

Capstone Project 1:

New York City Airbnb Price Prediction

Catherine Somers

PROBLEM STATEMENT

Airbnb is an American online marketplace and hospitality service for arranging lease or rent lodging in the short-term. The company does not own any of the real estate and instead a broker receives a percentage of the services fees with each booking.

Since 2008, Airbnb has become one of the largest marketplace for hospitality services with over 6 million listings in around 190 countries around the world.

An obstacle that both Airbnb hosts and Airbnb faces is determining an optimal rent price per night. As a host on the platform, charging too much causes the renter to choose more affordable options fit to their budget. Whereas, if they charge too low, they lose out on potential revenue on the property. If Airbnb could recommend an optimal listing price to hosts, it could help maximize revenue as well as service fees.

In this capstone project, we would like to be able to predict the listing price so that both the hosts and the company can optimally gain revenue.

CLIENT

The Airbnb hosts and Airbnb, the company would use price prediction for better pricing and to also maximize revenue.

DATASET

- The listing data for New York City can be found on Kaggle: [New York City Open Data](#)

- Home
- Compete
- Data**
- Notebooks
- Discuss
- Courses
- More

Dataset

New York City Airbnb Open Data

Airbnb listings and metrics in NYC, NY, USA (2019)

Dgomonov • updated 6 months ago (Version 3)

[Data](#)
[Tasks \(1\)](#)
[Kernels \(257\)](#)
[Discussion \(17\)](#)
[Activity](#)
[Metadata](#)
[Download \(7 MB\)](#)
[New Notebook](#)

Usability 10.0

License CC0: Public Domain

Tags business, internet, world, hotels and accommodations, vacation rentals and short-term stays and 2 more

Description

Context

Since 2008, guests and hosts have used Airbnb to expand on traveling possibilities and present more unique, personalized way of experiencing the world. This dataset describes the listing activity and metrics in NYC, NY for 2019.

Content

This data file includes all needed information to find out more about hosts, geographical availability, necessary metrics to make predictions and draw conclusions.



The data consists of 16 different columns:

Host descriptors:

- `host_id`: host ID
- `host_name`: host name

- `calculated_listings_count`: amount of listing per host

Listing descriptors:

- `id`: ID
- `name`: name of listing
- `room_type`: living space type
- `minimum_nights`: amount of nights minimum
- `availability_365`: number of days listing is available for booking
- `price`: price in dollars

Review descriptors:

- `number_of_reviews`: number of reviews
- `last_review`: latest review
- `reviews_per_month`: number of reviews per month

Location descriptors

- `neighbourhood_group`: location
- `neighbourhood`: area
- `latitude`: latitude coordinates
- `longitude`: longitude coordinates

DATA WRANGLING

1. Data Collection and Importing
 - a. Loaded in as a dataset. This dataset is publicly accessible and the `.csv` file contains summary information and metrics pertaining to Airbnb listing in New York City.
2. Selecting columns that may be useful in analysis
 - a. Visually inspect data and select columns useful for analysis
3. Drop duplicates
 - a. Removed repeat rows
4. Fix missing values
 - a. Drop rows that contain null values since

5. Consider the outlier data
 - a. In order for it not to skew the rest of the data we take the data between the 25th percentile and 75th percentile

