

Statistical Programming
for Data Science:
An investigation on the Airbnb price per night in Amsterdam

XXX

29 May 2024

Contents

1	Dataset	3
1.1	Introduction	3
1.2	Description of dataset	3
1.3	Three proposed research questions	4
1.3.1	Research question 1	4
1.3.2	Research question 2	4
1.3.3	Research question 3	4
2	Data Import and Cleaning	5
3	Data Analysis	7
3.1	Research question 1	7
3.2	Research question 2	8
3.3	Research question 3	9
4	Conclusion and Discussion	11
5	Appendix: Individual Assignment Coversheet	13
6	Appendix: R Environment	14

1 Dataset

1.1 Introduction

As a sharing economy product, Airbnb experiences a rapid development in recent decades. Airbnb is an alternative hotel business that shares the accommodation with others, short period or long period. There are lots of research on the price per night that Airbnb costs. In Chen Yong and Xie [2017], a wide array of utility-bearing attributes of Airbnb listings and the effects of these attributes on consumers' valuation in United States are measured. It provides a comprehensive study on the pricing of Airbnb listed properties and the results explain how the factors, i.e., listing functionality, attributes of hosts, customers reviews and market conditions affect the price. Another research in Cai et al. [2019] focuses on the market of Hong Kong. Five groups' variables were collected, i.e., listing attributes, host attributes, rental policies, listing reputation, and listing location to investigate the determinant of Airbnb price. Some use ordinary least square regression with geographically-weighted, which is introduced in Voltes-Dorta and Sánchez-Medina [2020], to study the factors that affect the price for different room types, i.e., entire room or private room.

In this project, we are going to build a suitable model that can explain the relationship between the rental price per night posted on Airbnb of apartment in the Netherlands, mainly in the city of Amsterdam, and several characteristics related to the apartment. Specially, to find the determinants of the price from the room features, e.g., number of bathrooms, bedrooms; host response rate, and the ratings received from the customers.

1.2 Description of dataset

The dataset to be analyzed is collected from <https://data.world/cannata/gairbnb> and is named "AirBNB.csv". In the raw dataset, there are 7833 observations on 41 variables. The selected variables to be analyzed are price, accommodates, bathrooms, bedrooms, room_type, host_response_rate and review_scores_rating. The description and type of each variable are listed as follows.

price: continuous variable, the price per night posted on the website.

accommodates: discrete variable, the number of guests that the property can accept.

bathrooms: continuous variable, the number of bathrooms the property has.

bedrooms: discrete variable, the number of bedrooms the property has.

room_type: nominal variable, the feature of the shared property, and there are three types, "Entire home/apt", "Private room" and "Shared room". host_response_rate: continuous variable, indicating the response frequency of the host when receiving message.

review_scores_rating: discrete variable, indicating the reputation of the shared property.

The screenshot of the dataset is displayed in Figure 1.

As for the data wrangling, we propose to filter out the observations related to property type "apartment" first, then select the necessary variables. At last removing the observations with missing values and changing the format or type of some variables.

