

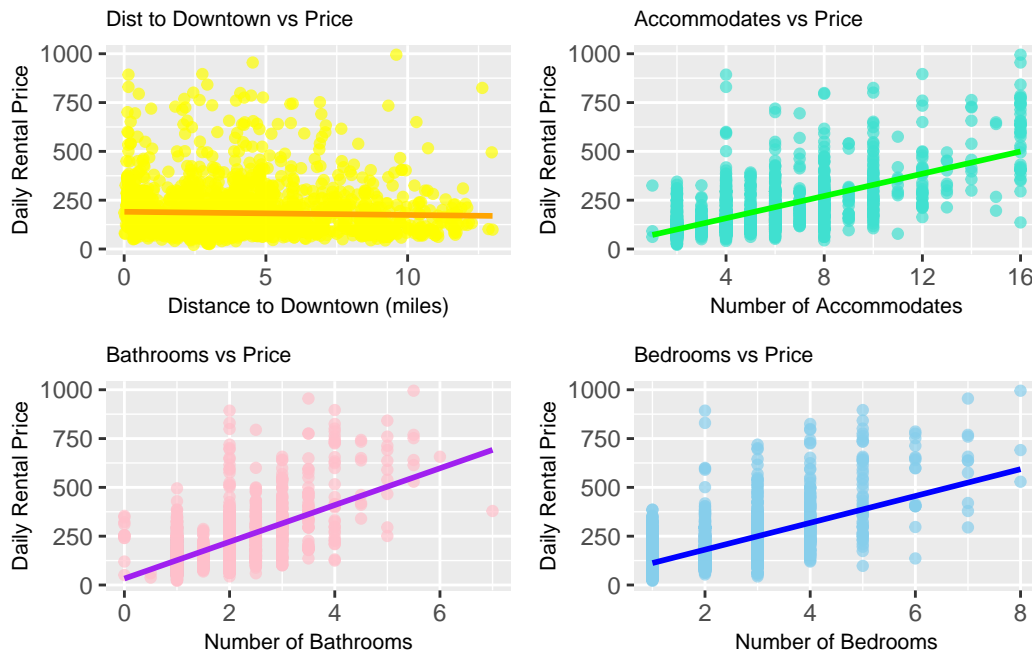
Airbnb Price Regression Modeling - Report for Airbnb executives

Introduction

The objective of this project is to generate a model that allows setting prices for Airbnb ads in Asheville, North Carolina. To achieve this goal, we have a database of various Airbnb rental listings in the city, which contains detailed information about the listings, such as price, number of rooms, amenities, number of bathrooms, property location, etc. This database was extracted in June 2023 by the company Inside Airbnb and contains information about 3,239 rental listings in Asheville, NC. Therefore, the conclusions drawn are based on this time period.

To develop a model that enables us to set prices for Airbnb ads in Asheville, we will use a linear regression approach. Linear regression is a statistical technique used to model the relationship between a dependent variable and one or more independent variables. In this context, the dependent variable would be the price of Airbnb ads, while independent variables may include features such as the number of rooms, amenities, number of bathrooms, and property location, among others.

In the following scatter plots, we observe positive linear relationships between '*Accommodates*', '*Bathrooms*' and '*Bedrooms*' and the dependent variable '*Price*'. In contrast, '*Dist to Downtown*' presents a negative relationship, indicating that homes closer to downtown tend to be more expensive. These patterns illustrate the goal of a regression model: to identify a line that minimizes the sum of squared errors between predicted and actual values



Methods

After data cleaning, variable consideration, and iterative model refinement, we developed a linear regression model for Airbnb pricing. Missing data were addressed through decision rules, resulting in a final database of 2,243 rows from the original 3,239.

Results

The final model obtained is the following:

$$\ln(\text{price}) = 3.74 + 0.13*(\text{is entire home apt}) + 0.08*(\text{n}^\circ \text{ bedrooms}) + 0.14*(\text{n}^\circ \text{ bathrooms}) - 0.05*(\text{dist to downtown}) + 0.19*(\text{has entertainment amenities}) + 0.44*(\text{has climate control}) + 0.07*(\text{n}^\circ \text{ accommodates})$$

The coefficient accompanying a variable in a linear regression is the amount by which the response variable (dependent variable) changes for each one-unit change in the predictor variable (independent variable), while keeping all other variables constant. For example, one would expect that, on average, for each additional room, the logarithm of the price would increase by 0.082.

If we were in the situation of defining the price of a new property, we should evaluate the variables of that property in the model. For example, let's assume we have a rental in Asheville, which is an entire home apartment (coded as 1), has 3 bedrooms, 1 bathroom, is located 3 miles from downtown, lacks a TV, Netflix, or any entertainment amenities (codes as 0), has air conditioning (coded as 1), and can accommodate 5 guests. Finally, we could predict a price as follows:

$$\ln(\text{price}) = 3.74 + 0.13*(1) + 0.08*(3) + 0.14*(1) - 0.05(3) + 0.19*(0) + 0.44*(1) + 0.07*(5) = 4.893$$

$$\text{price} = e^{4.893} = 133$$

So, a suggested price for this rental is \$133 per night

A measure obtained from our model is an R^2 of 0.52, which means that approximately 52% of the variability in the rental prices in Assville can be explained by the variables included in our model.

