

# Patterns of Airbnb Listings in NYC

Frances Hung, Yunran Chen, Keru Wu

## Abstract

Airbnb home rental listings vary in price and popularity; one question pertinent to hosts is which listings are most successful. We explore the relationships between certain rental characteristics (including neighbourhood location) and listing price/popularity in NYC.

## 1. Introduction

Airbnb is a platform provides more personalized home rentals for travelers compared with hotels. Our data observations consist of 48,895 individual Airbnb listings in New York City. Each listing observation contains the following variables: host ID, neighbourhood group, neighbourhood, longitude/latitude, available days of the listing in a year, room type, price, minimum nights required, number of reviews, and reviews per month.

From the perspective of a host, we are interested in exploring the patterns in price and popularity. Specifically, we are interested in (1) quantifying the influential factors in the price/popularity and evaluating their influence (2) exploring the most valuable neighborhoods with adjusting the influential factors (3) choose a location and set a price for the listing (4) name the listing.

## 2. Materials and Methods

Since the price and popularity are strongly related to the location of listings (Figure 16, 17) and neighborhoods provides a natural boundary for spatial characteristics of listings, we consider a multilevel conditional autoregressive Bayesian model (CARBayes)(Lee 2013) based on neighborhood units as follows:

$$Y_{kj} | \mu_{kj} \sim f(y_{kj} | \mu_{kj}, \nu^2), \quad k = \text{neighbourhood} = 1, \dots, K \\ j = \text{listings} = 1, \dots, m_k$$

$$g(\mu_{kj}) = x_{kj}^T \beta + \psi_{kj} \\ \psi_{kj} = \phi_k + \zeta_{kj}$$

, where  $\beta$  represents the potential effect of predictor  $x_{kj}$ , with a prior  $\beta \sim N(\mu_\beta, \Sigma_\beta)$ .  $\phi_k$  and  $\zeta_{kj}$  represents the neighborhoods' effect and individual effect. We consider a autoregressive prior for  $\phi_k$ :

$$\phi_k | \phi_{-k} \sim N\left(\frac{\rho \sum_{l=1}^K w_{kl} \phi_j}{\rho \sum_{j=1}^K w_{kl} + 1 - \rho}, \frac{\tau^2}{\rho \sum_{j=1}^K w_{kl} + 1 - \rho}\right)$$

where  $w_{kl}$  is known from data with  $w_{kl} = 1$  denotes neighbourhood  $k$  is adjacent to neighbourhood  $l$ , and 0 otherwise;  $\rho \sim U(0, 1)$  capture the relation between neighbourhood effects. This prior captures the spatial structure among neighborhoods for each neighborhoods' effect is centered at the weighted sum of the effects from its neighbors.

We consider `log(price)` and `log(1+review_per_month)`(popularity) as response variable and model them separately. We include room type, price, minimum nights required, price/popularity as predictors based on EDA results. Additionally, we incorporate the logarithm distance from a listing to the nearest metro station to account for the heterogeneity of individual spatial effect within the same neighborhood. Also, we extracted features from names of listings by applying Latent Dirichlet Allocation (LDA) (Blei, Ng, and Jordan 2003) model and introduced these features as predictors.

To carry out text analysis on names of listings, we first conducted a detailed text cleaning and applied Porter's stemmer algorithm to merge the words with the same root. Then we apply Latent Dirichlet Allocation (Blei, Ng, and Jordan 2003) to explore the latent topics. By assigning each word a weight of related topics (e.g. adjectives, locations), we extracted features from the listings' names and included them to multilevel CARBayes model. In addition, we conducted word frequency analysis for different boroughs, different levels of price and use wordcloud to visualize the results.

### 3. Results

#### 3.1 Exploratory Data Analysis

Figures 16,17,3 suggest a clear spatial structure for price, popularity and traffic. Most high-priced listings are located in midtown and downtown Manhattan, while some of them also lie in the part in Brooklyn that is closed to Manhattan. Similar pattern is discovered for traffic. Most popular neighborhood locates around LGA airport. Figures 4,5 demonstrate that the room type matters for price but not for popularity. We can also see the heterogeneity of room type across boroughs/neighborhoods as shown in Figure 6. A Pearson's Chi-squared test ( $p\text{-value} < 2.2\text{e-}16$ ) is conducted to further support the existence of heterogeneity. Figure ?? suggest a non-linear effect on price/popularity.

#### 3.2 Data Preprocessing

We remove 11 listings with price equal to 0 and impute 0 for `reviews_per_month` for 10052 listings with NA value. We consider the logarithm transformation for the response variable `price` and `reviews_per_month`. The choice of predictors are based on the results from EDA. Specifically, we categorize the `minimum_night` into 5 groups by the first three days and weeks accounting for the nonlinear association. To obtain the adjacency matrix of neighborhoods in New York In addition, we incorporate shape files for neighbourhoods in New York <sup>1</sup> and reallocate the listings' neighborhood based on the latitude and longitude. To account for the heterogeneity of spatial effects across listings within the same neighborhood, we introduce a new predictor which is the logarithm of distance from a listing to the closest metro<sup>2</sup>. For text cleaning, we first preprocess the listings' names by transforming them to lower case and removing non-informative characters such as punctuations, stopwords, whitespace, numbers. Then we apply Porter's stemmer algorithm (Porter 2001) for word normalization, allowing for extracting all the common roots of informative words.

#### 3.3 Main Results

According to model coefficient estimation (Fig 9), our multilevel CAR model on price demonstrates the following patterns. Numbers in brackets are median of corresponding coefficients. For room type, entire room (0) is more expensive than private ones (-0.7) and shared ones (-1.1), with shared room being the cheapest. Manhattan (0.57) stands out to be the most luxurious borough, and Bronx (0) has the lowest price. Availability (0.12) is positively related to price while reviews per month (-0.0) is negatively related. In addition, more strict requirement on minimum nights results in lower price, which aligns with our common sense. And longer distance to metro stations also reduces the price (-0.005). Table 1 shows the WAIC of different models, which allow us to choose the most influential factor for price(room type).

