

# Midterm Project

*Chuning Yuan*

*2019/12/5*

## I.Introduction

Nowadays, Airbnb has become very popular for the travelers due to its unique style and creative design and competitive prices and location than the hotels. Therefore I think it will be interesting to explore the data set of Airbnb, we can look at how are the price for a night stay for each room were affected by other variables such as reviews, number of bedrooms, minimum stays etc. This project will contain analysis on Airbnb data with EDA(Exploratory Data Analysis) and modeling. From this project, we will be able to have a general idea how to predict the price of rooms on Airbnb website and what's the most influential factor in predicting process.

Airbnb is a privately held global company headquartered in San Francisco that operates an online marketplace and hospitality service which is accessible via its websites and mobile apps. Members can use the service to arrange or offer lodging, primarily homestays, or tourism experiences. I choose the data of San Francisco is also because it is where Airbnb has grown around the world from, and I have the experience of searching Airbnb in bay area. First, I will read in the data and do some visualization to see which predictor will contribute more to the prediction of price. And then the modeling will be multi-level regression using room type and neighborhood and few other factors to predict the price.

## II. Data

### Data source:

The data set was extracted from the tomslee.net website– Airbnb Data Collection: Get the Data. There are zip file for many cities around the world. The zip file holds one or more csv files. Each csv file represents a single “survey” or “scrape” of the Airbnb web site for that city.

Specifically, the data used for this project was collected from the November 2013 to January 2017 in San Francisco. There are 9 variables that will be used in this project. Specifically, they are room id, host id, room type, neighborhood, number of reviews, overall satisfaction, number of accommodates, number of bedrooms and price. We are assuming the potential factors to influence the pricing are the room type, neighborhood, number of reviews, overall satisfaction, number of accommodates and number of bedrooms.

### Data Cleaning

For a large dataset, we want to extract the information as more useful and informative as possible, so we filter the observation with more than 100 reviews because these could be more representative in general, and extract neighborhood with more than 30 observations for the further modeling analysis. Below is the data overview after the cleaning processes.

### Overview of data:

Table 1. The head of the data:

room_id	host_id	room_type	neighborhood	reviews	overall_satisfaction	price
6910758	30920210	Shared room	South of Market	125	5.0	99
259622	329072	Shared room	Financial District	117	4.5	45
229240	329072	Shared room	Financial District	194	4.5	45
70753	329072	Shared room	Financial District	206	4.5	45
4518031	22931450	Shared room	North Beach	138	4.5	56
4519780	22931450	Shared room	North Beach	101	4.5	56

Table 2. Summary of the data:

```
##      room_id      host_id      room_type
##  Min.   : 5858   Min.    : 46   Entire home/apt:1794
## 1st Qu.: 678556 1st Qu.: 1032643 Private room :1955
## Median : 1827653 Median : 4393613 Shared room  : 92
## Mean   : 2689704 Mean    : 9009894
## 3rd Qu.: 4092288 3rd Qu.:11655078
## Max.   :13901641 Max.    :77807259
##
## borough      neighborhood      reviews
## Mode:logical Mission      : 497   Min.    :101.0
## NA's:3841     Castro/Upper Market: 416   1st Qu.:118.0
##              Western Addition  : 328   Median  :142.0
##              Bernal Heights    : 259   Mean    :162.5
##              Haight Ashbury    : 232   3rd Qu.:186.0
##              Noe Valley        : 219   Max.    :513.0
##              (Other)          :1890
## overall_satisfaction accommodates bedrooms      price
## Min.   :4.000      Min.    : 1.000   Min.    :0.000   Min.    : 35.0
## 1st Qu.:4.500      1st Qu.: 2.000   1st Qu.:1.000   1st Qu.: 95.0
## Median :5.000      Median : 2.000   Median :1.000   Median : 125.0
## Mean   :4.825      Mean    : 2.813   Mean    :1.088   Mean    : 140.6
## 3rd Qu.:5.000      3rd Qu.: 4.000   3rd Qu.:1.000   3rd Qu.: 160.0
## Max.   :5.000      Max.    :14.000   Max.    :4.000   Max.    :1000.0
##
## minstay
## Mode:logical
## NA's:3841
##
##
##
##
```

According to the summary we can delete the the column borough and minstay because there is no value in them. After eliminate all the NA, now we have the cleaned data we need, then we can start our EDA process.

