

BI2023_gr37_12238682_11711533

Philipp Rettig, 11711533, Person B*
Vienna University of Technology
Vienna, Austria
e117115332@student.tuwien.ac.at

Vladimir Panin, 12238682, Person A
Vienna University of Technology
Vienna, Austria
e12238682@student.tuwien.ac.at

CCS Concepts: • Data Mining; • Data Analytics; • Business Intelligence; • Machine Learning;

Keywords: Data Mining, Business Intelligence

ACM Reference Format:

Philipp Rettig, 11711533, Person B and Vladimir Panin, 12238682, Person A. 2023. BI2023_gr37_12238682_11711533. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 Business Understanding

In this class assignment, the objective is to address a data analytics problem using the CRISP-DM Process. The task involves navigating through the data analysis process within the constraints of a simulated, rather than a real-world, scenario. However, a general understanding for the CRISP-DM Process is to be obtained.

1.1 Scenario

The dataset deals with AirBnB data from Berlin in July 2021. The dataset was taken from "About Inside Airbnb", which is as described by themselves a "mission driven project that provides data and advocacy about AirBnB's impact on residential communities". The data is obtained by scraping the website ¹. An imagined scenario is that an AirBnB provider wants to set an appropriate price for the apartment being offered on the website.

1.2 Business Objectives

The business objective is to be able to predict the price based on the information given about this apartment. It has to be considered that for higher prices, the apartment could be available more frequently. This trade-off has to be considered.

*Both authors contributed equally to this research.

¹<http://insideairbnb.com/about/>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org. *Conference'17, July 2017, Washington, DC, USA*

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-XXXX-X/18/06

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1.3 Business Success Criteria

The Business Success Criteria are maximizing the annual revenue and find the ideal price. The goal might therefore lead to finding a balance between the price and the average vacant time.

1.4 Data Mining Goals

The data mining goals are, first of all, to predict the price based on the information given. Furthermore, variables which have an influence on the price should be identified. In our case these can be based on different groups, e.g. districts, which have higher or lower prices or another possibility, listings that are vacant more often than others. Finally, these goals should help in finding the right balance for predicting a price, especially considering availability and location.

1.5 Data Mining Success Criteria

The success criteria for the data mining goals encompass several key aspects. Firstly, to measure the effectiveness of the model that is going to be build in the process of this assignment, we will utilize metrics like Mean Absolute Error or Root Mean Squared Error. Since standalone values of these metrics will not be very meaningful, we aim to compare them against other implementations of price predictions on the same dataset ² and use this as a benchmark for our own implementation. All in all, the model should be able to give out fairly accurate price predictions and maybe could help us understand reasons or important variables for price differences.

1.6 AI risk aspects

One possible AI risk aspect might be concerning bias and fairness. As the names of hosts are provided in the dataset, the result might discriminate on the basis of names. A possible implication might be that non "German-sounding" names could get attributed a lower price or could implicate a less wealthy district. Additionally, when using AI in housing data, there's a risk that the system might unintentionally favor certain areas. For instance, it could perpetuate biases related to neighborhoods or historical socioeconomic factors. Another concern is that the AI might rely too much on outdated information, missing out on current trends or changes in the housing market, which especially in bigger urban areas change rapidly. This could lead to Historical Bias. To address

²<https://www.kaggle.com/code/lennarthaupt/airbnb-prices-in-berlin/notebook>

these risks, it's important to evaluate and update the produced regularly. Finally, another possible bias might be the Omitted Variable Bias as possible variables like for example, the average rating of a host, the host's age is not considered.

2 Data Understanding

This section outlines a comprehensive data analysis process, covering attribute types, statistical properties, data quality considerations (such as missing values, distributions, outliers, and provenance), visual exploration, ethical sensitivity assessment, identification of potential biases, and the formulation of necessary actions in data preparation based on the analysis results.

2.1 Attribute types and their semantics

The given attribute types and their semantics are described via :

1. id has a nominal attribute type: The Listening id of the housing offer
2. name has a nominal attribute type: The title of the offer
3. host_id has a nominal attribute type: The ID of the host
4. host_name has as nominal attribute type: The first name of the host
5. neighbourhood_group has a nominal attribute type: The district of the housing offered
6. neighbourhood has a nominal attribute type: A more precise description of the location inside the district
7. latitude has a ratio attribute type (although it could be debated if a true zero point exists): Latitude of the hosing offer
8. longitude has a ratio attribute type: Longitude of the housing offer
9. room_type has a nominal attribute type: Type of the housing offer
10. price has a ratio attribute type: Price of the housing offer per night
11. minimum_nights has a ratio attribute type: Amount of minimum nights to book this housing offer
12. number_of_reviews has a ratio attribute type: Number of total reviews per month for this housing offer
13. last_review has an interval attribute type: Last review for this housing offer
14. reviews_per_month has a ratio attribute type: Number of reviews per month for this housing offer
15. calculated_host_listings_count has a ratio attribute type: Amount of listenings per for host
16. availability_365 has a ratio attribute type: Amount of available days in a year