Model on popularity (9) has some similarity but is different as follows. Compared to other four boroughs, Queens borough (0.13) has the highest average reviews. Availability still has a positive effect (0.15) while price (-0.12) leads to a negative influence. Moreover, metro distance is no longer

---

<sup>1</sup><<https://data.cityofnewyork.us/City-Government/Neighborhood-Tabulation-Areas-NTA-/cpf4-rkhq>>

<sup>2</sup>Locations of metro stations: <<https://data.cityofnewyork.us/Transportation/Subway-Stations/ark3-7z49>>

significant for predicting popularity. Table 2 shows the WAIC of different models, which allow us to choose the most influential factor for popularity (availability & night).

Heterogeneity across neighbourhoods is shown in Fig 11 and 12. According to Fig 11, most neighbourhoods in Manhattan have higher average price, and their confidence interval is also narrower than others. Among all neighbourhoods, “Midtown South” in Manhattan turns out to be the most expensive one, while “New Drop-Midland Beach” in Staten Island becomes the one with lowest price. On the other hand, 12 indicates that “East Elmhurst” in Queen is the most popular neighbourhood, which makes sense since LaGuardia Airport is located here, and “Co-op City” is the most unpopular one. If we consider top 20 neighbourhoods for price and popularity separately, they have only one intersection at “Yorkville” in Manhattan.

Our text analysis (Fig 13, 14) indicates some critical words: luxury, manhattan, beautiful (Note that we use stemming algorithm so we get stem of words rather than words themselves). We further carry out LDA to find latent topics in listing names. We choose 4 topics which is not too complicated and has a reasonable result (Fig 15). The 4 topics can be categorized as adjectives, locations, Brooklyn related and Manhattan related. If we further add these 4 topics into our model (4 indicators), we conclude that Brooklyn and Manhattan has a positive significant coefficient, while the other two is significantly negative.

Wordcloud (Fig ??) implies some high-frequency words in high-priced listings: luxury, manhattan, apartment, etc.

### 3.4 Sensitivity Analysis

The availability\_365 variable has zero-valued observations which may correspond to hosts who temporarily take their listings off the market. Comparing the distribution of other variables for zero-valued vs. positive-valued availability\_365 observations suggests that the data may be missing at random because we don't see an obvious pattern in missingness. Using MICE (Buuren and Groothuis-Oudshoorn 2010), we impute the data, treating the zero-valued observations as missing values.

Our model using the imputed data had indistinguishable AIC with our model without imputed data. As a result, we choose to use the original dataset and in future work, explore missingness of availability\_365 further.

[The map for neighborhoods' effect cite these two:[1617](#)]

## 5. Discussion

Our multilevel CAR model incorporates both spatial information and individual random effects, which successfully discovers patterns across neighborhoods and account for individual effects. However, our model may lose information from exact locations (latitude, longitude) of listings and we didn't include temporal information in our model.

From this perspective, a direct extension of our model is using spatial temporal model that incorporates both information of locations as well as time (e.g. last review). Another promising direction is to use point reference models like Gaussian Process to include distance between each two listings.

We assume linear relationship between other predictors and outcome, which might not be the case. Therefore considering nonlinear model using spline regression such as GAM would be more reasonable. Another critical part is how to impute missing data. Although MICE doesn't perform better than imputing with 0, discovering other imputation methods could be helpful. Moreover, different hosts have different number of listings, we can further try approaches that accounts for their influence (e.g. random effects).