### III. EDA

Figure 1. Distribution of room price

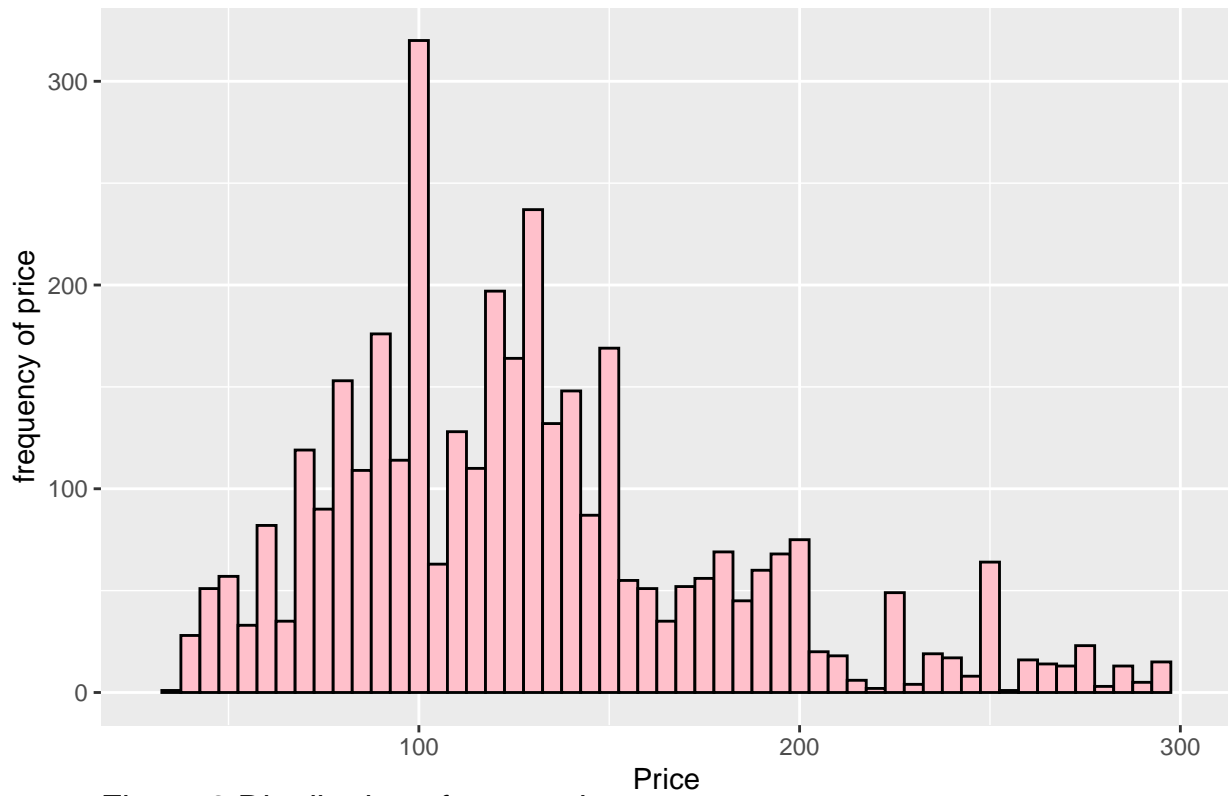
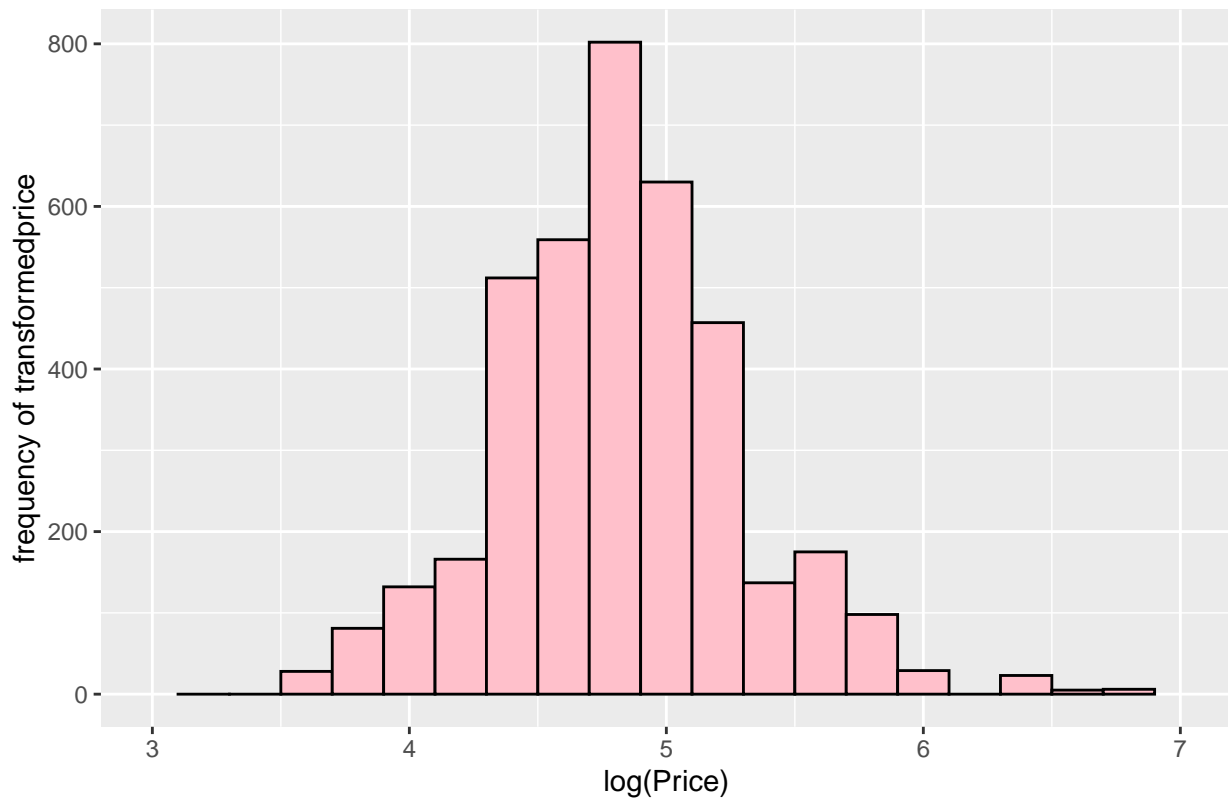


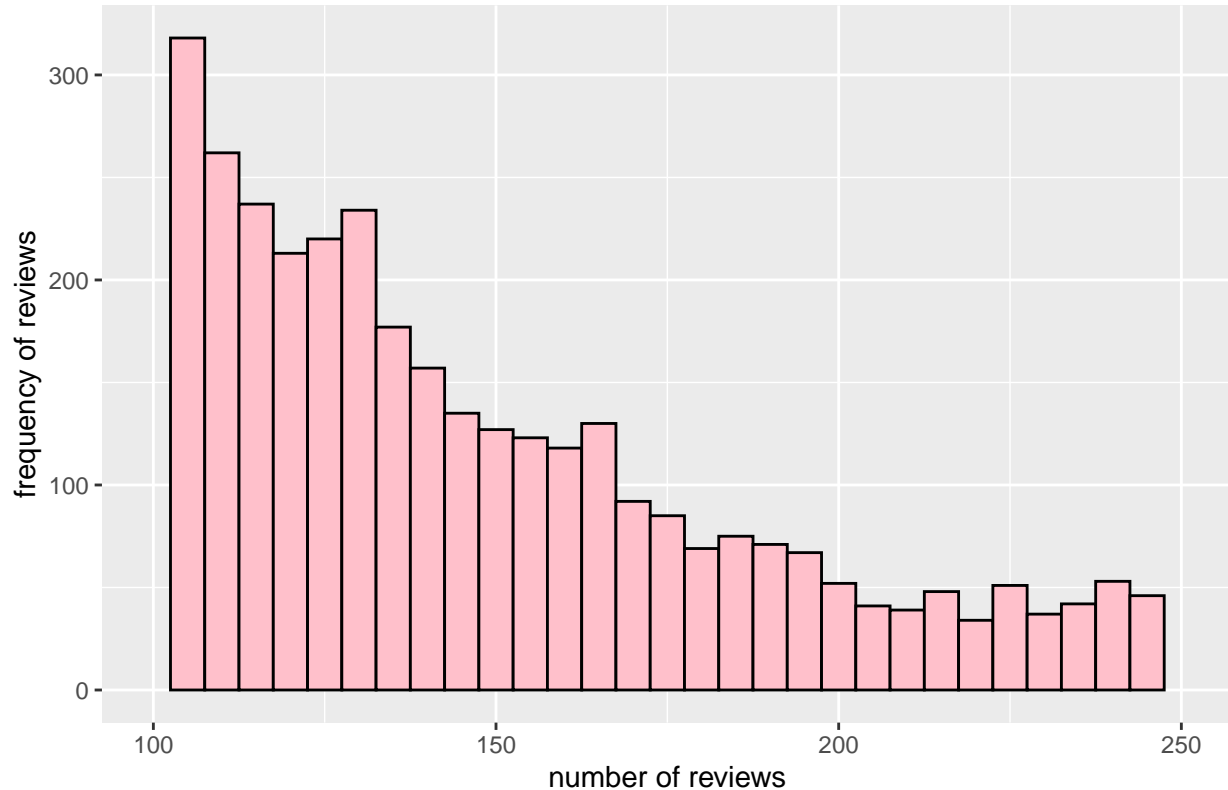
Figure 2. Distribution of room price



Distribution of the room price:

