

Airbnbs in New York City

Team lol: Tamsin Connerly, Hannah Lee, Jasmine Xiang

2025-03-17

Introduction

The rise of short-term rental platforms, particularly Airbnb, has significantly disrupted the traditional hospitality industry and transformed urban housing markets worldwide. In New York City, one of the world's most popular tourist destinations, the impact of Airbnb has been particularly pronounced, raising questions about its effects on local communities, housing affordability, and the broader urban economy.

Previous research has identified several factors that impact Airbnb pricing. One study found that host attributes, site and property attributes, amenities and services, rental rules, and online review ratings all play significant roles in determining listing prices [wang2017]. Furthermore, recent studies have provided evidence of Airbnb's influence on housing markets. Another study found that a 1% increase in Airbnb listings leads to a 0.018% increase in rents and a 0.026% increase in house prices [barron2018]. This effect is more pronounced in areas with a lower share of owner-occupiers, suggesting that non-owner-occupiers are more likely to reallocate their properties from long-term to short-term rentals.

Our research question is: "How do various factors, such as bedroom number, room type, review scores, and neighborhood, influence the price of Airbnb listings in New York City?"

Understanding the determinants of Airbnb pricing in New York City is crucial for several reasons. Firstly, it can provide valuable insights for policymakers grappling with the challenges posed by the growth of short-term rentals, including potential impacts on housing affordability and neighborhood character [toader2021]. Secondly, it can help hosts make more informed pricing decisions, potentially leading to more efficient market outcomes.

Based on existing literature and our understanding of the New York City housing market, we hypothesize that:

- Listings with more bedrooms will command higher prices, reflecting the premium placed on space in urban environments.

- The type of room (entire home/apartment vs. private room) will significantly impact pricing, with entire homes/apartments having a higher price.
- Higher review scores will be associated with higher prices, as positive feedback may justify premium pricing.
- Properties in more affluent neighborhoods like Manhattan will have higher prices compared to less affluent ones like the Bronx because of real estate price differences in each borough.

Exploratory Data Analysis

The Airbnb dataset that we are utilizing can be found on Inside Airbnb (<https://insideairbnb.com/>). Inside Airbnb has collected data on dozens of countries and cities, but we decided to focus on New York City. The data was sourced from publicly available data on the Airbnb website on March 1, 2025.

Each row in the dataset represents a unique Airbnb listing in New York City. Each of these correspond to individual properties available for rental on the platform and have many (58) variables such as name of the listing, latitude and longitude, room type, price, minimum number of nights required for booking, total number of reviews the listing has reviewed, and more. We are particularly interested in the following explanatory variables that we believe could impact the price of an Airbnb listing:

- Number of bedrooms
 - These could give insights into the size and comfort level of the Airbnb, likely affecting the price.
- Room type
 - Type of room (whether the listing is an entire home/apartment or a private room) can impact pricing
- Review scores (overall rating)
 - Listings with better or higher number of reviews could be priced higher because of higher perceived value and trustworthiness
- Neighborhood
 - Which New York City Neighborhood the Airbnb is located in. More affluent neighborhoods like Manhattan may have higher prices than neighborhoods like the Bronx.

