

# Airbnb Price Prediction

Team 1: Veronika Junková, Filip Rott, Lucie Pinterová, Daniel Borner

Course: Data X

## 1 Data Understanding

For predicting Airbnb prices in Prague, we use the listings dataset.

### **Listings.csv.gz**

This dataset contains records of individual Airbnb listings in Prague. It includes the price, ratings, maximum occupancy, and many other accommodation-related details.

The dataset contains 8366 rows and 75 columns. It encompasses text, numerical, categorical, boolean, date, ID, json and null columns. Among these, 26 columns (34,6% ) contain some null values, while 3 columns (4 %) solely consist of null values. There are in total 7,82 % missing cells in Listings dataset.

Text Columns: listing\_url, name, neighborhood\_overview, picture\_url, host\_url, host\_name, host\_location, host\_about, host\_thumbnail\_url, host\_picture\_url, host\_neighbourhood, neighbourhood, property\_type, bathrooms\_text, description

Numeric Columns: host\_response\_rate, host\_acceptance\_rate, host\_listings\_count, host\_total\_listings\_count, latitude, longitude, accommodates, beds, price, minimum\_nights, maximum\_nights, minimum\_minimum\_nights, maximum\_minimum\_nights, minimum\_maximum\_nights, maximum\_maximum\_nights, minimum\_nights\_avg\_ntm, maximum\_nights\_avg\_ntm, availability\_30, availability\_60, availability\_90, availability\_365, number\_of\_reviews, number\_of\_reviews\_ltm, number\_of\_reviews\_l30d, review\_scores\_rating, review\_scores\_accuracy, review\_scores\_cleanliness, review\_scores\_checkin, review\_scores\_communication, review\_scores\_location, review\_scores\_value, calculated\_host\_listings\_count, calculated\_host\_listings\_count\_entire\_homes, calculated\_host\_listings\_count\_private\_rooms, calculated\_host\_listings\_count\_shared\_rooms, reviews\_per\_month, bathrooms, bedrooms

Categorical Columns: source, host\_response\_time, host\_verifications, neighbourhood\_cleansed, room\_type

Null columns: neighbourhood\_group\_cleansed, calendar\_updated, license

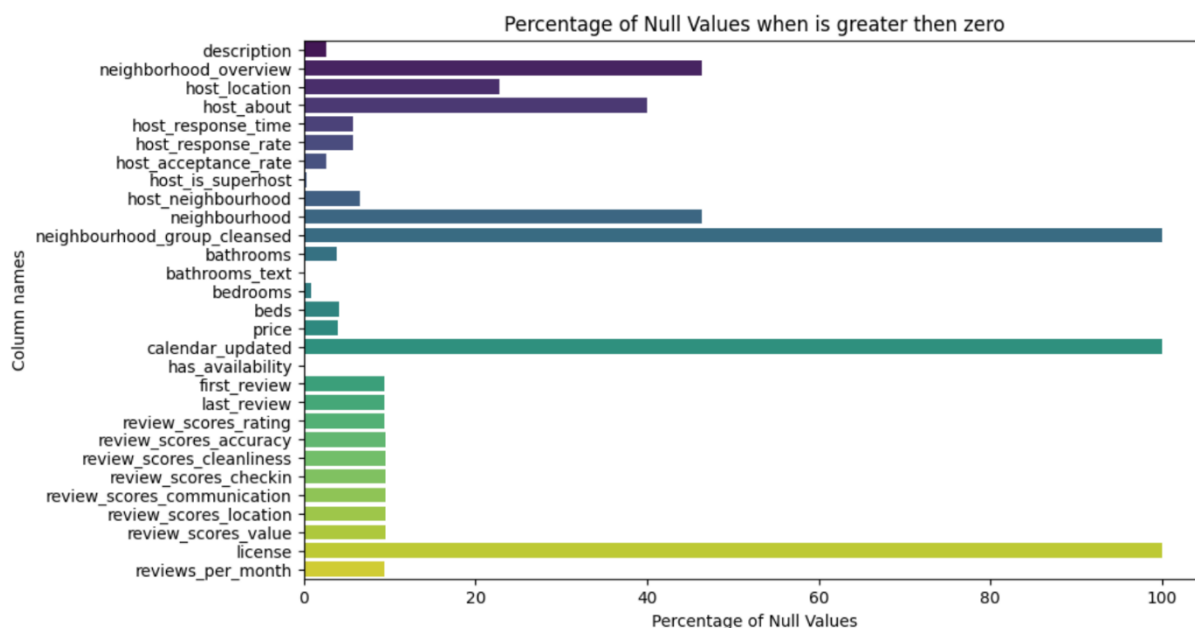
Datetime columns: last\_scraped, host\_since, calendar\_last\_scraped, first\_review, last\_review

Id columns: id, host\_id, scrape\_id

Boolean columns: host\_is\_superhost, host\_has\_profile\_pic, host\_identity\_verified, has\_availability, instant\_bookable

JSON columns: amenities

The chart above illustrates the percentages of null values present in each column.