## Appendix

### EDA

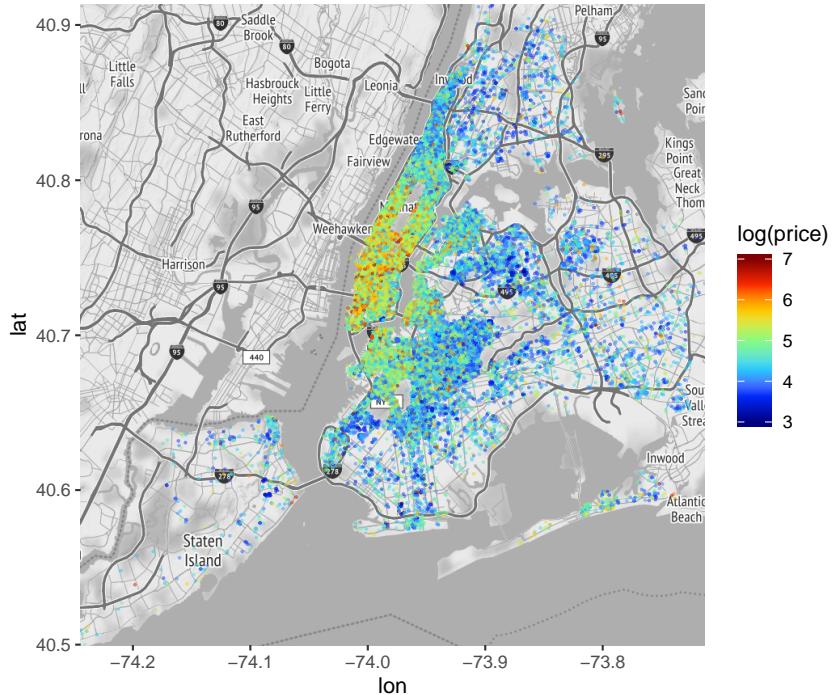


Figure 1: Distribution of  $\log(\text{price})$

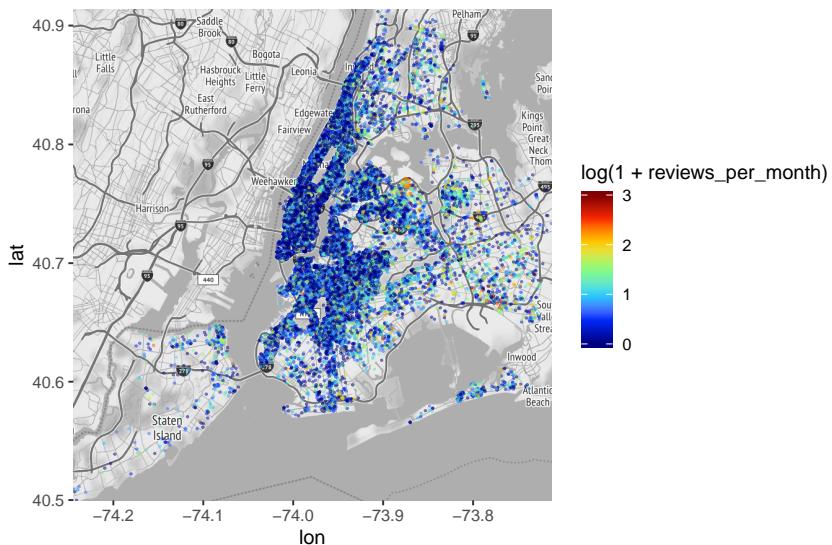


Figure 2: Distribution of  $\log(1 + \text{reviews/mon})$

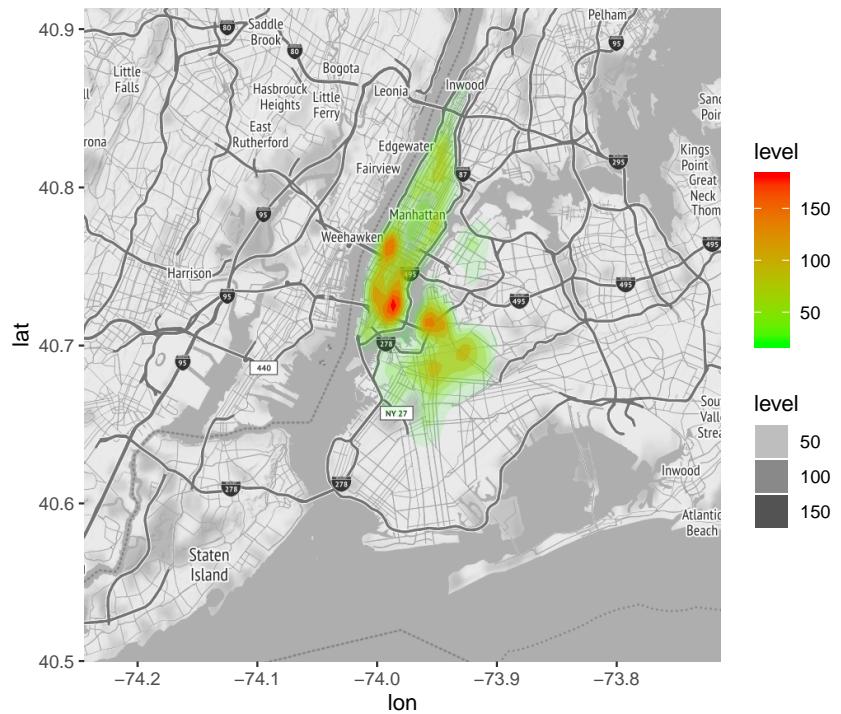


Figure 3: 2D-Density estimation

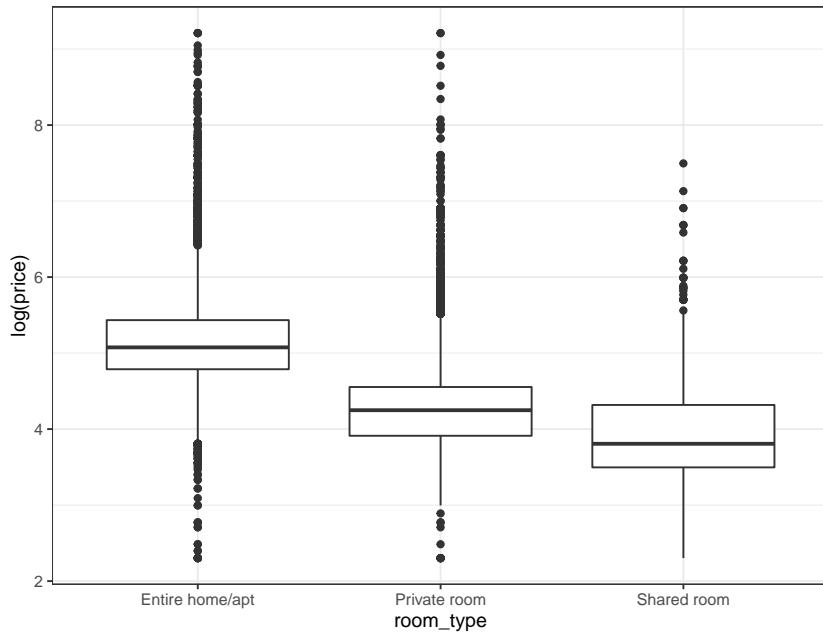


Figure 4: Association between price and room type

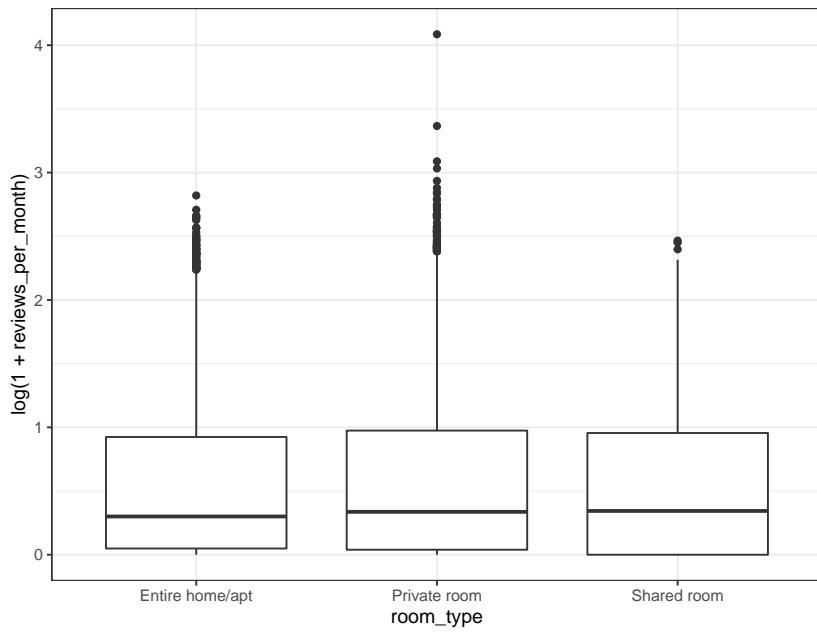


Figure 5: Association between review/mon and room type

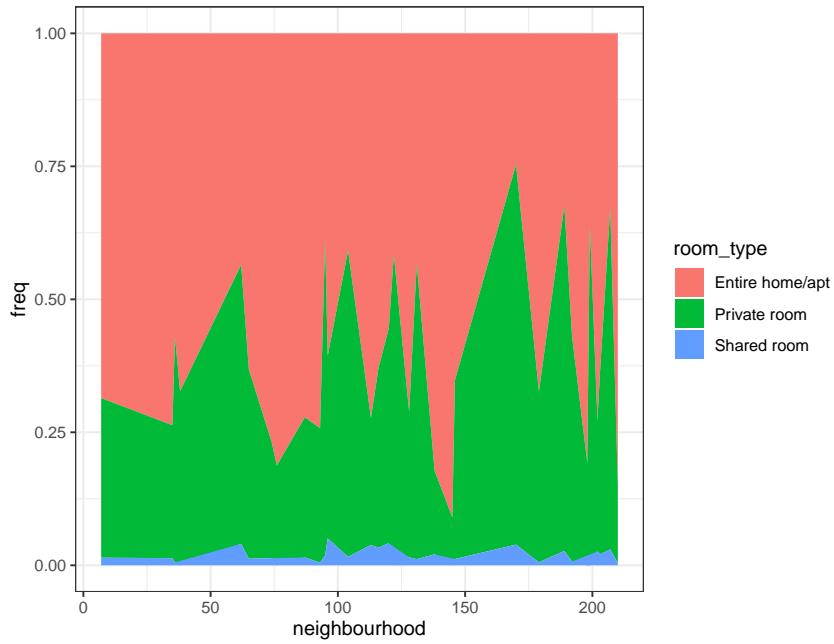


Figure 6: Heterogeneity of Room Type Across Neighborhoods (Manhattan)

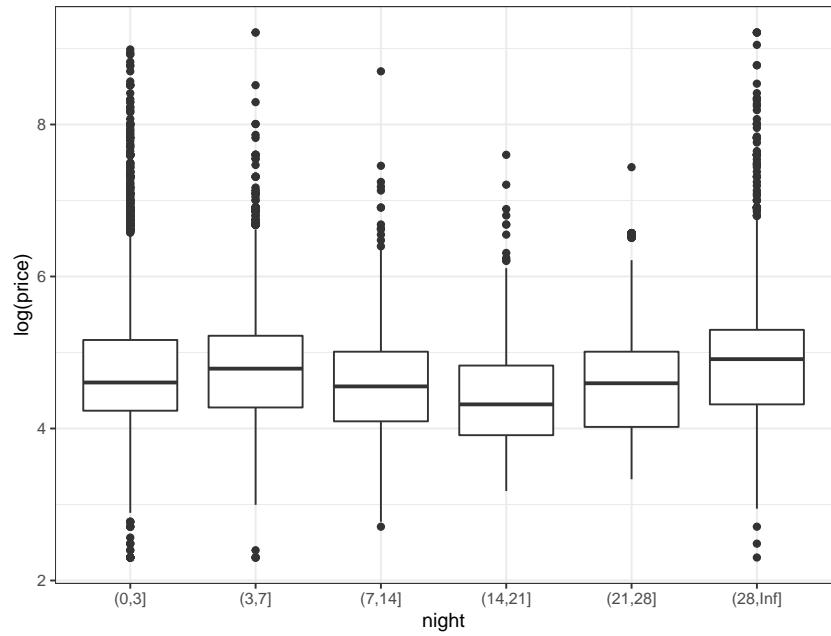


Figure 7: Association between price and minimum night

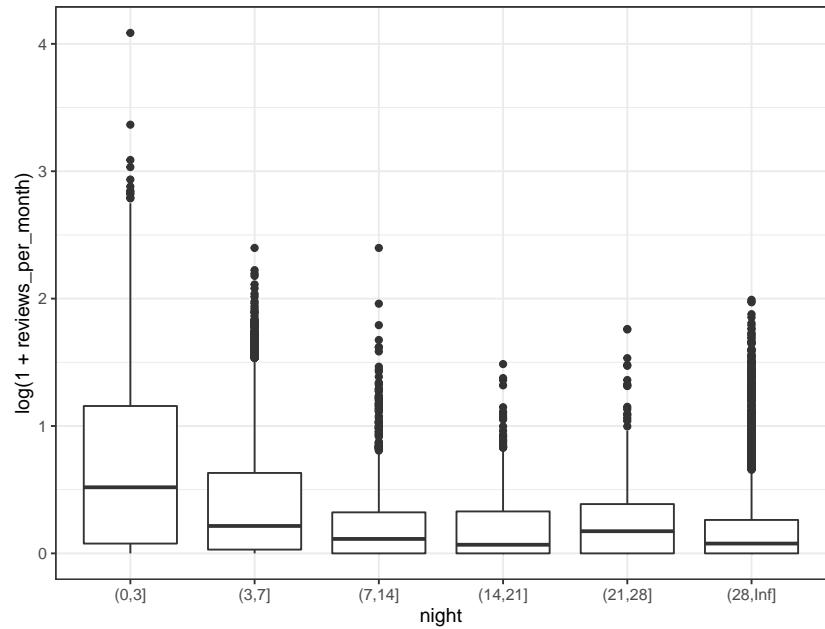


Figure 8: Association between review/month and minimum night