host_id	host_name	host_since_year	host_since_anniversary	Customer Since	Age in years	id	neighbourhood_cleaned	city	city_translated	state	state_translated	zipcode	country	latitude	longitude	property_type	room_type	accommodates	bathrooms	bedrooms	beds	bed_type
1662	Chloé	2008	8/11	8/11/08	8.93	304658	Westerpark	Amsterdam	Amsterdam	North Holland	North Holland	1053	Netherlands	52.37320064	4.868460823	Apartment	Entire home/apt	4	2	2	2	Real Bed
3159	Daniel	2008	8/24	9/24/08	8.80	2818	Oostelijk Havengebied - Indische Buurt	Amsterdam	Amsterdam	North Holland	North Holland		Netherlands	52.36575451	4.941418235	Apartment	Private room	2	1	1	2	Real Bed
3718	Britta	2008	10/19	10/19/08	8.74	103026	De Baarsjes - Oud-West	Amsterdam	Amsterdam	Noord-Holland	North Holland	1053	Netherlands	52.36938767	4.866872319	Apartment	Entire home/apt	4	1	1	1	Real Bed
4716	Stefan	2008	11/20	11/20/08	8.62	550017	Centrum-Oost	Amsterdam	Amsterdam	North Holland	North Holland	1017	Netherlands	52.36190508	4.888503037	Apartment	Entire home/apt	2	1	1	1	Real Bed
5271	Tyler	2008	12/17	12/17/08	8.57	4728389	Centrum-West	Amsterdam	Amsterdam	Noord-Holland	North Holland	1016 AM	Netherlands	52.37153345	4.887057291	Apartment	Entire home/apt	6	1	2	2	Real Bed
5271	Tyler	2008	12/17	12/17/08	8.57	5009954	Centrum-West	Amsterdam	Amsterdam	NH	North Holland	1016 AM	Netherlands	52.37153392	4.886072287	Apartment	Private room	4	1	1	1	Real Bed
5271	Tyler	2008	12/17	12/17/08	8.57	6181918	Centrum-West	Amsterdam	Amsterdam	Noord-Holland	North Holland	1016 AM	Netherlands	52.3704458	4.880966479	Apartment	Private room	2	1	1	1	Futon
5968	Ramona	2009	1/4	1/4/09	8.53	2774524	Zuid	Amsterdam	Amsterdam	North Holland	North Holland	1071 VV	Netherlands	52.35564811	4.893534819	House	Private room	2	1	1	1	Real Bed
9616	Laura	2009	3/9	3/9/09	8.35	23651	De Pijp - Rivierenbuurt	Amsterdam	Amsterdam	North Holland	North Holland	1078	Netherlands	52.34591098	4.891982805	Apartment	Private room	3	1	1	1	Real Bed
14589	Rutger	2009	4/23	4/23/09	8.23	738545	Centrum-West	Amsterdam	Amsterdam	North Holland	North Holland	1015	Netherlands	52.37935439	4.883276386	House	Entire home/apt	2	1	1	1	Real Bed
15618	Shelly	2009	5/2	5/2/09	8.20	51969	De Pijp - Rivierenbuurt	De Pijp	De Pijp	North Holland	North Holland	1072	Netherlands	52.35748276	4.887099693	Apartment	Entire home/apt	3	1.5	2	2	Real Bed
21669	Mark	2009	6/15	6/15/09	8.08	8061	De Baarsjes - Oud-West	Amsterdam	Amsterdam	Noord-Holland	North Holland	1056 TM	Netherlands	52.371207	4.857291017	Apartment	Entire home/apt	3	1	2	2	Real Bed
28919	Hugo	2009	7/22	7/22/09	7.98	98558	Centrum-Oost	Amsterdam	Amsterdam	North Holland	North Holland	1011 JX	Netherlands	52.36959599	4.890890308	Apartment	Entire home/apt	2	1	1	1	Real Bed
32366	Sabine & Sander	2009	8/18	8/18/09	7.91	9693	Centrum-West	Amsterdam	Amsterdam	North Holland	North Holland	1013	Netherlands	52.37891663	4.892703442	Apartment	Entire home/apt	3	1.5	1	1	Real Bed
36701	Levin	2009	9/7	9/7/09	7.85	1832819	Box en Lommer	Amsterdam	Amsterdam	North Holland	North Holland	1059DP	Netherlands	52.38141023	4.892742701	Apartment	Entire home/apt	2	1	1	1	Real Bed
42212	Miguel	2009	9/29	9/29/09	7.79	280106	Centrum-West	Amsterdam	Amsterdam	North Holland	North Holland	1013	Netherlands	52.38209988	4.886143665	Apartment	Entire home/apt	4	1	0	2	Real Bed
42212	Miguel	2009	9/29	9/29/09	7.79	3527892	Centrum-West	Amsterdam	Amsterdam	North Holland	North Holland	1013HE	Netherlands	52.38147315	4.886808875	Loft	Shared room	1	1	1	1	Real Bed
42725	Marco	2009	10/1	10/1/09	7.79	933385	De Baarsjes - Oud-West	Amsterdam	Amsterdam	North Holland	North Holland	1053	Netherlands	52.36761407	4.8868895471	Apartment	Private room	2	1	1	2	Real Bed
46431	Jennifer & Michiel	2009	10/17	10/17/09	7.74	1182306	Zuid	Amsterdam	Amsterdam	North Holland	North Holland	1059	Netherlands	52.34658737	4.84919711	Apartment	Private room	2	1	1	1	Real Bed
47817	Gert	2009	10/21	10/21/09	7.73	3047091	Watergraafmeer	Amsterdam	Amsterdam	North Holland	North Holland	1087 AM	Netherlands	52.35342049	4.92442006	Apartment	Entire home/apt	2	1	1	1	Real Bed
50517	Sanne	2009	11/1	11/1/09	7.70	4003922	Centrum-Oost	Amsterdam	Amsterdam	North Holland	North Holland	1018	Netherlands	52.36928364	4.909938668	Apartment	Entire home/apt	4	1	1	2	Real Bed
56142	Joan	2009	11/20	11/20/09	7.65	1003965	De Baarsjes - Oud-West	Amsterdam	Amsterdam	North Holland	North Holland	1053 LB	Netherlands	52.36975078	4.871953549	Apartment	Entire home/apt	4	1	1	2	Real Bed
56142	Joan	2009	11/20	11/20/09	7.65	25628	Centrum-West	Amsterdam	Amsterdam	North Holland	North Holland	1016	Netherlands	52.3731584	4.882339196	Apartment	Entire home/apt	3	1	1	1	Real Bed
59058	Marius	2009	12/1	12/1/09	7.62	75963	Stoetvaart	Amsterdam	Amsterdam	North Holland	North Holland	1056	Netherlands	52.36520971	4.836338494	Apartment	Private room	4	1.5	1	2	Real Bed
89297	Jan	2009	12/2	12/2/09	7.62	15061	Westerpark	Amsterdam	Amsterdam	North Holland	North Holland	1052	Netherlands	52.38268456	4.876129664	Apartment	Private room	4		1	2	Real Bed
95844	Alex	2009	12/2	12/2/09	7.62	20168	Centrum-Oost	Amsterdam	Amsterdam	North Holland	North Holland	1017	Netherlands	52.365088703	4.893541008	House	Private room	2	1	1	1	Real Bed