2.2 Statistical properties describing the dataset including correlations

For the numerical variables the statistics for the quantiles, min and max, the count and the standard deviation are shown in the following three tables: [3a](#), [3b](#) and [3c](#). It is to note, that prices can get quite high up to 8000 €. It has to be checked, whether these offers are plausible. Furthermore, there is at least one room that is vacant for the whole year, 365 days. Also, for the maximum minimum_nights, there is an instance with 1124, which is probably an error, as this would be more of a rental offer with a time span of more than three years. This will be investigated further in the data preparation steps. The other variables based on the summary seem to follow a reasonable distribution.

2.3 Data Quality

There is comparably few missing data. In the dataset there are 19095 rows and 16 columns. The following table shows the missing data [1](#). From the table we can derive, that for a sizeable portion of the instances the information about the reviews is lacking. Hence, this might make the predictions more difficult, as the variables "last_review" and "reviews_per_month".

Variable	Number of NaNs
last_review	4155
reviews_per_month	4155
name	30
host_name	12

Table 1. Number of NaNs in Each Variable

The outliers were discussed in the section [2.2](#)

2.4 Visual exploration

During the exploration of the data, many insightful plots were created, of which a handful are going to be shown and explored. Figure [1](#) shows the median prices for the various neighborhoods. The x-axis represents the neighborhoods, while the y-axis indicates the mean price scale. It is evident that Neukölln and Reinickendorf offer the most affordable accommodations, whereas Charlottenburg-Wilmersdorf and Mitte host the highest-priced listenings. Additionally, it is worth mentioning the significant price variability observed in the Spandau and Marzahn-Hellersdorf neighborhoods.

Figure [2](#) shows the correlation matrix for the numerical variables. One can observe that there is a high positive correlation for the three review variables, which was to be expected. However, the price doesn't show high correlation with any of the numerical variables. The highest correlation can be observed with the availability_365 variable.

Finally, it can be pointed out that some numerical variables follow a right skewed distribution. An example is given by

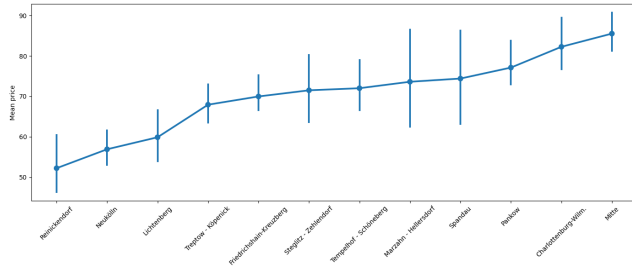


Figure 1. Mean price by neighborhood

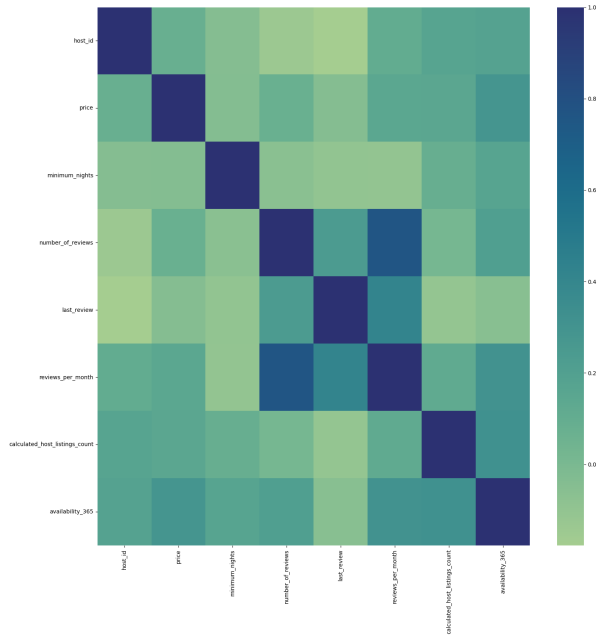


Figure 2. Correlation matrix

the variable "availability_365" following plot 3. The skewness will be further analyzed in the data preparation section.

