



# Predicting Airbnb Listing Prices in New York City Using Regression Models

*Midterm Presentation*

Peter Huang

Brown University

10/22/2019

<https://github.com/huang960404/data1030-project>

# Introduction

- Goal: predict the price of future Airbnb listings as accurately as possible
- Why: Airbnb is becoming more and more popular these days, and customers want to understand the pattern of pricing and what characteristics impact the pricing decision the most.
- Data set: New York City Airbnb Open Data
  - 48,895 observations, 16 variables, 782,320 data points
  - 11 numeric variables, 5 categorical variables
- Source: Kaggle
  - <https://www.kaggle.com/dgomonov/new-york-city-airbnb-open-data>
- Type: regression model
  - Target variable: price



# Preprocessing - Missing Values

- id, host\_id: dropped

number_of_reviews	last_review	reviews_per_month
0		

- reviews\_per\_month: set missing values to 0
- last\_review
  - Days between dates (reference date: 7/9/2019)
  - Set missing values to 0

# Preprocessing - Feature Engineering

- name
  - Length
    - “!” : 1 (contains “!”) or 0 (does not contain “!”)
- host\_name
  - Length: 1 ( $\geq 3$ ) or 0 ( $\leq 2$ )

# Preprocessing - Encoders

- OneHotEncoder
  - neighbourhood\_group: Bronx, Brooklyn, Manhattan, Queens, Staten Island
  - neighbourhood: Kensington, Midtown, Harlem, Hell's Kitchen, Chinatown, etc.
- OrdinalEncoder
  - room\_type: 0 (Shared room), 1 (Private room), 2 (Entire home/apt)