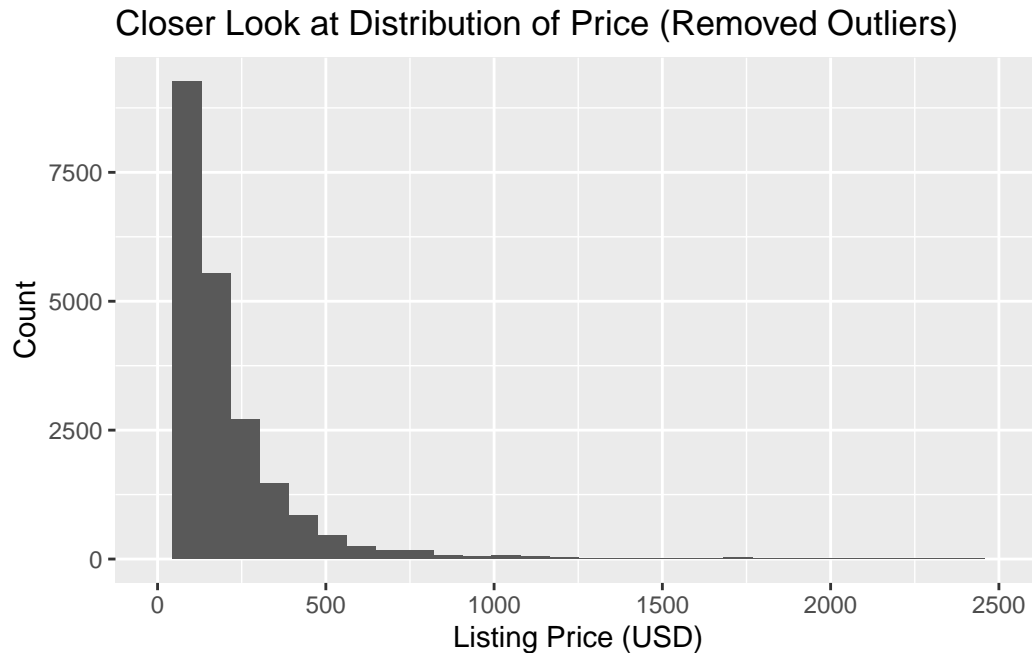
Univariate Exploratory Data Analysis

Response Variable - Price



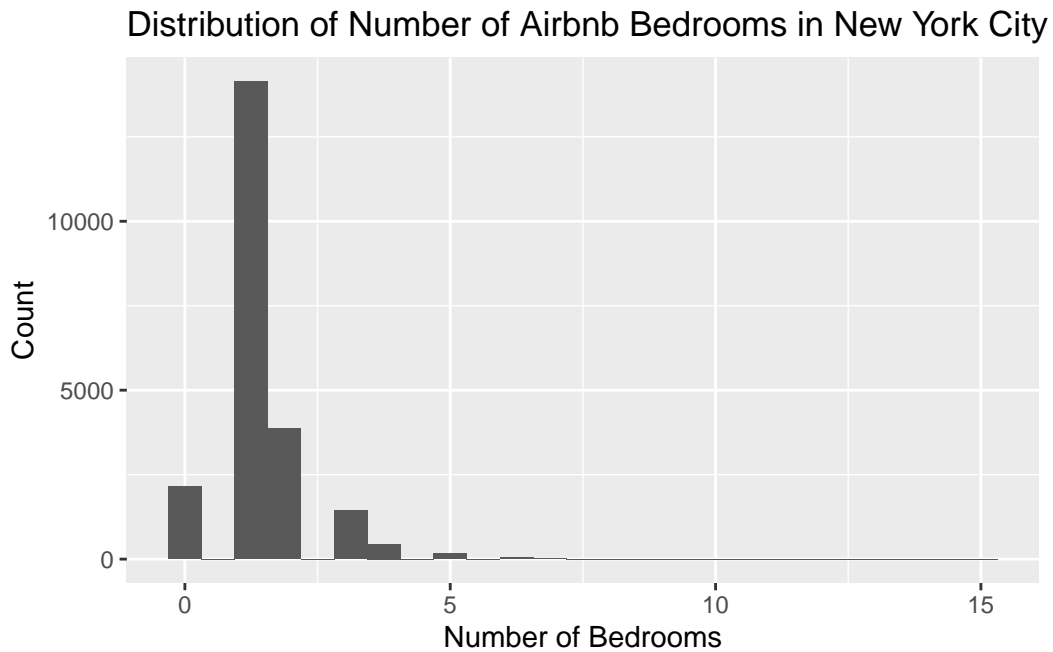
```
# A tibble: 6 x 2
  Statistic Value
  <chr>      <table>
1 Min.          7.0000
2 1st Qu.       85.0000
3 Median       140.0000
4 Mean        213.8352
5 3rd Qu.      240.0000
6 Max.      20000.0000
```

The distribution is pretty heavily right skewed, as can be seen from both of the histograms. It is difficult to analyze the distribution in the first histogram because there is an outlier at \$20,000 and makes the bins and binwidth very narrow and zoomed out (since the range of the data is too large). It is also clear that this outlier impacts the mean, since the median of \$140 is quite a bit less than the mean of around \$213.84, and the mean is roughly equal to the 3rd quartile which is also around \$240. We plan to remove this outlier as a result when doing our analysis, and we will go into more depth later and check for additional outliers.



From the initial histogram, since the majority of the listings seem to fall between 0 and 2500 dollars, we visualized the distribution of the listings between this price range to get a better view of it. We can see that the distribution is still right skewed, and the vast majority of the listings seem to cost between \$50-\$200. There are a few more points visible that were not in the initial histogram that appear to be quite far from where the majority of the listings are clustered, so we will go back and check to see if they are greater than $1.5 \times \text{IQR}$ and can be classified as outliers. We also plan to apply log transformation to this variable to address the skew of the response variable.