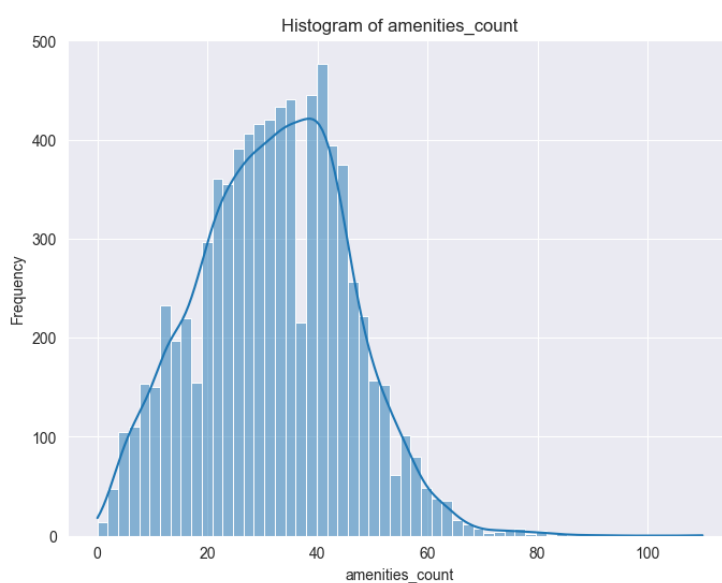


## 2 Data Preparation

### Added Columns:

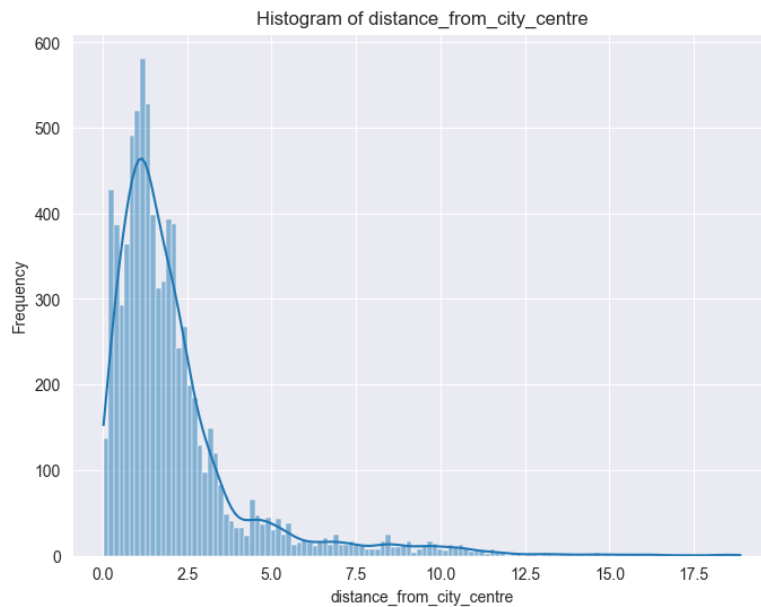
We created some other columns, which we thought would be useful for the modelling.

amenities\_count – after parsing the amenities column we added the column amenities\_count, which represents how many amenities the Airbnb has



count\_verifications – instead of having one column (list of verifications), we decided to have a number of how many verifications host has in another column

distance\_from\_city\_centre – as we all know, when booking an apartment, we all consider how far is the Airbnb from the city centre, so we took the location of Old Town Square as a city centre and calculated thanks to longitude and longitude of the apartment distance in kilometres



season, seasonal\_availability – thanks to availability\_30, availability\_60, availability\_90, availability\_365 we added a column which represents the seasonal availability

