "</li></ul>")
36
37   html = HTML(extern)
38
39 })
40
41
42   ## Panel 1 (Find number of clusters)
43
44   output$ctest_text <- renderText({
45     "Before clustering the airbnb properties, we need to find out what is

```

```

    the optimal number of clusters.

41   We use the elbow method, which means that from plot showing the
42     percentage of the total variance
43     explained for that number of clusters we choose the number before the
44     increase becomes too flat."
45 }

46 tot_var_explained <- eventReactive(input$ctest_start, {

47   clusters = seq(input$ctest[1], input$ctest[2], 1)

48   cluster_count = input$ctest[2] - input$ctest[1] + 1

49   tve = data.frame(clusters = clusters,
50                     tve       = rep(NA, cluster_count))

51

52   clk = list()

53

54   n = 1

55

56   for (k in seq(input$ctest[1], input$ctest[2], 1)) {

57

58     clk[[k-1]] = kmeans(listings_scaled, centers = k, iter.max = 20)
59     tve$tve[k-1] = 1 - clk[[k-1]]$tot.withinss / clk[[k-1]]$totss

60

61     updateProgressBar(session = session, id = "ctest_progress",
62                       value = (n / cluster_count * 100), total = 100,
63                       title = paste(n, "of", cluster_count))

64     n = n + 1
65   }

66

67   return(tve)
68 }

69

70 }

71

72

73 output$tve = renderPlot({
74

75   # Dependency on input$ctest_start
76   if (input$ctest_start == 0) return()
77

78   tot_var_explained() %>%

```

```

79     ggplot(aes(x = clusters, y = tve)) +
80       geom_line(color = "grey") +
81       geom_point(color = "red") +
82       scale_x_continuous(breaks = seq(input$ctest[1], input$ctest[2], 2))
83       +
84       theme_bw()
85   })
86
87 ## Panel 2 (Actual clustering)
88
89 output$cnum_ui <- renderUI({
90   text <- "Select the optimal number of clusters you found in the
91         previous tab and run the clustering
92         function. In this case we use k-means, which splits the data into the
93         given number of clusters
94         in such a way as to minimize the within-cluster sum of squares. For
95         detailed information on
96         this and other clustering algorithms you can read here:"
97   url <- a("Data Clustering: A Review",
98           href="http://delivery.acm.org/10.1145/340000/331504/p264-jain
99           .pdf?ip=141.20.217.40&id=331504&acc=PUBLIC&key=2
100          BA2C432AB83DA15%2E04C410D892D772C0%2E4D4702B0C3E38B35%2
101          E4D4702B0C3E38B35&__acm__=1551719946_9
102          d82c9a89fd8cd69cfa3d310ea27bc7f")
103   tagList(text, url)
104 })
105
106 # REACTIVE ELEMENTS
107
108 # k-means
109 k_means <- eventReactive(input$cnum_start, {
110
111   kmeans(listings_scaled, centers = input$cnum, iter.max = 20)
112
113 })
114
115 # clusters
116 clusters <- eventReactive(input$cnum_start, {
117
118   k_means()$cluster %>% as.factor()
119
120 })
121
122 
```

```

112
113 })
114
115 # Plot of cluster on coordinates
116 output$clustercoord = renderPlot({
117
118 # Dependency on input$ctest_start
119 if (input$cnum_start == 0) { return() }
120
121 ggplot() +
122   geom_sf(data = berlin_sf %>% dplyr::filter(view == "Districts"),
123           show.legend = FALSE, color = "black") +
124   coord_sf(datum = NA) +
125   geom_point(data = listings, aes(x = long, y = lat, color = clusters
126     ), alpha = 0.5) +
127   theme(plot.title = element_text(hjust = 0.5)) +
128   labs(color = "clusters") +
129   theme_void()
130
131 })
132
133 # Plot of clusters on Principal Components
134 output$clusterpc = renderPlot({
135
136 # Dependency on input$ctest_start
137 if (input$cnum_start == 0) { return() }
138
139 ggplot2::autoplot(listings_pc, data = k_means(), colour = "cluster") +
140   labs(color = "clusters") +
141   theme_bw()
142
143 })
144 }
```

**Listing 50:** |app.R|

```

1 # Install and load needed packages
2 needed_packages <- c("shiny",
3                         "shinydashboard",
4                         "leaflet",
5                         "shinyWidgets",
```

```

6      "colorspace",
7      "htmltools",
8      "rstudioapi",
9      "devtools",
10     "lindia",
11     "tidyverse",
12     "Jmisc",
13     "stringr")
14 for (package in needed_packages) {
15   if (!require(package, character.only=TRUE)) {
16     install.packages(package, character.only=TRUE)}
17   library(package, character.only=TRUE)
18 }
19 devtools::install_github('andrewsali/shinycssloaders')
20 library("shinycssloaders")
21 rm("needed_packages", "package")
22
23 # Set working directory to the one where the file is located
24 # setwd(dirname(sys.frame(1)$ofile)) # This works when sourcing
25 setwd(dirname(rstudioapi::getActiveDocumentContext()$path)) # This works
when running the code directly
26
27 # Load helper functions and other scripts
28 load("workspace_for_shinyapp.RData")
29 Jmisc::sourceAll(file.path(getwd(), "Helpers", fsep="/"))
30 Jmisc::sourceAll(file.path(getwd(), "Server", fsep="/"))
31 source("ui.R")
32 source("server.R")
33
34 # Run app
35 shinyApp(ui      = ui,
36           server = server)

```

## **Declaration of Authorship**

I hereby confirm that I have authored this seminar paper independently and without use of others than the indicated sources. All passages which are literally or in general matter taken out of publications or other sources are marked as such.

Berlin, March 15, 2019

Silvia Ventoruzzo