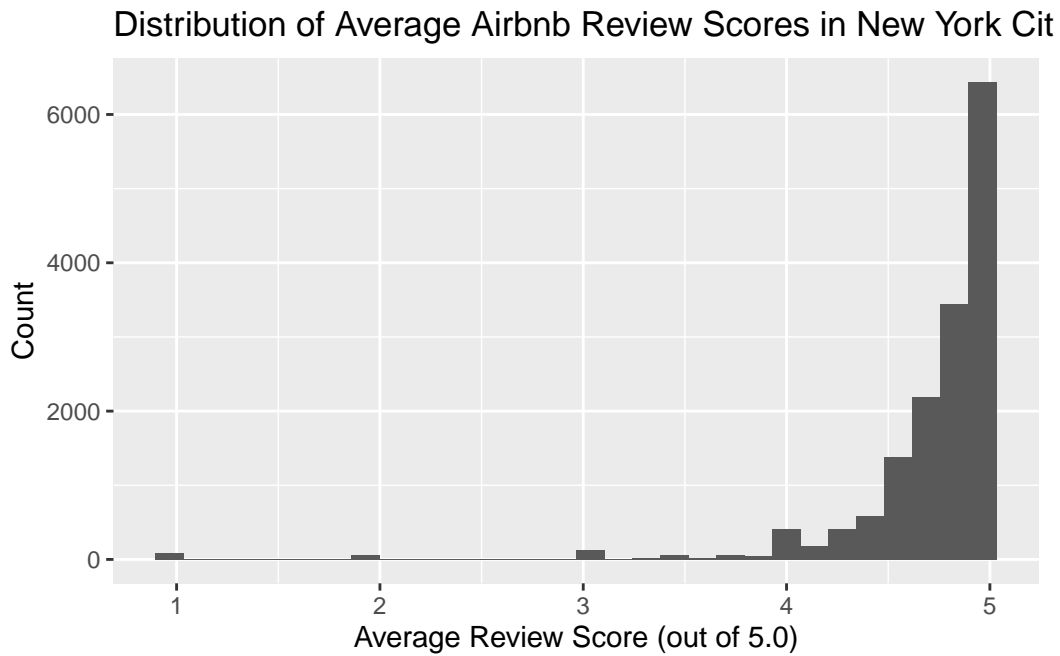
Predictor Variable - Bedrooms



The distribution of the number of bedrooms for Airbnb listings in New York City is skewed to the right. Most properties are one- or two-bedroom units. This aligns with the typical demand for smaller accommodations, which are more suited for individuals or small groups that might be more common. Listings with more bedrooms, such as 5 or more, are far less common. Larger apartments or homes, typically with 3 or more bedrooms, are less frequent, but they most likely consist of listings that may be priced much higher.

```
# A tibble: 7 x 2
  Statistic Value
  <chr>      <table>
1 Min.      0.000000
2 1st Qu.   1.000000
3 Median    1.000000
4 Mean      1.315468
5 3rd Qu.   2.000000
6 Max.      15.000000
7 NA's      49.000000
```

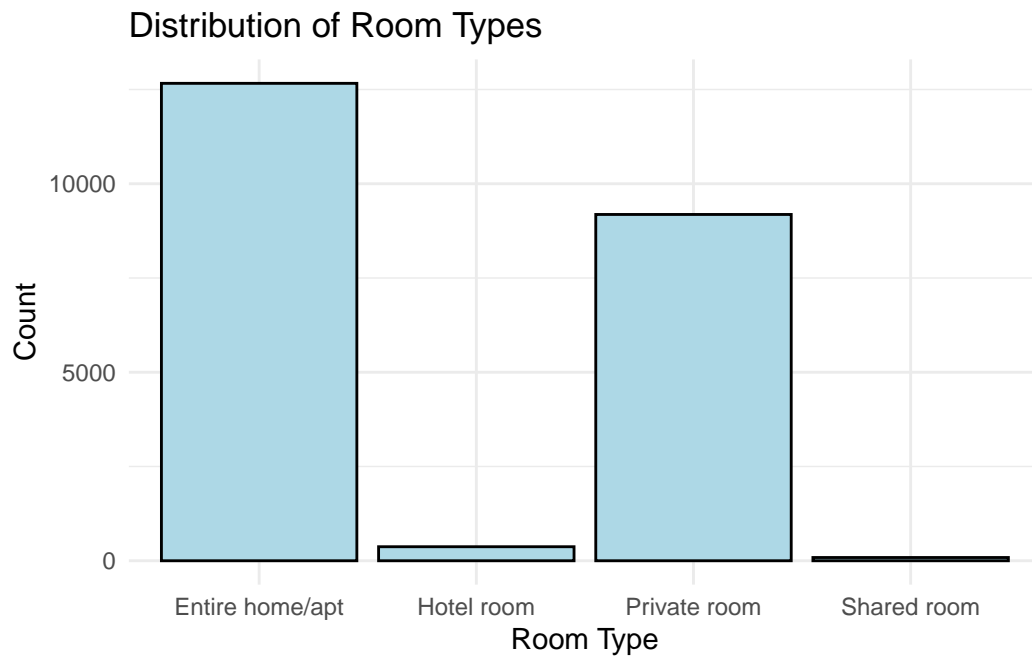
Predictor Variable - Review Scores



```
# A tibble: 7 x 2
  Statistic Value
  <chr>      <table>
1 Min.      1.000000
2 1st Qu.   4.662500
3 Median    4.850000
4 Mean      4.724384
5 3rd Qu.   5.000000
6 Max.      5.000000
7 NA's      6798.000000
```

The distribution of the predictor variable “review score” is skewed left, with the median of 4.85 being slightly higher than the mean of 4.72. The majority of the reviews are around 5 (over 6000 of them), and the vast majority of the observations are between 4 and 5, with very few of them being below 3. Additionally, there are quite a few NA values for this variable (around a third of the rows have NAs).