minimal\_rating, maximum\_rating

#### **Modified Columns:**

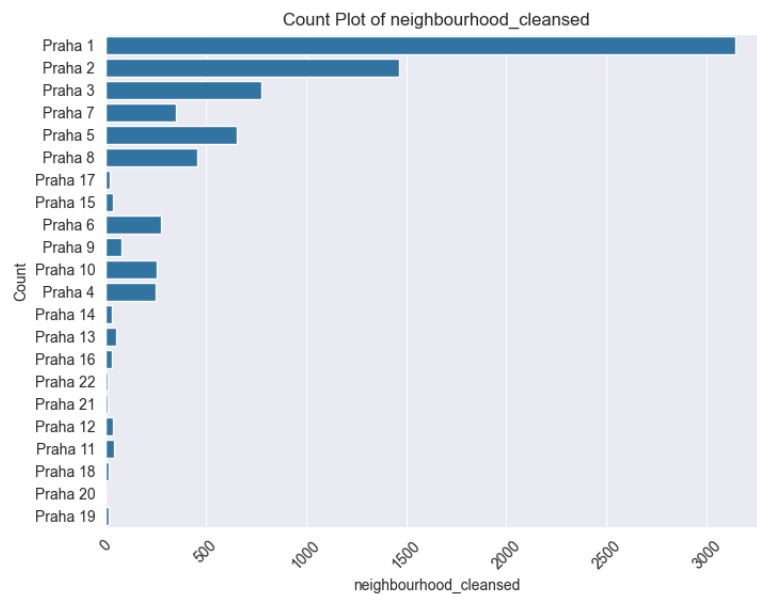
bathrooms\_text: text descriptions of bathrooms were extracted into numeric values for easier analysis and aggregation (2 baths into 2)

has\_availability, host\_has\_profile\_pic, host\_identity\_verified, host\_is\_superhost, instant\_bookable: boolean values were converted to 0 and 1, which allows simpler logical operations and integration into mathematical models

host\_acceptance\_rate, host\_response\_rate: Percentage values were extracted from the text format and converted to a numeric format

price: We removed the currency symbols and commas, converted it into a numeric format

neighbourhood\_cleansed: we unified the neighbourhoods



### Removed Columns:

scrape\_id, calendar\_last\_scraped, last\_scraped: columns about data collection, not useful for our analytical models

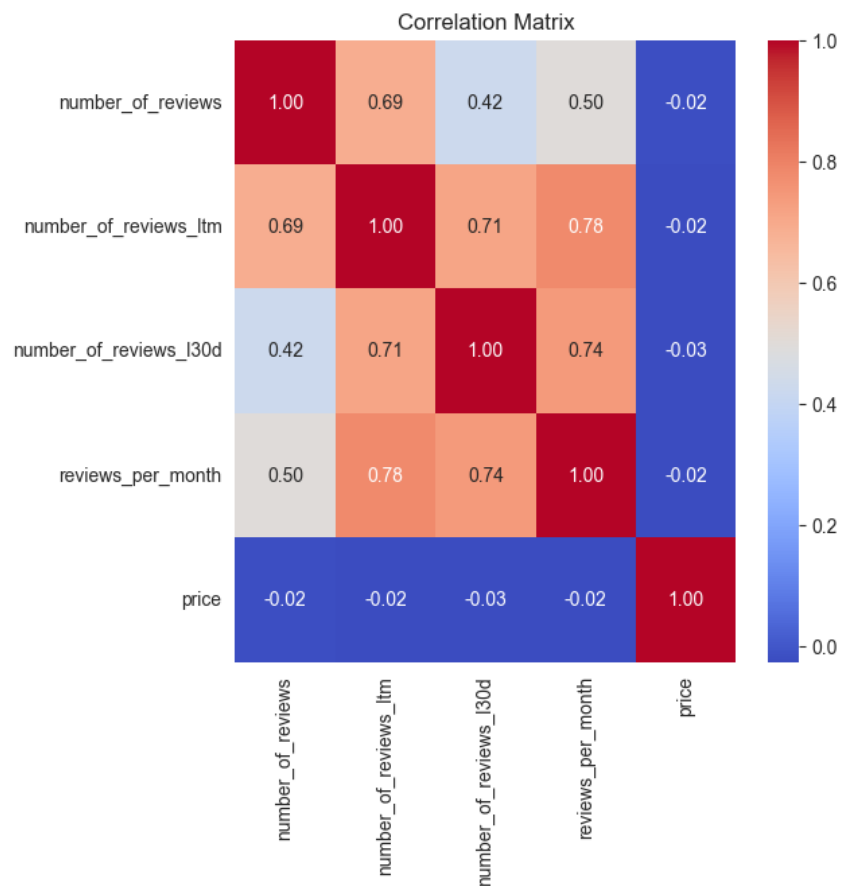
latitude, longitude, neighbourhood, neighbourhood\_group\_cleansed: we decided to use the information from these columns into new columns neighbourhood\_cleansed, distance\_from\_city\_centre, which are more useful for modelling, so we decided to remove the additional columns

