

Milestone 03 - Group 09

Daniel Hadley, Kristina Wright

March 14, 2020

Airbnb Listings for Barcelona

Introduction

Airbnb, Inc. is a company founded in 2008 that offers an online marketplace connecting people who offer lodging with people who require accommodations in that locale. The company does not own any of the listed properties and operates as a broker, collecting commissions once a lodging is booked. As a direct competitor to hotels, we are interested in how the users listing properties determine the price they charge.

When accommodations are offered through Airbnb, the person listing the property is called a host, and they must provide a variety of information about the listing including price, neighborhood, type of accommodations offered, and the minimum number of nights a guest must stay if they want to make a booking. In addition to information provided by the host, Airbnb collects and disseminates information about the listing which we use to perform our analysis.

The data is collected using public information compiled from the Airbnb website. Specific collection techniques are not specified, though the Inside Airbnb website states that it uses Open Source technologies such as D3, Bootstrap, jQuery, etc. to collect the data and much code was “copied and pasted” from the internet. A major contributor to this code, Tom Slee, described it as a “scrape” of the Airbnb website for each city.

Data Description

The dataset used in this analysis is collected and offered by Inside Airbnb, an independent, non-commercial project started by Murray Cox and John Morris. Their goal is to allow people to see how Airbnb might be affecting the residential housing market. We use the summary data for listings, since it includes the data we are interested in exploring and is more manageable, size-wise, than the detailed listings data.

The data used in this analysis was compiled on November 9, 2019 and includes 20,428 Airbnb listings that travellers see when using the Airbnb website to find accommodations in Barcelona, Spain. The table below describes the available data for each listing in the dataset.

Variable Name	Column Name	Type of Data	Description
Listing ID	<code>id</code>	Categorical/Numeric	Numeric identifier unique to each listing
Name	<code>name</code>	Character	Short title for the listing provided by the host
Host ID	<code>host_id</code>	Categorical/Numeric	Numeric identifier for the host of the listing
Host Name	<code>host_name</code>	Categorical/String	Name of the host or hosts of the listing provided by the host(s) to Airbnb

Variable Name	Column Name	Type of Data	Description
Neighbourhood Group	neighbourhood_group	Categorical/String	Districts of Barcelona as determined by the coordinates of the listing and the city's definition of its districts; this data is not the data provided by the host
Neighbourhood	neighbourhood	Categorical/String	Neighbourhoods of Barcelona are smaller geographical areas than districts and are determined by the coordinates of the listing and compared to the city's boundaries of its neighbourhoods; this data is not the neighbourhood provided by the host
Latitude	latitude	Numeric	Latitude coordinates of the listing
Longitude	longitude	Numeric	Longitude coordinates of the listing
Type of Accommodation	room_type	Categorical/String	Type of accommodations specify whether the listing is for an entire home or apartment, a private room in a shared home or apartment, a hotel room, or a shared room
Price	price	Numeric	The price per night, in euros, to book a listing
Minimum Stay	minimum_nights	Numeric	The minimum number of nights that a guest must reserve in order to book a listing
Number of Reviews	number_of_reviews	Numeric	The number of reviews left by guests after their stay
Last Review	last_review	Date	The date of the last review left by a guest
Reviews per Month	reviews_per_month	Numeric	The number of reviews left by guests of a listing divided by the number of months the listing has been active
Number of Listings by Host	calculated_host_listings_count	Numeric	A count of the number of listings under the same Host Name

Variable Name	Column Name	Type of Data	Description
Availability	<code>availability_365</code>	Numeric	The number of days over the next 365 days that the listing can be booked by guests; calculated as 365 minus booked days minus days listing is unavailable as per the host

Exploring the Dataset

Remove Unwanted Data

In this section, we remove columns from the dataset that should have no fundamental influence on listing price. This includes the short title of the listing (`name`), the name of the host(s) (`host_name`), and the availability of the listing over the next 365 days (`availability_365`). While there might end up being a relation between availability and price, since cheap listings for a desirable neighbourhood are likely to be booked, this relationship is backwards; we want to find factors that affect the listing price, not factors affected by the listing price.