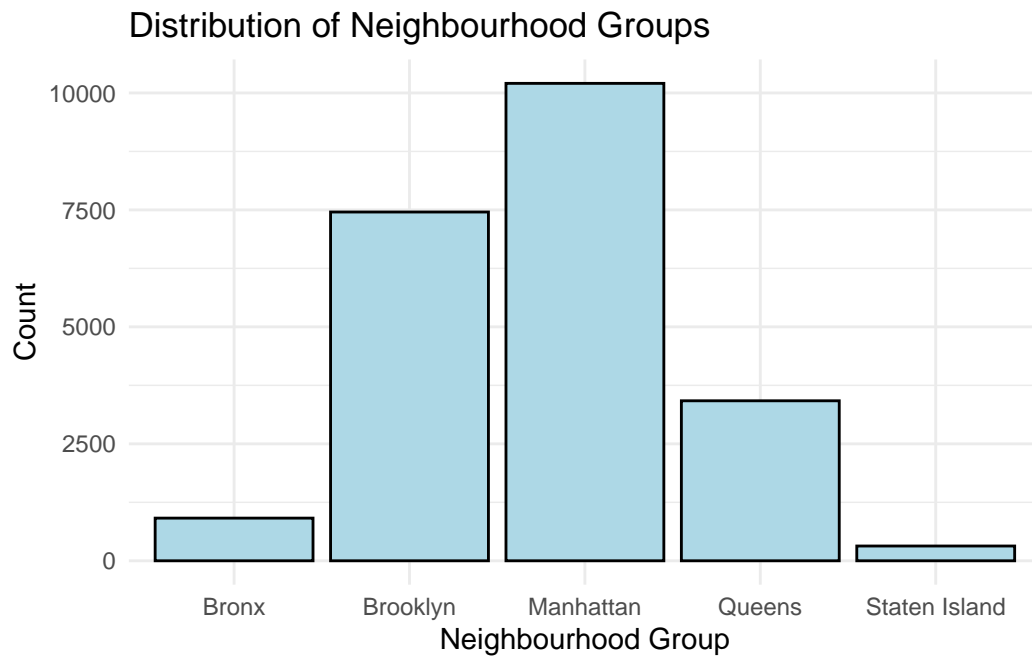
Predictor Variable - Room Type



Var1	Freq
Entire home/apt	12664
Hotel room	372
Private room	9186
Shared room	86

The mode, or most frequent room type in this dataset is Entire home/apt, followed by private room; these may be more popular and sought out, and might tend to be more expensive than the hotel or shared rooms. There are very few hotel rooms and even fewer shared rooms.