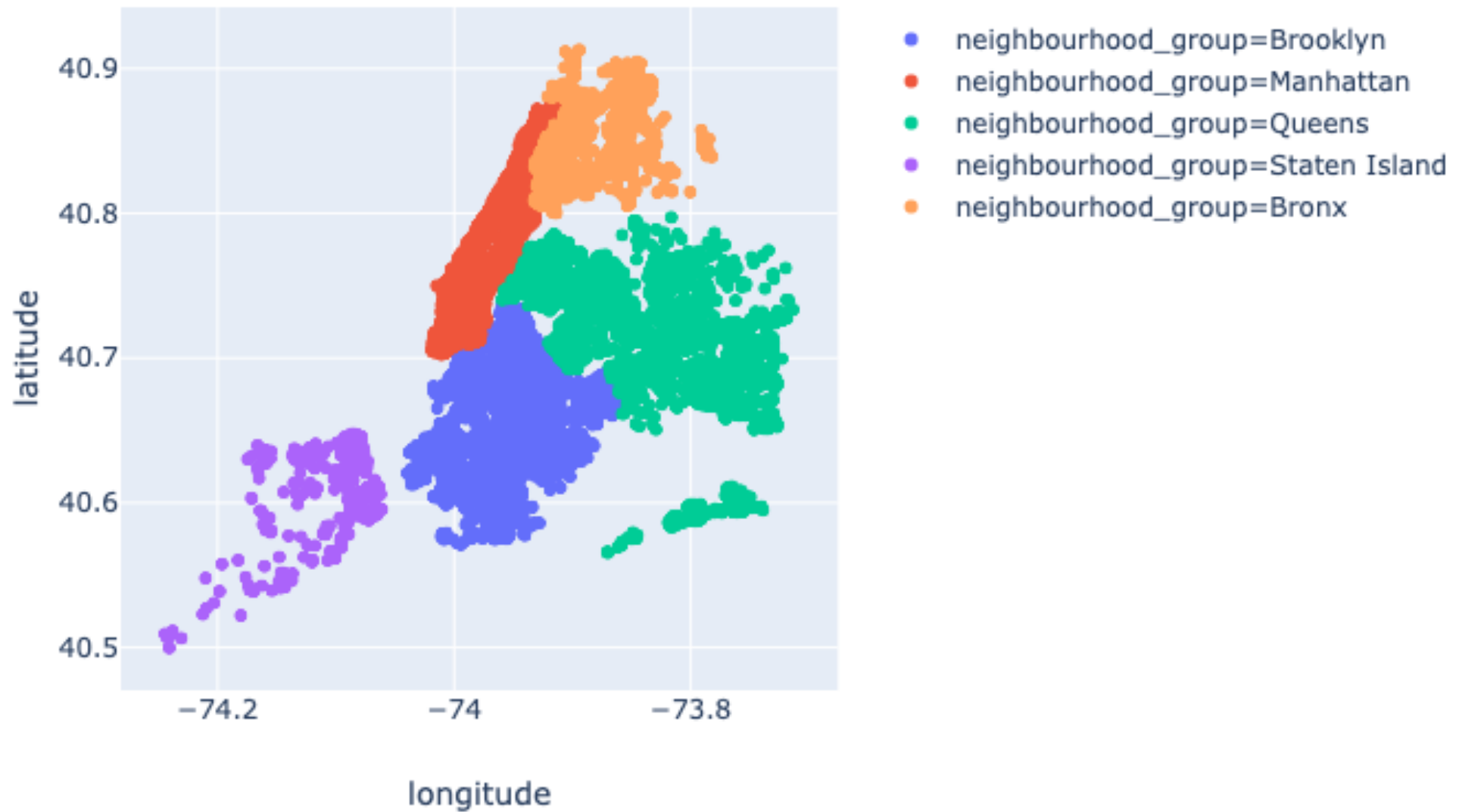


# Preprocessing - Scalers

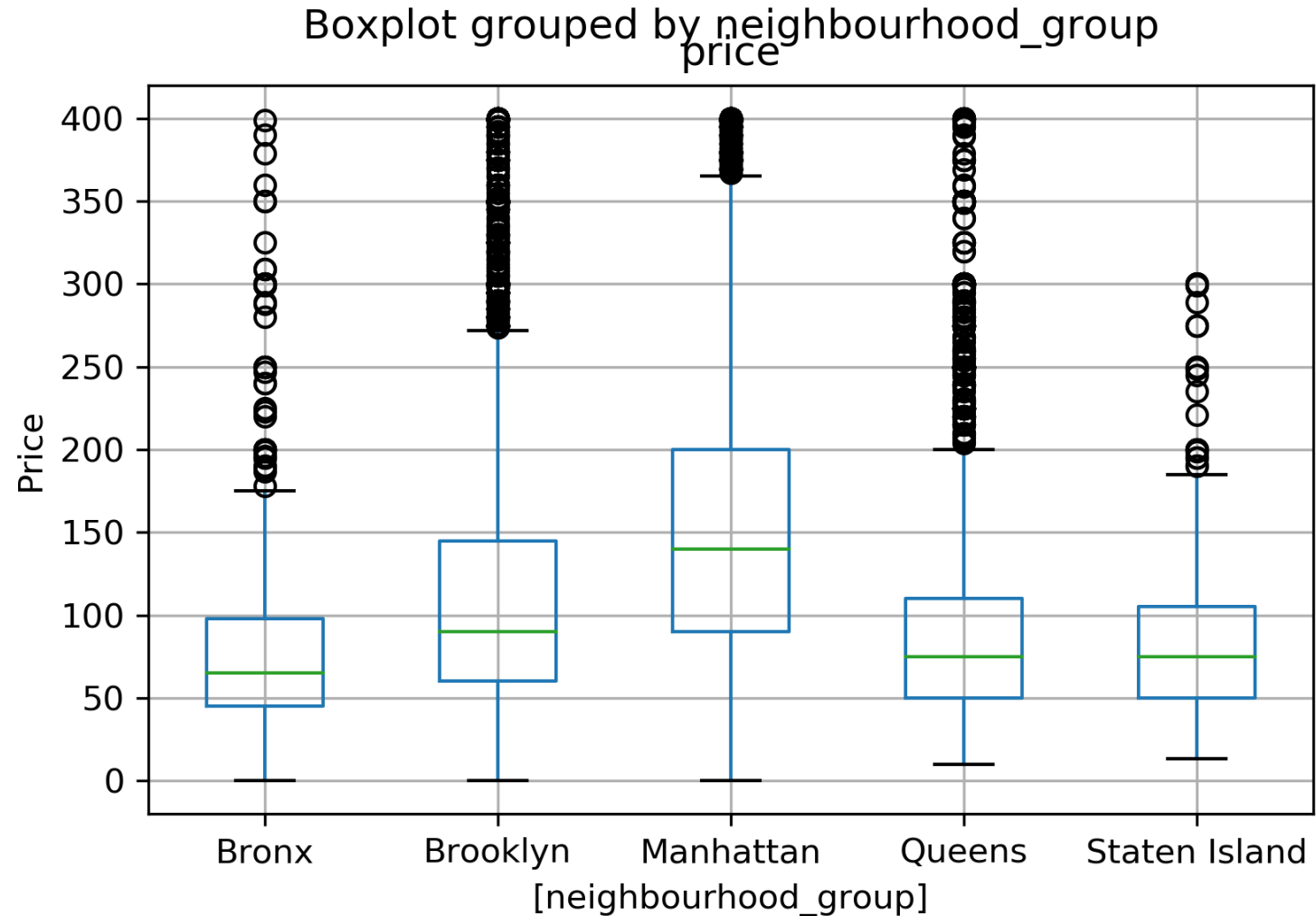
- MinMaxScaler
  - latitude, longitude, availability\_365
- StandardScaler
  - name, minimum\_nights, number\_of\_reviews, last\_review, reviews\_per\_month, calculated\_host\_listings\_count
- 239 features, 11,685,905 data points