Figure 1: Screenshot of the dataset

1.3 Three proposed research questions

1.3.1 Research question 1

The first proposed question: “Are the average prices per night the same for different room type?”

1.3.2 Research question 2

The second proposed question: “Is there a linear relationship between price per night and accommodates, if there is, how to interpret the relationship?”

1.3.3 Research question 3

The third proposed question: “Is there a multiple linear relationship between price per night and a group of predictors. i.e., accommodates, bathrooms, bedrooms, room_type, host_response_rate and review_scores_rating?”

2 Data Import and Cleaning

In this section, we are going to prepare a well-structured dataset for subsequent analytics, which contains three steps, import, cleaning and tidy. Since the raw data is in “csv” format. We use the function “read.csv” to import the data. Then we filter out the observations related to property type “Apartment”.

```
# import dataset and filter out apartment
tb<-read.csv("AirBnb.csv") %>%
  filter(property_type=="Apartment")
```

There are lots of variables that we don't use, we select some specified columns by name to make the dataset more concise.

```
# select the necessary variables
tb.selected<-tb %>%
  select(price,accommodates,bathrooms,bedrooms,room_type,
         host_response_rate,review_scores_rating)

str(tb.selected)
```

```
## 'data.frame':    6265 obs. of  7 variables:
## $ price          : chr  " $130 " " $59 " " $95 " " $100 " ...
## $ accommodates   : int   4 2 4 2 6 4 2 3 3 3 ...
## $ bathrooms      : num   2 1 1 1 1 1 1 1 1.5 1 ...
## $ bedrooms       : int   2 1 1 1 2 1 1 1 2 2 ...
## $ room_type      : chr   "Entire home/apt" "Private room" "Entire home/apt" "Entire
## $ host_response_rate : chr  "0.8" "1" "1" "1" ...
## $ review_scores_rating: int   98 97 92 97 100 NA 95 96 95 100 ...
```

Given a view of the selected variables, we found that there are some variables in wrong type. For example, the variable price and host_response_rate should be numerical but they were in character type.

```
# change the type of some variables
tb.selected$price<-parse_number(tb.selected$price)
tb.selected$host_response_rate<-as.numeric(tb.selected$host_response_rate)
```

The last step is to tidy the missing values. Since the sample size is large, we are going to remove the observations with missing value in any column.

```
# remove the observations having missing values
tb.clean<-tb.selected %>% na.omit()
```