Predictor Variable - Neighborhood

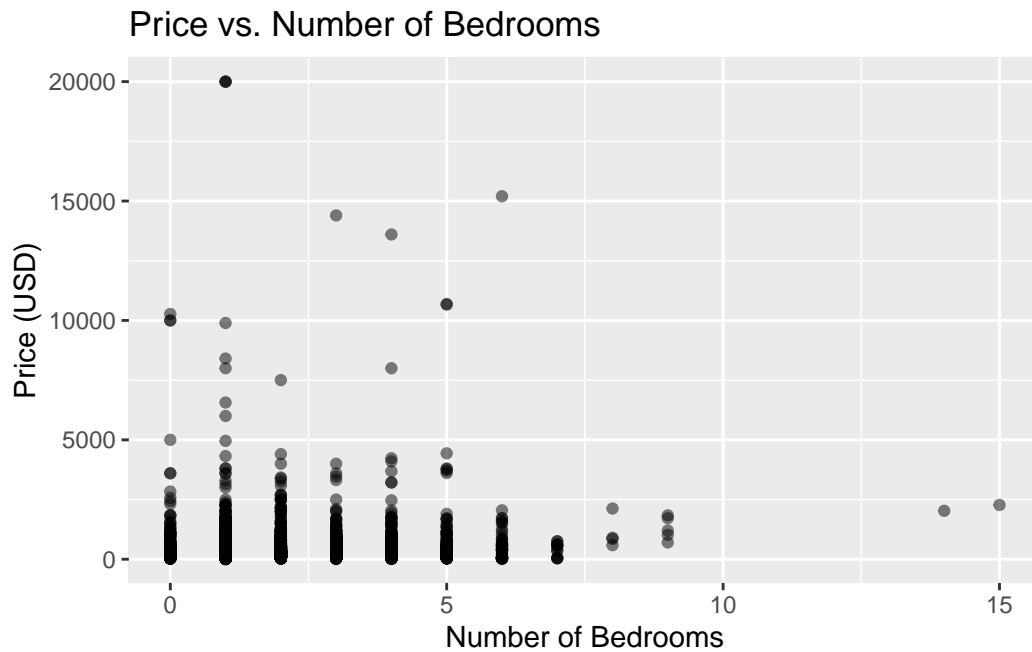


Var1	Freq
Bronx	912
Brooklyn	7455
Manhattan	10205
Queens	3420
Staten Island	316

The majority of the airbnb listings are in Manhattan (just over 10,000), followed by Brooklyn (around 7500), and Queens (a bit over 2500). A few of them are in Bronx and even fewer in Staten Island. This is to be expected, as Manhattan and Brooklyn are prime areas for tourism and business, while other areas might be less popular for short-term rentals.

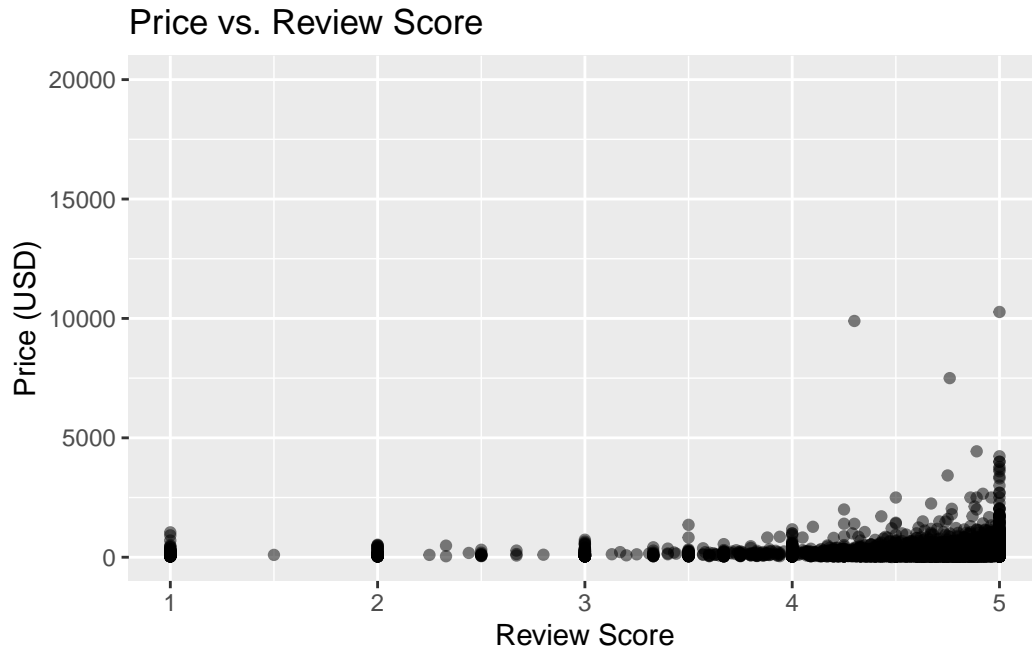
Bivariate Exploratory Data Analysis

Response (Price) vs Predictor Variable (number of prices)

