

Airbnb Optimizer

Smart Pricing and Investment Advisory Platform
based on the Madrid Airbnb Dataset



Agenda

1. Introduction
2. Motivation
3. Dataset overview
4. Feature visualisation
5. Predictive models
6. Classification models



Introduction

Problem & Value Proposition:

- Development of an AI-powered platform that helps property owners and investors maximize their Airbnb rental potential in Madrid
- **Provided data-driven insights for:**
 1. Optimal pricing strategies
 2. Neighborhood performance analysis
 3. Potential return on investment
 4. Competitive positioning



Dataset overview

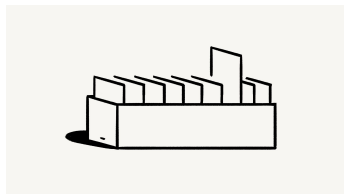
- **Madrid Airbnb Data**
- Information about Airbnb listings in Madrid, Spain.
- Relatively up-to-date data:
11 September, 2024
- source: Kaggle platform
- It consists of 6 csv files that contain data about airbnb places in the city

kaggle

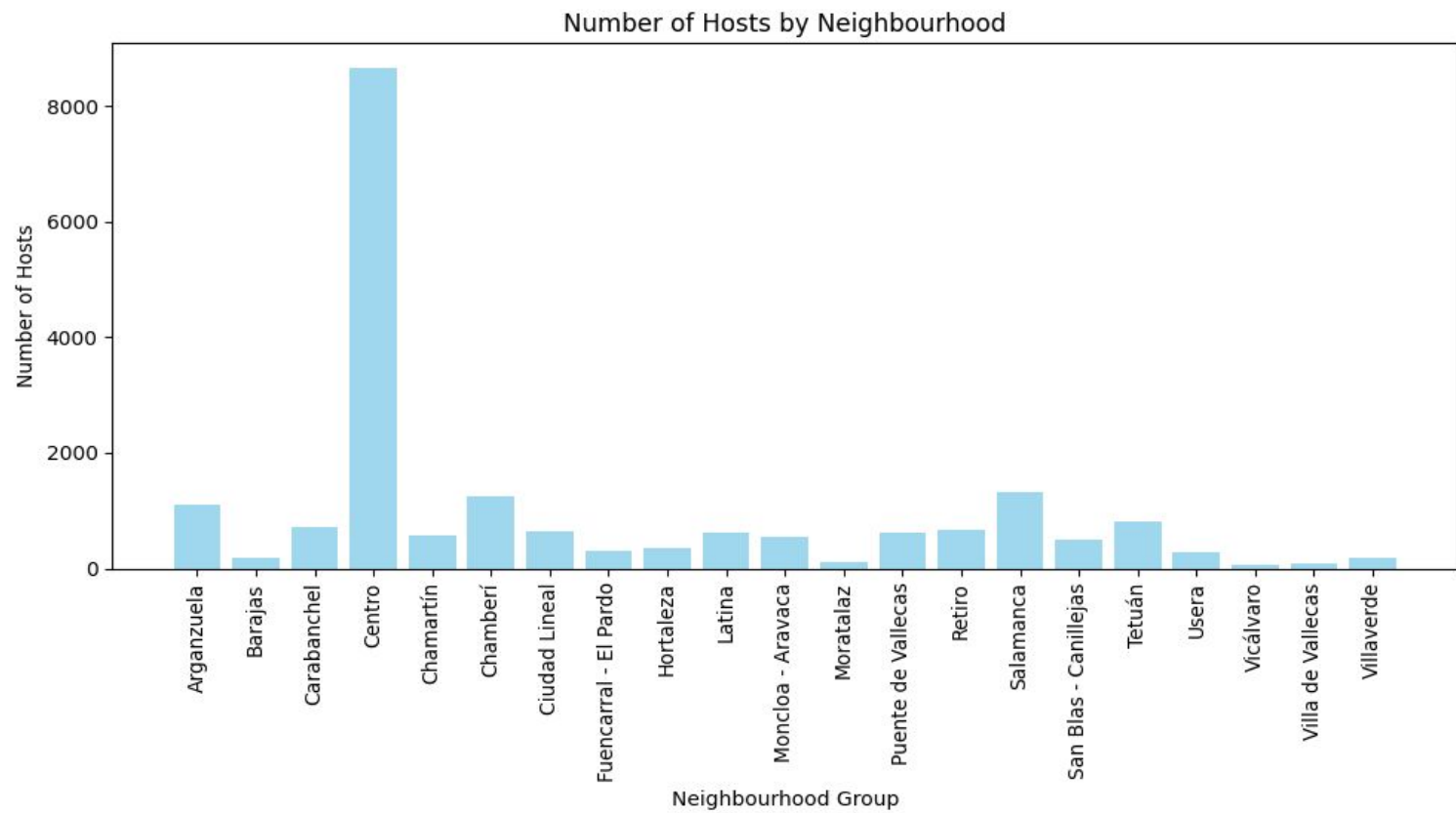


Contents – the same, probably not enough time

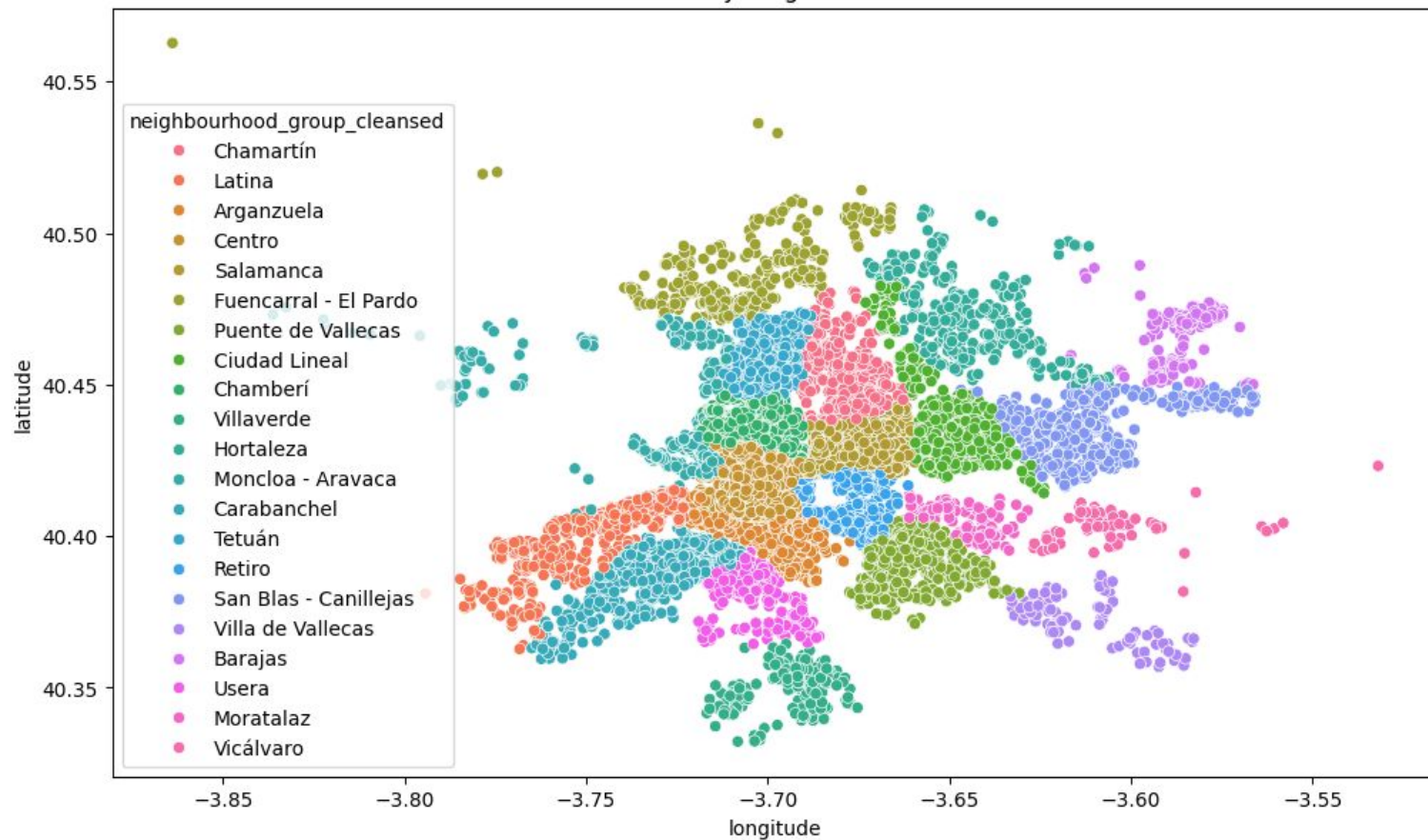
- **calendar.csv** - availability of places depending on the dates of arrival, as well as the prices and number of min and max nights that are possible to book.
- **reviews.csv** - dates of the collected reviews
- **reviews_detailed.csv** - extension of the previous file, with added information about the author of the review and the review content
- **neighbourhoods.csv** - assembly of the particular neighbourhoods
- **listings.csv** - information such as name, hostname, neighbourhood, location, type of room/flat, price, min nights, availability and reviews about a particular place
- **listings_detailed.csv** - extension of the previous file with added links to the websites of such place and its description. On top of that we can obtain some statistical data describing number of reviews or room scores.