Rename Columns

Some of the column names are a little long, so we perform the following renamings:

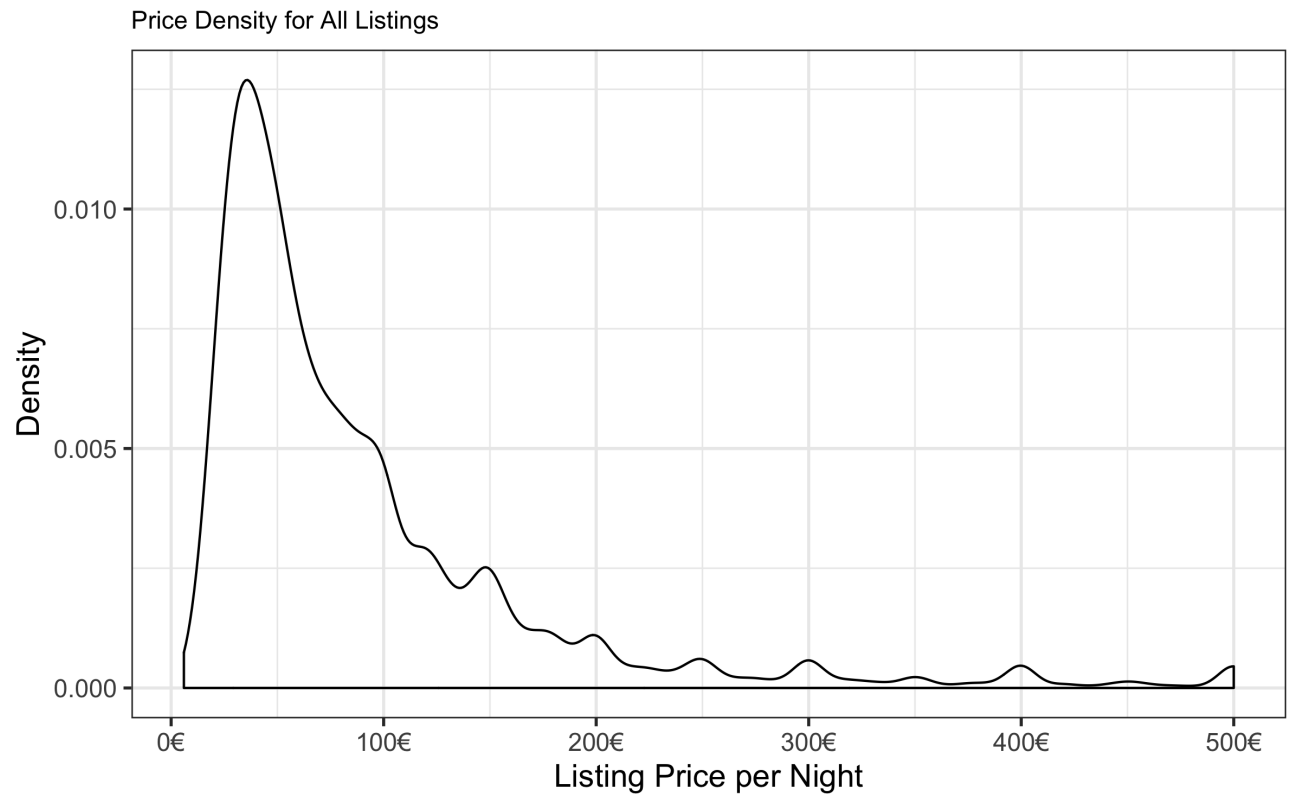
- `neighbourhood_group` is renamed to `district`
- `minimum_nights` is renamed to `min_stay`
- `number_of_reviews` is renamed to `reviews`
- `calculated_host_listings_count` is renamed to `host_listings`

Filter Data

A few extreme outliers skew the density of the price of listings to the right. As a result, we exclude the top 2.5% of listings. Then, we exclude listings with a minimum stay over 5 nights. This should help to limit listings that are catered to tourists by eliminating listings that are better classified as short-term rentals.