# Exploratory Data Analysis - Map

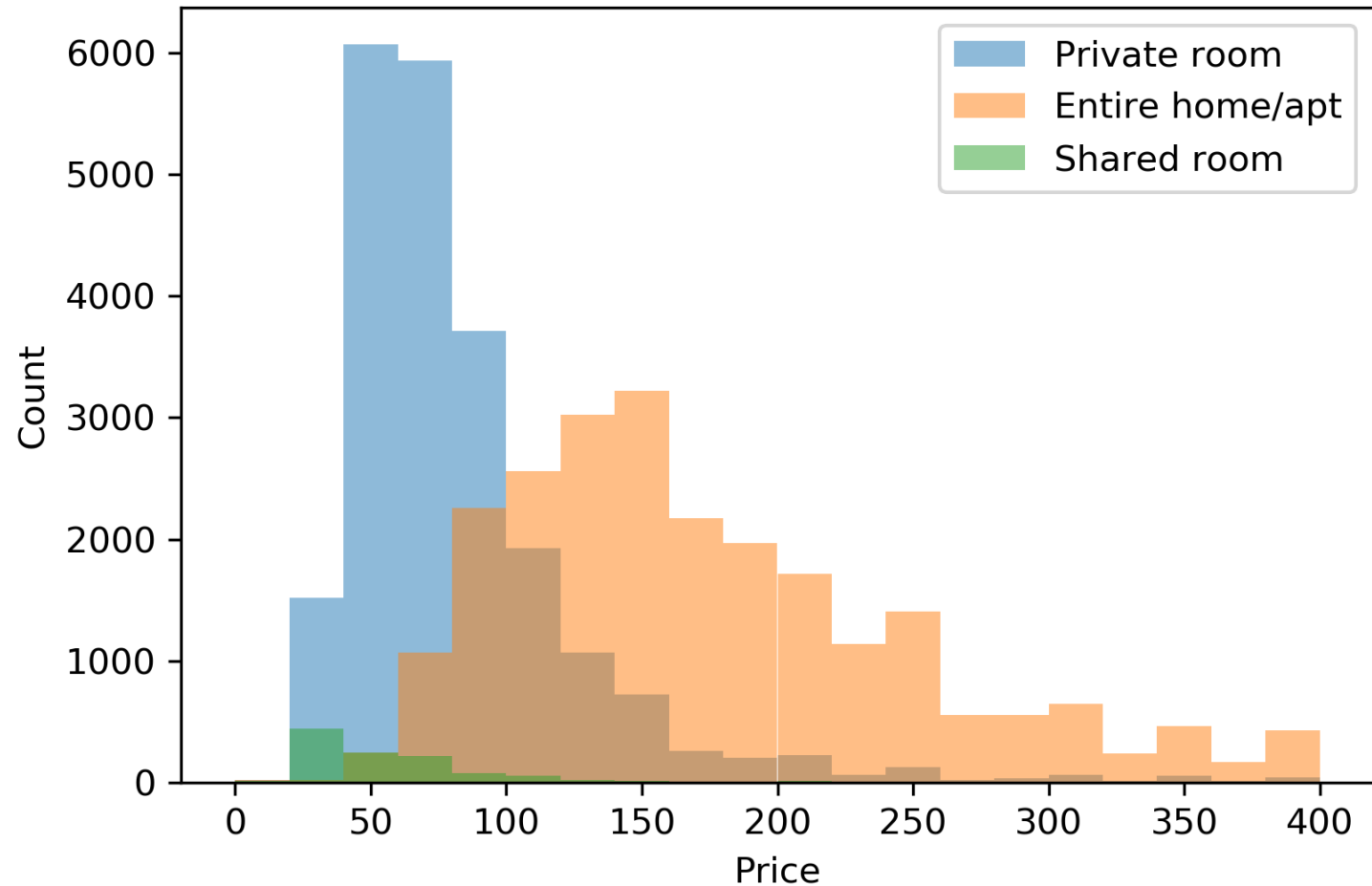


# EDA - Price vs. Region

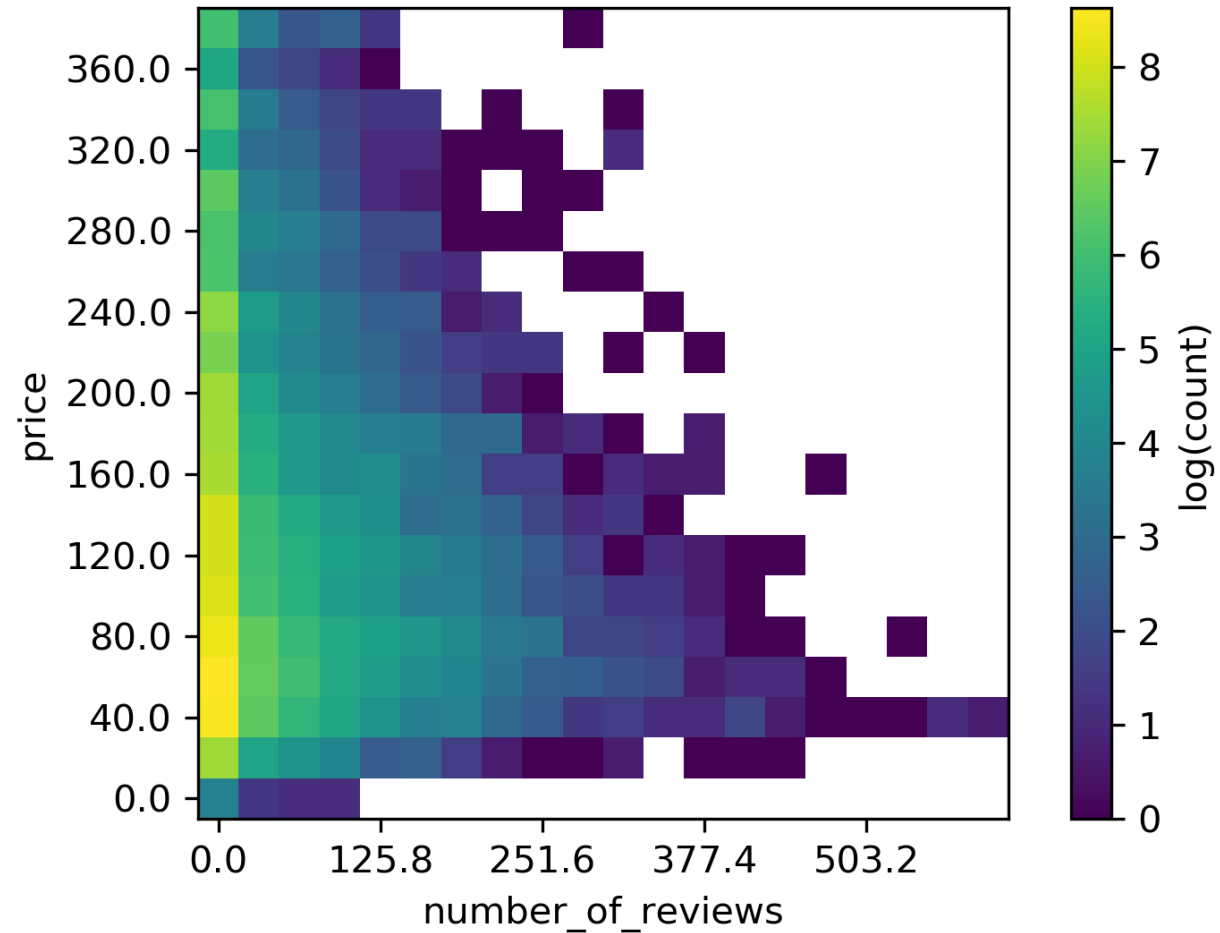




