# Exploratory Data Analysis of Airbnb Listings in Berlin

Silvia Ventoruzzo

Ladislaus von Bortkiewicz Chair of Statistics
Humboldt–Universität zu Berlin
http://lvb.wiwi.hu-berlin.de

# Outline

# Objectives

- ⊡ Airbnb is a website where private people, and commercial entities, can rent apartments, or part of them.
- ⊡ Aim of the project:
  - ▶ Answer the following questions:
    - • How are the Airbnb properties distributed in Berlin?
    - • Which variables drive the property price?
  - ▶ Display the information through a ShinyApp

# Creating Polygons for the districts and neighbourhoods

- ⊡ Loading the shapefile of Berlin neighbourhoods: (downloaded from www.statistik-berlin-brandenburg.de)
- ⊡ For the neighbourhoods we just need to transform the SpatialPolygonsDataFrame into an sf object.

```
berlin_neighbourhood_sf <- berlin %>%
  sf::st_as_sf() %>%
  sf::st_set_crs(proj4string(berlin)) %>%
  dplyr::rename(id    = Name,
                group = BEZNAME) %>%
  dplyr::select(id, group, geometry) %>%
  dplyr::arrange(group) %>%
  dplyr::mutate(view = "Neighbourhoods")
```

- ⊡ The neighbourhood Buckow is composed of two separate parts, which we unite by grouping the neighbourhoods by their id.

```
1  berlin_neighbourhood_singlebuckow_sf <-
     berlin_neighbourhood_sf %>%
2    dplyr::group_by(id, group, view) %>%
3    dplyr::summarize(do_union = TRUE) %>%
4    dplyr::arrange(group)
```

- ⊡ To get the districts we unite the neighbourhoods by their district (group).

```
1  berlin_district_sf <- berlin_neighbourhood_sf %>%
2    dplyr::group_by(group, view) %>%
3    dplyr::summarize(do_union = TRUE) %>%
4    dplyr::mutate(id    = group,
5                  view  = "Districts")
```
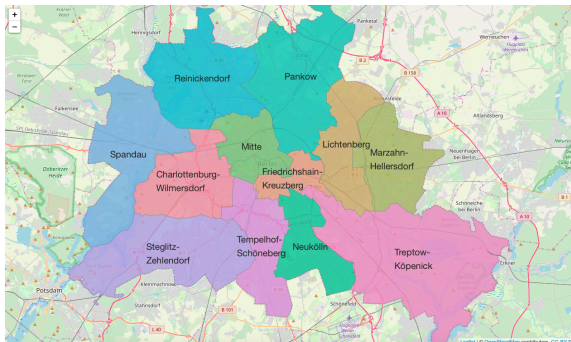
# Mapping the polygons

☐ Districts:



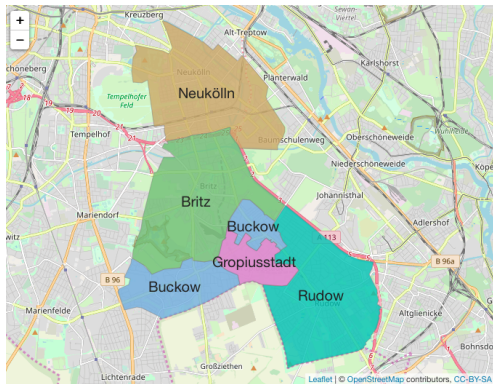Figure 1: Leaflet map of Berlin Districts

▣ Neighbourhoods:



Figure 2: Leaflet map of Neukölln's Neighbourhoods

# Creating Polygons of the VBB Areas

⊡ Loading the shapefile of the stations and stops in Berlin: (downloaded from www.geofabrik.de)