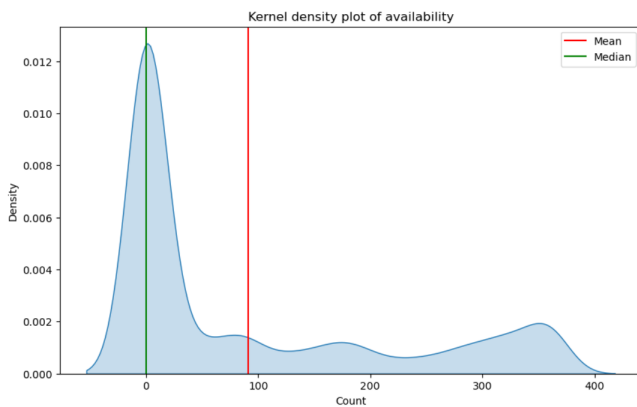


Figure 3. Kernel Density for variable availability_365

Table 2. Frequency of Neighborhoods in Berlin

Neighborhood	Count
Friedrichshain-Kreuzberg	4197
Mitte	4173
Pankow	2981
Neukölln	2608
Charlottenburg-Wilm.	1567
Tempelhof - Schöneberg	1371
Treptow - Köpenick	620
Lichtenberg	612
Steglitz - Zehlendorf	415
Reinickendorf	271
Marzahn - Hellersdorf	142
Spandau	138

2.5 Ethically sensitive information

Ethically sensitive information in this dataset, could be the distribution of names of the host. It is to say that the top 10 most occurring feature hosts are Anna, Michael, Julia, David, Baharbin, Daniel, Flo, Martin, Laura and Jan. Therefore, the classes with not classically German names could use a special sampling strategy. Apart from that, not ethically sensitive information is stored.

2.6 Risk and Bias

It has to be pointed out that some districts have much more entries than other ones. This can be seen from the counts of the different neighbourhoods 2. While Friedrichshain-Kreuzberg and Mitte have the most entries, different neighborhoods have much less offers. Therefore, this could skew the model and appropriate techniques such as under- or oversampling could be used.

An expert could answer the questions whether different housing options like hotels for example also have the same distribution around the neighborhoods. However, it would also make sense to compare to the population of the districts and the youth population, to obtain insights about the distribution itself.

2.7 Actions in data preparation

As stated before in the section 2.2 a few variables had unreasonable ranges. These outliers have to be further investigated and dealt with by for example dropping or imputing with reasonable values. It has to be considered that machine learning methods usually can't deal with categories. Therefore they will be converted to a one hot encoded representation. Possibly, connections between the "host_name" and the price could be found. Additionally, the variables have to be converted to reasonable types. Going on, missing values have to be dealt with by for example imputation.

3 Data preparation

This section comprises the data preparation phase, including analysis of the possible and necessary steps as well as incorporation of additional external data sources.

3.1 Necessary actions

Using the information gathered in the data understanding section, a few necessary actions have to be performed.

3.1.1 Missing values. Using the information layed out in section 2.3, we decided to remove the entire columns for name and host_name, taking into consideration the high number of unique values in them. After some further investigation, we figured that the missing values in last_review and reviews_per_month always occurred when a listing had zero reviews. Thus, the strategy was not to drop the missing values, but to impute them with a numerical value. For the reviews_per_month we replaced the NaNs with zeros and for the last_review we first transformed the date to "days since last review" and then replaced the NaNs with the highest value.

3.1.2 Outliers. Some listings in the dataset did show a price of zero. We considered this an error in the data and removed the corresponding rows. Other outliers in this column included listings with a price of 4000 or 8000. Such high values are not very frequent, but still we thought it reasonable that high valued listing like these could exist and decided to keep them unchanged. The minimum nights attribute includes some high values up to 1124 nights for a few listings, see figure 4, which could be valid, since there probably exists some hosts who focus on long term rentals. However, for these few values to not have a strong impact on the models performance, we decided to only focus on short term rentals and removed all rows with "minimum nights" > 30.

3.1.3 One hot encoding. For our machine learning model to be able to handle the categorical input features, some variables are one to n encoded. These comprises the attributes neighbourhood group and room type. It has to be noted that this procedure increases the dimensionality of our dataset by the number of categories -1. Considering the high number of records (19000) this should have a sizeable influence.