description, neighbourhood\_overview, amenities: these textual columns were not useful for other modelling, thus they were removed

availability\_30, availability\_60, availability\_90, availability\_365: instead of using these columns, column season and seasonal\_availability replaced them

For outlier detection, we used the Z-score method to identify unusual data points. These outliers were visualized using box plots and removed.

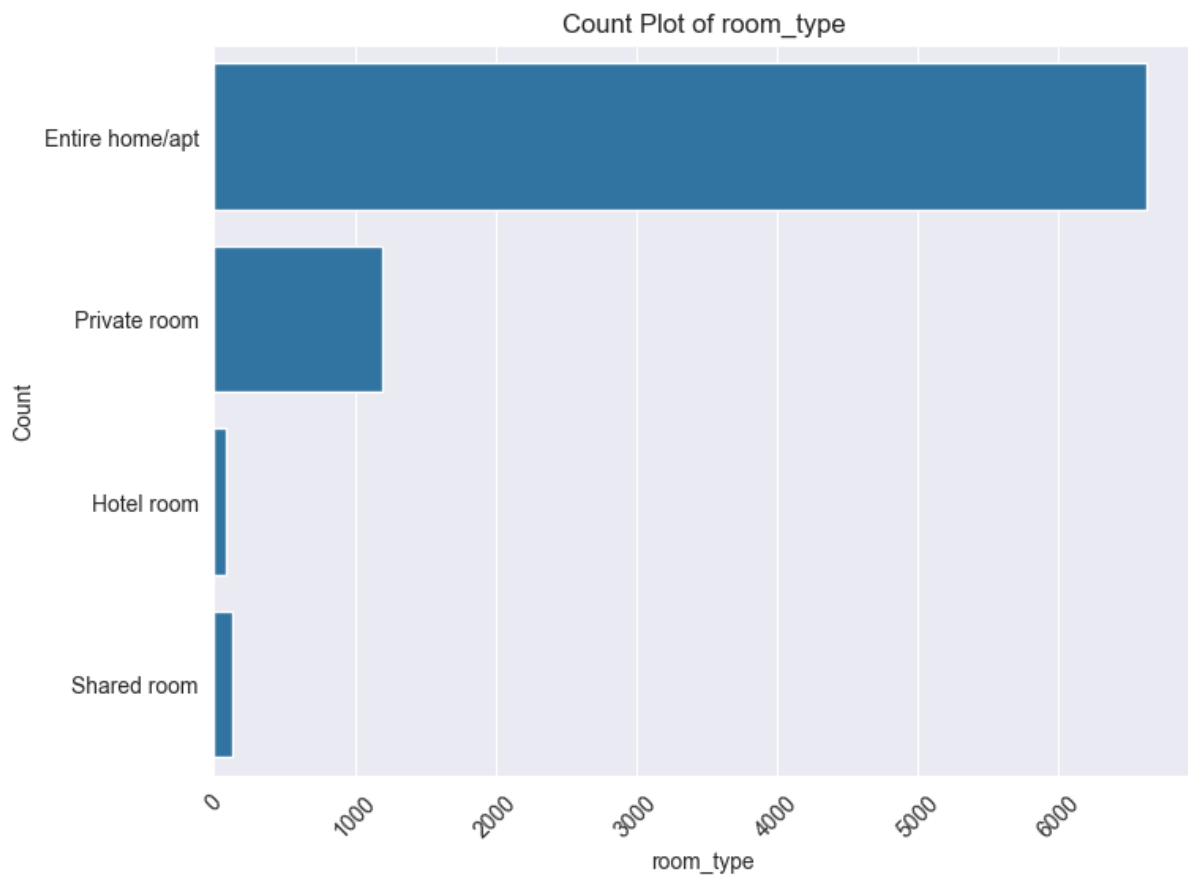
number\_of\_reviews, number\_of\_reviews\_ltm, review\_scores\_rating – since we decided to just use minimal\_rating and maximum\_rating and don't do a sentiment analysis these columns were dropped