Price Density

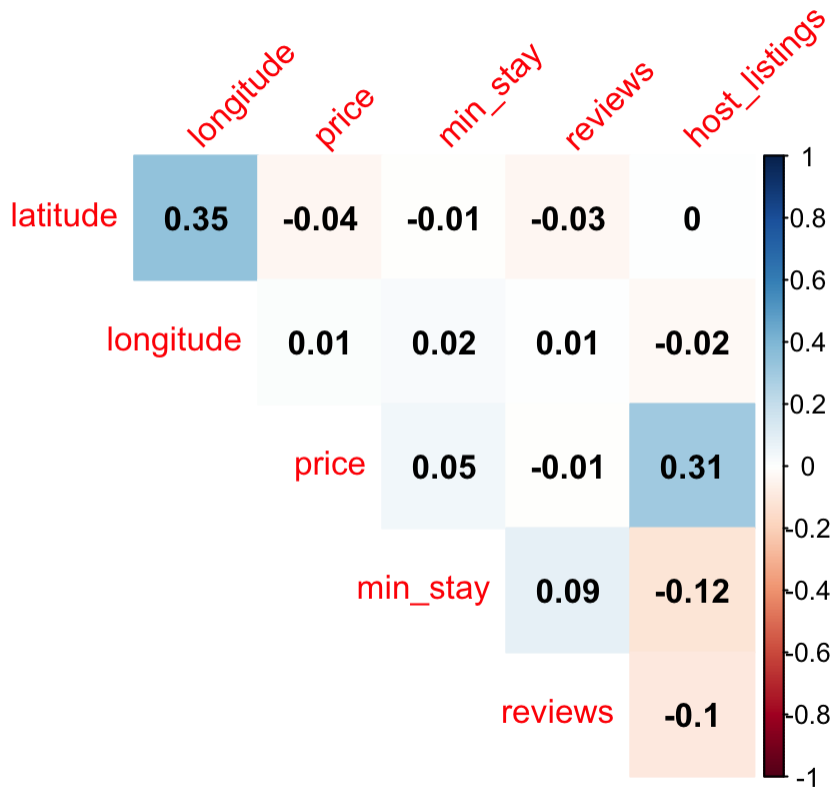
A kernel density plot is presented for listing prices. An interesting observation from the price density is the tendency for people to price their listings in increments of 50 Euros. For example, the Density Plot, we see multi-modes, where each mode after the largest mode occurs at every 50 Euro increment along the x-axis



Correlogram

Based on the correlogram shown below there is little correlation between the 6 numerical variables presented. All positive correlations are in blue, and all negative correlations are in red.