Here we take a view of the prepared dataset with dimension 4568 * 7. It is demonstrated that each variable have its own column, each observation has its own row and each value has its own cell. The three rules of tidy data is satisfied.

```
kable(head(tb.clean))
```

	price	accommodates	bathrooms	bedrooms	room_type	host_response_rate	review_scores_rating
1	130	4	2	2	Entire home/apt	0.80	98
2	59	2	1	1	Private room	1.00	97
3	95	4	1	1	Entire home/apt	1.00	92
4	100	2	1	1	Entire home/apt	1.00	97
5	250	6	1	2	Entire home/apt	0.89	100
7	115	2	1	1	Private room	0.89	95

3 Data Analysis

3.1 Research question 1

The objective is to analyze whether the Airbnb posted price per night of apartment are different among different room types. Since the room type is a categorical variable, a one-way ANOVA approach is suitable. Before conducting any statistical analysis, a descriptive summary for the price is tabulated in Table 1. It is found that the average price 129 for the Entire home or entire apartment is much higher than that for private room 69 and shared room 56. Meanwhile, the variability of the price for entire home is also the highest.

Table 2: Summary statistics for price per night for different room types

room_type	average	SD
Entire home/apt	128.63	60.62
Private room	68.76	28.86
Shared room	55.96	28.34

The boxplot displayed in Figure 2 gives a direct comparison of the distribution of price for each room type. The median for price of an entire home/apt is 110, which is much higher than the upper quartile for price of private room and shared room, which are 80 and 67, respectively. The comparison results are as expected because the entire room provides more private space and usable area, and more costs are charged accordingly. Moreover, there are lots of extreme high outliers for the price of entire room. The range is quite large, which indicates there are lots of luxury property you can rent from Airbnb.

The ANOVA analysis yields the p-value is below 0.05, which means that the average prices are significantly different among different room types.