*First, we want to take a look at the room price distribution, and it is obvious that the most popular price are around 100, and between 100-150.*

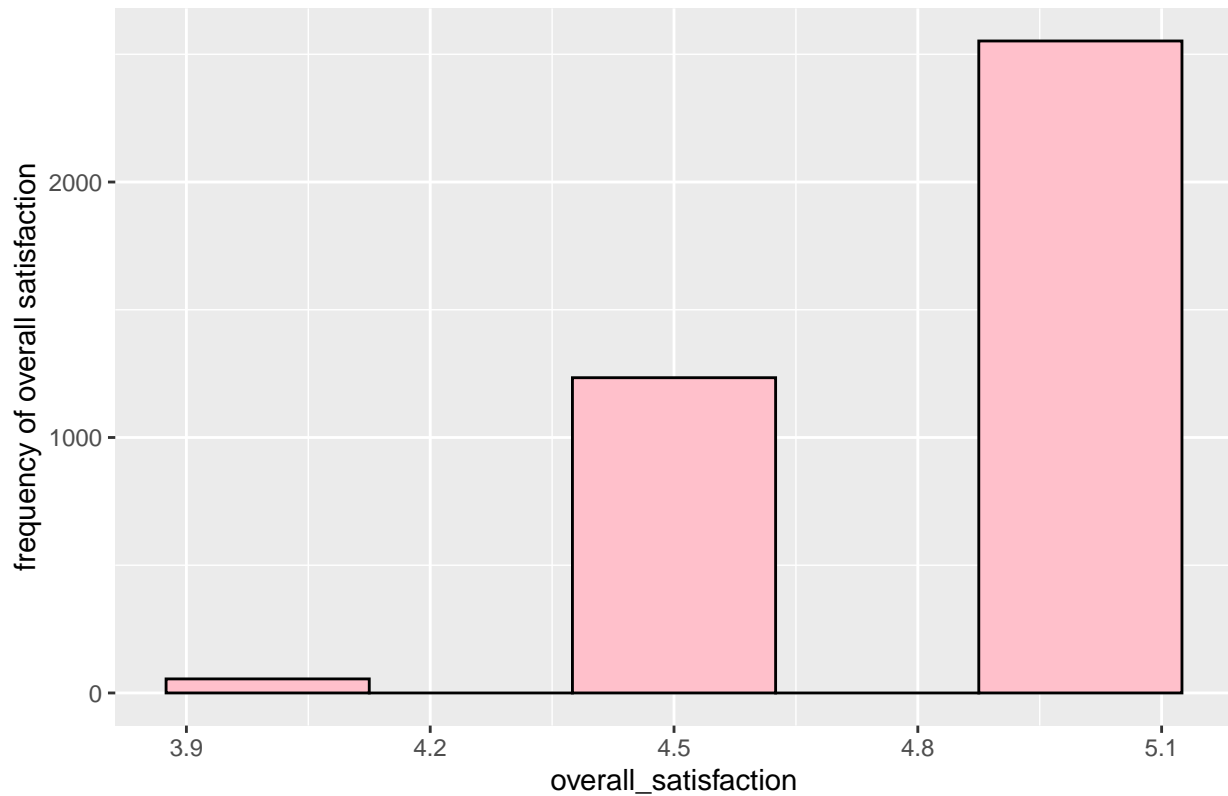
**Figure 3.distribution of number of reviews**



Distribution of number of reviews:

*Since we have already filter the reviews that is less than 100, here the distribution plot can tell us that most of the reviews are around 100 to 175, most of the reviews are under 200. This provide us some general ideas of number of reviews for the San Francisco Airbnb, so we can know how much reviews to expect when we are choosing the Airbnb from the website based on the reviews.*

Figure 4.distribution of overall satisfaction



Distribution of overall satisfaction:

*In this histogram plot, most of overall rating is around 4.5 and 5. There are also a small portion of people rate the room 3.9 to 4 star. Overall, customers are satisfied with most of rooms in San Francisco area.*

Figure 5.Average number of reviews per neighbor

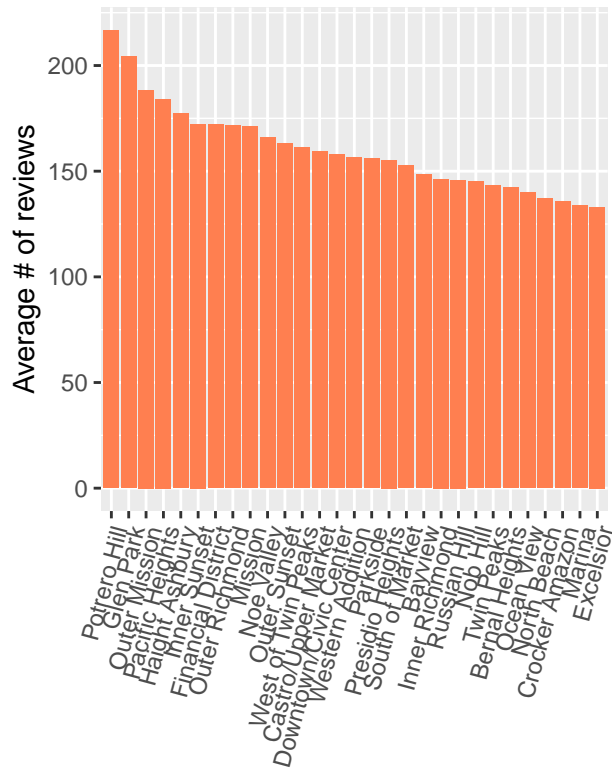
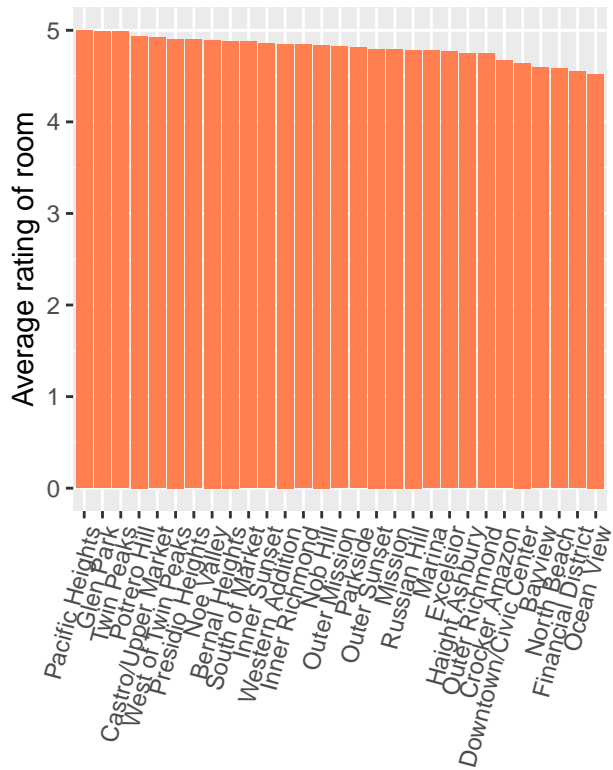