⊡ Create dataframe with the information for train/subway stations.

```
1  bahn_stations_df <- stations %>%
2     base::as.data.frame() %>%
3     dplyr::filter(fclass %like% "railway") %>%
4     dplyr::rename(id   = name,
5                   long = coords.x1,
6                   lat  = coords.x2) %>%
```

```r
dplyr::mutate(id      = gsub("Berlin ", "", id),
              id      = gsub("Berlin-", "", id),
              id      = gsub(" *\\(.*?\\) *", "",
                            id),
              s_bahn = ifelse(startsWith(id,
                              "S "), TRUE, FALSE),
              u_bahn = ifelse(startsWith(id,
                              "U "), TRUE, FALSE),
              id      = gsub("S ", "", id),
              id      = gsub("U ", "", id)) %>%
dplyr::group_by(id) %>%
dplyr::summarize(long   = mean(long),
                 lat    = mean(lat),
                 s_bahn = any(s_bahn),
                 u_bahn = any(u_bahn))
```

⊡ Create dataframe with names and order position of stations that form the Ringbahn (VBB Area A).

```
1  ringbahn_names_df <- base::data.frame(
2      id = c("Südkreuz", "Schöneberg",
3          "Innsbrucker Platz", "Bundesplatz",
4          "Heidelberger Platz", "Hohenzollerndamm",
5          "Halensee", "Westkreuz", "Messe Nord/ICC",
6          "Westend", "Jungfernheide",
7          "Beusselstraße", "Westhafen", "Wedding",
8          "Gesundbrunnen", "Schönhauser Allee",
9          "Prenzlauer Allee", "Greifswalder Straße",
10         "Landsberger Allee", "Storkower Straße",
11         "Frankfurter Allee", "Ostkreuz",
12         "Treptower Park", "Sonnenallee",
13         "Neukölln", "Hermannstraße",
14         "Tempelhof", "Südkreuz"),
15         stringsAsFactors = FALSE) %>%
```

```
1       tibble :: rownames_to_column (var = "order") %>%
2       dplyr :: mutate (order = as.numeric (order))
```

⊡ Create the sf object of the VBB Areas by binding the Area A
  with the complete-Berlin polygon and finding their interaction.

```
1  vbb_AB_sf <- bahn_stations_df %>%
2    dplyr :: right_join (ringbahn_names_df ,
3                       by = "id") %>%
4    dplyr :: arrange (order) %>%
5    dplyr :: select (long , lat) %>%
6    base :: as.matrix () %>%
7    base :: list () %>%
8    sf :: st_polygon () %>%
9    sf :: st_sfc () %>%
10   sf :: st_sf (crs = proj4string (berlin)) %>%
```

```
1   base::rbind(berlin_all_sf) %>%
2   sf::st_intersection() %>%
```

- ⊡ Area A: interaction (n.overlaps $> 1$)
- ⊡ Area B: no interaction (n.overlaps $= 1$)

```
1   dplyr::mutate(id      = ifelse(n.overlaps > 1,
2                                  "A", "B"),
3                 view    = "VBB Areas",
4                 group   = id) %>%
5   dplyr::select(-n.overlaps, -origins) %>%
6   dplyr::arrange(desc(id))
```
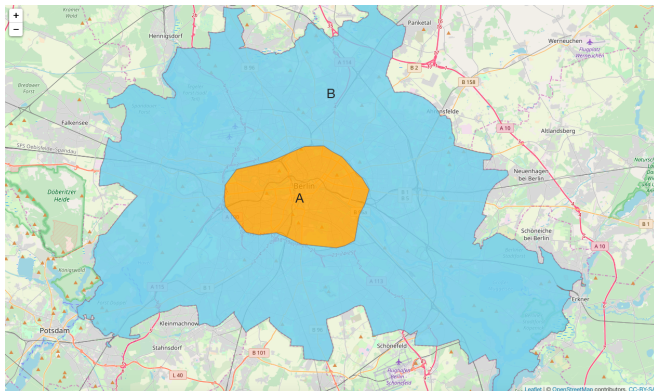
# Mapping the VBB Areas



