# Exploratory Analysis of Data for Airbnb Listings in NYC

Youngsoo Baek, Irene Yi Ji, Phuc Nguyen
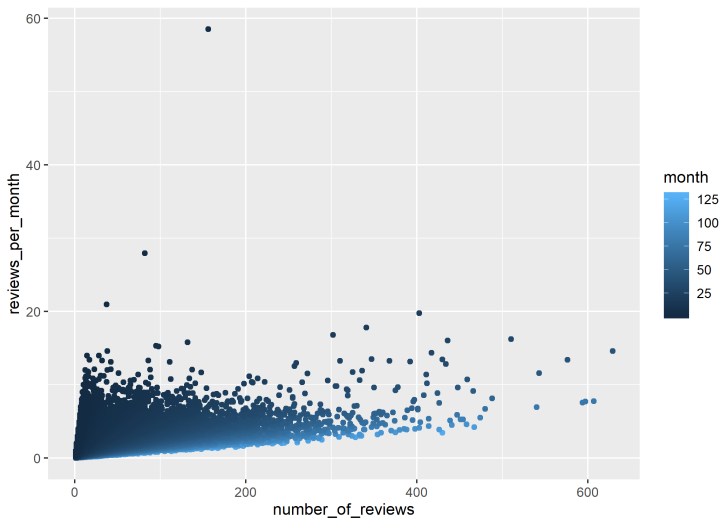
# Introduction

- Data: Airbnb New York City open data collected from 2019, with 48,895 listings and 16 variables.
- Goals:
  - Identify most influential factors for price/popularity
  - Examine heterogeneity across boroughs and neighbourhoods
  - Recommend best location and name for airbnb

# Data Processing

- Remove 14 observations with *minimum_nights* $> 365$
- *Price*: the lowest non-zero value is 10, added 5 to 0's
- *Reviews per Month*: missing values are set to 0 (last review dates are missing and total number of reviews are 0)
- *Last Review*: group by years from 2019 (e.g. 2019 -> 0; 2018 -> 1, etc.)
- *availability_365*: create a new variable *available_spec* to indicate whether the value is 0
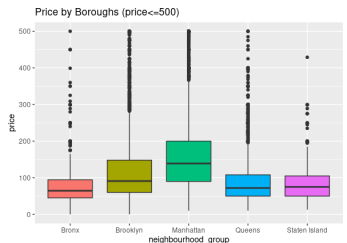
# What is a Valid Metric for Popularity?



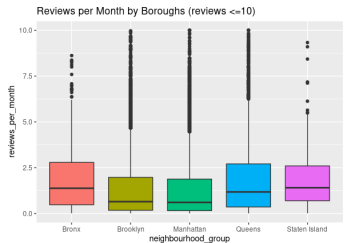▶ **Monthly reviews** adjusts for the history of a listing (albeit not perfectly)

# Heterogeneity of Price / Popularity across Boroughs

- Create new variables "Price Level" and "Popularity Level":
  - "Below Q1" for values < 25th Percentile
  - "Between Q1 and Q3" for values from 25th to 75th Percentile
  - "Above Q3" for values > 75th Percentile
- Create contingency table and conduct Chi-squared Test for Homogeneity

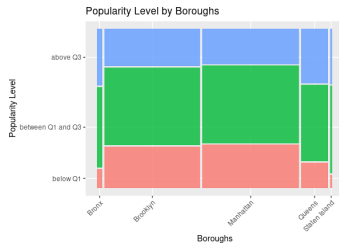# Heterogeneity of Price / Popularity across Boroughs
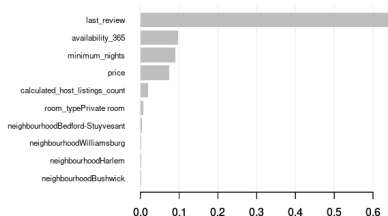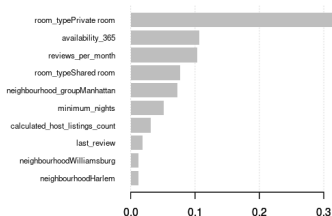


▶ Small p-value suggests heterogeneity across boroughs.

# Heterogeneity of Room type across Boroughs



▶ Small p-value suggests heterogeneity across boroughs.

# Price: XGBoost for Important Variables
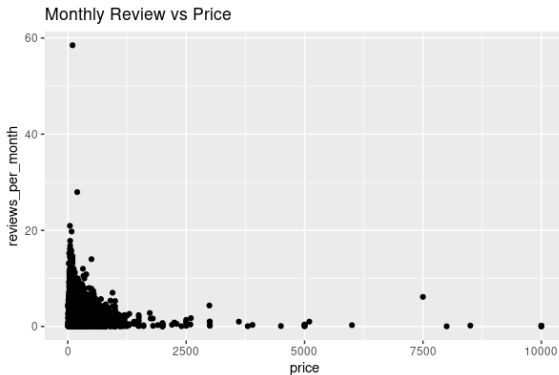


▶ The most influential factors for price of airbnb include: room type, availability, monthly reviews, boroughs, etc.