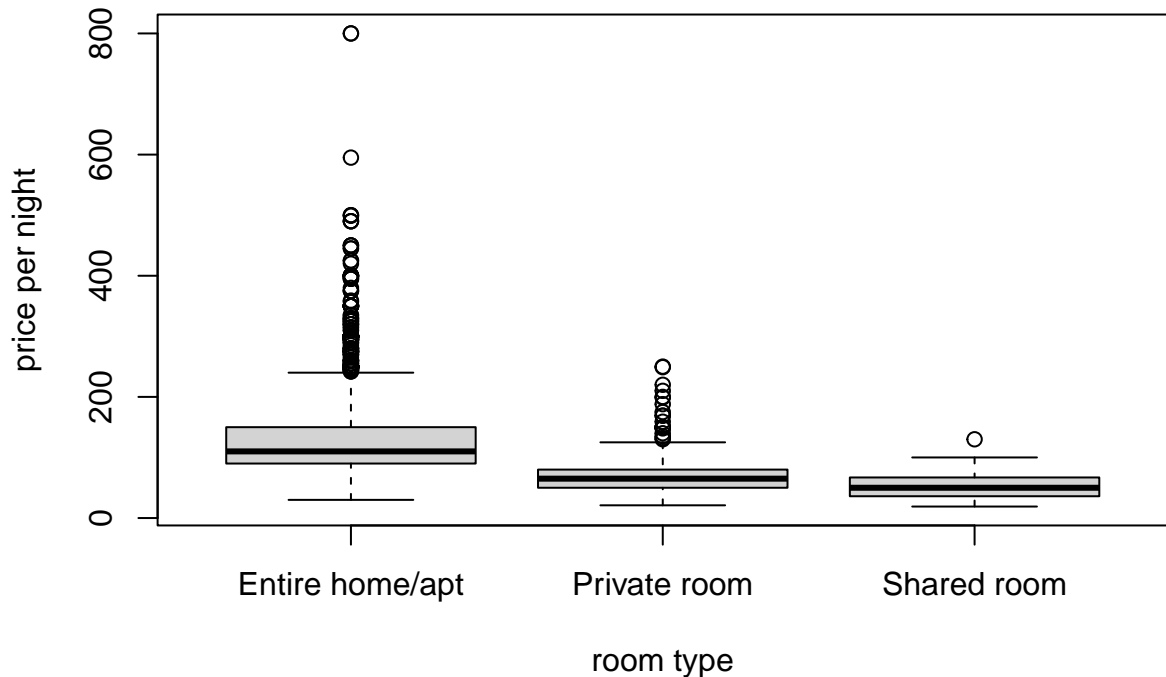


Figure 2: Boxplot for price per night

3.2 Research question 2

The object is to find whether there is a linear relationship between price per night and the number of guests that the shared property can hold. Figure 3 displays the scatter plot between the two variables. It is noted that there is an increasing trend for the price when accommodates value increases. And it seems the relationship is linear. Simple linear regression is a model that describes the relationship between one dependent and one independent variable using a straight line. Through the fitted model, the estimated coefficient on each variable indicates the association between response variable price and the predictor. The regression model is

$$price = \beta_0 + \beta_1 * accommodates + \epsilon \quad (model\ 1)$$

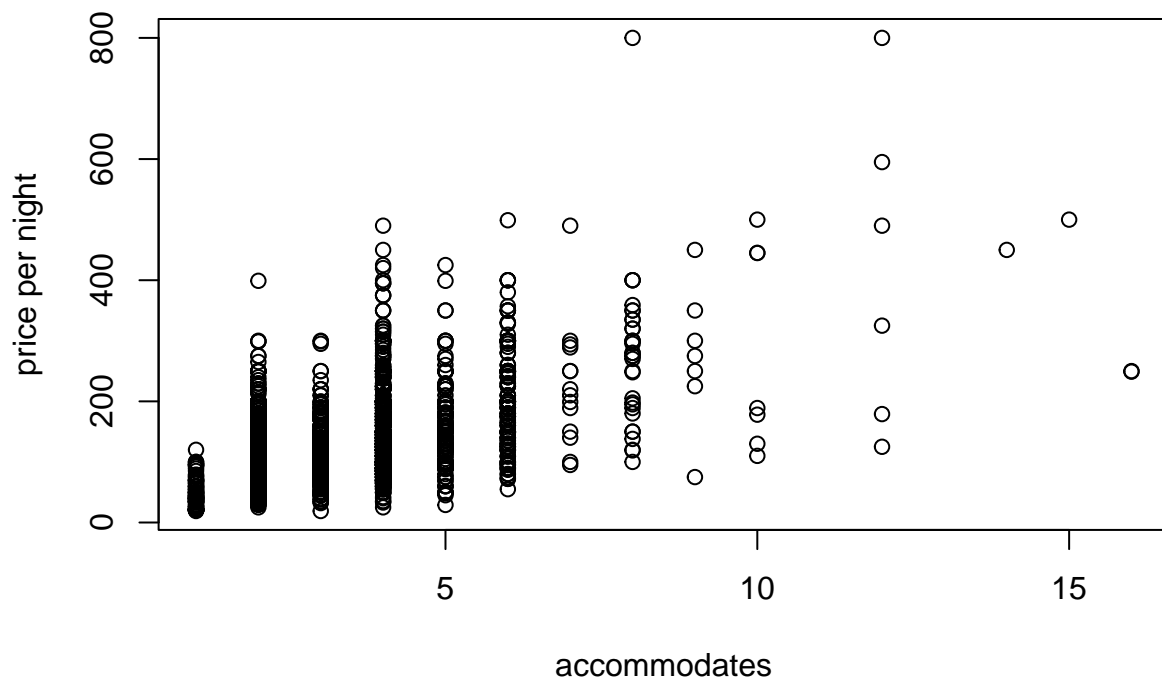


Figure 3: Scatter plot between price and accommodates

Table 3: Regression summary for model 1

term	estimate	std.error	statistic	p.value
(Intercept)	46.601	1.772	26.304	0
accommodates	24.605	0.541	45.496	0

Table 2 lists the regression summary. It is found that the accommodates has significant effect on the performance of price. The estimated coefficient on variable accommodates is 24.605, which is positive, and it indicates when the number of accommodates increase one person, the price per night is expected to increase 24.605 dollars. The coefficient of determination is around 0.31, that about 31% variation of the price can be explained by variable accommodates.

