

MA678 Midterm Project –

Yunyi Zhang

27/11/2019

1: Introduction

All datasets are downloaded from <http://tomslee.net/airbnb-data-collection-get-the-data>, datasets are separated by month and location. Each dataset contains information such as room_id, room_type, city, neighborhood, accommodates, number of bedrooms. In this project I am going to first inspect each important variable, and then do the regression on different variables against the price.

2: Load and merge the data

After merging all datasets, there are in total 7112 rows with 19 columns

3. Quick summary of my dataset

These quick summary shows that there are 13 different neighborhoods, 3 different room types, 15 different accommodates and 6 number of bedroom. I also included a detailed summary below.

```
## [1] 13
```

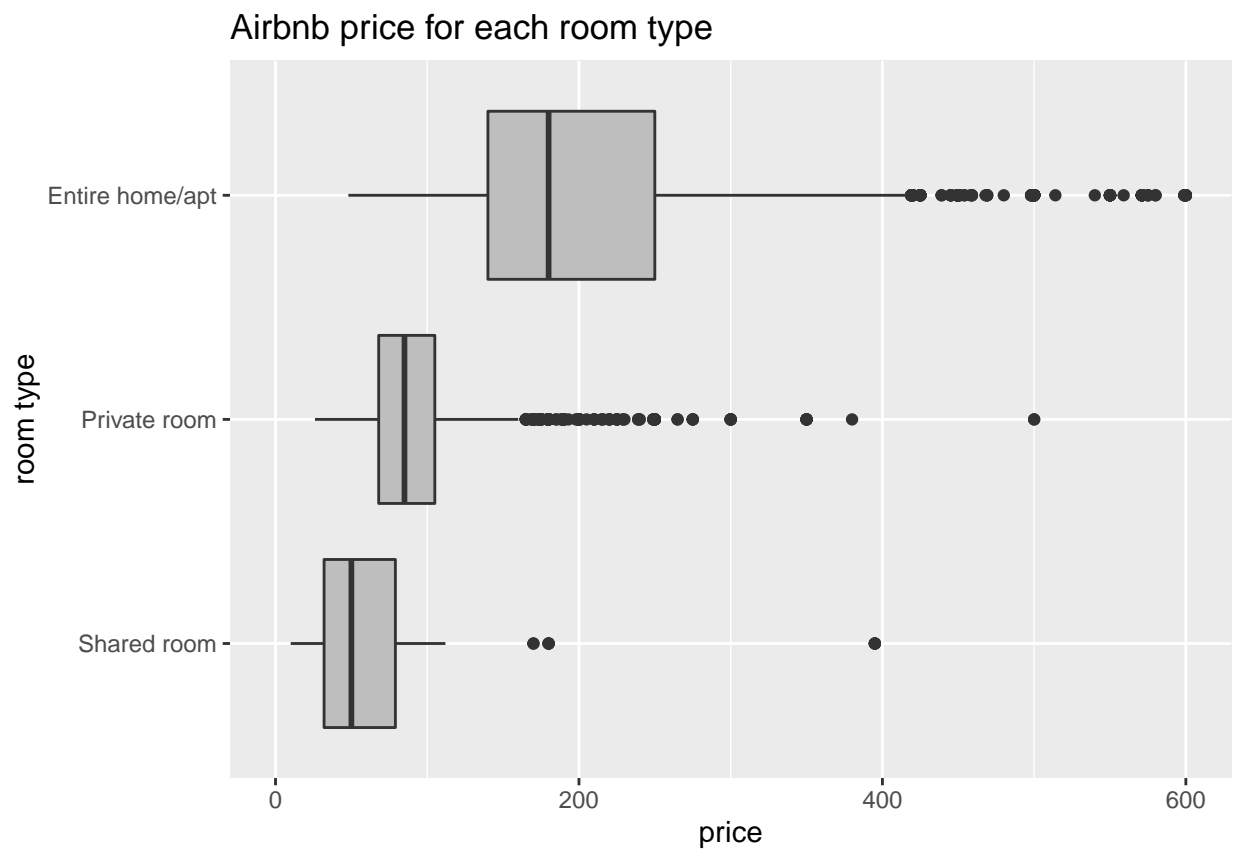
```
## [1] 3
```

```
## [1] 15
```

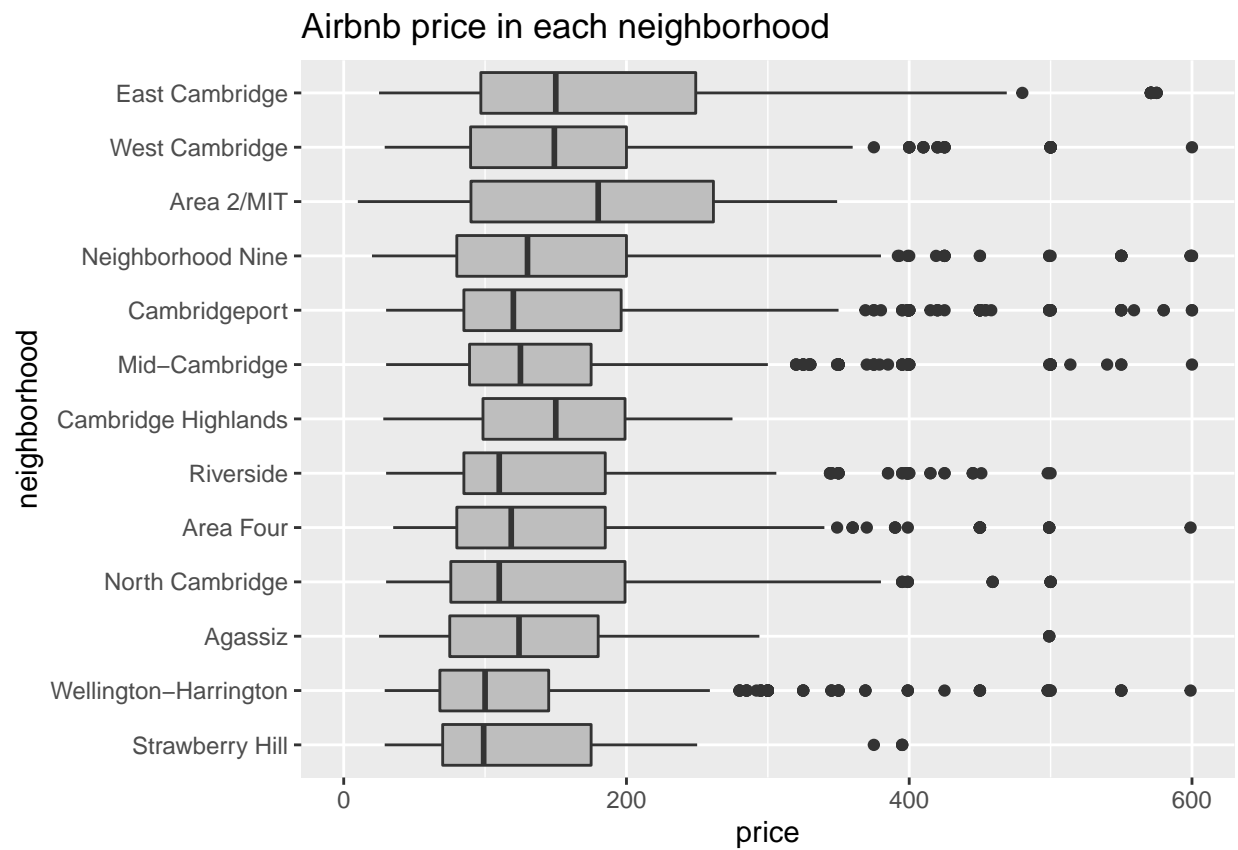
```