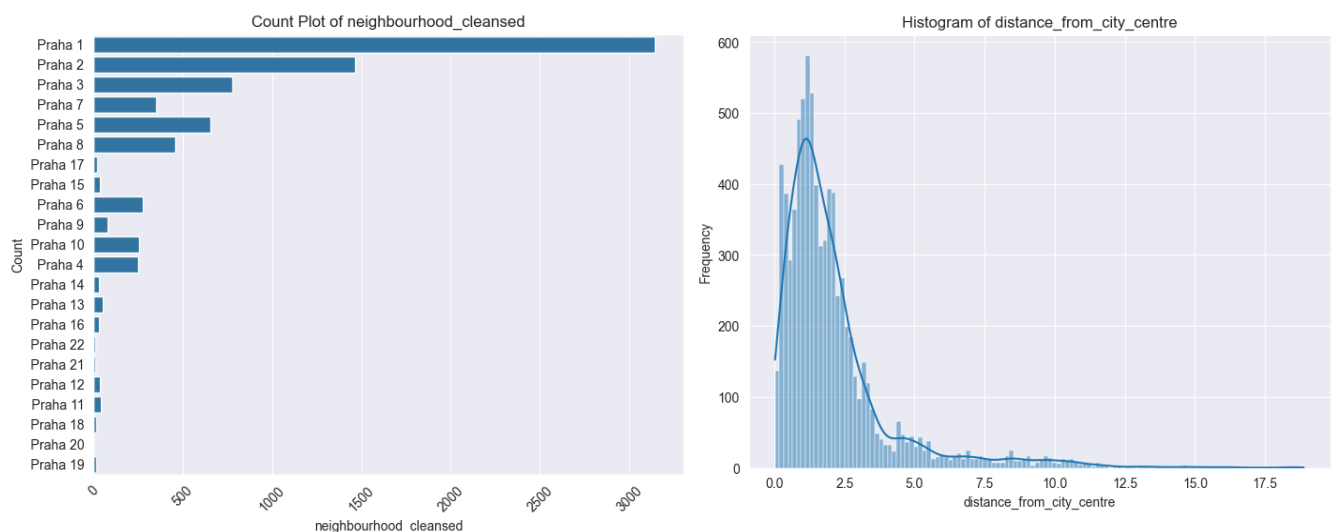
minimum\_minimum\_nights, maximum\_minimum\_nights, minimum\_maximum\_nights, maximum\_maximum\_nights, minimum\_nights\_avg\_ntm, maximum\_nights\_avg\_ntm – we thought that enough information is included in minimum\_nights and maximum\_nights so the others were removed