**CARBayes**

	Median	2.5%	97.5%
(Intercept)	4.8153	4.7443	4.8862
room_typePrivate room	-0.7238	-0.7322	-0.7142
room_typeShared room	-1.1091	-1.1379	-1.0836
neighbourhood_groupBrooklyn	0.1874	0.1089	0.2657
neighbourhood_groupManhattan	0.5775	0.4893	0.6526
neighbourhood_groupQueens	0.0964	0.0280	0.1787
neighbourhood_groupStaten Island	0.0404	-0.0698	0.1578
availability_365	0.1174	0.1129	0.1222
log(1 + reviews_per_month)	-0.0919	-0.1008	-0.0835
night(3,7]	-0.0758	-0.0871	-0.0646
night(7,14]	-0.2247	-0.2490	-0.2005
night(14,21]	-0.2865	-0.3193	-0.2503
night(21,28]	-0.2536	-0.3088	-0.2053
night(28,Inf]	-0.3288	-0.3452	-0.3141
metrodist	-0.0054	-0.0124	0.0017
topic1TRUE	-0.0655	-0.0767	-0.0532
topic2TRUE	0.0434	0.0270	0.0608
topic3TRUE	-0.0164	-0.0270	-0.0063
topic4TRUE	0.0283	0.0175	0.0391