## [1] 6
```

4. Boxplots for different variables vs rating

4.1 Room Type and Price



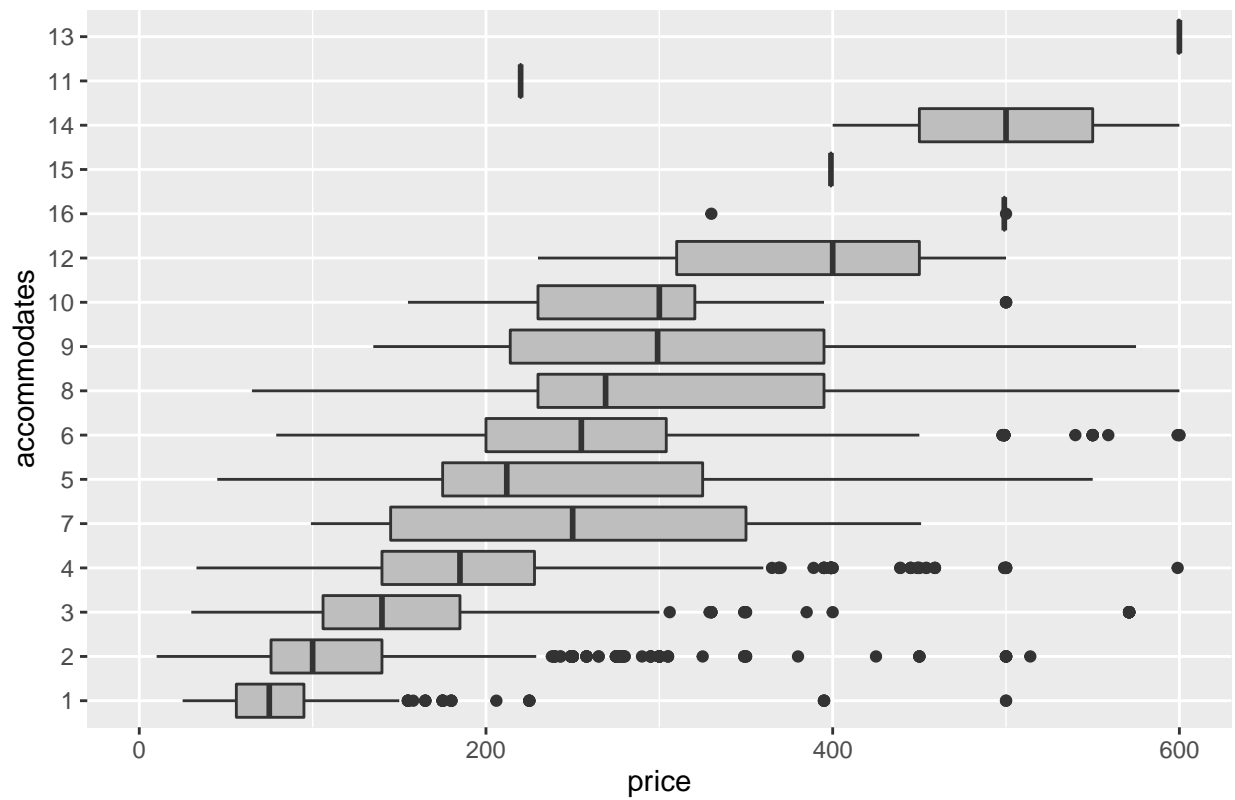
4.2 Neighborhood and Price



Neighborhood	average price
Agassiz	143.15
Area 2/MIT	172.16
Area Four	146.28
Cambridge Highlands	150.49
Cambridgeport	159.56
East Cambridge	191.00
Mid-Cambridge	151.24
Neighborhood Nine	169.90
North Cambridge	144.85
Riverside	146.50
Strawberry Hill	127.29
Wellington-Harrington	129.05
West Cambridge	176.91