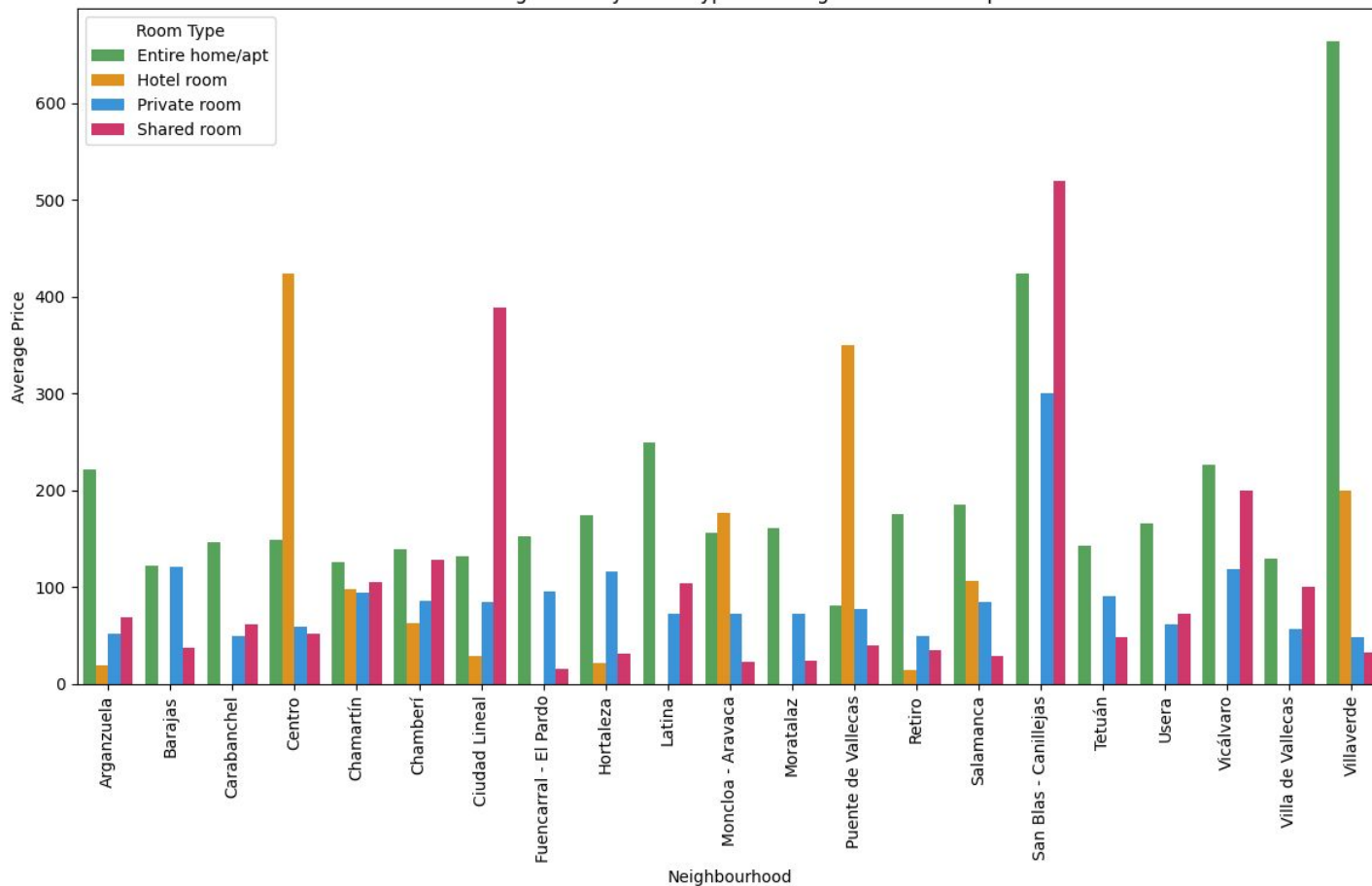
Data Visualization



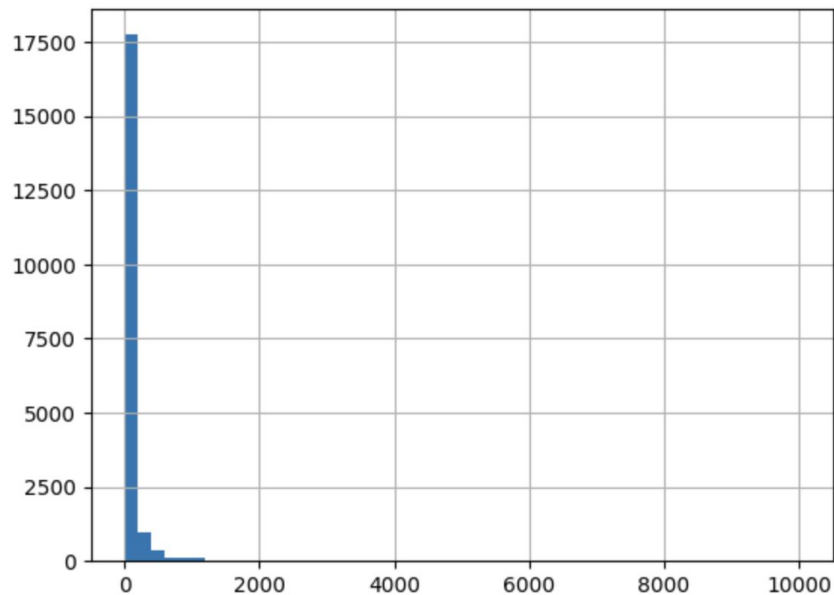
Airbnb by Neighbourhood



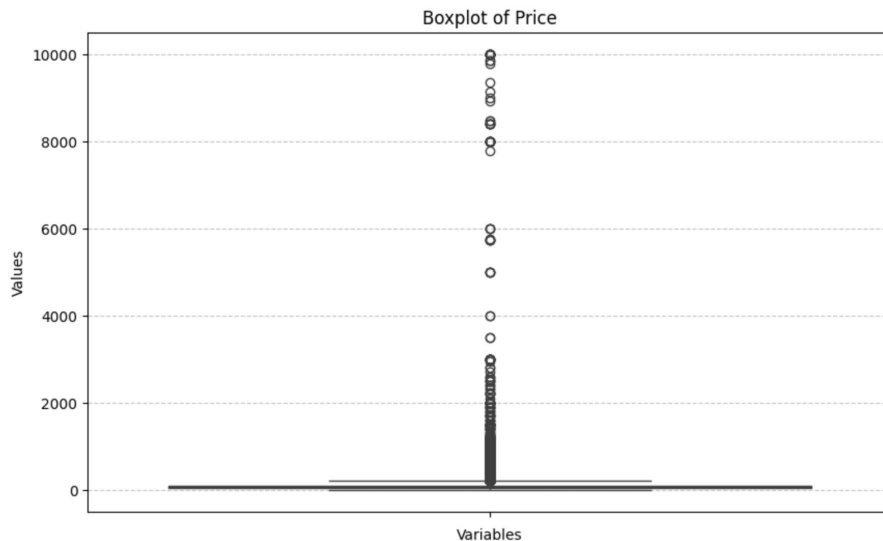
Average Price by Room Type and Neighbourhood Group



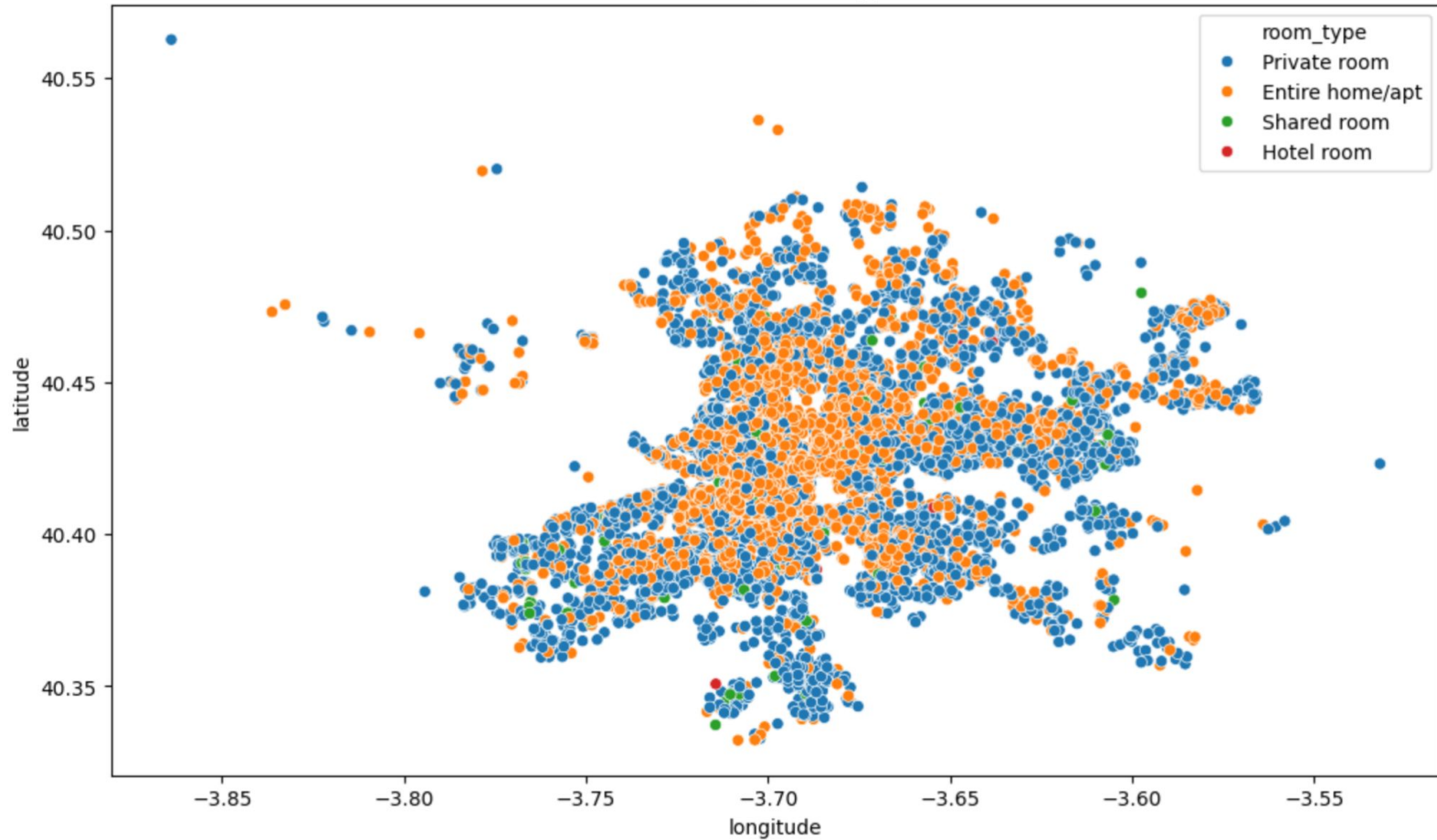
Distribution of the target feature – price (raw data)