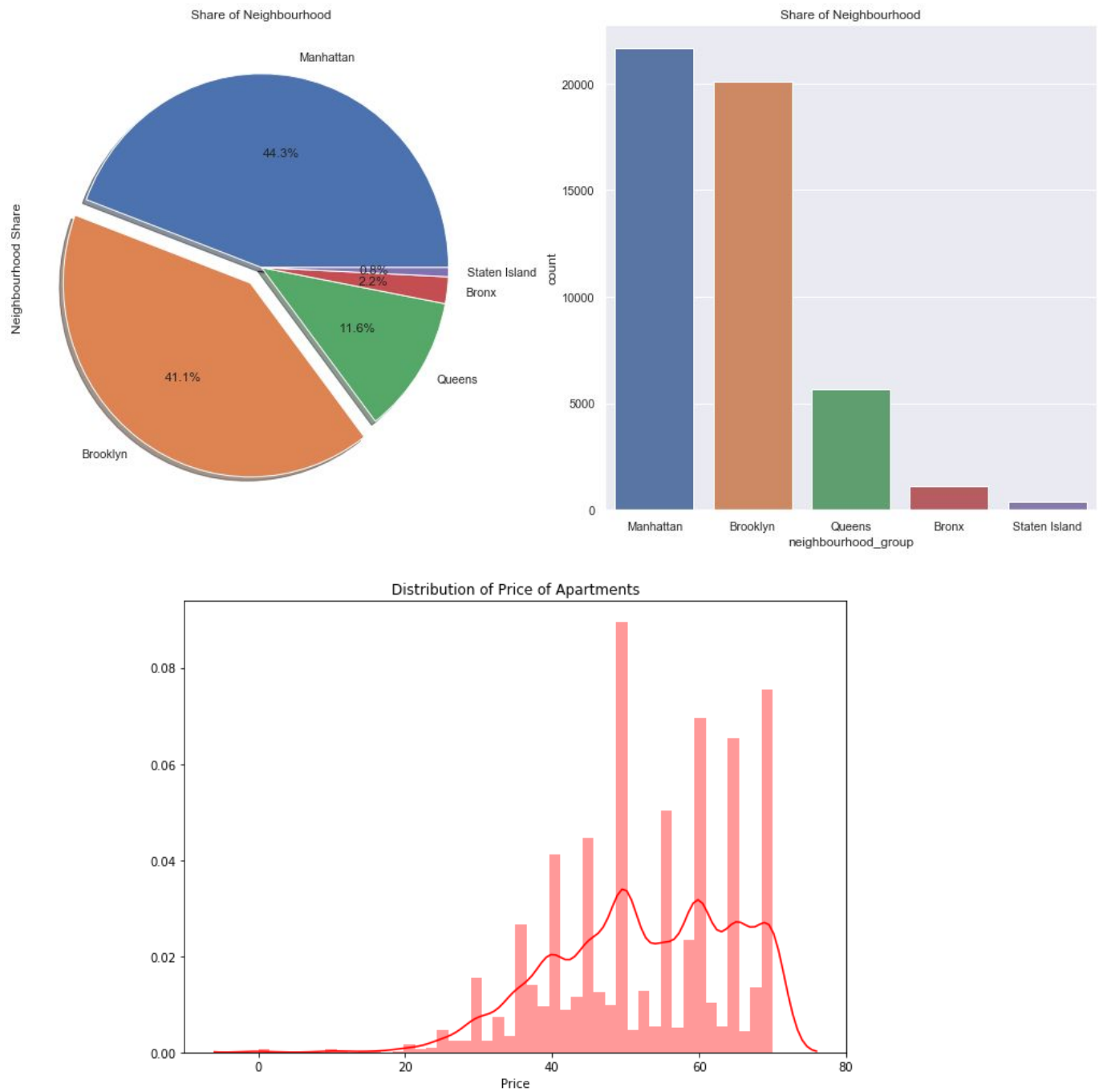
INITIAL FINDINGS

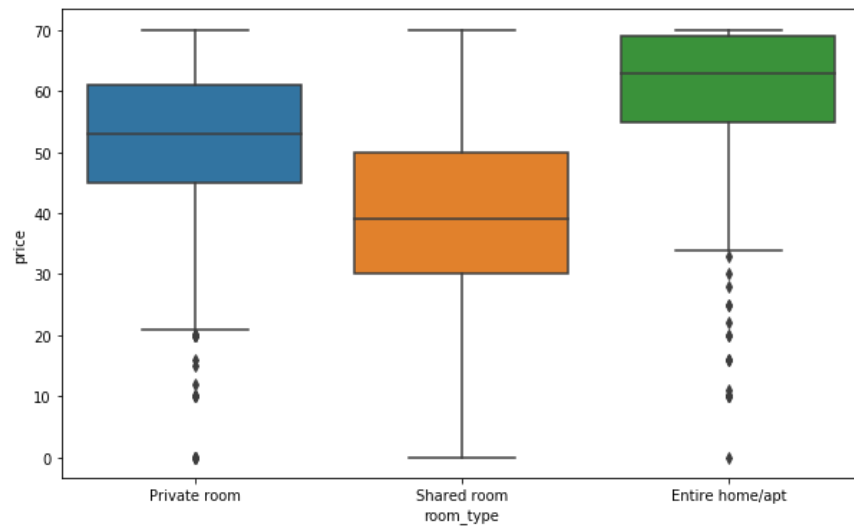
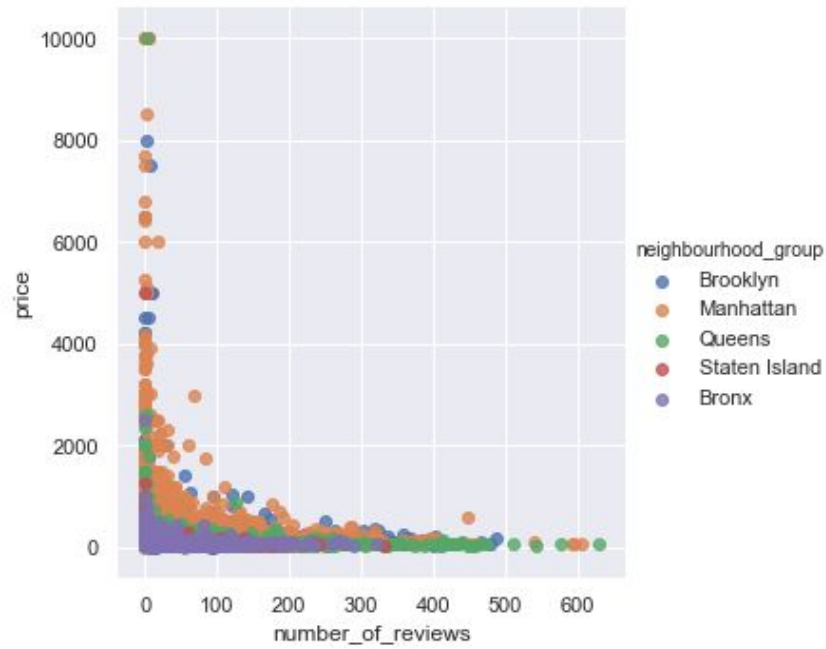
During exploratory data analysis, we asked the following questions:

1. Which neighbourhood group are these Airbnbs mostly located?
2. What is the price distribution based on the number of reviews?
3. What is the relationship between price and room type across neighbourhood and neighbourhood group?