### 3 Data Visualization

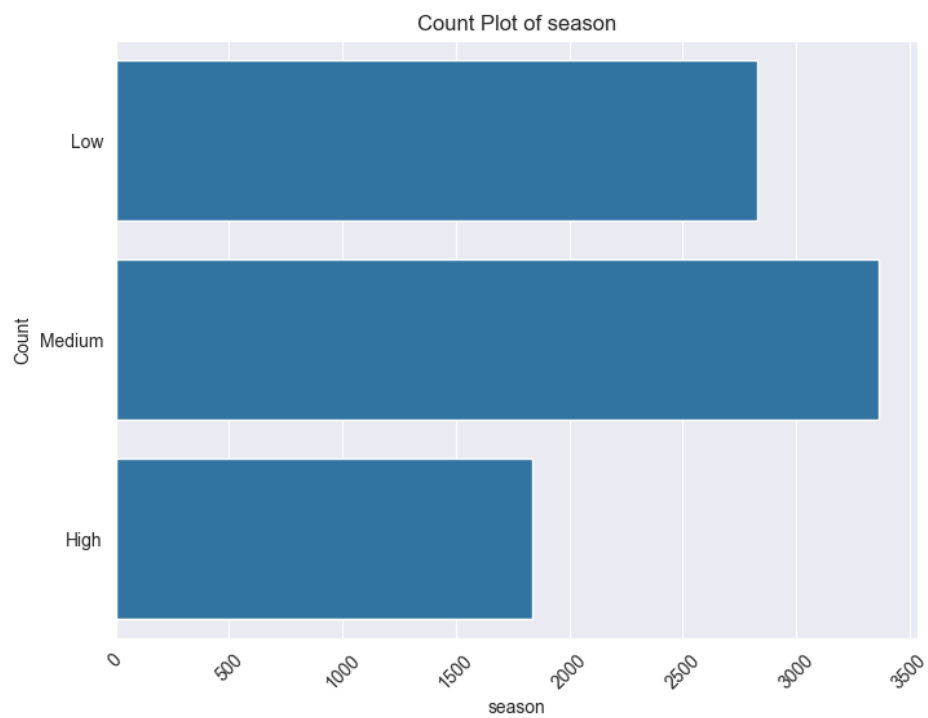
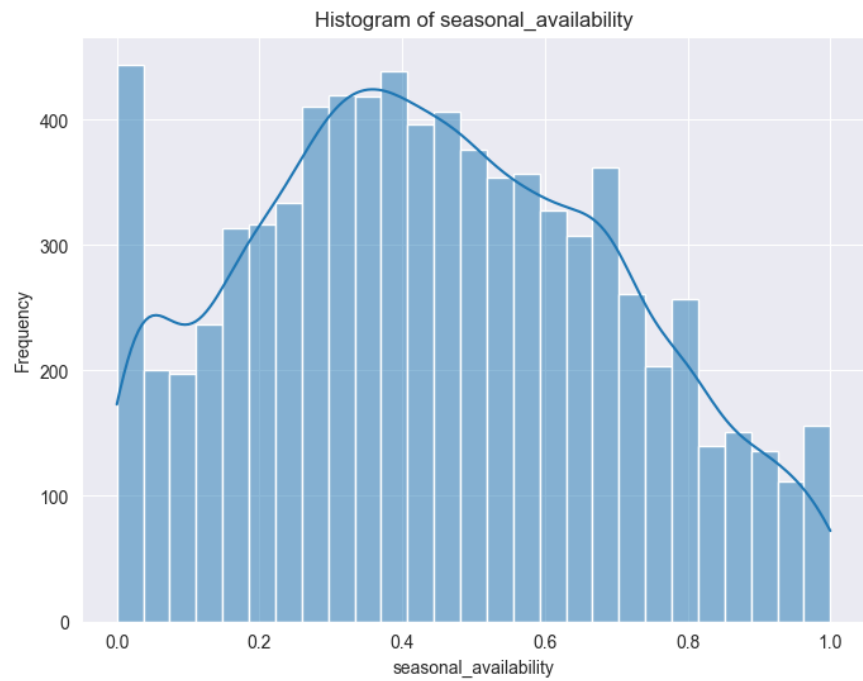
Here we see that most of the Airbnbs are Entire homes/apartments.



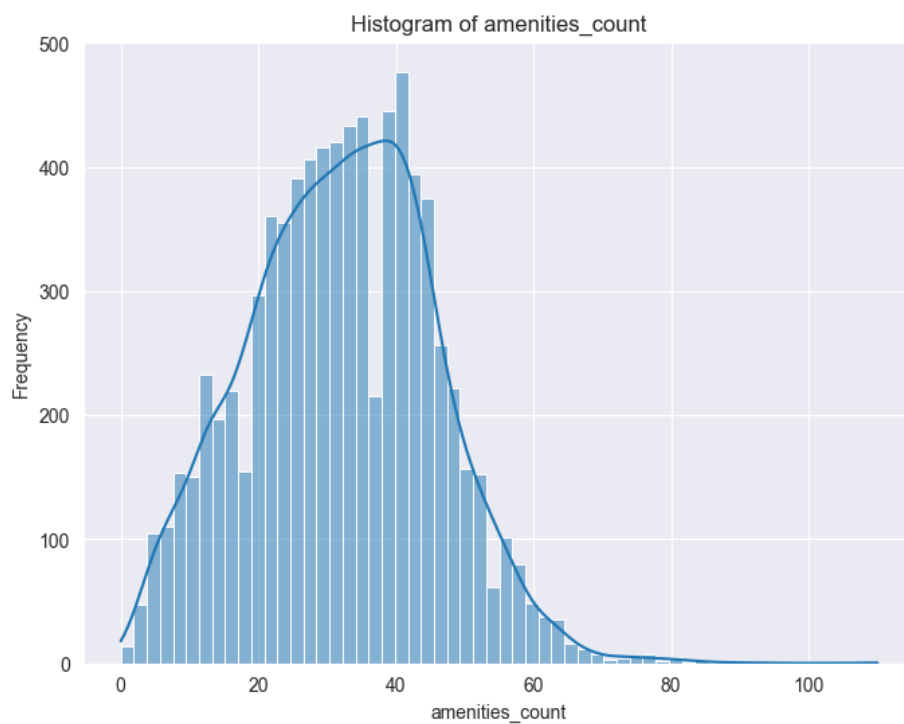
Most of the being in the city centre in Prague1, Prague 2 or Prague 3.