Conclusion

With the database, we were able to create a model to determine the rental price in Assville, as seen in the previous example. However, it is important to consider that other factors may be explaining the rental price, and here we are bound by our original database. With our model, we were only able to explain 52% of the variance, However, it can still be useful as a guide if our goal is to determine the listing price for a new rental on Airbnb.

Airbnb Price Regression Modeling - Report for Data Science Team

Introduction

The objective of this project is to generate a linear regression model that allows setting prices for Airbnb ads in Asheville, North Carolina.

Dataset

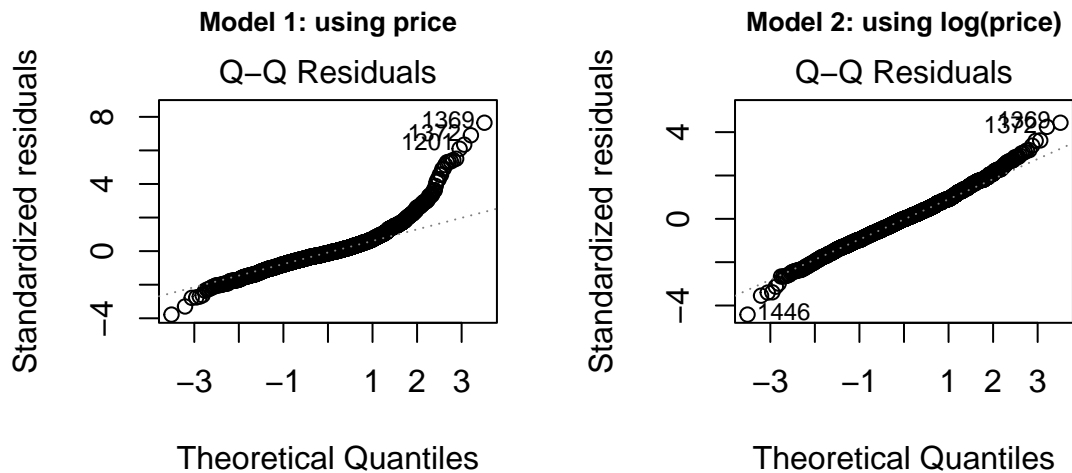
We have a database of various Airbnb rental listings in the city, obtained through web scraping, which contains detailed information about the listing. The original database includes **75 variables**. This dataset was extracted in June 2023 by the company Inside Airbnb and contains information about **3,239** rental listings in Asheville, North Carolina. Therefore, the conclusions drawn are based on this time period. The database and dictionary can be downloaded from [listings.csv](#) and [Inside Airbnb Data Dictionary](#)

Data cleaning

- First, the database was cleaned to be used for a linear regression model. Variables that did not contribute information to the model were excluded, such as: *“id”, “listing_url”, “scrape_id”, “last_scraped”, “source”, “name”, “description”, “neighborhood_overview”, “picture_url”, “host_id”, etc.*
- Secondly, certain transformations were performed on some variables to handle them within the model, such as *“host_response_time”, “host_response_rate”, “host_acceptance_rate”, “host_is_superhost”, “neighbourhood”, “latitude”, “longitude”, “property_type”, “room_type”, “bathrooms_text”, “amenities” and “price”.*
- It is worth mentioning that the **‘amenities’** variable considered combinations of 2,115 different features, so the following transformation was performed: new binary variables were generated to group the most frequent amenities. For example, **‘security’** considers *‘smoke and fire alarms’, ‘fire extinguisher’, and ‘first aid kit’*. **‘kitchen_amenities’** includes *‘dishes and silverware’, ‘microwave’, ‘kitchen’, ‘refrigerator’, ‘cooking basics’, ‘coffee’, ‘freezer’, ‘wine glasses’, ‘oven’, ‘toaster’, and ‘dining table’.*
- Certain variables were handled for missing values, such as in the case of the ‘beds’ variable. For records with NA in the ‘beds’ field, we observed that the majority represent campsites; therefore, we will consider ‘beds = 0’ for these instances. In other cases, rows with missing values were simply removed. The final model, as a result, incorporated 2,243 out of the original 3,239 rows.