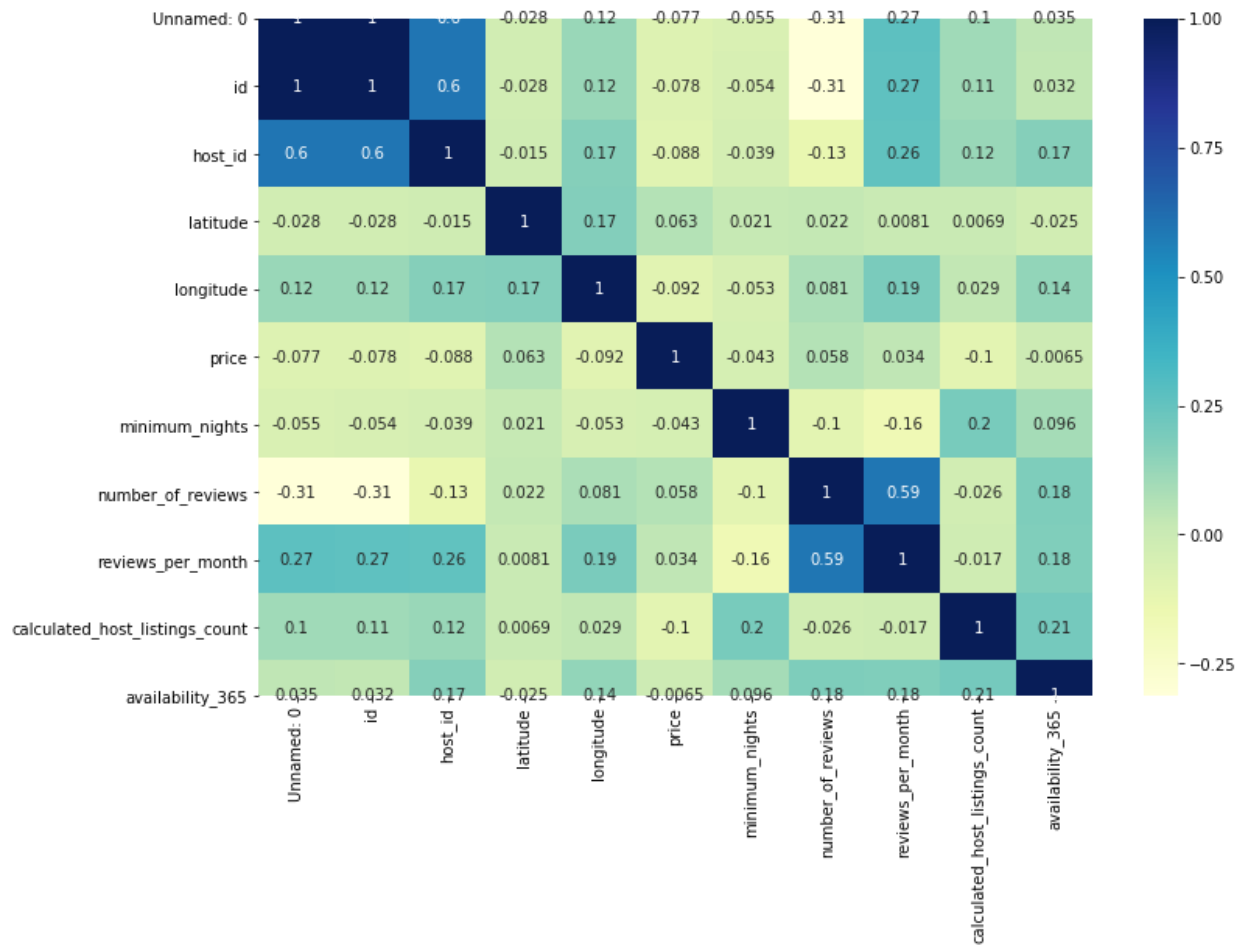
Initial findings are:

1. Most Airbnbs in New York City are located in the Manhattan and Brooklyn neighbourhood groups.
 - a. 44.3% of Airbnbs located in Manhattan
 - b. 41.1% of Airbnbs located in Brooklyn
 - c. We can also see this evidence of both being the majority in the bar graph.
2. Among the number of reviews, price distribution of apartments are more concentrated around the Manhattan, Queens, and Bronx neighbourhood groups.
 - a. Based on how we see the data spread out in the graph below
3. Price of the property is dependent on the room type across the neighbourhood and neighbourhood group.
 - a. The mean price across all room types is not equal.
 - b. Price is dependent on the neighbourhood group home is available in.
 - c. Between the room type and neighbourhood group, the proportion of availability of a particular room type is dependent on the neighbourhood group searched for.





In



In the heatmap seen above, the highest correlation exists between `host_id` and `id` (0.6) closely followed `number_of_reviews` and `reviews_per_month`, while the lowest correlation exist between `longitude` and `price` (-.092).