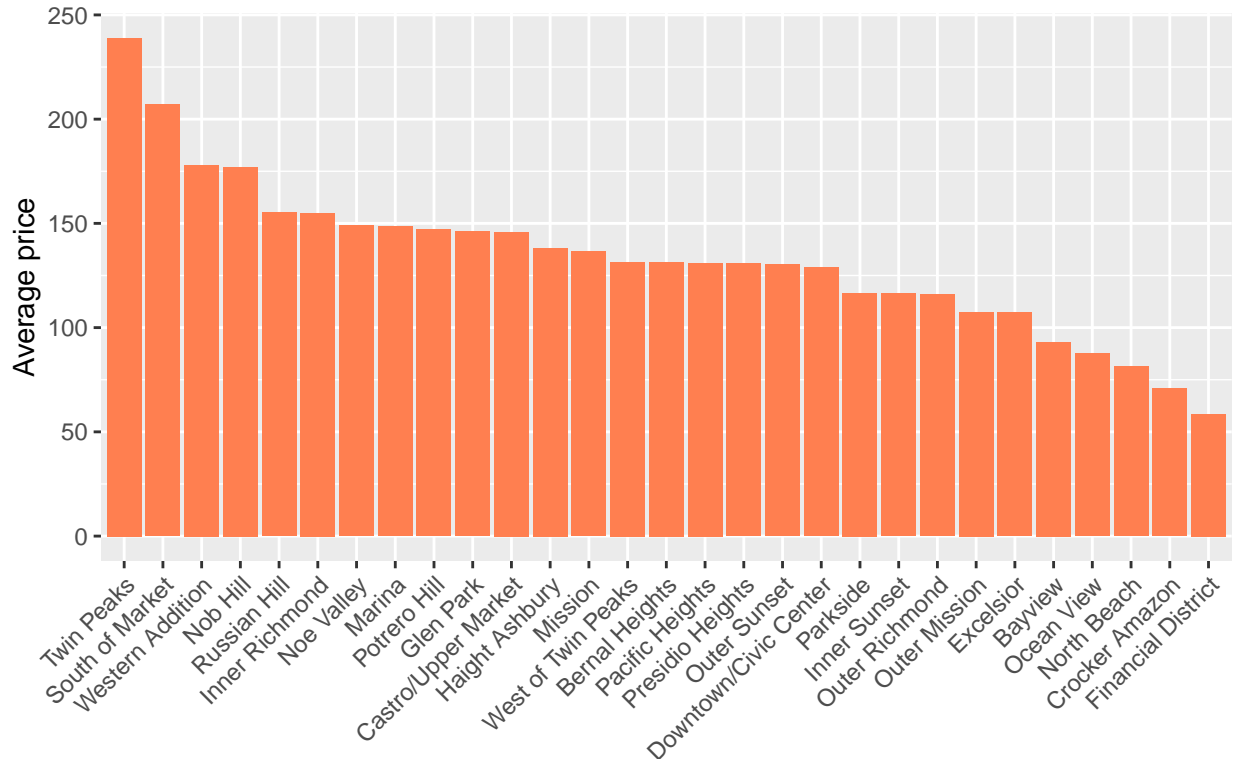
Figure 6.Average rating of airbnb rooms per neighbor



reorder(neighborhood, -reviews)

reorder(neighborhood, -Avg\_rating)

Figure 7.Average price per neighborhood

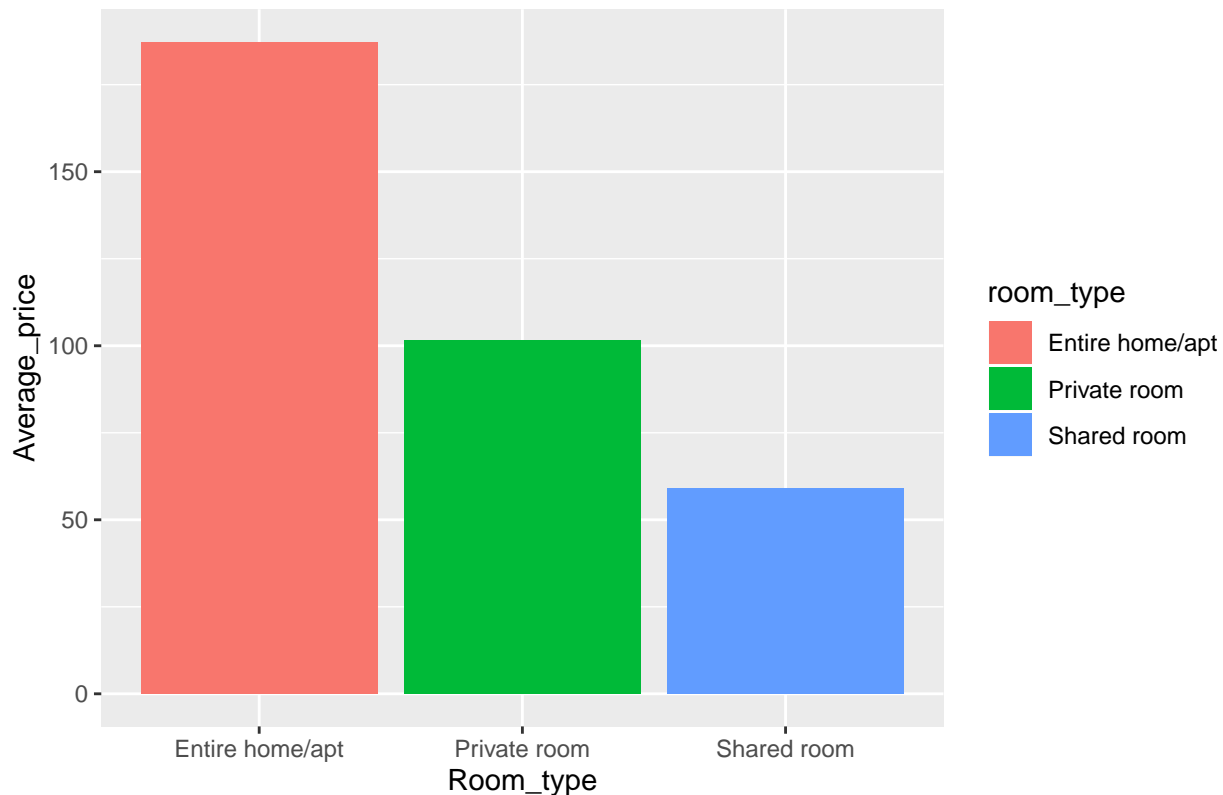


reorder(neighborhood, -Avg\_price)