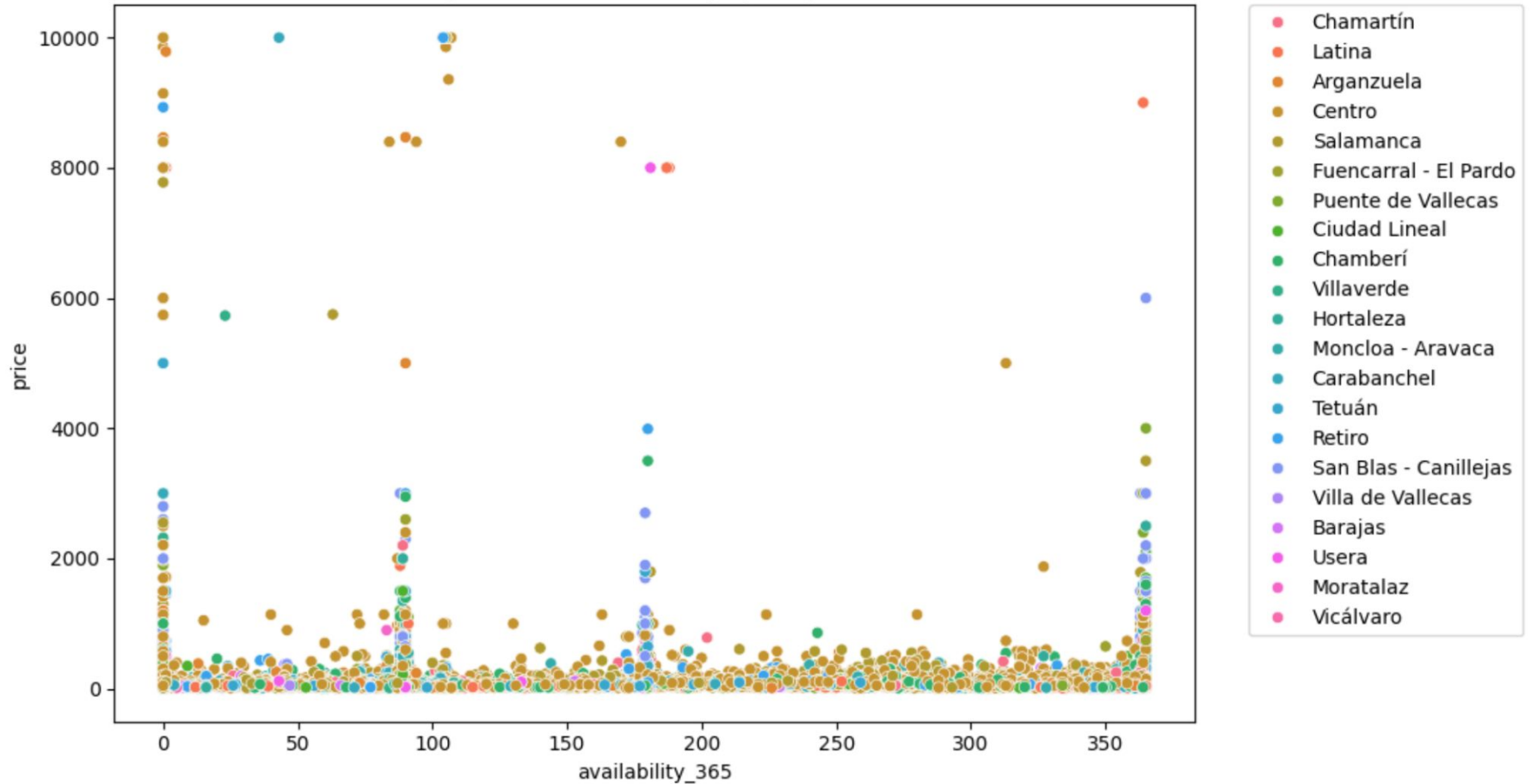
min	0.000000
25%	35.000000
50%	58.000000
75%	100.000000
max	9999.000000