Figure 9: CAR Model on price - Model Summary

	Median	2.5%	97.5%
(Intercept)	1.2715	1.1960	1.3567
room_typePrivate room	-0.1425	-0.1538	-0.1303
room_typeShared room	-0.2697	-0.3008	-0.2335
neighbourhood_groupBrooklyn	0.0065	-0.0621	0.0646
neighbourhood_groupManhattan	0.0026	-0.0746	0.0660
neighbourhood_groupQueens	0.1284	0.0507	0.1882
neighbourhood_groupStaten Island	0.0491	-0.0545	0.1513
availability_365	0.1508	0.1457	0.1553
log(price)	-0.1153	-0.1255	-0.1062
night(3,7]	-0.2439	-0.2560	-0.2314
night(7,14]	-0.4046	-0.4324	-0.3760
night(14,21]	-0.4521	-0.4925	-0.4141
night(21,28]	-0.4421	-0.5008	-0.3836
night(28,Inf]	-0.6003	-0.6175	-0.5833
metrodist	0.0007	-0.0071	0.0081
topic1TRUE	-0.0537	-0.0678	-0.0401
topic2TRUE	0.0099	-0.0112	0.0298
topic3TRUE	0.0006	-0.0115	0.0115
topic4TRUE	0.0318	0.0201	0.0447

Figure 10: CAR Model on popularity - Model Summary

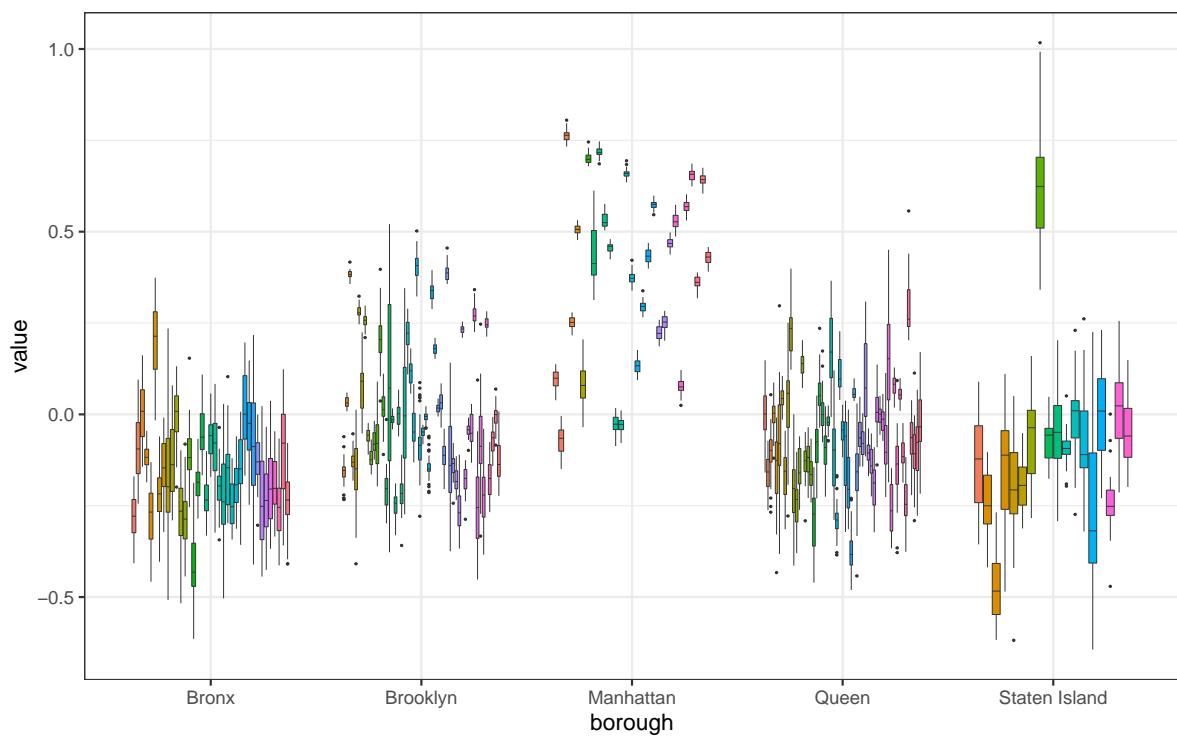


Figure 11: CAR Model on price - Neighbourhoods

Model	All var	Room type	Availability	Reviews	Night	neighborhood
WAIC	63998	85372	66426	64501	66023	70860

Table 1: WAIC for model on price: without 1 variable

Model	All var	Room type	Availability	Price	Night	neighborhood
WAIC	74803	75370	78011	75297	80749	75881