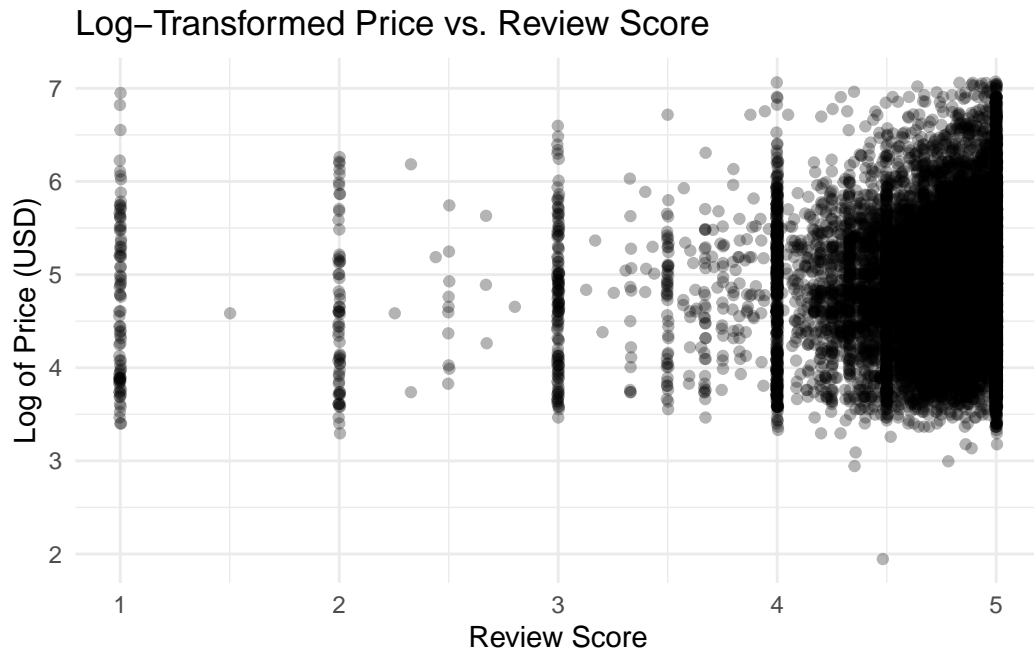


This scatter plot visualizes the relationship between the number of bedrooms and the price of Airbnb listings in New York City. There is not a clear trend. However, listings with 0 to 5 bedrooms exhibit a wide range of prices, with some listings priced significantly higher than the majority. Some outliers even exceed \$10,000 per night, likely representing luxury or highly unique accommodations. This indicates that other factors such as location, room type, and review scores may have a stronger influence on pricing.

Response (Price) vs Predictor Variable (review scores)

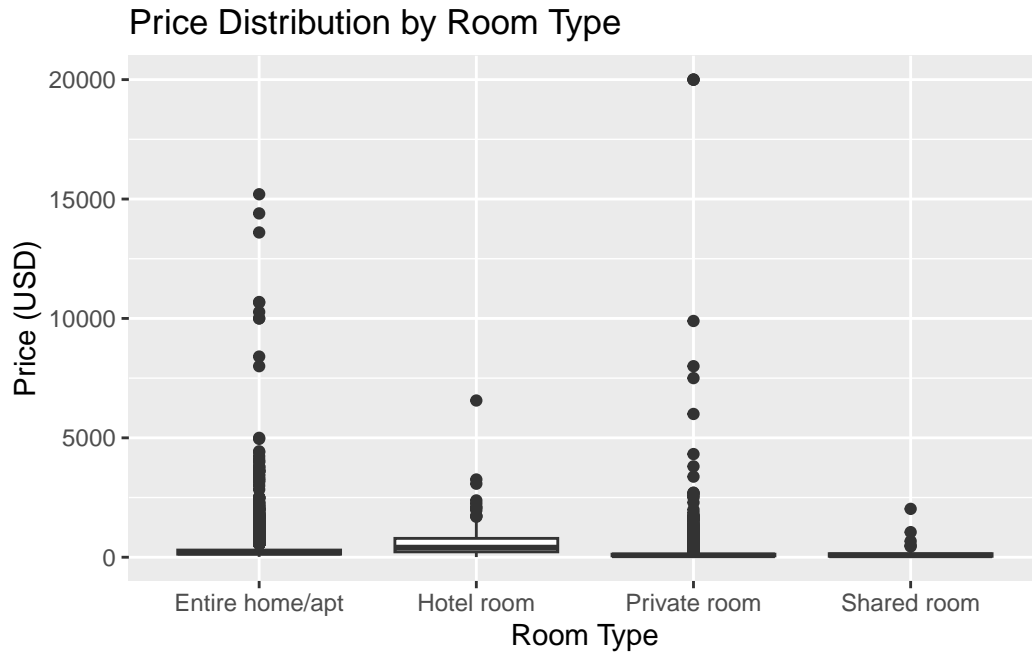


This scatter plot shows the relationship between price and review scores. However, it looks a little problematic and hard to interpret because high data density at certain score levels, especially between 4 and 5. Since the price variable is highly skewed, we applied a log transformation to try to help spread out values and make trends more visible.