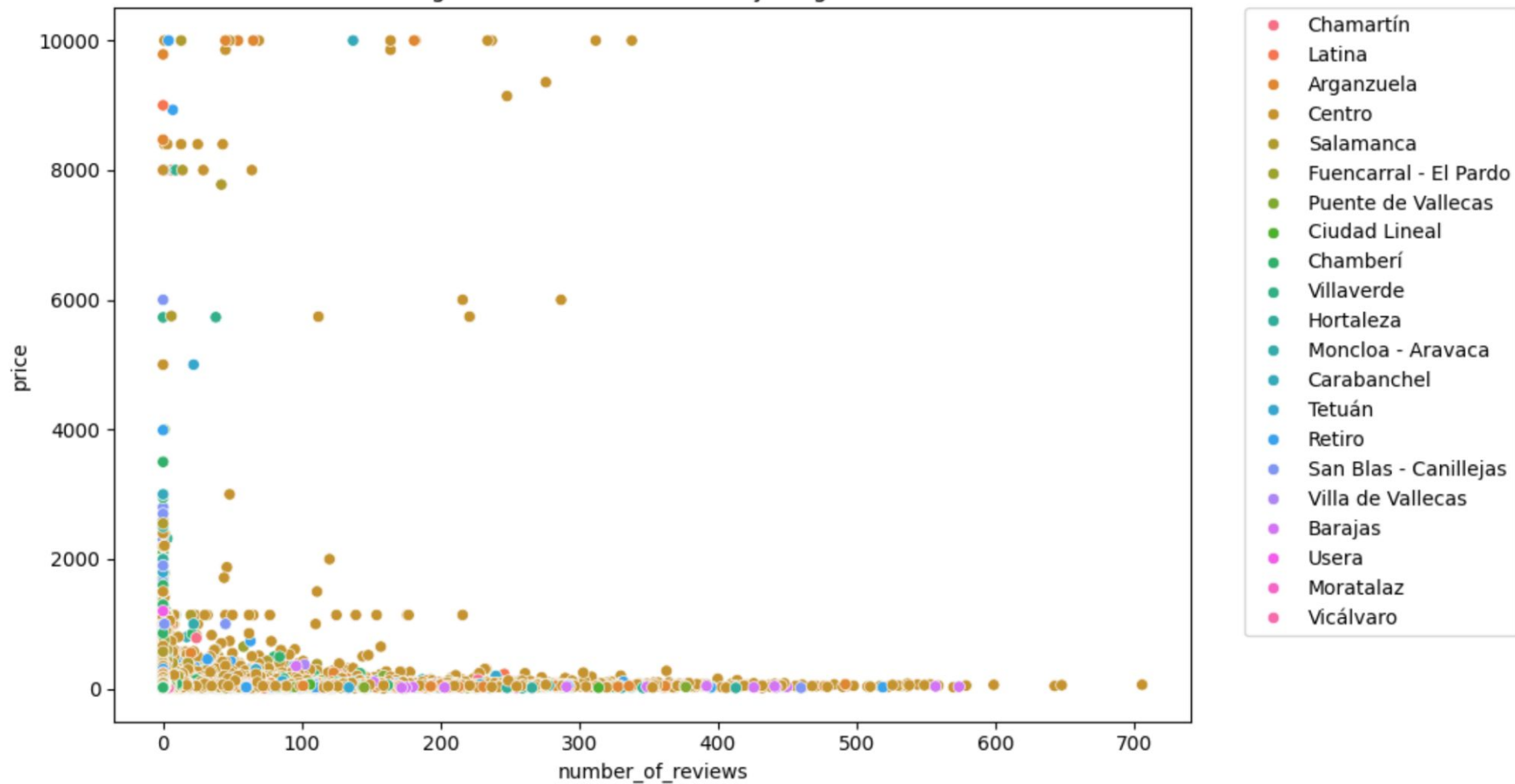
Airbnb by Room Type

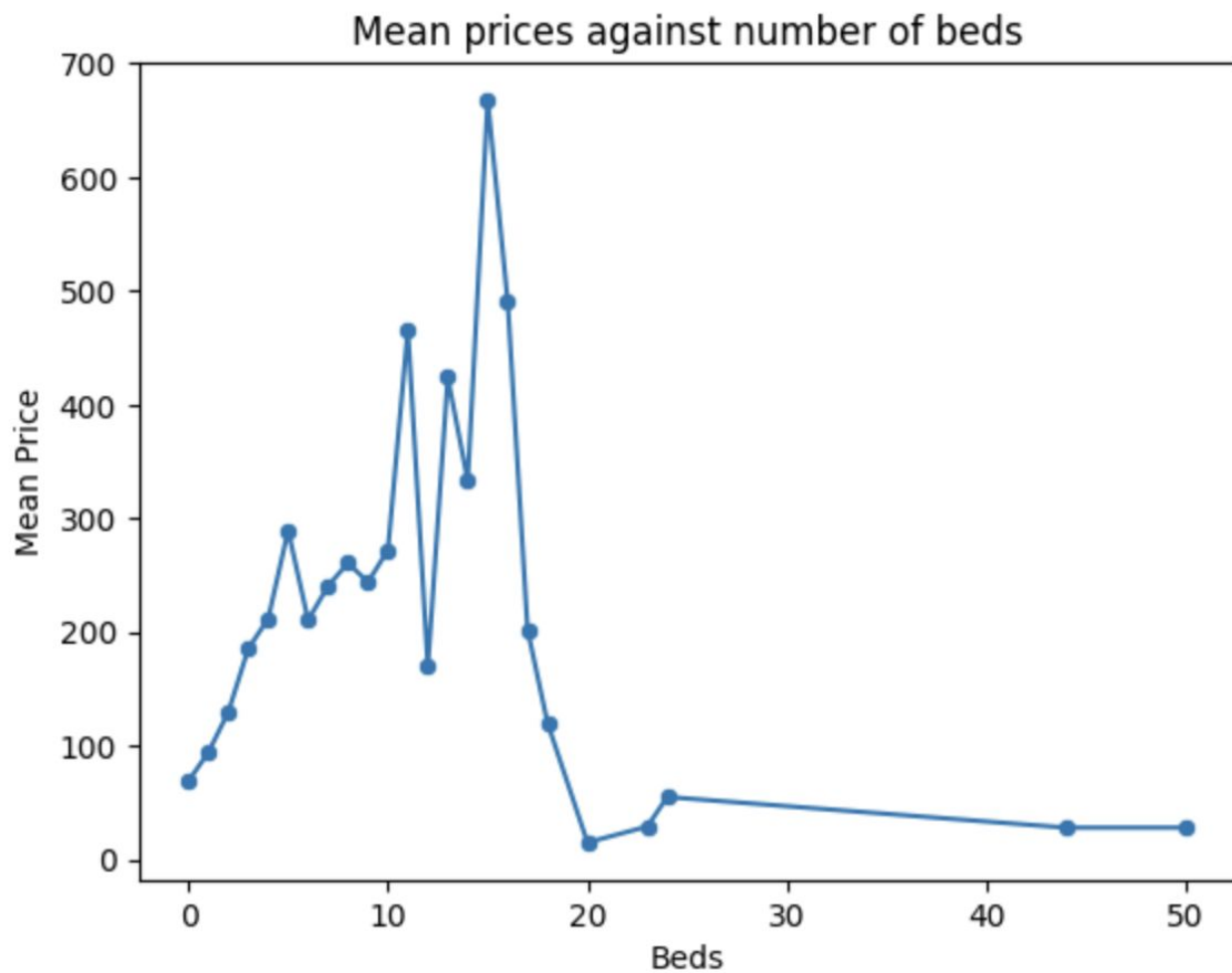


Price against room availability by Neighbourhood