Table 2: WAIC for model on popularity: without 1 variable

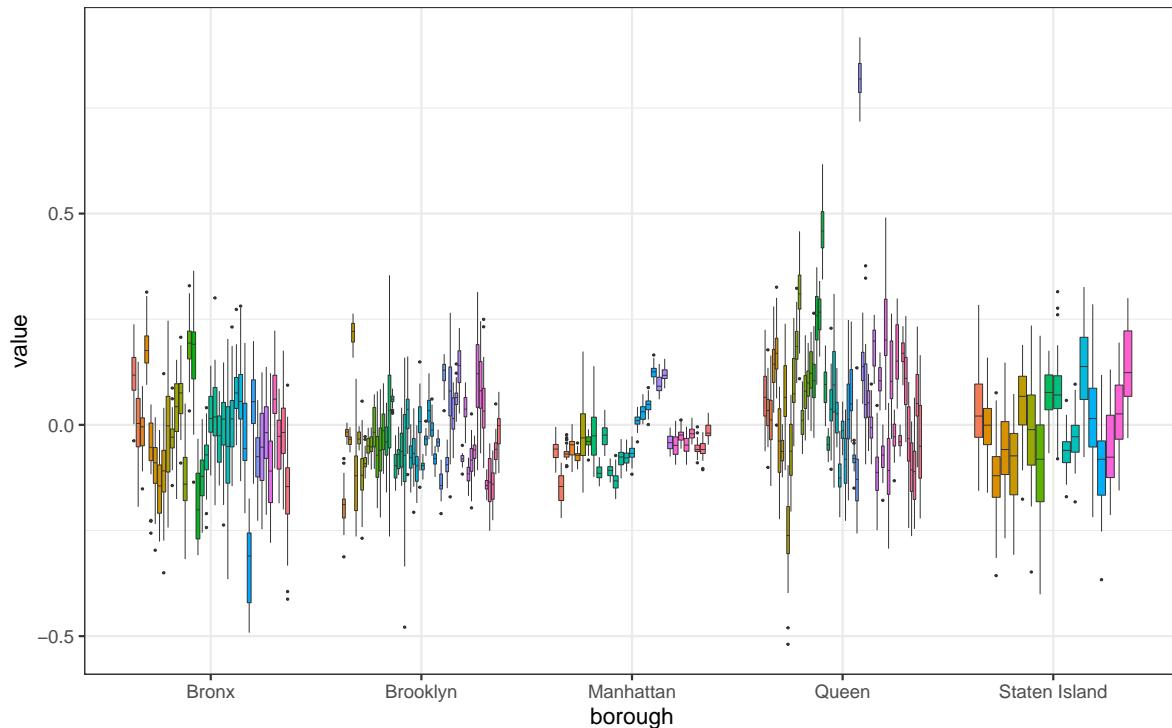


Figure 12: CAR Model on popularity - Neighbourhoods

### Latent Dirichlet Allocation

- Terms:
  - Corpus  $D = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_M\}$
  - Document  $\mathbf{w} = \{w_1, w_2, \dots, w_N\}$
  - Word  $w_i \in \{1, \dots, V\}$ ,  $V$  is total number of unique words.
- LDA Model:  
For all document  $\mathbf{w}$  in  $D$ :
  1.  $N \sim \text{Poisson}(\xi)$
  2.  $\theta \sim \text{Dir}(\alpha)$
  3. For word  $w_n$  ( $n = 1, \dots, N$ )

- (a) choose a topic  $z_n|\theta \sim \text{Multinomial}(\theta)$   
(b) choose a word  $w_n|z_n, \beta \sim \text{Multinomial}(\beta_{z_n})$

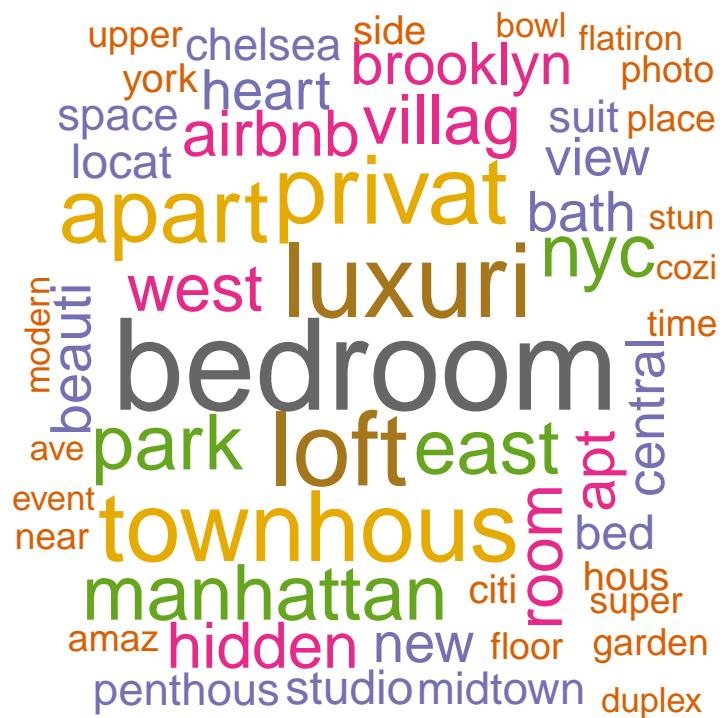
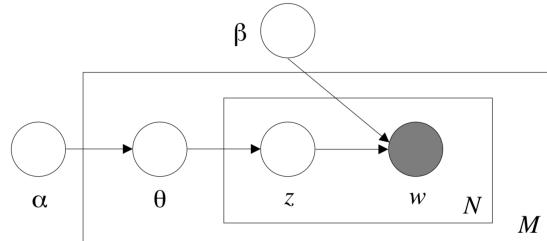