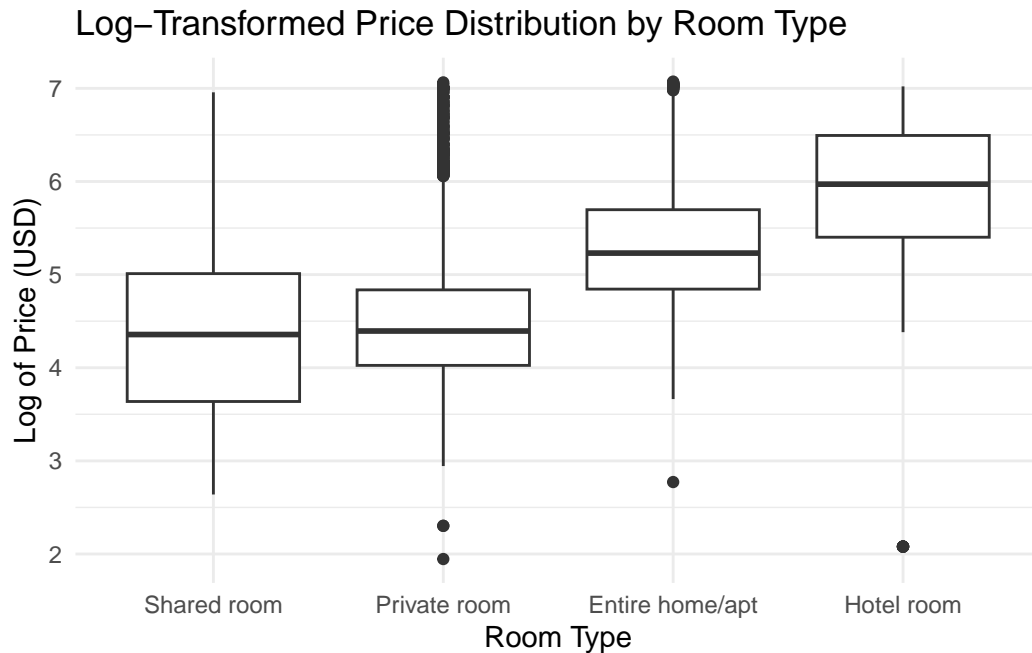


The majority of listings have review scores between 4 and 5, this shows that most listings have scores within this range. Higher-rated listings tend to have slightly higher prices, but the effect is weak.

Response (Price) vs Predictor Variable (Room Type)



This box plot tries to visualize the distribution of prices across different room types, but it appears compressed due to the presence of extreme outliers. This makes it difficult to interpret the main trends effectively. To improve the visualization, We applied a log transformation, which may help spread out the values and make the differences between room types clearer, because prices are highly skewed.

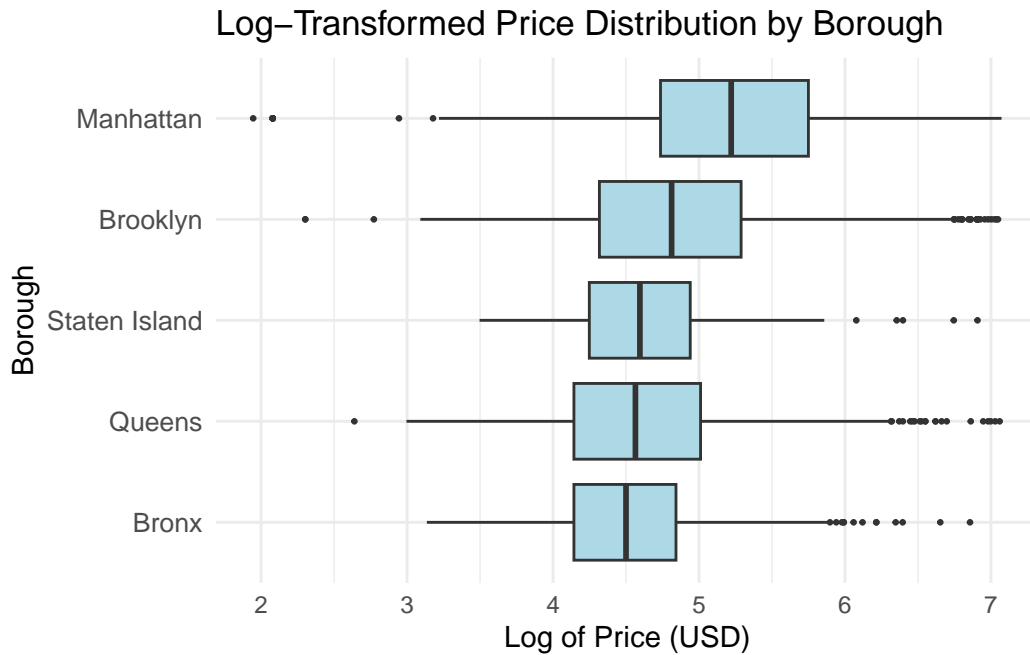


This box plot provides a clearer comparison of price distributions across different room types. According to this plot, hotel rooms have the highest median price and general highest price among all room types, suggesting that they are generally priced higher than other Airbnb listings such as shared room and private room, etc. Also, the IQR for hotel rooms and entire homes/apartments is larger compared to private and shared rooms, indicating greater variation in pricing. Even after log transformation, private rooms and entire homes/apartments still exhibit a notable number of higher-end outliers. However, the log transformation has greatly reduced the impact of extreme values, making the distribution more interpretable.