▶ The most influential factors for popularity of airbnb include: last review, availability, minimum nights, price, etc.
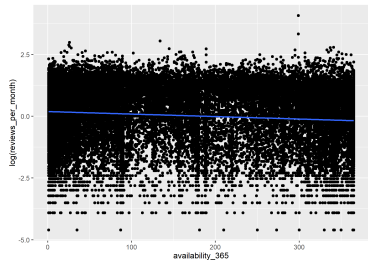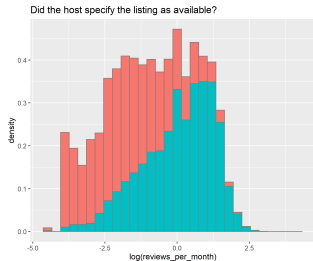
# EDA - Price and Popularity

- ▶ From XGBoost outputs, price and popularity are closely related, both being an important variable of the other.

- ▶ The plot below shows a negative correlation between them:



Monthly Review vs Price

- ▶ We may consider model them as bivariate reponse.

# Possibly Unreliable Predictors

# Modeling: Bivariate Mixed Effects Regression
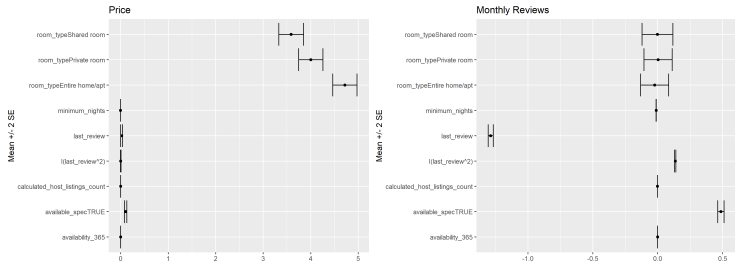
▶ Varying intercept model: Random effects for each neighbourhood, and each borough

For the $i-$th observation in neighbourhood $j$, in borough $k$,

$$\begin{pmatrix} \text{Price}_{k[j[i]]} \\ \text{Monthly review}_{k[j[i]]} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\beta}_1^T \mathbf{X}_i \\ \boldsymbol{\beta}_2^T \mathbf{X}_i \end{pmatrix} + \boldsymbol{\eta}_{k[j]} + \boldsymbol{\theta}_j + \boldsymbol{\epsilon}_{k[j[i]]},$$

$$\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 I_2).$$

▶ Quadratic term for how "old" a listing is included
▶ Observations with no reviews excluded (21% of the data)

# What Are the Important Predictors for Price/Popularity?



- ▶ In terms of magnitude, not significance, **room type** for price, and **last review** for popularity

- ▶ Apartments > Pvt room > Shared room for price, and more popular if the listing is young
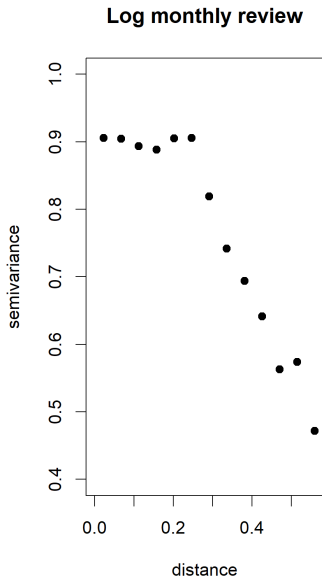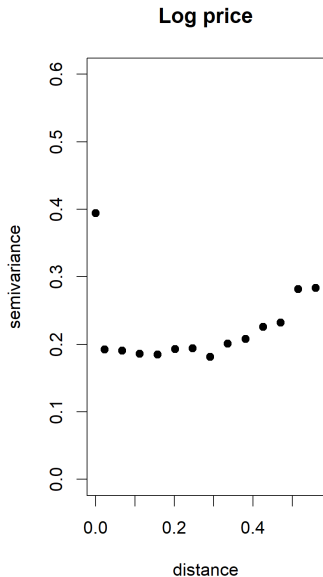
# Estimates for Group Heterogeneities

|  | variableprice | variablereviews_per_month |
|---|---|---|
| variableprice | 0.0285710 | -0.0073307 |
| variablereviews_per_month | -0.0073307 | 0.0369894 |

|  | variableprice | variablereviews_per_month |
|---|---|---|
| variableprice | 0.0796294 | -0.0238452 |
| variablereviews_per_month | -0.0238452 | 0.0124883 |

▶ Many coefficients for significant predictors (adjusted for other variables) are swamped by the variability within/between different neighborhoods and boroughs

▶ Negative correlation between coefficients for price and popularity

# Examining Spatial Correlation of the Residuals

# Possible Insights

▶ When two listings are very close (identical coordinates), the market effect takes sway over all others. One potential customer is being sapped away from one listing to another.

▶ As a result, closer things have more dissimilar popularity measures. As distance increases, however, the effect becomes less severe and association between a listing's features and sales becomes noticeable.

▶ However, price is relatively "inelastic"; unless two listings are extremely close to each other, the hosts' pricing policy remains relatively indifferent to their neighbors, adjusted for other features of a listing.

▶ Hence, we observe no evidence of spatial correlation, conditional on what neighborhood a listing belongs to, except in extreme proximity (high semivariogram for price).

# Text Analysis for Listing Names

(. . . Phuc's analysis. . . )

# Limitations and Further Work

- Including varying slopes calls for strong shrinkage

- Care is needed for spatial covariance models: "soft" adjacency matrix for neighborhoods/boroughs, negative autocorrelation, etc.

- Missing data/latent space model for `availability_365`

- Nonparametric approach for bivariate model