Pricing, reviews and average rating with different neighborhood:

From this plot, the Potrero Hill area has the highest number of review (over 200), while Excelsior has the lowest average number of reviews beside those less than 100 reviews. We can observe that the average number of review do vary a lot by neighborhood. Although there is not much diffence of rating among different neighborhood, still the neighborhood of the Airbnb room could be an influence predictor based on figure 6. We need to include this predictor in the model to see whether rating is a significant for predicting price of rooms.

Figure 8. Average price with different room types



Average rating and price with different room types:

This result is consistent with our common sense and meaning the pricing of Airbnb in San Francisco are reasonable as the bigger the room type has the higher average price.

Testing other predictors

The accomodates and bedrooms could be two correlated terms in the model, because the number of bedroom will limit the number of customers served. So the correlation test will be conducted in next step.

```
##
## Pearson's product-moment correlation
##
## data: SFAirbnb$accommodates and SFAirbnb$bedrooms
## t = 43.983, df = 3839, p-value < 2.2e-16
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## 0.5574291 0.5995021
## sample estimates:
## cor
## 0.5788508
```

The accomodates and bedrooms could be two correlated terms in the model, because the number of bedroom is related to the number of customers served for the room. So we can do a correlation test for this two predictors.

The  $p$ -value of this test is  $2.2e-16$ . Reject the null hypothesis. So correlation between the those two variables is significant. Therefore we can add the correlation term into the model to test whether this influence term is significant.

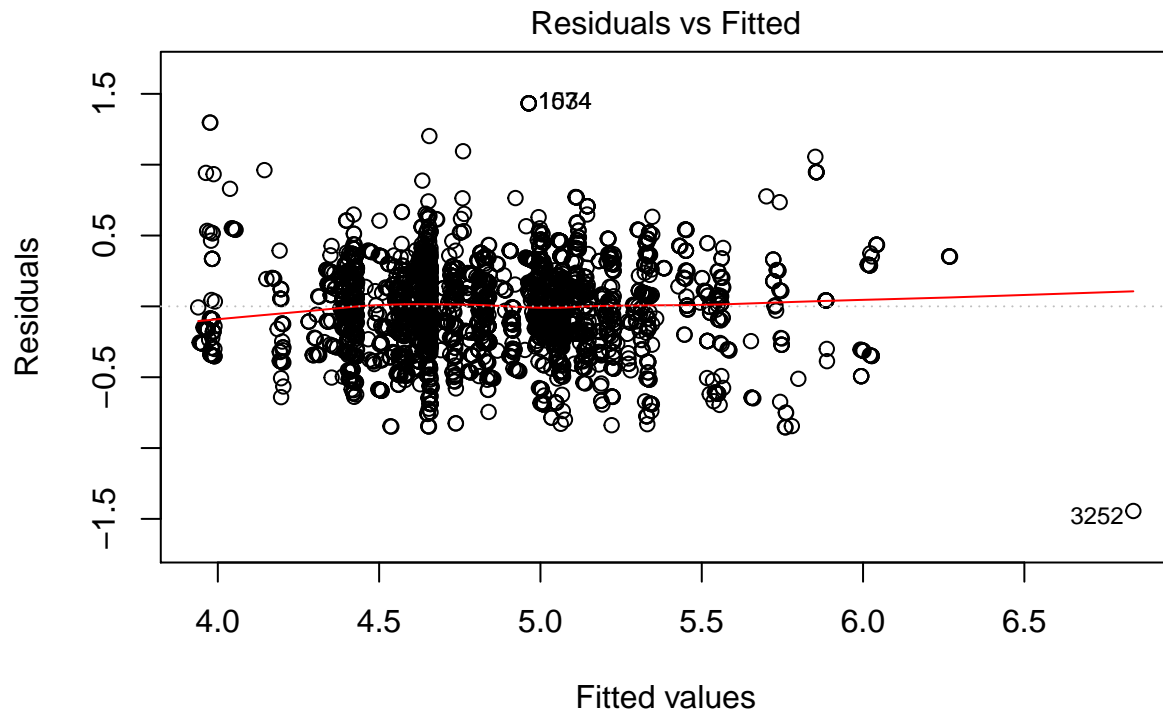
## IV. Modelling:

### Model1: Simple linear regression:

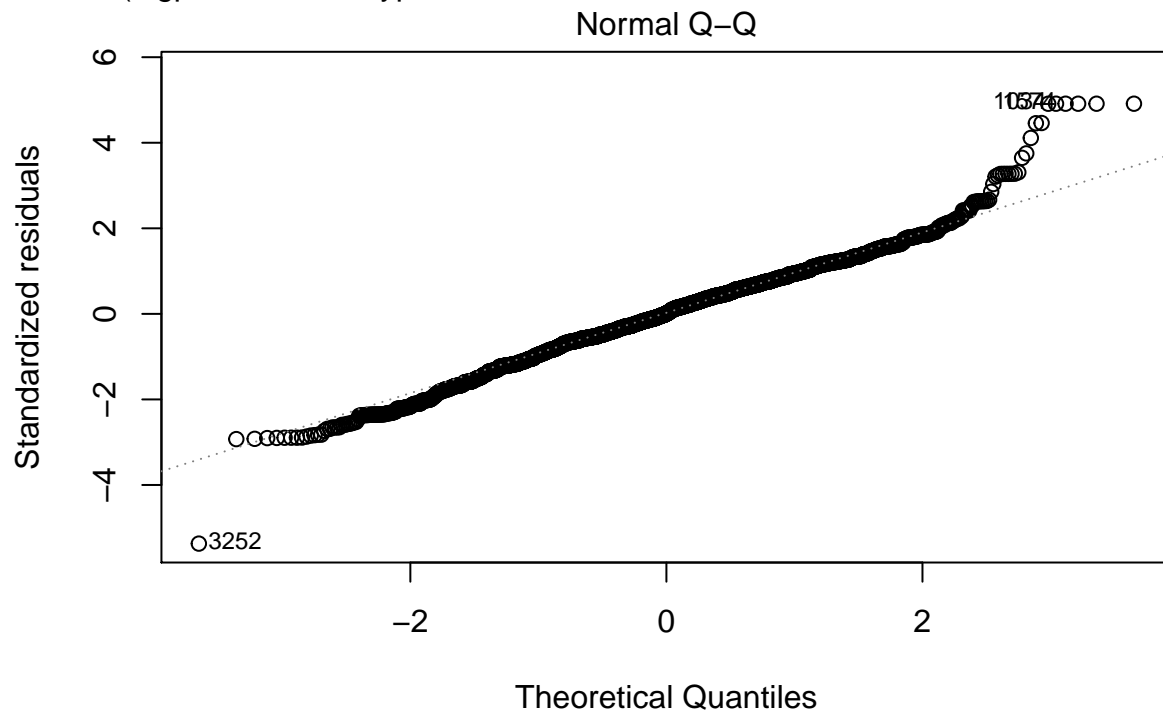
$$\log(\text{price}) = \alpha + \beta_1 x_{\text{roomtype}} + \beta_2 x_{\text{reviews}} + \beta_3 x_{\text{rating}} + \beta_4 x_{\text{accommodates} * \text{bedrooms}} + \beta_5 x_{\text{accommodates}} + \beta_6 x_{\text{bedrooms}}$$

```
##
## Call:
## lm(formula = logprice ~ room_type + reviews + overall_satisfaction +
##      accommodates + accommodates * bedrooms, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.44414 -0.17748  0.00088  0.19270  1.43330
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      2.662e+00  9.788e-02  27.196 < 2e-16 ***
## room_typePrivate room -4.109e-01  1.083e-02 -37.936 < 2e-16 ***
## room_typeShared room -1.009e+00  3.160e-02 -31.926 < 2e-16 ***
## reviews          -2.824e-04  7.562e-05  -3.734 0.000191 ***
## overall_satisfaction  4.565e-01  1.894e-02  24.103 < 2e-16 ***
## accommodates        4.791e-02  7.338e-03   6.529 7.48e-11 ***
## bedrooms          -2.776e-03  2.270e-02  -0.122 0.902668
## accommodates:bedrooms  3.073e-02  4.411e-03   6.967 3.79e-12 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2925 on 3833 degrees of freedom
## Multiple R-squared:  0.6137, Adjusted R-squared:  0.613
## F-statistic: 869.9 on 7 and 3833 DF,  p-value: < 2.2e-16
```





lm(logprice ~ room\_type + reviews + overall\_satisfaction + accommodates + a ...



lm(logprice ~ room\_type + reviews + overall\_satisfaction + accommodates + a ...

From this simple linear regression, most of the predictors are significant. From the summary of the model we notice that the predictor bedrooms is not significant to price. We may want to eliminate those predictors in the multilevel models. Also the R-square in the model is 0.613, so the model is not well fitted. However, in the residual plot, there are some points having big residuals: 1674, 3252. Those might be the prices are very high that lead to big residuals. The rest of points are symmetrically distributed around the line  $\mu = 0$ . In the QQ plot, we can see most dots in the middle fall on the line. However, the data have more extreme values on the tail of the distribution. Now let's expand simple linear model to multilevel linear model.

## Model2: Multilevel linear model with random intercept:

$$\log(\text{price}) = \alpha_i + \beta_1 x_{\text{roomtype}} + \beta_2 x_{\text{reviews}} + \beta_3 x_{\text{overall_satisfaction}} + \beta_4 x_{\text{accommodates} * \text{bedrooms}} + \beta_5 x_{\text{accommodates}}$$

```
## lmer(formula = log(price) ~ room_type + reviews + overall_satisfaction +  
##     accommodates + accommodates * bedrooms + (1 | neighborhood) -  
##     1, data = df)
```

```
##               coef.est coef.se  
## room_typeEntire home/apt 3.02    0.10  
## room_typePrivate room   2.61    0.10  
## room_typeShared room    2.02    0.10  
## reviews                  0.00    0.00  
## overall_satisfaction     0.36    0.02  
## accommodates             0.07    0.01  
## bedrooms                 0.05    0.02  
## accommodates:bedrooms    0.02    0.00  
##  
## Error terms:  
## Groups      Name      Std.Dev.  
## neighborhood (Intercept) 0.19  
## Residual              0.26  
## ---  
## number of obs: 3841, groups: neighborhood, 29  
## AIC = 672.1, DIC = 518.4  
## deviance = 585.3  
  
## Computing profile confidence intervals ...  
  
##               2.5 %      97.5 %  
## .sig01          0.1429034414 0.2444516626  
## .sigma          0.2516234569 0.2631786577  
## room_typeEntire home/apt 2.8260397676 3.2066884226  
## room_typePrivate room   2.4232281236 2.8030337606  
## room_typeShared room    1.8157086836 2.2263119360  
## reviews             -0.0005326705 -0.0002619241  
## overall_satisfaction    0.3228740310 0.3925677058  
## accommodates          0.0618700559 0.0880053494  
## bedrooms             0.0067551356 0.0871202090  
## accommodates:bedrooms   0.0112923800 0.0268447733
```

*This model gets rid of non-significant terms bedrooms. The facotor of room type plays the most important part in the model. The coeifficent of reviews is zero so we dont need to keep it for the next model. Besides, the overall satisfaction is the second influential term in this model*

### Model3 : Multilevel linear model with random slope:

$$\log(\text{price}) = \alpha_i + \beta_1 x_{\text{roomtype}} + \beta_2 x_{\text{overall_satisfaction}} + \beta_3 x_{\text{accommodates} * \text{bedrooms}} + \beta_4 x_{\text{accommodates}}$$

```
## lmer(formula = logprice ~ room_type + overall_satisfaction +
##      accommodates + accommodates * bedrooms + (0 + overall_satisfaction |
##      neighborhood) - 1, data = df)
##               coef.est coef.se
## room_typeEntire home/apt 2.99    0.09
## room_typePrivate room    2.58    0.09
## room_typeShared room     1.98    0.10
## overall_satisfaction     0.35    0.02
## accommodates              0.08    0.01
## bedrooms                  0.05    0.02
## accommodates:bedrooms     0.02    0.00
##
## Error terms:
##   Groups      Name              Std.Dev.
## neighborhood overall_satisfaction 0.04
## Residual                      0.26
## ---
## number of obs: 3841, groups: neighborhood, 29
## AIC = 687, DIC = 570
## deviance = 619.5
## Computing profile confidence intervals ...
##
##               2.5 %    97.5 %
## .sig01          0.029772925 0.05095609
## .sigma          0.252750959 0.26435799
## room_typeEntire home/apt 2.809961398 3.16559622
## room_typePrivate room    2.401449234 2.75589727
## room_typeShared room     1.789302925 2.17917757
## overall_satisfaction     0.311936110 0.38720061
## accommodates           0.062808603 0.08903638
## bedrooms              0.008721213 0.08948098
## accommodates:bedrooms    0.011030174 0.02664824
```