3.3 Research question 3

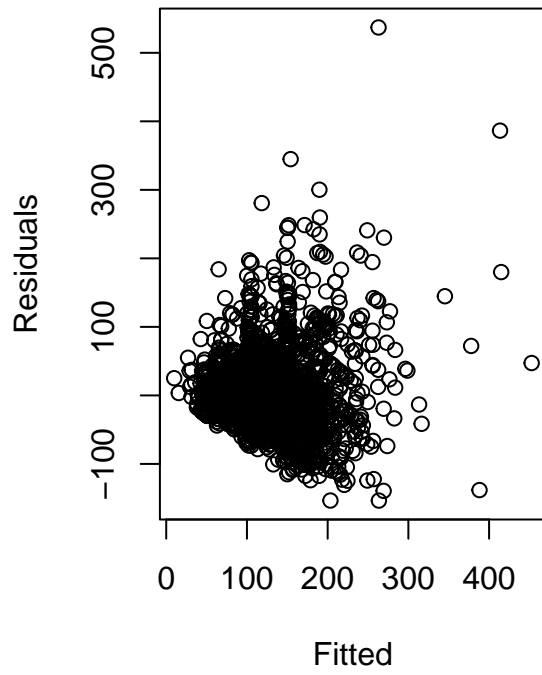
Table 4: Regression summary for model 2

term	estimate	std.error	statistic	p.value
(Intercept)	-18.095	10.000	-1.809	0.070
accommodates	13.249	0.675	19.630	0.000
bathrooms	38.141	2.811	13.568	0.000
bedrooms	20.127	1.344	14.976	0.000
room_typePrivate room	-37.469	2.041	-18.360	0.000
room_typeShared room	-47.818	9.664	-4.948	0.000
host_response_rate	7.209	4.505	1.600	0.110
review_scores_rating	0.317	0.096	3.302	0.001

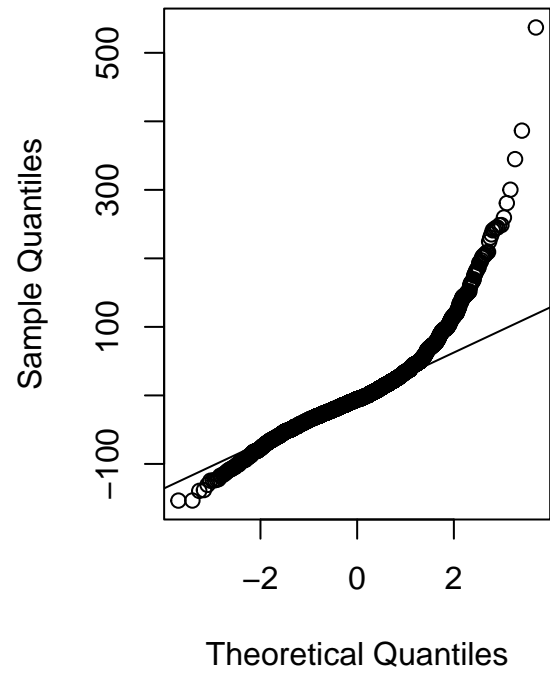
The purpose is to add more explanatory variables to the simple model, and to find whether there is a multiple linear relationship between dependent variable price and all the predictors. Table 3 presents the regression summary. We can notice that except the variable host response rate, all the other estimated coefficients are statistically significant at 5% level. In general, accommodates, bathrooms, bedrooms and review scores rating have positive effect on the prices. Compared to model 1, the estimated coefficient on accommodates decreases to 13.249, because other factors also have effects on the price. As for the categorical variable room type, there are three levels. In the model, entire room serves as the baseline level, therefore, the negative coefficients on private room and shared room means holding other variables constant, the entire rooms cost the highest price per night. The coefficient of determination increases to 0.31, hence the model with more predictors can explain more variability of the dependent variable. Model 2 has a better goodness of fit.

For the assumptions' assessment for the linear regression, Figure 4 displays the residuals diagnostics. The left panel shows that the residuals display a fanning pattern with increasing spread. And in the right panel, the QQ plot tells the majority of the points are aligned with the diagonal line, but there are some deviations on both tails. Considering the large sample size, the normality assumption is considered moderately hold.