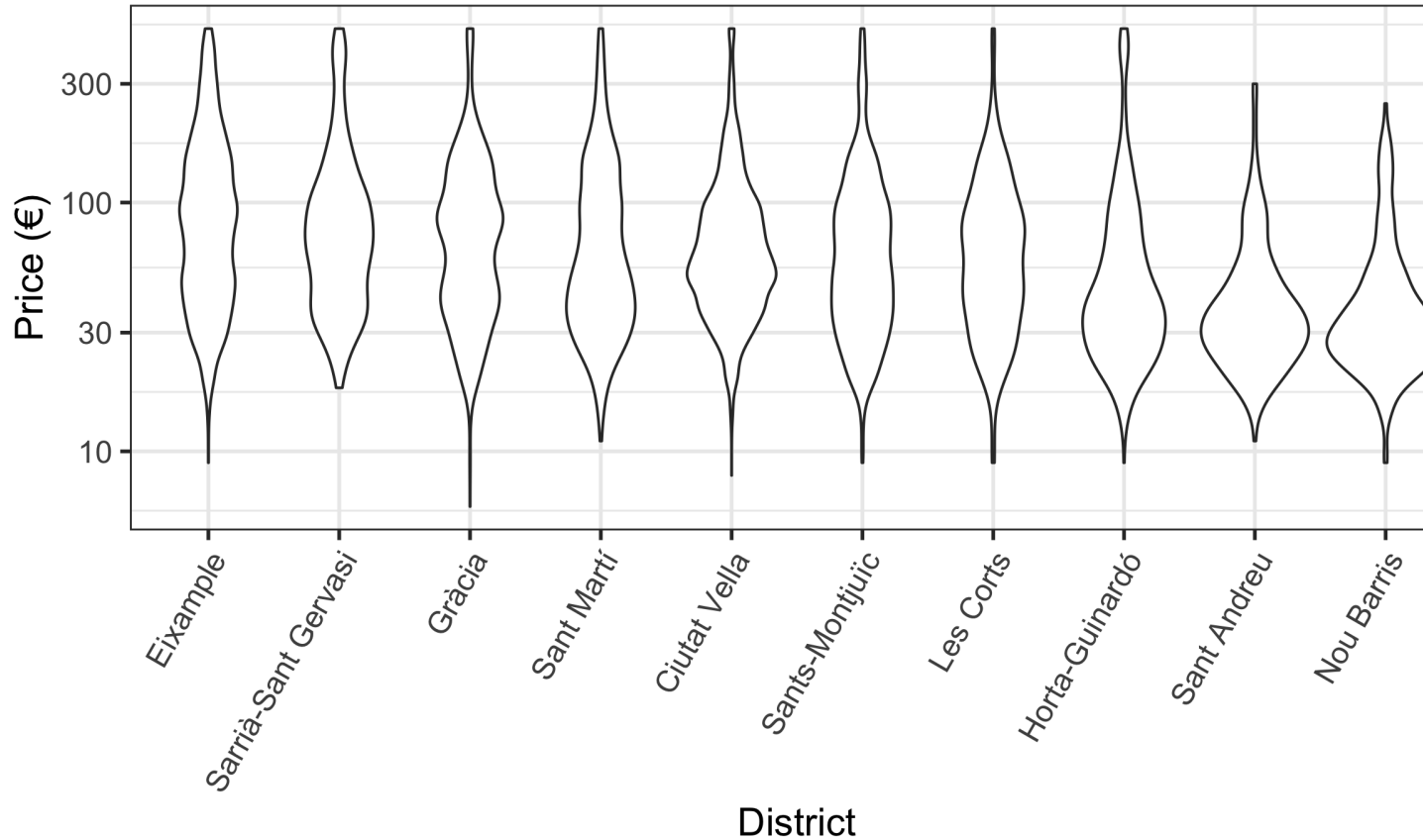
Correlation of Some Columns



Violin Plot

The violin plot below shows the distribution of price in log10 scale for each district in descending order of average price. Based on the plot, Example has the highest priced and Nou Barris has the lowest priced listings.

Distribution of Price for Each Barcelona District



Research Question

We are interested in the linear relationship between an Airbnb listing's district, type of room, reviews left per month, and distance from city center to the listing's price. The results of the linear model are presented below.

```
lm.1 <- readRDS(file=here::here("data", "lm_results"))
tidy(lm.1)
```

```
## # A tibble: 15 x 5
##   term                                estimate std.error statistic  p.value
##   <chr>                                <dbl>     <dbl>     <dbl>    <dbl>
## 1 (Intercept)                        154.         1.76      87.6     0.
## 2 districtEixample                     13.7         1.84       7.44 1.05e-13
## 3 districtGràcia                       -0.959        2.80     -0.342 7.32e- 1
## 4 districtHorta-Guinardó                4.49         3.92       1.14 2.53e- 1
## 5 districtLes Corts                     0.938         5.09       0.184 8.54e- 1
## 6 districtNou Barris                    4.50         6.08       0.741 4.59e- 1
## 7 districtSant Andreu                   -3.51         5.11     -0.687 4.92e- 1
## 8 districtSant Martí                     4.21         2.49       1.69 9.02e- 2
## 9 districtSants-Montjuïc                -2.76         2.67     -1.03 3.02e- 1
## 10 districtSarrià-Sant Gervasi          18.2         4.20       4.33 1.48e- 5
## 11 room_typeHotel room                   30.7         3.58       8.58 1.03e-17
```

```
## 12 room_typePrivate room      -89.9      1.20    -74.7    0.
## 13 room_typeShared room      -92.5      6.85    -13.5    3.19e-41
## 14 distance                -489.      68.2     -7.17    7.80e-13
## 15 reviews_per_month        -3.13      0.320     -9.80    1.40e-22
```