*This is a multilevel linear model with random slope. From the confidence interval we can see all the predictors are significant.*

## model4 : Multilevel linear model with random slope and random intercept:

```

log(price) =  $\alpha_i + \beta_1 x_{roomtype} + \beta_2[i] x_{overallsatisfaction} + \beta_3 x_{accommodates*bedrooms} + \beta_4 x_{accommodates}$ 

## Warning in checkConv(attr(opt, "derivs"), opt$par, ctrl =
## control$checkConv, : Model failed to converge with max|grad| = 0.00411589
## (tol = 0.002, component 1)

## lmer(formula = logprice ~ room_type + overall_satisfaction +
##       accommodates + accommodates * bedrooms + (1 + overall_satisfaction |
##       neighborhood) - 1, data = df)
##               coef.est coef.se
## room_typeEntire home/apt 2.95    0.29
## room_typePrivate room    2.54    0.29
## room_typeShared room     1.99    0.29
## overall_satisfaction     0.36    0.06
## accommodates              0.07    0.01
## bedrooms                  0.04    0.02
## accommodates:bedrooms     0.02    0.00
##
## Error terms:
## Groups      Name                      Std.Dev. Corr
## neighborhood (Intercept)              1.37
##               overall_satisfaction 0.29    -0.99
## Residual                               0.25
## ---
## number of obs: 3841, groups: neighborhood, 29
## AIC = 605.3, DIC = 489.3
## deviance = 536.3

## Computing profile confidence intervals ...

## Warning in nextpar(mat, cc, i, delta, lowcut, upcut): unexpected decrease
## in profile: using minstep

## Warning in FUN(X[[i]], ...): non-monotonic profile for .sig03

## Warning in confint.thpr(pp, level = level, zeta = zeta): bad spline fit
## for .sig03: falling back to linear interpolation

##               2.5 %      97.5 %
## .sig01          0.9424766290  1.88414758
## .sig02         -0.9926812541 -0.97714462
## .sig03          0.1965420753  0.39785803
## .sigma          0.2481679760  0.25961447
## room_typeEntire home/apt 2.3685574472  3.52199751
## room_typePrivate room    1.9637948284  3.11674585
## room_typeShared room     1.4070261997  2.56680488
## overall_satisfaction     0.2437658766  0.48248678
## accommodates           0.0612101788  0.08740651
## bedrooms             0.0001592506  0.08100124
## accommodates:bedrooms    0.0115477115  0.02717413

```

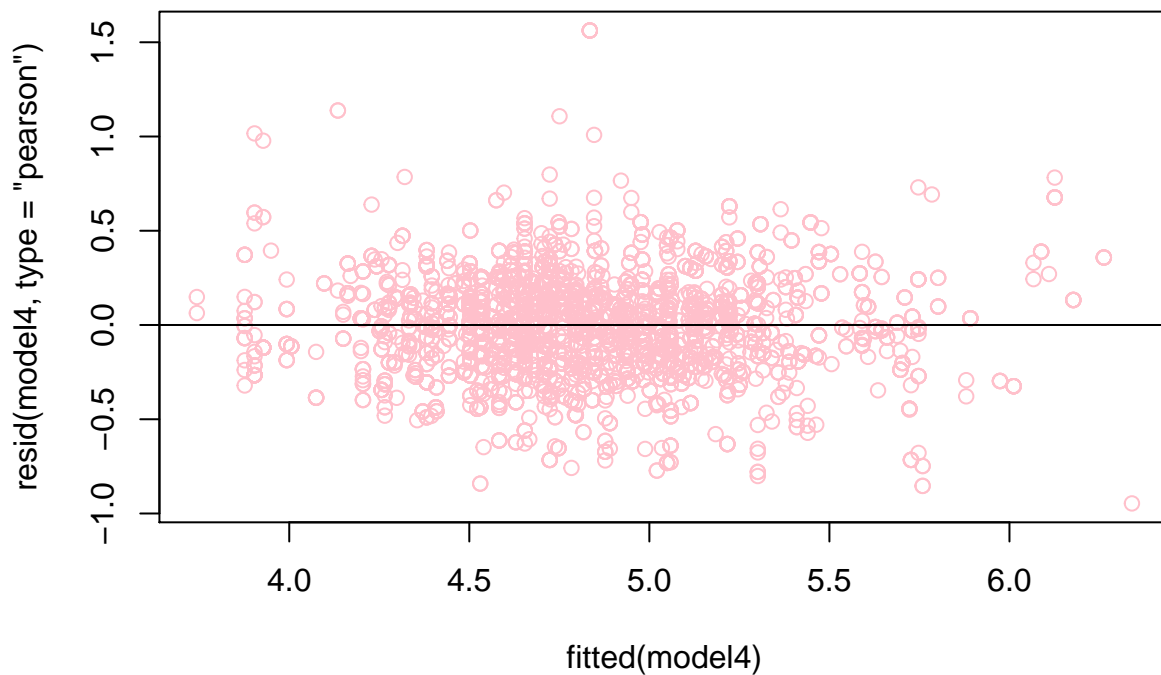
*This is the model is based on model2 and model3 with random slope as well random intercept. According to the confidence interval and coefficient, all the predictors are significant.*

## V. Result:

Since we have 3 multilevel models with similar structures, we want to run ANOVA test to test whether there's any difference among the models and which model has best goodness of fit.

```
## Data: df
## Models:
## model3: logprice ~ room_type + overall_satisfaction + accommodates +
## model3: accommodates * bedrooms + (0 + overall_satisfaction | neighborhood) -
## model3: 1
## model2: log(price) ~ room_type + reviews + overall_satisfaction + accommodates +
## model2: accommodates * bedrooms + (1 | neighborhood) - 1
## model4: logprice ~ room_type + overall_satisfaction + accommodates +
## model4: accommodates * bedrooms + (1 + overall_satisfaction | neighborhood) -
## model4: 1
##      Df    AIC    BIC logLik deviance  Chisq Chi Df Pr(>Chisq)
## model3  9 686.99 743.28 -334.50   668.99
## model2 10 672.14 734.68 -326.07   652.14 16.851    1 4.042e-05 ***
## model4 11 605.33 674.12 -291.67   583.33 68.809    1 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

**Figure 9. residual plot for model4**