Residual vs Fitted



Normal Q-Q Plot



4 Conclusion and Discussion

In this report, we have built a linear regression model with price as a response, room features, host services and accommodates ratings as the predictors. Regarding the estimated coefficients, number of accommodates, bathroom, bedrooms, host response rate and review scores rating have positive effect on the price per night, which is consistent with the results of the peer-reviewed articles. The coefficients on room type are negative, considering the base group is Entire room/apt, we can conclude that the entire room/apartment costs higher average price than private or shared room. The limitations of analysis lie in the assessment of assumptions for the regression model. The residuals are not randomly scattered round the horizontal zero line, which may cause the incorrect or misleading of the inference analysis using the estimated model. For the future research, it is suggested to do some data transformation if the distribution of the dependent variable is not a normal one, logarithm transform is a good way and a log-linear model could be fitted.

References

- Yuan Cai, Yongbo Zhou, Noel Scott, et al. Price determinants of airbnb listings: evidence from hong kong. *Tourism Analysis*, 24(2):227–242, 2019.
- Chen Yong Chen Yong and K Xie. Consumer valuation of airbnb listings: a hedonic pricing approach. 2017.
- David Robinson, Alex Hayes, and Simon Couch. *broom: Convert Statistical Objects into Tidy Tibbles*, 2023. URL <https://CRAN.R-project.org/package=broom>. R package version 1.0.5.
- Augusto Voltes-Dorta and Agustín Sánchez-Medina. Drivers of airbnb prices according to property/room type, season and location: A regression approach. *Journal of Hospitality and Tourism Management*, 45:266–275, 2020.
- Hadley Wickham, Romain François, Lionel Henry, Kirill Müller, and Davis Vaughan. *dplyr: A Grammar of Data Manipulation*, 2023a. URL <https://CRAN.R-project.org/package=dplyr>. R package version 1.1.4.
- Hadley Wickham, Jim Hester, and Jennifer Bryan. *readr: Read Rectangular Text Data*, 2023b. URL <https://CRAN.R-project.org/package=readr>. R package version 2.1.4.
- Yihui Xie. knitr: A comprehensive tool for reproducible research in R. In Victoria Stodden, Friedrich Leisch, and Roger D. Peng, editors, *Implementing Reproducible Computational Research*. Chapman and Hall/CRC, 2014. ISBN 978-1466561595.

5 Appendix: Individual Assignment Coversheet


INDIVIDUAL ASSESSMENT COVER SHEET Faculty of Design and Creative Technologies				 TE WĀNANGA ARONUI O TĀMAKI MAKĀU RAU	
---	--	--	--	--	--

First Name		Family Name		Student ID No	
Paper Name		Paper Code:		Assignment Due Date	
Lecturer:		Tutorial Day		Date Submitted	
Tutor:		Tutorial Time		No.Words/Pages	

In order to ensure fair and honest assessment results for all students, it is a requirement that the work that you hand in for assessment is your own work. If you are uncertain about any of these matters then please discuss them with your lecturer.

Plagiarism and Dishonesty are methods of cheating for the purposes of General Academic Regulations (GAR)
<http://www.aut.ac.nz/calendar>

Assignments will not be accepted if this section is not completed and signed.

Please read the following and **tick**  to indicate your understanding:

1. I understand it is my responsibility to keep a copy of my assignment.	<input type="checkbox"/> Yes	<input type="checkbox"/> No
2. I have signed and read the Student's Statement below .	<input type="checkbox"/> Yes	<input type="checkbox"/> No
3. I understand that a software programme (Turnitin) that detects plagiarism and copying may be used on my assignment.	<input type="checkbox"/> Yes	<input type="checkbox"/> No

Student's Statement:

This assessment is entirely my own work and has not been submitted in any other course of study. I have submitted a copy of this assessment to Turnitin, if required.

In this assessment I have acknowledged, to the best of my ability:

- The source of direct quotes from the work of others.
- The ideas of others (includes work from private or professional services, past assessments, other students, books, journals, cut/paste from internet sites and/or other materials).
- The source of diagrams or visual images.

Student's Signature: _____ **Date:** _____

The information on this form is collected for the primary purpose of submitting your assignment for assessment. Other purposes of collection include receiving your acknowledgement of plagiarism policies and attending to administrative matters. If you choose not to complete all questions on this form, it may not be possible for the Faculty of Design and Creative Technologies to accept your assignment.

6 Appendix: R Environment

All the statistics were done using R 4.3.2 (2023-10-31), the packages readr Wickham et al. [2023b], dplyr Wickham et al. [2023a], knitr Xie [2014] and broom Robinson et al. [2023] were used.