Price against number of reviews by Neighbourhood



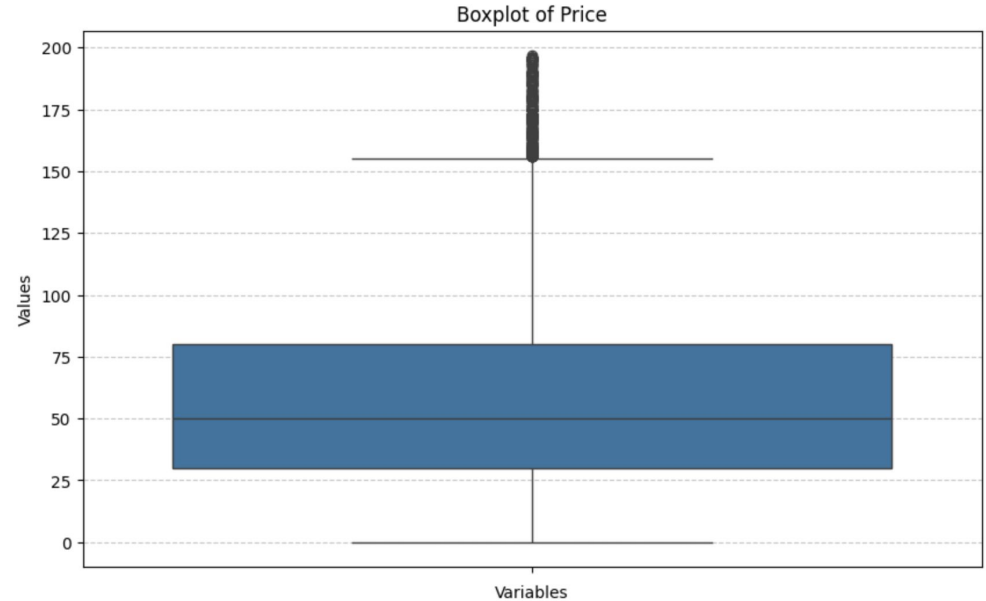
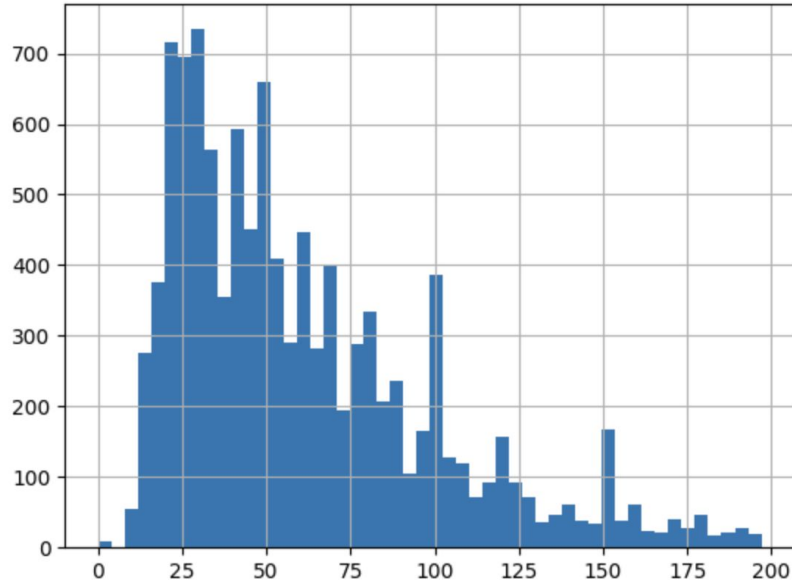