Response (Price) vs Predictor Variable (Neighborhood)



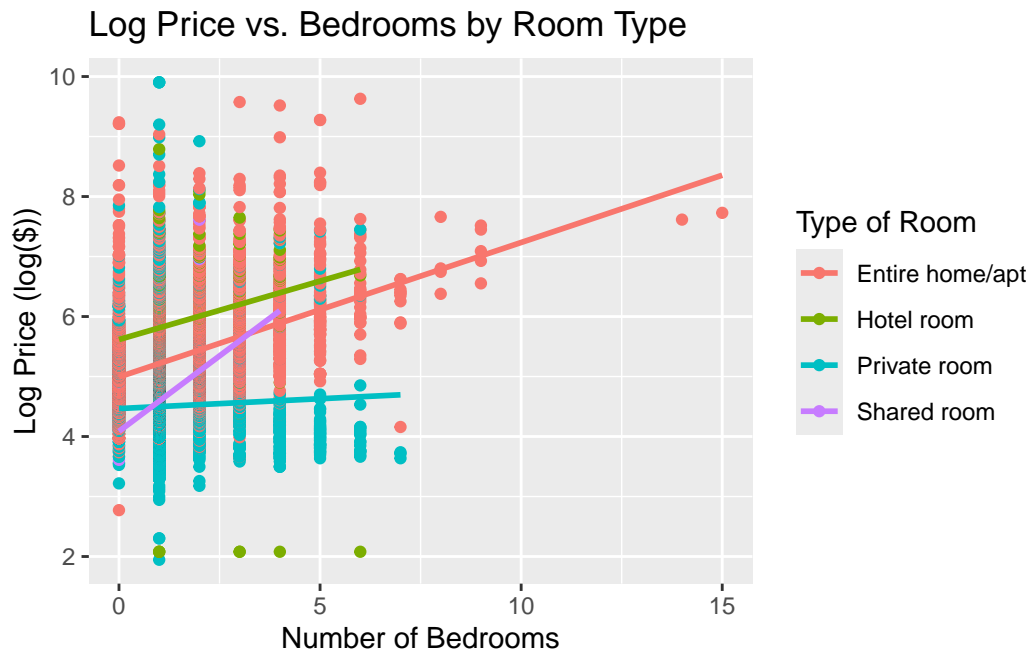
The current plot has similar issues as plots that are shown above, which is difficult to interpret due to issues like extreme outliers and others. We removed the extreme values and also did log transformation.



According to this plot, we can see that Manhattan has the highest median price, showing that it is the most expensive borough for Airbnb listings. It also exhibits the widest IQR, suggesting a high variation in listing prices. The median of Brooklyn follows Manhattan, with a slightly lower median price but still a wider spread. There are still some outliers shown in the plot, but the interpretability is much better.

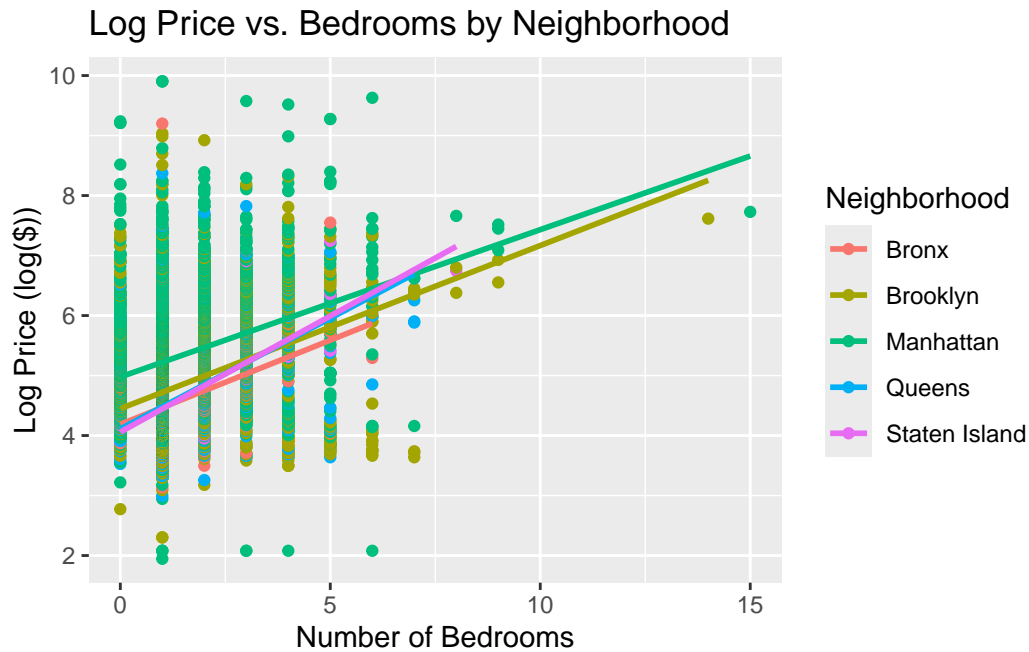
Interaction Effects

Bedrooms and Room Type



Based on the graph, it appears that the rate at which price increases per number of bedrooms varies across room types. The slope of the shared room especially seems to differ from the others. Thus, there may be an interaction effect here.

Bedrooms and Neighborhood



Based on the graph visually, it appears that the rate at which price increases per bedroom does not greatly vary based on neighborhood, especially compared to the previous graph by room type. Staten Island does have a noticeably different slope, though. Thus, there is a potential for an interaction effect between bedroom number and neighborhood, but not as much as the previous graph with bedroom and room type.