Most of the Airbnb is only available between 40 % and 80 % of the time. Most of the Airbnbs are medium-low available during the time.



The number of amenities listed.



The median price is very low compared to the maximum, most prices are in the lower range. The outliers may affect the modelling.





## 4 Modelling

### Setting train and test set

In the first step of modelling, we split our data into two groups: a test set and a training set. We divided the data into 20% for the test set and 80% for the training set. The training set is used for training the model, and the test set is used for comparing the models with each other. All models are trained on the same dataset. We used seed 30.

### Model selection

For predicting Airbnb prices, we tried using the models Linear Regression, Random Forest, CatBoost, and XGBoost. For Random Forest and XGBoost, we tuned the hyperparameters using cross-validation. We call these improved models Optimized Random Forest and Optimized XGBoost.

### Metrics

To compare the models, we chose the R-Squared value. This value, ranging from 0 to 1, shows how well the data fit the regression model.

### Linear Regression

R-squared = 0.018

Linear regression with its R-squared value performed clearly the worst.

### XGBoost

R-squared = 0.23

XGBoost performed slightly better than linear regression. Even though we tried different hyperparameter settings, we couldn't achieve an optimal model configuration.

### CatBoost

R-squared = 0.36

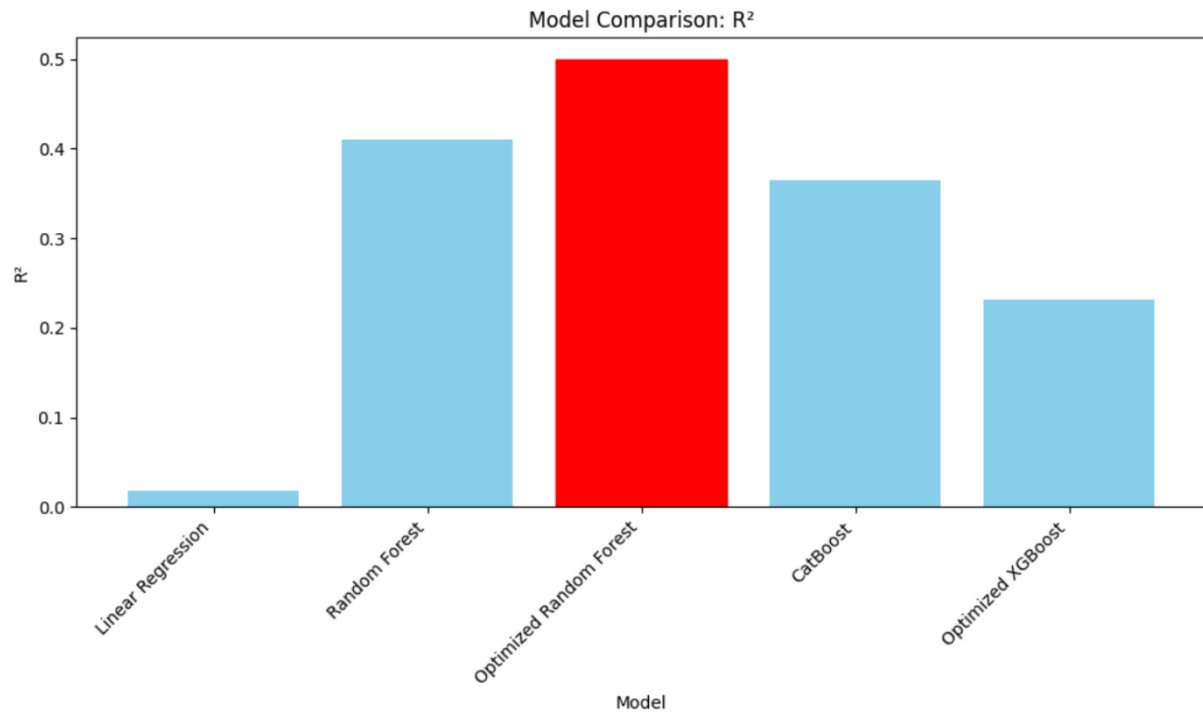
CatBoost performed better compared to XGBoost and linear regression, but its R-squared results were still not the highest.