Figure 13: Wordcloud for listings with price > 2000

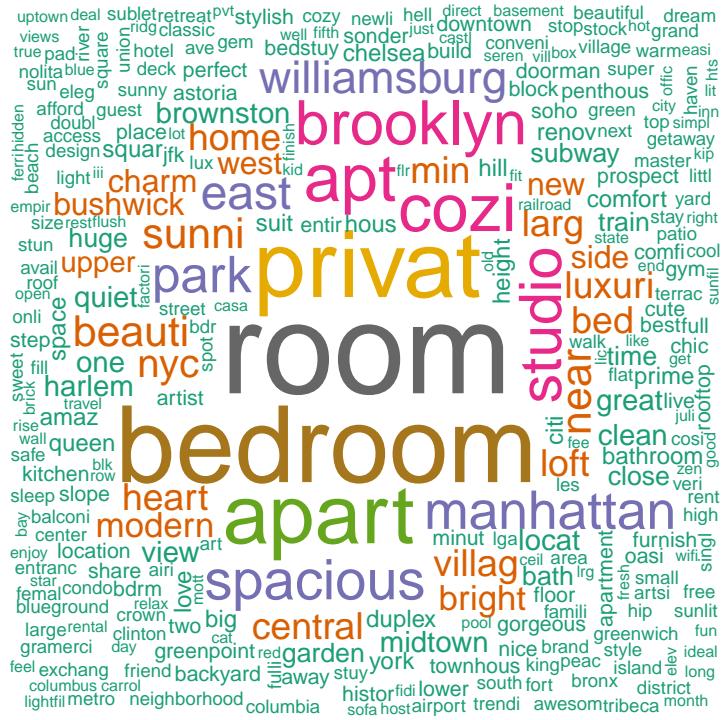


Figure 14: Wordcloud for listings

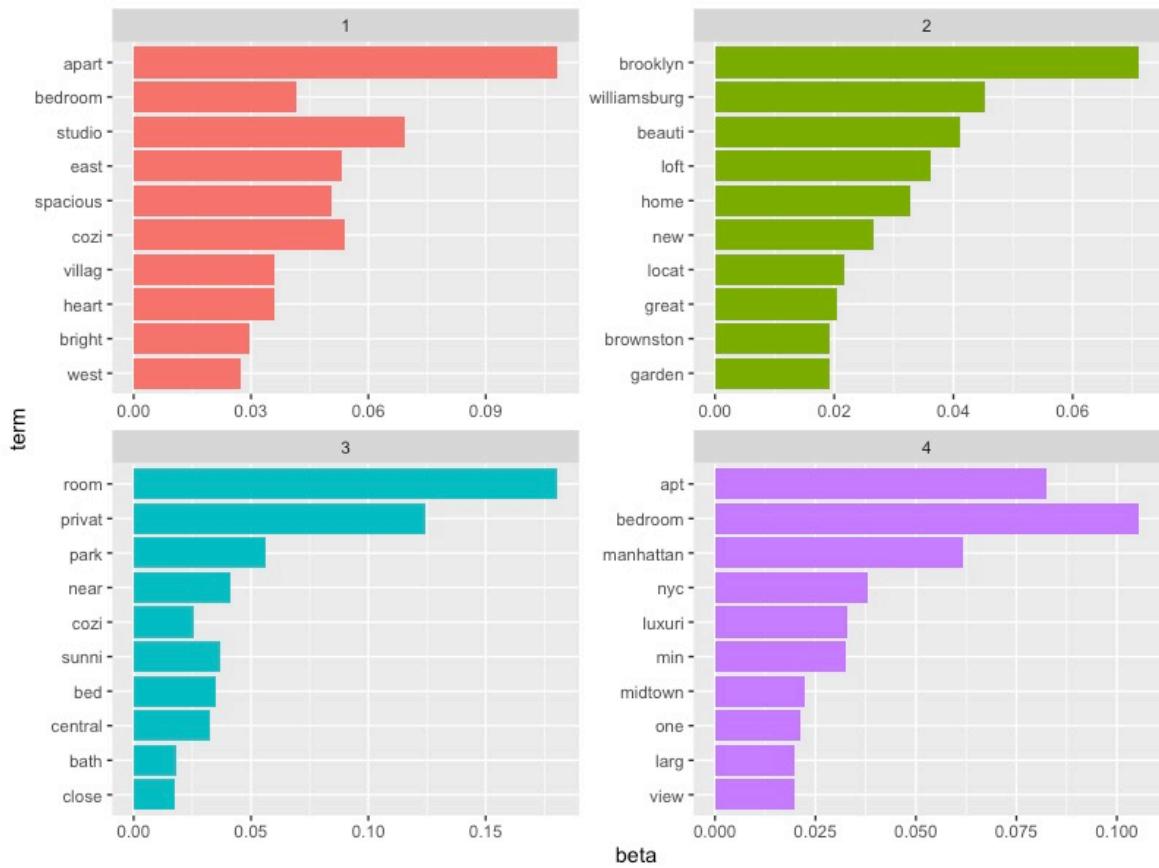


Figure 15: LDA: Top 10 words in each topic

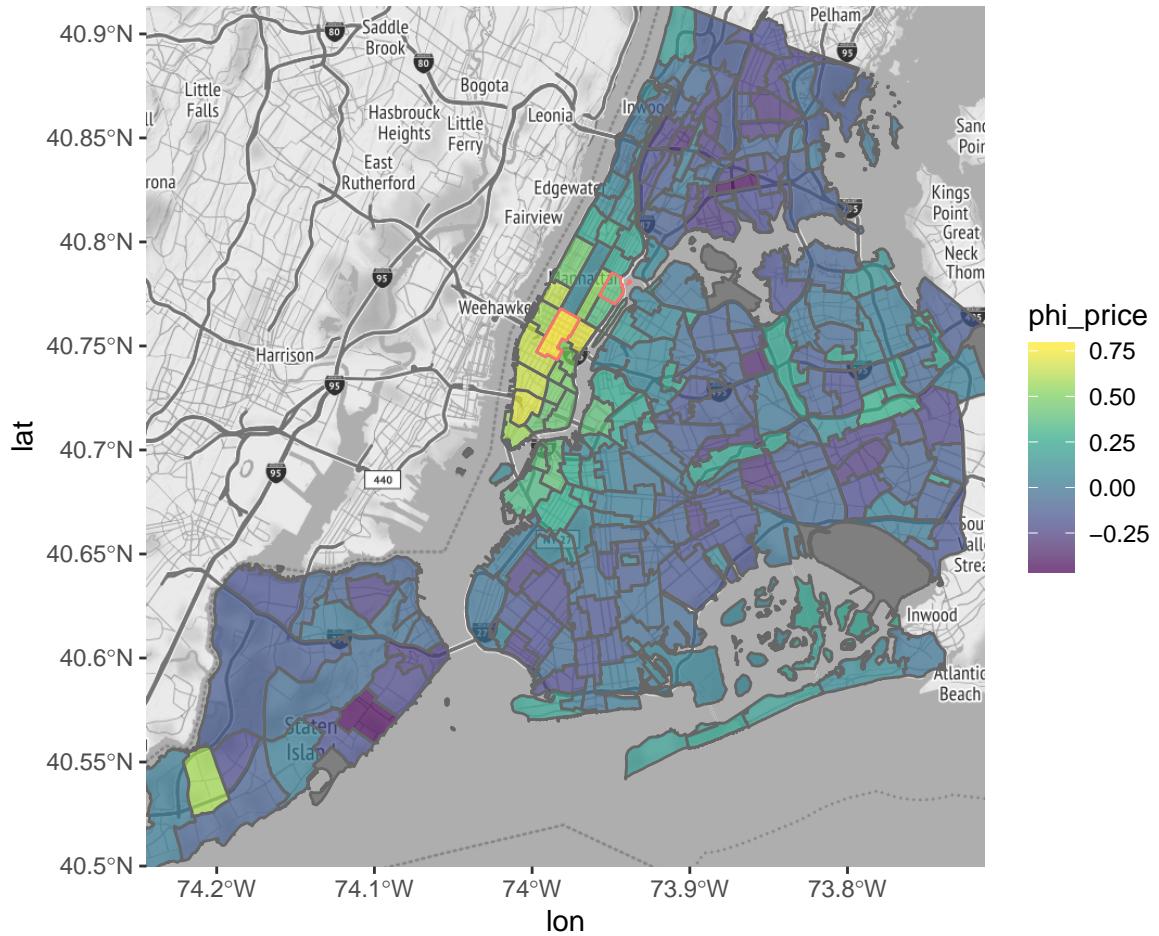


Figure 16: Neighborhoods' effects for price

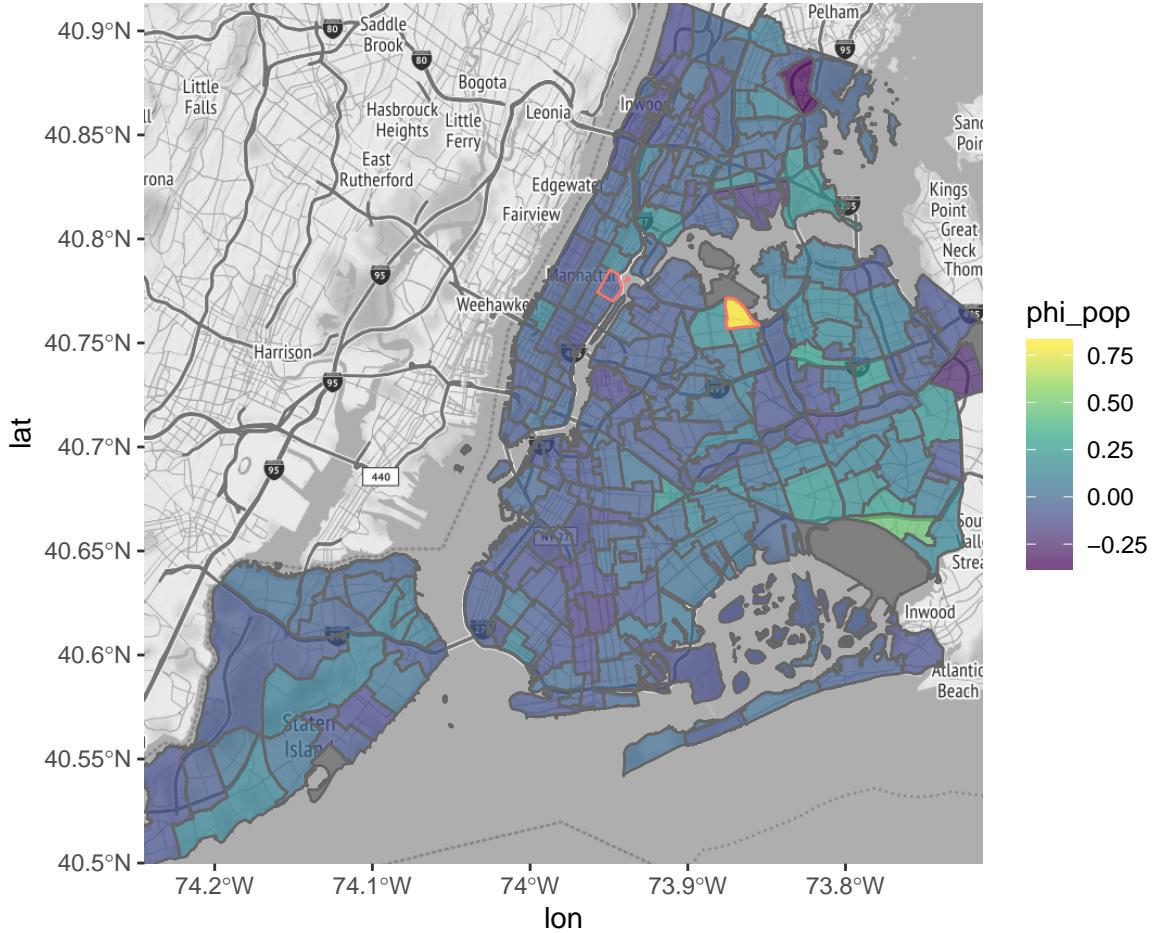


Figure 17: Neighborhoods' effects for popularity

## References

- Blei, David M, Andrew Y Ng, and Michael I Jordan. 2003. “Latent Dirichlet Allocation.” *Journal of Machine Learning Research* 3 (Jan): 993–1022.
- Buuren, S van, and Karin Groothuis-Oudshoorn. 2010. “Mice: Multivariate Imputation by Chained Equations in R.” *Journal of Statistical Software*. University of California, Los Angeles, 1–68.
- Lee, Duncan. 2013. “CARBayes: An R Package for Bayesian Spatial Modeling with Conditional Autoregressive Priors.” *Journal of Statistical Software* 55 (13). American Statistical Association: 1–24.
- Porter, Martin F. 2001. “Snowball: A Language for Stemming Algorithms.”