Preprocessing

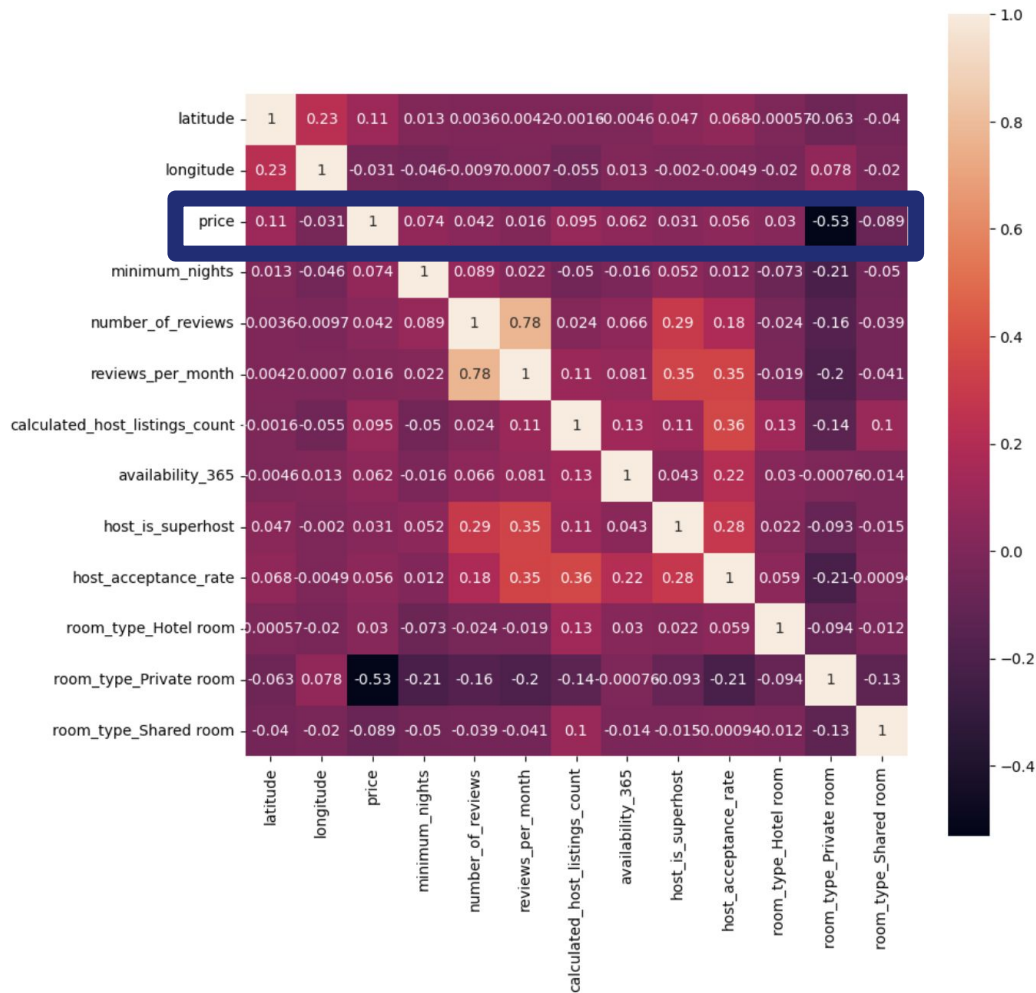
Preprocessing steps

- 1- Feature Engineering: Remove redundant columns (id, host_id, neighbourhood_group and neighbourhood [as the location given as longitude and latitude]) & Add new helpful columns (host_is_superhost, host_acceptance_rate merged from the other dataset on column "id")
 - 2- Handling missing values: fill with False & 0
 - 3- Handling outliers: Identify and remove outliers beyond the Interquartile Range threshold
 - 4- Dealing with categorical data: Apply one hot encoding
 - 5- Preparing and splitting data: Organize data into features and targets and split the data into training and test sets
 - 6- Normalization: Power Transformer to minimize the skewness of the data and improve the performance of the models
- 

Distribution of the target feature – price (after preprocessing)



Correlations



Predictive Modeling – Regression

Price Prediction

Models:

1. Linear Regression
2. Decision Tree (CART)
3. Random Forest
4. Gradient Boosting
5. Multilayer Perceptron



Price Prediction

Models:

1. Linear Regression
2. Decision Tree (CART)
3. Random Forest
4. Gradient Boosting
5. Multilayer Perceptron

Training and Evaluation (optimized with Gridsearch using 5-fold cross validation)

- performance optimization
- avoiding overfitting



Evaluation Metrics

Mean Absolute Error (MAE)

Root Mean Squared Error (RMSE)

R-Squared



Results

Linear Regression:

MAE: 23.04€
RMSE: 31.82€
 R^2 : 0.34

Decision Tree

MAE: 23.15€
RMSE: 31.97€
 R^2 : 0.34

Random Forest

MAE: 21.68€
RMSE: 30.31€
 R^2 : 0.40

Gradient Boosting

MAE: 21.71€
RMSE: 30.43€
 R^2 : 0.40

MLP

MAE: 22.52€
RMSE: 31.22€
 R^2 : 0.37

Results

Linear Regression:

MAE: 23.04€
RMSE: 31.82€
 R^2 : 0.34

Decision Tree

MAE: 23.15€
RMSE: 31.97€
 R^2 : 0.34



Random Forest

MAE: 21.68€
RMSE: 30.31€
 R^2 : 0.40

Gradient Boosting

MAE: 21.71€
RMSE: 30.43€
 R^2 : 0.40

MLP

MAE: 22.52€
RMSE: 31.22€
 R^2 : 0.37

Results

Linear Regression:

MAE: 23.04€
RMSE: 31.82€
 R^2 : 0.34

Decision Tree

MAE: 23.15€
RMSE: 31.97€
 R^2 : 0.34



Random Forest

MAE: 21.68€
RMSE: 30.31€
 R^2 : 0.40



Gradient Boosting

MAE: 21.71€
RMSE: 30.43€
 R^2 : 0.40

MLP

MAE: 22.52€
RMSE: 31.22€
 R^2 : 0.37

Results

Linear Regression:

MAE: 23.04€
RMSE: 31.82€
 R^2 : 0.34

Decision Tree

MAE: 23.15€
RMSE: 31.97€
 R^2 : 0.34

1

Random Forest

MAE: 21.68€
RMSE: 30.31€
 R^2 : 0.40

2

Gradient Boosting

MAE: 21.71€
RMSE: 30.43€
 R^2 : 0.40

3

MLP

MAE: 22.52€
RMSE: 31.22€
 R^2 : 0.37

Results

4

Linear Regression:

MAE: 23.04€
RMSE: 31.82€
 R^2 : 0.34

Decision Tree

MAE: 23.15€
RMSE: 31.97€
 R^2 : 0.34

1

Random Forest

MAE: 21.68€
RMSE: 30.31€
 R^2 : 0.40

2

Gradient Boosting

MAE: 21.71€
RMSE: 30.43€
 R^2 : 0.40

3

MLP

MAE: 22.52€
RMSE: 31.22€
 R^2 : 0.37

Results

4

Linear Regression:

MAE: 23.04€
RMSE: 31.82€
 R^2 : 0.34

5

Decision Tree

MAE: 23.15€
RMSE: 31.97€
 R^2 : 0.34

1

Random Forest

MAE: 21.68€
RMSE: 30.31€
 R^2 : 0.40

2

Gradient Boosting

MAE: 21.71€
RMSE: 30.43€
 R^2 : 0.40

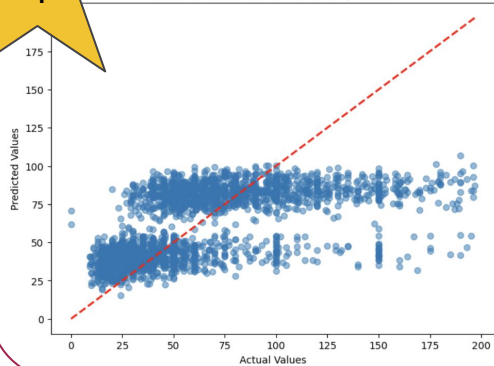
3

MLP

MAE: 22.52€
RMSE: 31.22€
 R^2 : 0.37

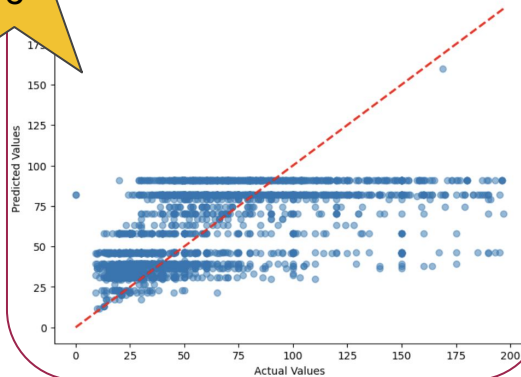
4

Actual vs Predicted for Linear Regression



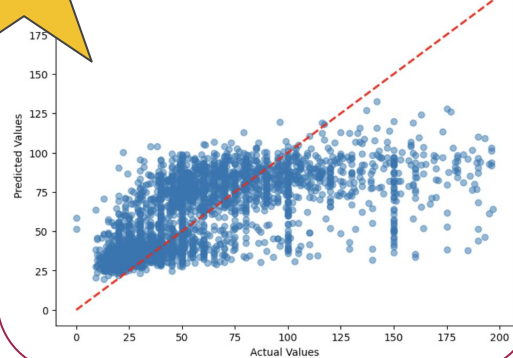
5

Actual vs Predicted for Decision Tree



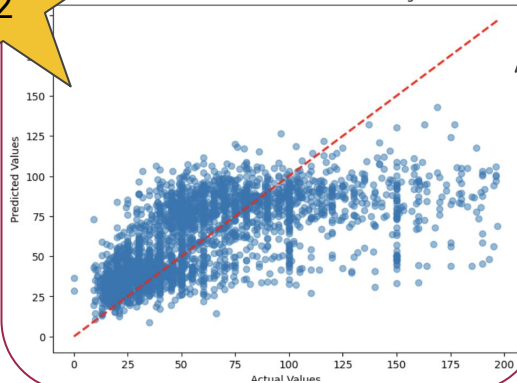
1

Actual vs Predicted for Random Forest



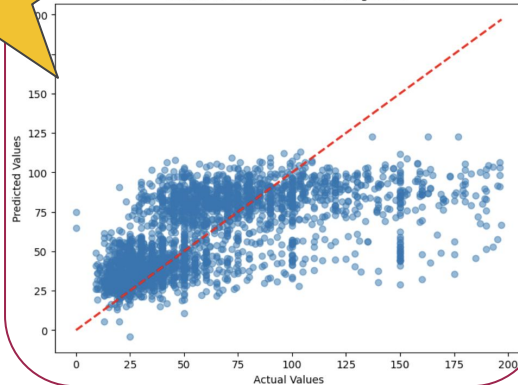
2

Actual vs Predicted for Gradient Boosting



3

Actual vs Predicted for MLP Regressor



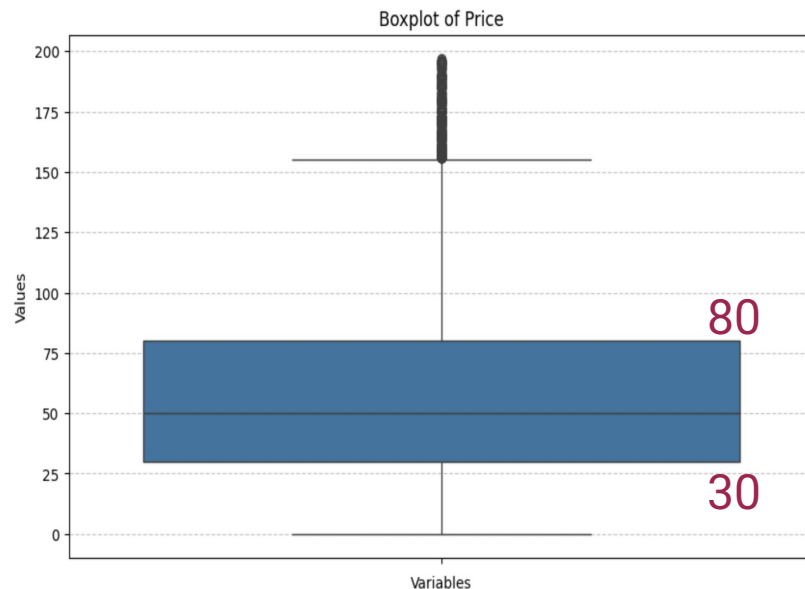


Predictive Modeling – Classification

Price Tier Prediction

Defining Price tiers according to the distribution of the data

- Low (under 30€)
- Medium (between 30€ and 80€)
- High (above 80€)

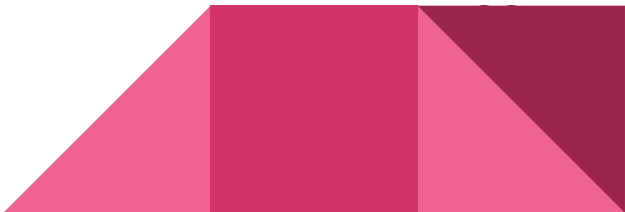


Price Tier Prediction

Defining Price tiers according to the distribution of the data

- Low (under 30€)
- Medium (between 30€ and 80€)
- High (above 80€)

Models:

- Naive Bayes
 - Decision Tree
 - Neural Network (MLP classifier)
 - Support Vector Machine (SVM)
 - k-nearest neighbours (kNN)
- 

Price Tier Prediction

Defining Price tiers according to the distribution of the data

- Low (under 30€)
- Medium (between 30€ and 80€)
- High (above 80€)

Models:

- Naive Bayes
- Decision Tree
- Neural Network (MLP classifier)
- Support Vector Machine (SVM)
- k-nearest neighbours (kNN)