Figure 3: Leaflet map of VBB Areas

# Initial variables

- ⊡ Loading the following files (dowloaded from www.insideairbnb.com):
  - ▶ *listings.csv.gz:* Detailed Listings data for Berlin
  - ▶ *listings.csv:* Summary information and metrics for listings in Berlin
  - ▶ *calendar.csv.gz:* Detailed Calendar Data for listings in Berlin
- ⊡ Last update: 7th November 2018
- ⊡ Keeping the following variables:
  - ▶ *id:* Reference to the listing
  - ▶ *long, lat:* Coordinates of the listing
  - ▶ *price:* Price per night in dollars
  - ▶ *security_deposit_yn, cleaning_fee_yn:* If extra money will be asked as deposit and/or for cleaning
  - ▶ *property_type, room_type:* Type of accommodation

- ▶ *accommodates, bedrooms, beds:* Important for the number of people who could stay in the property
- ▶ *minimum_nights:* How many night at least a person has to book
- ▶ *reviewed_yn, number_of_reviews, review_scores_rating:* If a listing has been reviewed, by how many peope and what its average scoring is
- ▶ *availability_30, availability_60, availability_90, availability_365:* Availability in days in the next 30, 60, 90 and 365 days
- ▶ *listing_url:* Website of the specific listing

## New variables

⊡ district, neighbourhood, vbb_area:

```
1 listings <- point_in_polygons(
2                points_df = listings,
3                polys_sf  = berlin_neighbourhood_sf,
4                var_name  = "neighbourhood")
```

⊡ station_count, station_dist, attraction_count, attraction_dist:

```
1 listings <- distance_count(
2                main      = listings,
3                reference = bahn_stations_df,
4                var_name  = "station",
5                distance  = 1000)
```

# Functions

⊡ point_in_polygons:

```
1  point_in_polygons <-
2      function(points_df, polys_sf, var_name)  {
3          is_in <- data.frame(
4                      matrix(ncol = nrow(polys_sf),
5                             nrow = nrow(points_df)))
6          is_in[,var_name] <- NA
7          coordinates <- as.data.frame(st_coordinates(
8            polys_sf))
8          name <- polys_sf$id
```

```
1      for (k in 1:nrow(polys_sf)) {
2          is_in[,k] <- sp::point.in.polygon(
3                      point.x = points_df$long,
4                      point.y = points_df$lat,
5                      pol.x   = coordinates$X
6                                [coordinates$L2 == k],
7                      pol.y   = coordinates$Y
8                                [coordinates$L2 == k])
9          is_in[,var_name][is_in[,k] == 1] <- name[k]
10     }
11     is_in <- is_in %>%
12     dplyr::select(var_name) %>%
13     dplyr::mutate(id = points_df$id)
14     points_df <- dplyr::full_join(points_df, is_in,
15                                   by = "id")
16     return(points_df)
17   }
```

⊡ distance_count:

```
1  distance_count <- function(main, reference,
2                             var_name, distance) {
3    var_name_count <- paste(var_name, "count",
4                            sep = "_")
5    var_name_dist <- paste(var_name, "dist",
6                           sep = "_")
7    point_distance <- geosphere::distm(
8                      x    = main %>%
9                             dplyr::select(long,
10                                           lat),
11                     y    = reference %>%
12                            dplyr::select(long,
13                                          lat),
14                     fun = distHaversine) %>%
15     as.data.frame() %>%
16     data.table::setnames(as.character(reference$id))
```

```
1   point_distance [, var_name_count] <- rowSums (
      point_distance <= distance)
2   point_distance [, var_name_dist] <- apply (
3         point_distance [, -ncol (point_distance)],
4         MARGIN = 1, FUN = min) %>% round (0)
5   main <- point_distance %>%
6     dplyr :: mutate (id = main$id) %>%
7     dplyr :: select (id,
8                     var_name_count , var_name_dist) %>%
9     dplyr :: right_join (main, by = "id")
10   return (main)
11 }
```