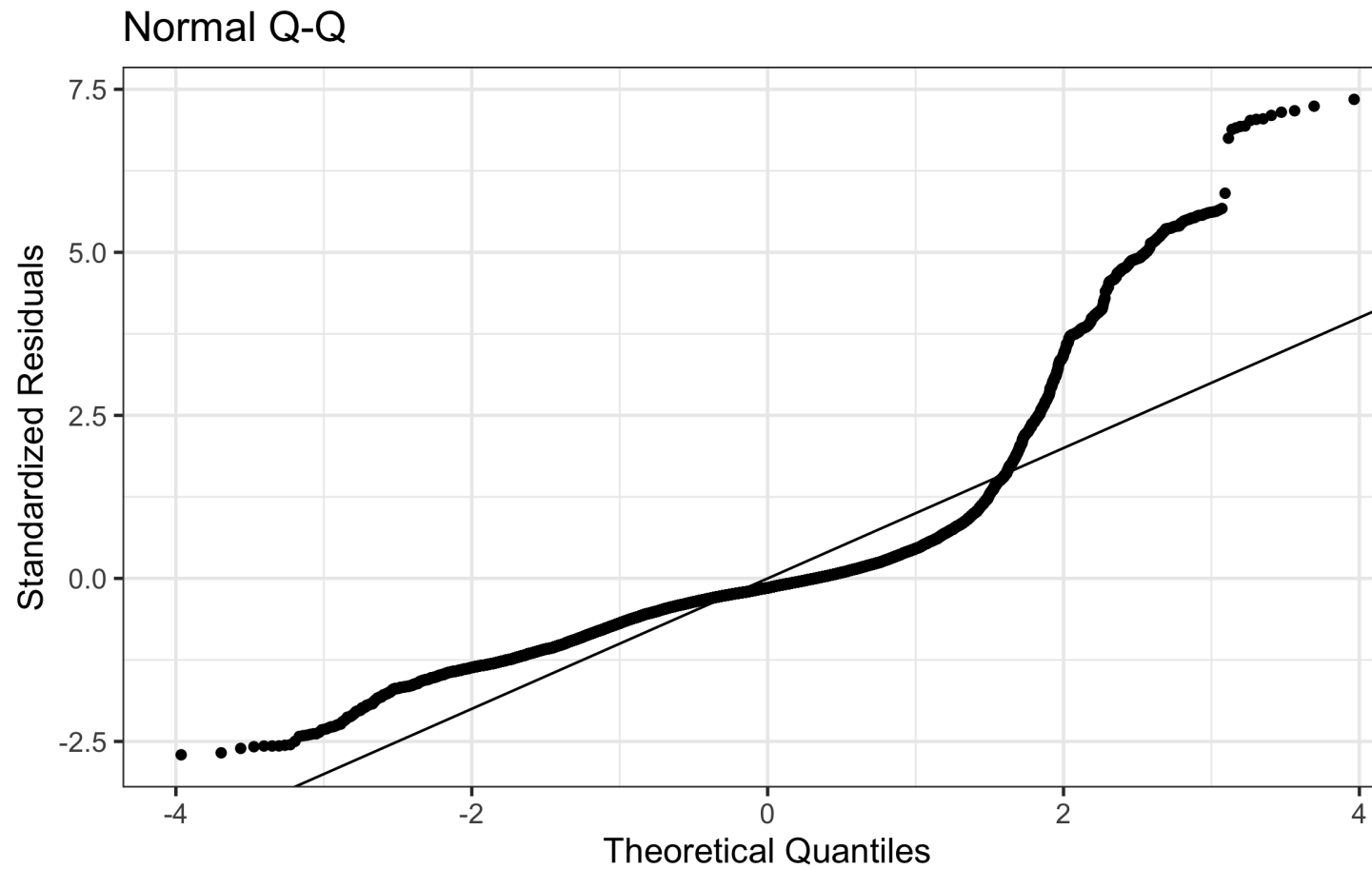
```
augment(lm.1)
```

```
## # A tibble: 13,590 x 13
##   .rownames price district room_type distance reviews_per_mon~ .fitted .se.fit
##   <chr>      <dbl> <chr>   <chr>      <dbl>      <dbl>      <dbl>      <dbl>
## 1 1          130 Sant Ma~ Entire h~  0.0253      0.02      146.      1.92
## 2 2           60 Eixample Entire h~  0.0224      0.25      156.      1.25
## 3 3          210 Sant Ma~ Entire h~  0.0482      0.48      133.      2.36
## 4 4           32 Gràcia  Private ~  0.0312      2.38       40.2      2.01
## 5 5           60 Gràcia  Entire h~  0.0343      1.71      131.      2.02
## 6 6           70 Gràcia  Entire h~  0.0330      0.84      134.      2.03
## 7 7          140 Gràcia  Entire h~  0.0244      0.580     139.      2.09
## 8 8          100 Ciutat ~ Private ~  0.00820     0.07       59.6      1.51
## 9 9          250 Ciutat ~ Entire h~  0.00708     1.28      146.      1.55
## 10 10         40 Eixample Private ~  0.0301      2.97       53.5      1.25
## # ... with 13,580 more rows, and 5 more variables: .resid <dbl>, .hat <dbl>,
## #   .sigma <dbl>, .cooksad <dbl>, .std.resid <dbl>
```

```
glance(lm.1)
```

```
## # A tibble: 1 x 11
##   r.squared adj.r.squared sigma statistic p.value    df logLik   AIC    BIC
##   <dbl>      <dbl> <dbl>      <dbl>   <dbl> <int>  <dbl>  <dbl>  <dbl>
## 1    0.360    0.359  65.3       545.     0     15 -76071. 1.52e5 1.52e5
## # ... with 2 more variables: deviance <dbl>, df.residual <int>
```

The QQ-plot is presented to assess whether the assumption of normally distributed residuals is reasonable. We can see that the residuals have an extremely heavy upper tail and a light lower tail. This plot shows pretty convincing evidence that the normality assumption is not appropriate. For a better model, perhaps we should try transformations of the variables such as a log transformation of price.



Plan of Action

With our research question, the first goal is to determine which factors are most important to explaining list price and perform a linear regression analysis. This may require some data transformation, handling or removal of outliers, and removing incomplete observations.

Methods

Results

Discussion

References

Airbnb dataset