Used grid search with cross validation for classification models too



Results – Evaluation Metric Accuracy



Naive Bayes
Accuracy: 0.48



Decision Tree
Accuracy: 0.59



Neural Network
Accuracy: 0.61



SVM
Accuracy: 0.60

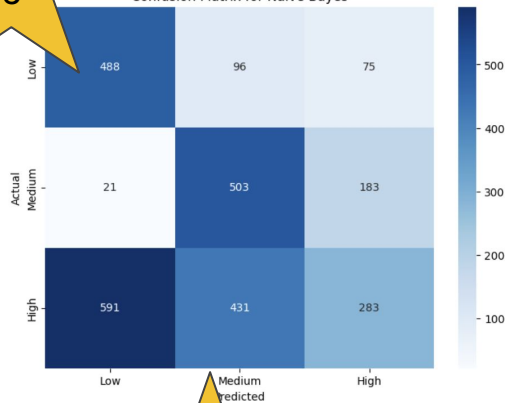


kNN
Accuracy: 0.59



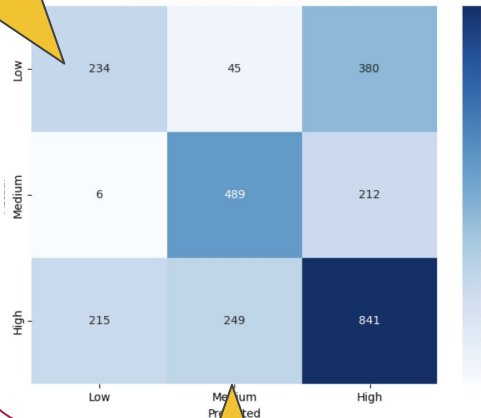
5

Confusion Matrix for Naive Bayes



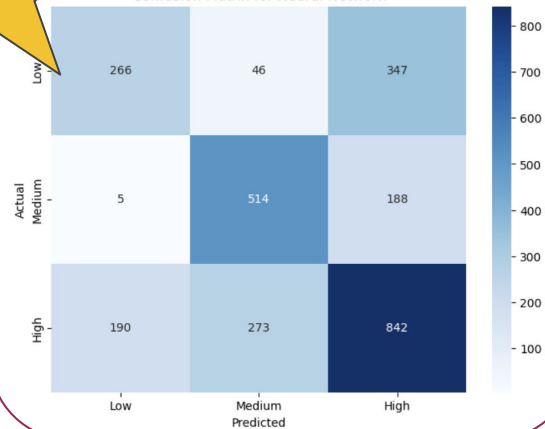
4

Confusion Matrix for Decision Tree



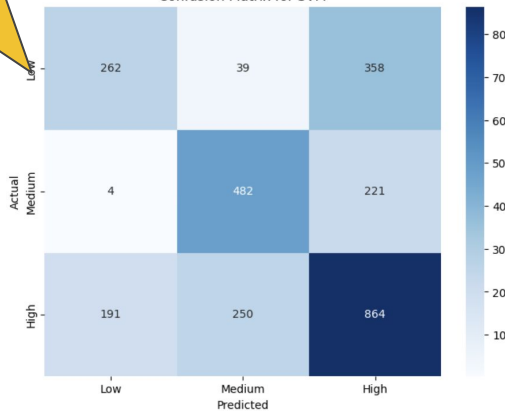
1

Confusion Matrix for Neural Network



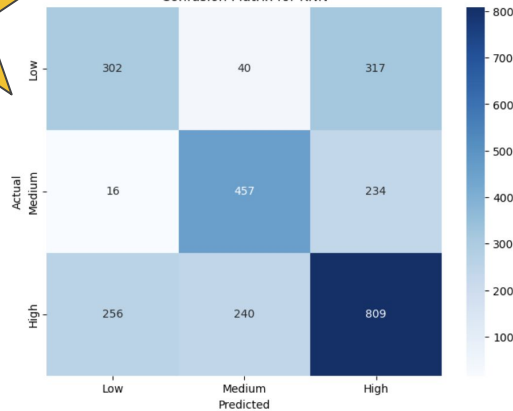
2

Confusion Matrix for SVM

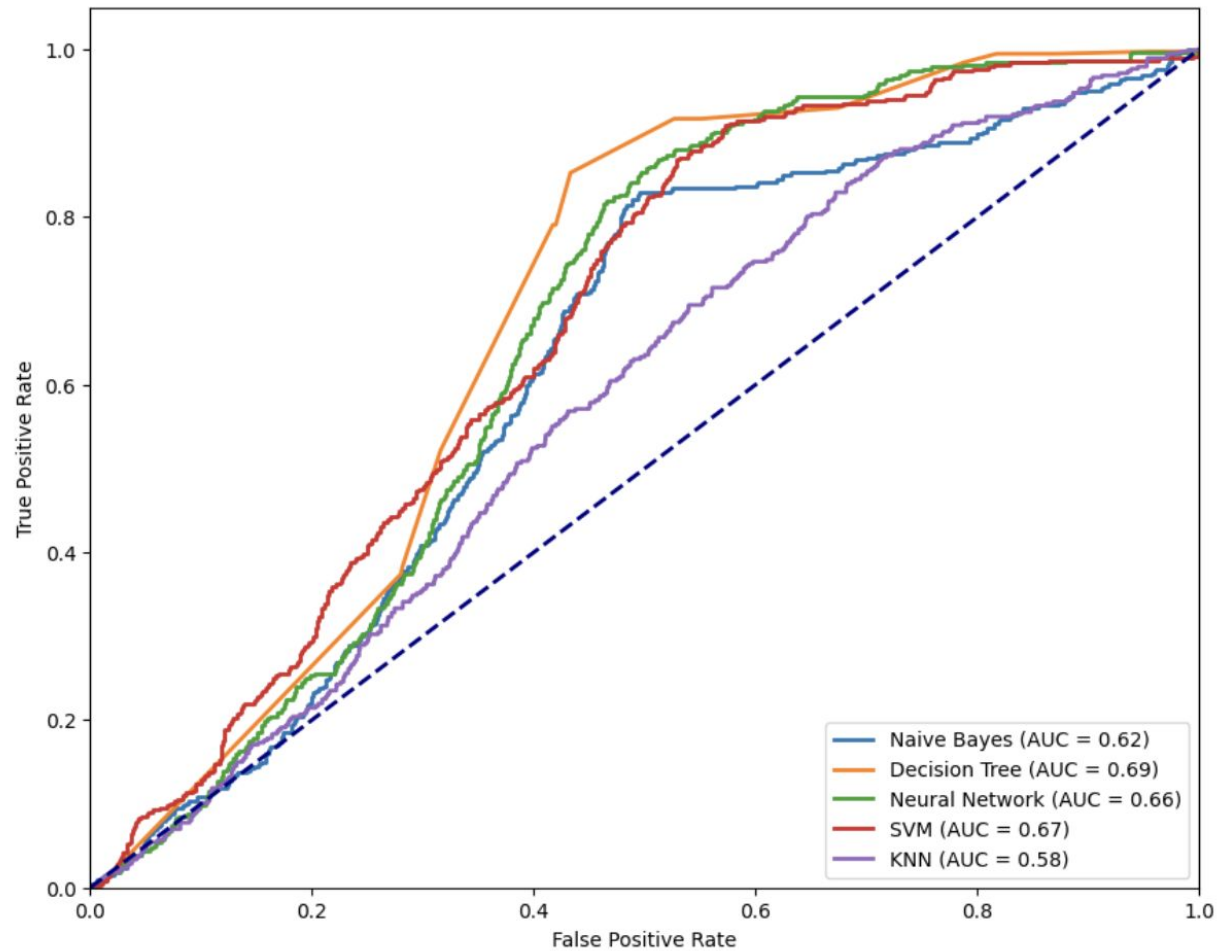


3

Confusion Matrix for KNN



ROC Curve for All Models



Sources

- <https://www.kaggle.com/datasets/rusiano/madrid-airbnb-data/data>
- <https://insideairbnb.com/get-the-data/>
- <https://insideairbnb.com/madrid/>