### Random Forest

Default random forest = 0.41

Optimized random forest = 0.50

The model with the highest R-squared is Random Forest. We first tried this model in its basic form, where it surprisingly performed the best. After tuning the hyperparameters, we managed to achieve an R-squared of 0.50. For selecting the best hyperparameters, we chose the GridSearchCV method, which, with the help of cross-validation, selected their optimal values. We selected hyperparameters within a reasonable range based on our knowledge from the lectures.



The image shows a comparison of the different models.

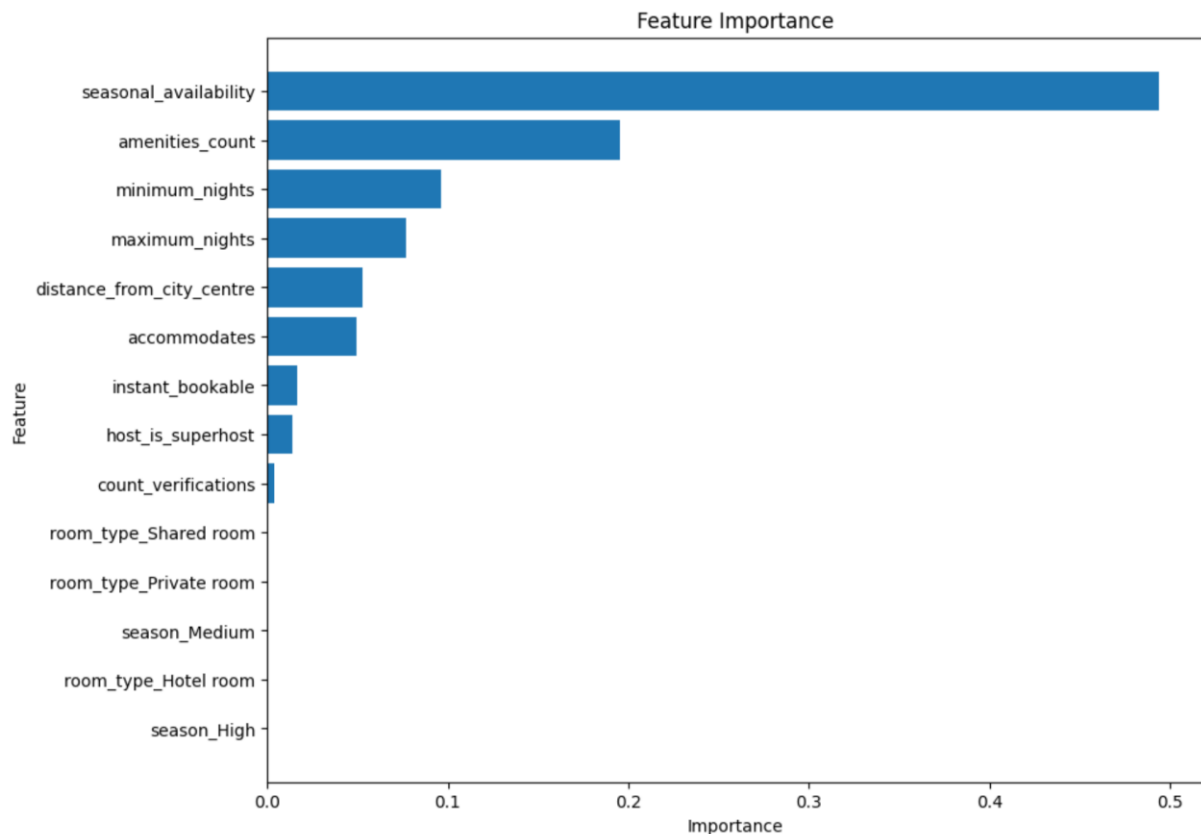
### Model limitation

A significant limitation of the models we tested is data processing. Although we spent a considerable amount of time on data manipulation and cleaning, we believe that adding additional information and utilizing current dataset data more comprehensively could lead to better results. The second limitation is the choice of hyperparameters. With greater computational power, we could explore a wider range of hyperparameters and thus refine the results of our model.

## 5 Model interpretation

### Feature importance

The top predictors of property prices in the Optimized Random Forest model are seasonal availability, distance from the city center, and amenities count. While seasonal availability indicates significant price fluctuations, amenities count suggests that more amenities tend to lead to higher prices, alongside proximity to the city center.



*Feature importance of Optimized random forest model.*