```
## [1] 2.816359
```

According to the output of the test, we can see that model 4 with random intercept and random slope is the best fit among three multilevel models. It has lowest deviance with 583.33. The second plot is a residual plot form model 4. As we can see the plot the points are symmetrically distributed around the line  $h = 0$ . The neighborhood with the maximum intercept is Parkside with intercept 2.816

From the analysis above we found the best model among is the Model4:

$$\log(\text{price}) = \alpha_i + \beta_1 x_{\text{roomtype}} + \beta_2 x_{\text{overall satisfaction}} + \beta_3 x_{\text{accommodates} * \text{bedrooms}} + \beta_4 x_{\text{accommodates}}$$

*From the model, we can find that the most influential term is the factor of room type. The second influential term is the district of neighborhood. Also, more higher overall satisfaction will lead to a higher price. The ideas of building these model is to predicting models by different levels of neighborhood. The model in each level will have a unique intercept and a unique slope for the predictor overall satisfaction. In this way it will lead the last model minimize the deviance comparing to the previous two multilevel linear model.*

## VI. Discussion:

The result turns out to be not very surprising, although I thought the reviews would so how affect the price of Airbnb in San Francisco. The reality here is that the factor of room type will determine the price the most, which is reasonable because from what my knowledge, San Francisco is one of the most expensive living environment in the U.S. especially the bay area has the highest housing price. Therefore it is a place the every inch is like the price as gold. The room type of Airbnb is more like the different room size for the hotel. The entire home/apt tends to serve more people and have bigger space at one time. So the price of entire home/apt should be higher than the other two room type. Also, the second term neighborhood is obvious to be significant. This is because based on the living environment in San Francisco there are many tech company and university in the city, as well as the neighborhood around union square, twin peak is close to downtown area and there are many sight-seeing spots for tourists. We can conclude that our findings are reasonable.

However, this analysis is definitely not perfect. There are only 4 predictors in the multilevel model. Those predictors are the most significant variables we found in the dataset. This is just a very basic analysis of the Airbnb data, there are many other factors that can take in to account, such as the academic institution in the neighborhood or the transportation condition, and also the geographic difference in each country will be interesting to explore in the future. That maybe the future direction of this project.\*

## VII. Reference

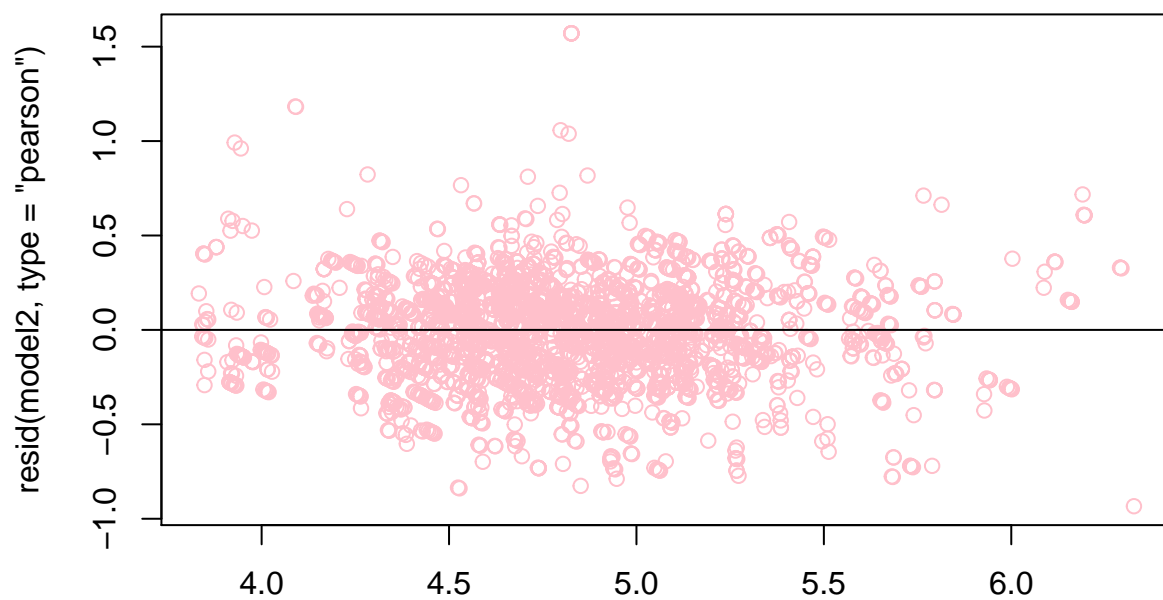
<http://tomslee.net>

<https://en.wikipedia.org/wiki/Airbnb>

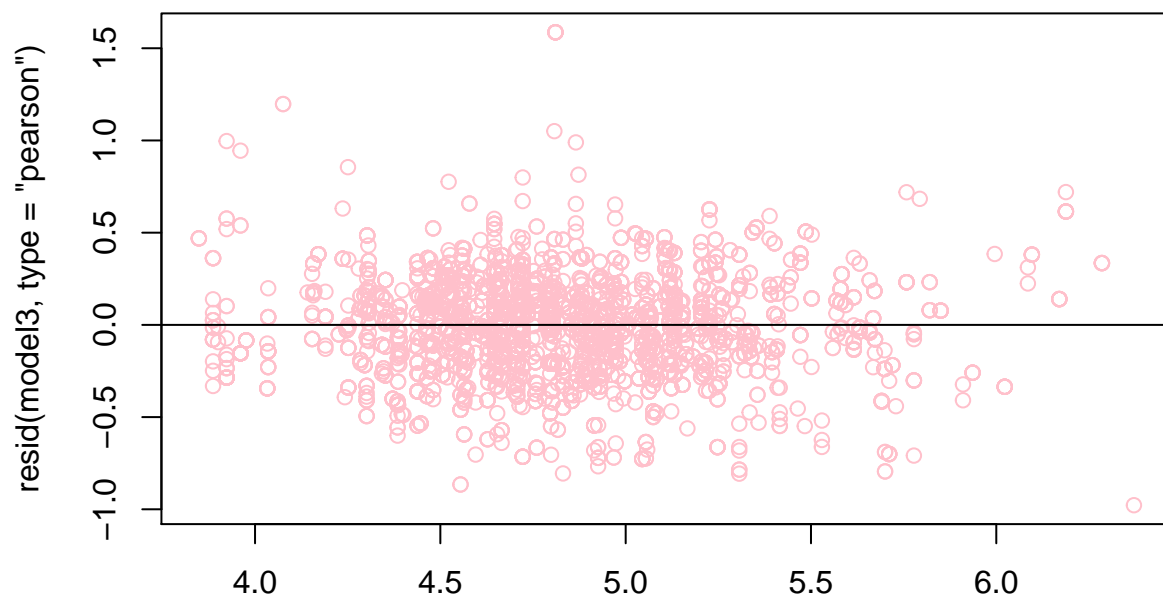
<https://www.airbnb.com/>

## VIII. Appendix

**residual plot for model2**



**residual plot for model3**



Residual for model 2&3