Methods

- After data cleaning, the following explanatory variables were considered as a foundation:
 - Room type (Transformed into the binary variable ‘**entire_home_aprt**’)
 - Number of bedrooms.
 - Distance to downtown.
 - Amenities, which, as mentioned earlier, were categorized into binary variables.
- Additional variables were considered for the model, focusing on factors relevant to new hosts for price determination. Variables like “*review scores*” or “*superhost status*” were omitted, given the emphasis on new, unrated listings or hosts new to the platform.
- Also, interactions were not considered in the model to enhance its interpretability.
- Examined correlations using plots; removed highly correlated variables. For instance, “accommodates” strongly correlates with “n° of beds”; only “accommodates” remained.
- The binary amenities variables “*security*”, “*essentials*”, “*bedroom amenities*”, “*convenience features*”, “*kitchen amenities*” and “*toiletries*” proved to be less significant in predicting the price, and thus, they were removed.
- In the final model, logarithmic transformation of the price was performed, as the QQ plot exhibited a systematic deviation between the theoretical distribution and the actual distribution of the data. This could result from both the non-normality of the residuals and heteroscedasticity.



After several iterations, the finally selected model was:

$$\ln(\text{price}) = 3.74 + 0.13*(\text{is entire home apt}) + 0.08*(\text{n}^\circ \text{ bedrooms}) + 0.14*(\text{n}^\circ \text{ bathrooms}) - 0.05*(\text{dist to downtown}) + 0.19*(\text{has entertainment amenities}) + 0.44*(\text{has climate control}) + 0.07*(\text{n}^\circ \text{ accommodates})$$

The details of the coefficients and statistics obtained with this model are as follows:

Predictors	log_price Estimates	CI	p
(Intercept)	3.74	3.57 – 3.92	< 0.001
entire home aptTRUE	0.13	0.01 – 0.25	0.035
bedrooms	0.08	0.05 – 0.12	< 0.001
bathrooms numeric	0.14	0.11 – 0.18	< 0.001
dist to dt	-0.05	-0.05 – -0.04	< 0.001
entertainmentTRUE	0.19	0.12 – 0.26	< 0.001
climate controlTRUE	0.44	0.30 – 0.58	< 0.001
accommodates	0.07	0.06 – 0.08	< 0.001
Observations	2243		
R ² / R ² adjusted	0.524 / 0.523		

F-statistic: 351.5 on 7 and 2235 DF, p-value: < 2.2e-16

Mean Squared Error: 0.1494199

- In general, the model explains approximately 52% of the variability in the logarithm of the price. The F-statistic and p-value suggest that the model is statistically significant. Additionally, the coefficients have small p-values, indicating that they are statistically significant.
- The Mean Squared Error (MSE) of 0.15 suggests that, on average, the predictions of the regression model exhibit a relatively low mean squared error, indicating a reasonable level of accuracy. In the context of model evaluation, this MSE is considered relatively good, signifying that the model performs well in minimizing prediction errors.
- If we calculate the Variance Inflation Factor (VIF) for our model, we obtain the following:

entire_home_apt	bedrooms	bathrooms_numeric	dist_to_dt
1.019634	5.926032	3.363982	1.080871
entertainment	climate_control	accommodates	
1.067202	1.043949	5.006517	

- VIF measures the degree to which the variance of an estimated regression coefficient is inflated due to multicollinearity with other predictors. VIF values close to 1 indicate minimal multicollinearity, while values exceeding 5 suggest a moderate to high level. In this context, variables such as “*entire home apt*”, “*dist to dt*”, “*entertainment*”, and “*climate_control*” exhibit minimal multicollinearity with VIFs less than 5, implying stable regression coefficients. Conversely, variables like “*bedrooms*” and “*accommodates*” show moderate multicollinearity. However, the decision was made to retain both variables as both are theoretically expected to influence the rental price.

Conclusion

- Considering the task of helping new hosts set prices for Airbnb listings in Asheville, NC, falls more into the category of a **Prediction Problem**. The primary goal is to build a model that can accurately predict or generate prices for Airbnb listings based on various factors. We want to provide hosts with a tool that can predict the optimal price for their listings rather than drawing in-depth inferences about the underlying relationships between variables. However, this model is also useful for explaining some of the variables that have a significant impact on the determination of the price, even though its primary goal is predictive.
- The price is transformed using a logarithm due to a systematic deviation observed in the QQ plot between the theoretical distribution and the actual distribution of the data.
- Examined correlations using plots; removed highly correlated variables. For instance, “*accommodates*” strongly correlates with “*n° of beds*”; only “*accommodates*” remained. Finally, an acceptable VIF was obtained where the values of your variables were close to 5 or lower.
- A key metric for evaluating the model’s performance is the Mean Squared Error (MSE). In summary, the model demonstrates a satisfactory level of accuracy, evidenced by a relatively low MSE. The successful minimization of prediction errors highlights commendable performance, aligning with expectations. According to R², the model explains approximately 52% of the variability in the logarithm of the price.
- Generally, we can consider this to be a valid model for predicting or determining the price of a home in Asheville, given its low multicollinearity, statistical significance of the model, significance of its variables, and the associated MSE (mean squared error). However, It is important to consider that other factors may be explaining the rental price, and here we are bound by our original database. With our model, we were only able to explain 52% of the variance, However, it can still be useful as a guide if our goal is to determine the listing price for a new rental on Airbnb.