3.1.4 Attribute removal. As a final necessary step we considered to remove some attributes, which might not be that important for the data mining goal or simply not useful for the model:

- id and host_id: not useful for the model, only unique values
- name and host_name: too many unique values and difficult to group them further
- latitude and longitude: instead we create a new feature (see section 3.2)

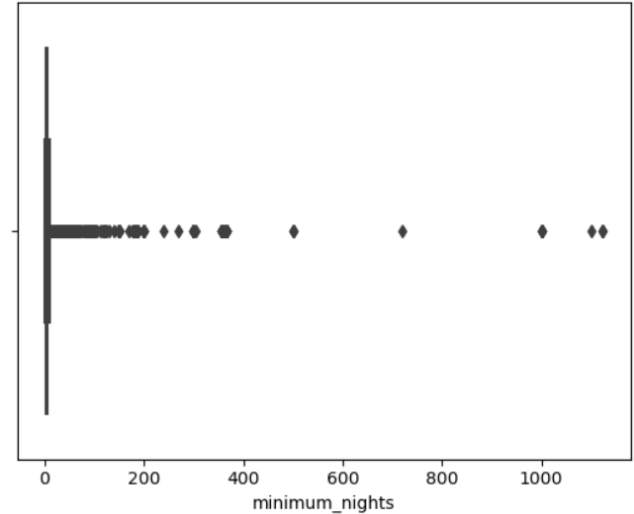


Figure 4. Boxplot minimum nights

- neighbourhood: too many categories (138), instead use the neighbourhood group

3.2 Derived attributes

As mentioned in section 3.1.1 the attribute "last_review" is transformed to a new attribute, which represents the days since the last review instead of the date of the last review. The numerical representation should aid in better understandability for the machine learning model. Additionally, the latitude and longitude information is used to create a new attribute, which represents the distance to the berlin city center. The chosen coordinates are subjective: (52.51638889, 13.37888889). The distance is calculated using the hamming function and the resulting values are given in kilometers.

Other potentials include the transformation of the host_name to a binary category indicating, whether the name is German sounding or not. We speculate that this might influence some Airbnb guests on choosing a listing. However, as this would take a very long time to manually annotate, we decided not to perform these steps.

3.3 External data sources

The following section will discuss options about additional data sources and attributes.

To enhance the dataset one option is to consider **demographics** about the neighbourhood, like population size and density, average income, proximities to public transport or age distribution. These could provide some further insights into the different regions and why prices differ.

Especially useful for airbnb listings are probably **events** happening in the near area. Hence, the data could be enriched with information about local events, to provide better understanding of price differences at different times.

As already mentioned, it would be ideal to make use of the host name and classify it as German sounding or not. Ideally, the data could be enriched with an extensive list comprising **German names**, which could be used to perform the aforementioned binary classification.

Lastly, a more extensive analysis using **additional months** would be preferred. This would allow use to consider temporal trends and perform monthly comparisons.

3.4 Other preprocessing steps

In this section we outline some additional preprocessing steps we considered, but didn't deem imperative for a successful machine learning model.

3.4.1 Scaling. Most of the numerical values in the dataset, are simply count based, which means they are in the same range of magnitude. However, the attributes "distance to city center" and "reviews per month", are not simply counts, but feature a different range that might not be easily comparable. To avoid for an explanatory variable to have an unreasonable impact in the model, we decided to align the scale for each using a min max scaler. It has to be noted though that we will lose some sense of interpretability using this procedure.

3.4.2 Log transformation. By utilizing a kernel density plot we did get some idea of the distribution of the numerical attributes. Most of them did show a slightly right skewed distribution. However, the median and mean seem to not deviate too much in most cases. Therefore, we are not sure if log transforming these attributes will have a significant impact in the performance of the model, especially considering that we will also lose some sense of explainability. After log transforming the number of reviews and reviews per month, the distribution seems to be more normally distributed. Eventually, we decided against transforming the attributes, but keep this option in mind in case the results are not satisfactory.

4 Appendix

Variable	count	mean	std
latitude	19095	52.510215	0.032391
longitude	19095	13.404654	0.062953
price	19095	73.303221	136.249622
minimum_nights	19095	9.105944	33.635956
number_of_reviews	19095	21.637078	48.670427
reviews_per_month	14940	0.718274	1.445272
host_listings_count	19095	3.135847	7.773246
availability_365	19095	91.271694	127.645330

(a) Part 1

Variable	min	25%	50%
latitude	52.340070	52.489710	52.509950
longitude	13.097150	13.367160	13.414090
price	0	35	52
minimum_nights	1	2	3
number_of_reviews	0	1	4
reviews_per_month	0.01	0.09	0.27
host_listings_count	1	1	1
availability_365	0	0	0

(b) Part 2

Variable	75%	max
latitude	52.533320	52.656110
longitude	13.438900	13.757370
price	81	8000
minimum_nights	5	1124
number_of_reviews	17	620
reviews_per_month	0.83	94.35
host_listings_count	2	76
availability_365	175	365

(c) Part 3

Table 3. Numeric Summary