# EDA - Price vs. Room Type

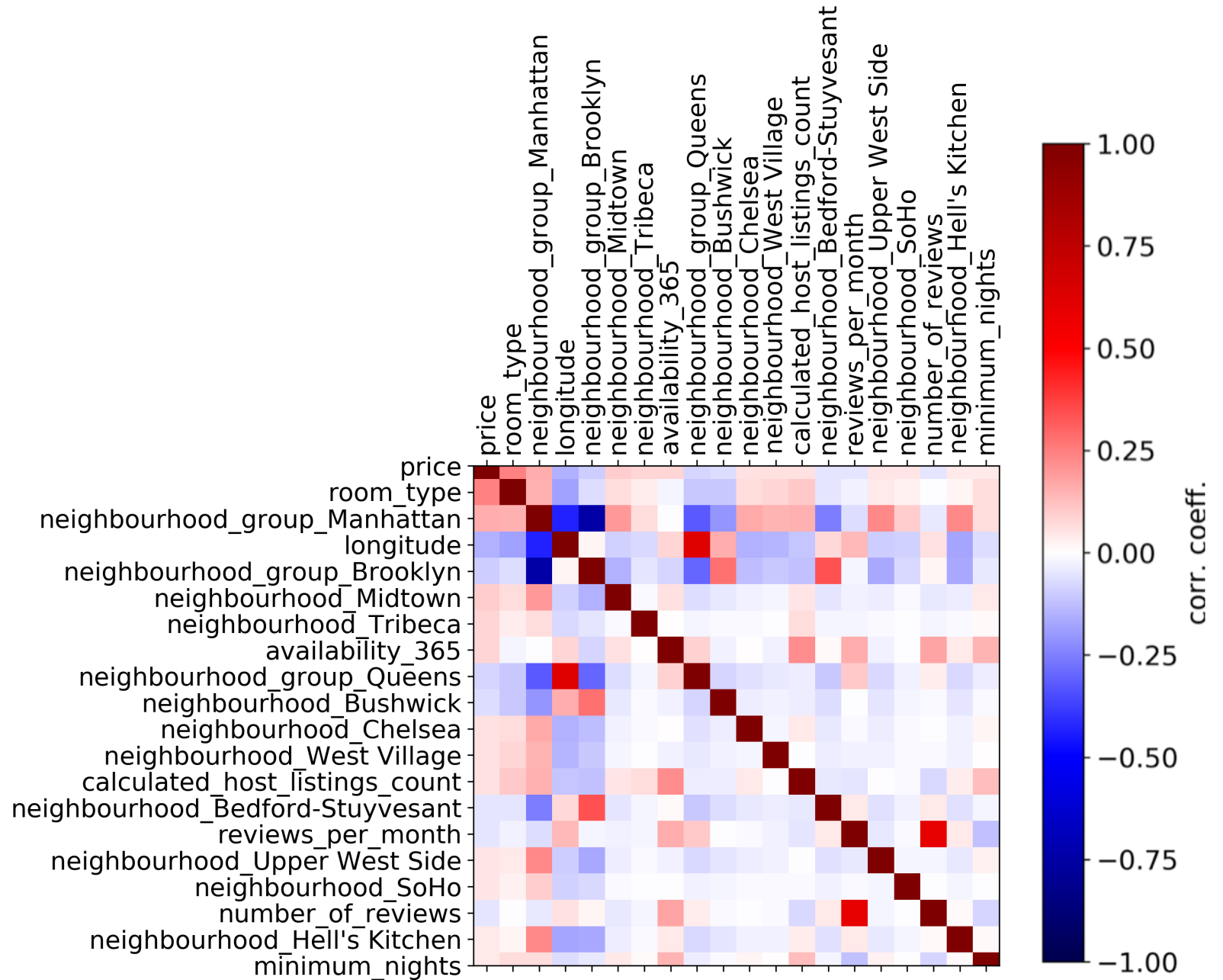


# EDA - Price vs. Number of Reviews





# EDA - The Correlation Matrix



**Thank you!**

**Q & A**