# Correlation with price

⊡ Correlation of each variable with price.
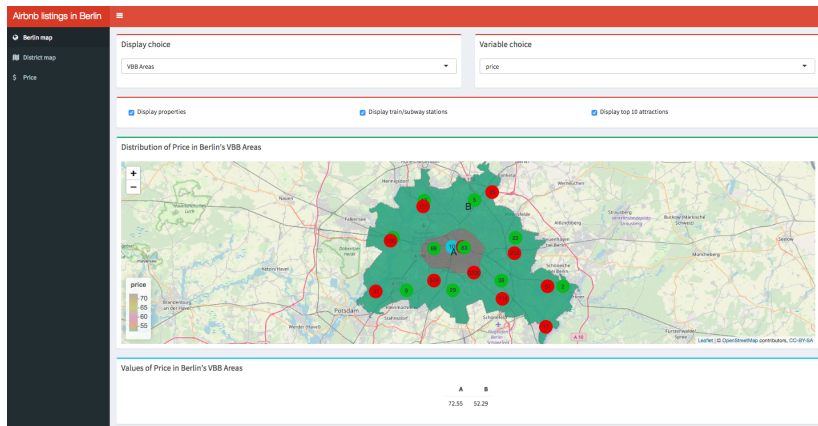


Figure 4: Correlation plot
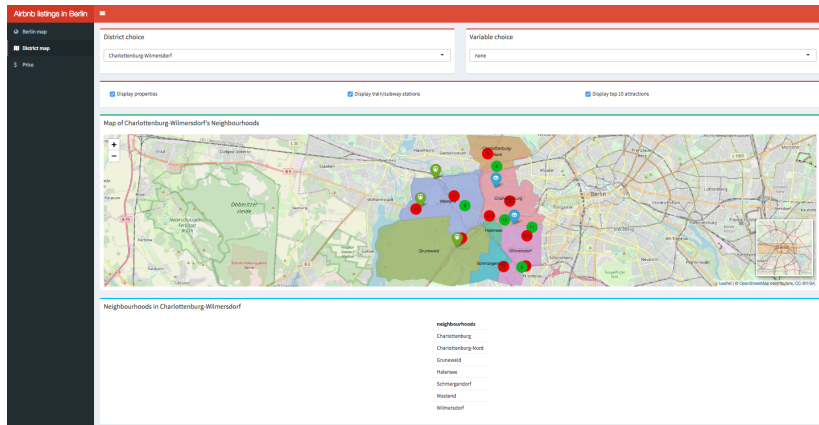
# Tab Berlin map



Figure 5: Tab Berlin map

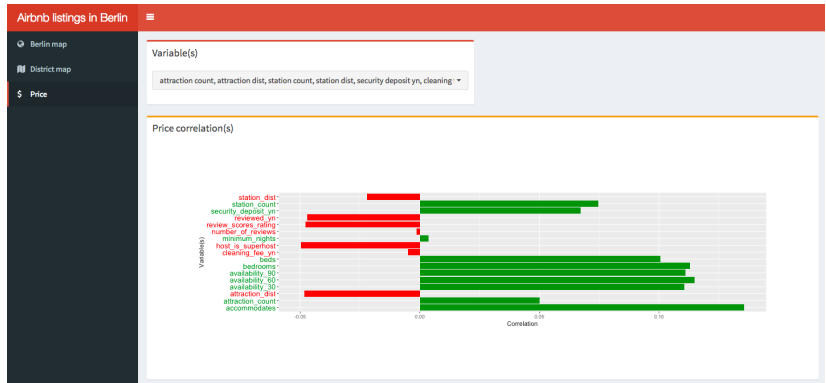# Tab District map



Figure 6: Tab District map

# Tab Price



Figure 7: Tab Price