4.3 Accommodate and Price

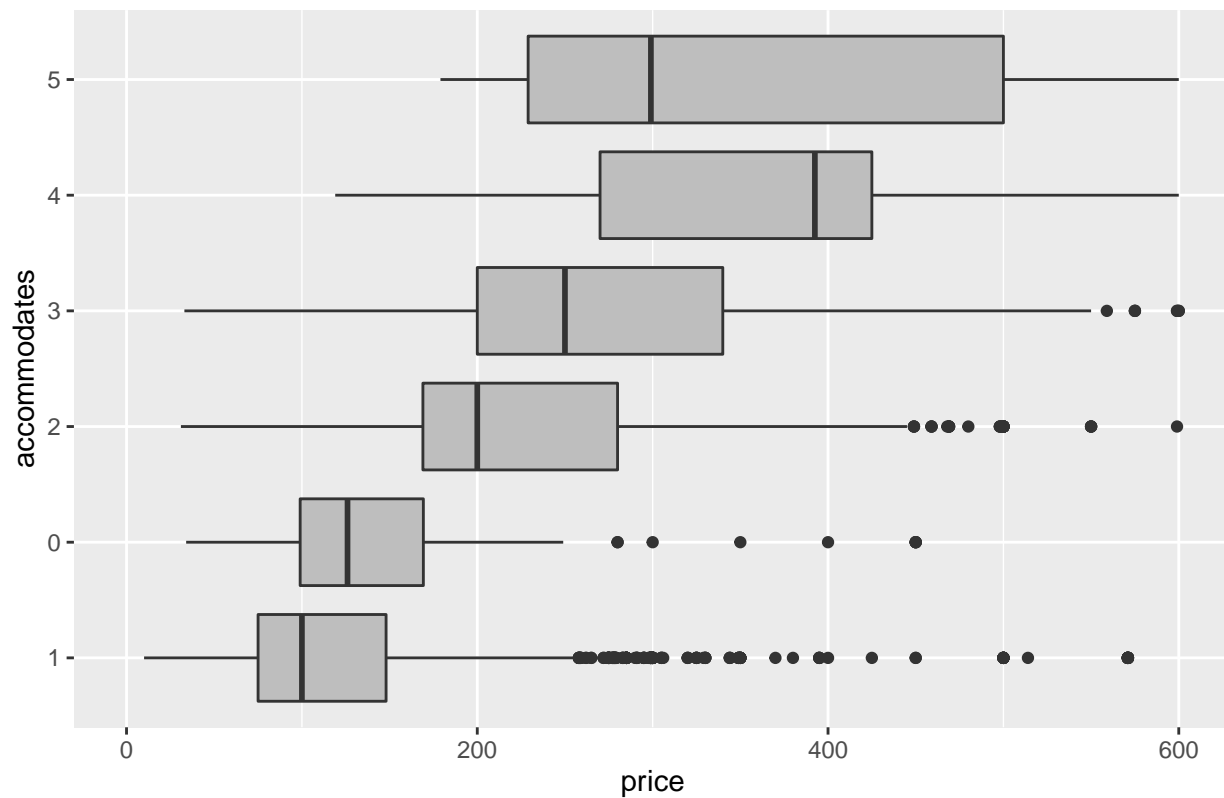
Airbnb price for each accommodates



Accommodates	average price
1	79.90
2	114.08
3	162.51
4	195.88
5	255.95
6	282.82
7	251.33
8	341.92
9	356.82
10	386.53
11	584.00
12	436.15
13	599.67
14	531.25
15	469.25
16	465.40

4.4 Bedrooms and Price

Airbnb price for each accommodates



Bedrooms	average price
0	138.20
1	116.17
2	234.23
3	299.10
4	397.35
5	477.68

5 Multilevel Model analysis

5.1 Regress price on Neighborhood and Bedrooms, no between group

```
## lmer(formula = price ~ neighborhood + bedrooms + (1 | room_type) +  
##       (1 | accommodates), data = CAMB_main, REML = FALSE)  
##               coef.est coef.se  
## (Intercept)      175.18    33.61  
## neighborhoodArea 2/MIT      29.46    13.12  
## neighborhoodArea Four       4.17     5.63  
## neighborhoodCambridge Highlands -25.86    13.22  
## neighborhoodCambridgeport      9.13     5.24  
## neighborhoodEast Cambridge     35.02     5.33  
## neighborhoodMid-Cambridge     10.53     5.10  
## neighborhoodNeighborhood Nine   10.63     5.81
```

```
## neighborhoodNorth Cambridge      -4.87      5.87
## neighborhoodRiverside             5.85      5.45
## neighborhoodStrawberry Hill      -26.67      9.17
## neighborhoodWellington-Harrington -9.86      5.83
## neighborhoodWest Cambridge       -7.11      6.37
## bedrooms                          40.07      2.01
##
## Error terms:
## Groups      Name      Std.Dev.
## accommodates (Intercept) 93.16
## room_type    (Intercept) 38.59
## Residual                        81.15
## ---
## number of obs: 7112, groups: accommodates, 16; room_type, 3
## AIC = 82834.4, DIC = 82800.4
## deviance = 82800.4
```

5.2 Regress price on Neighborhood, Bedrooms, Room Types, with one between-group of bedrooms and accommodates.

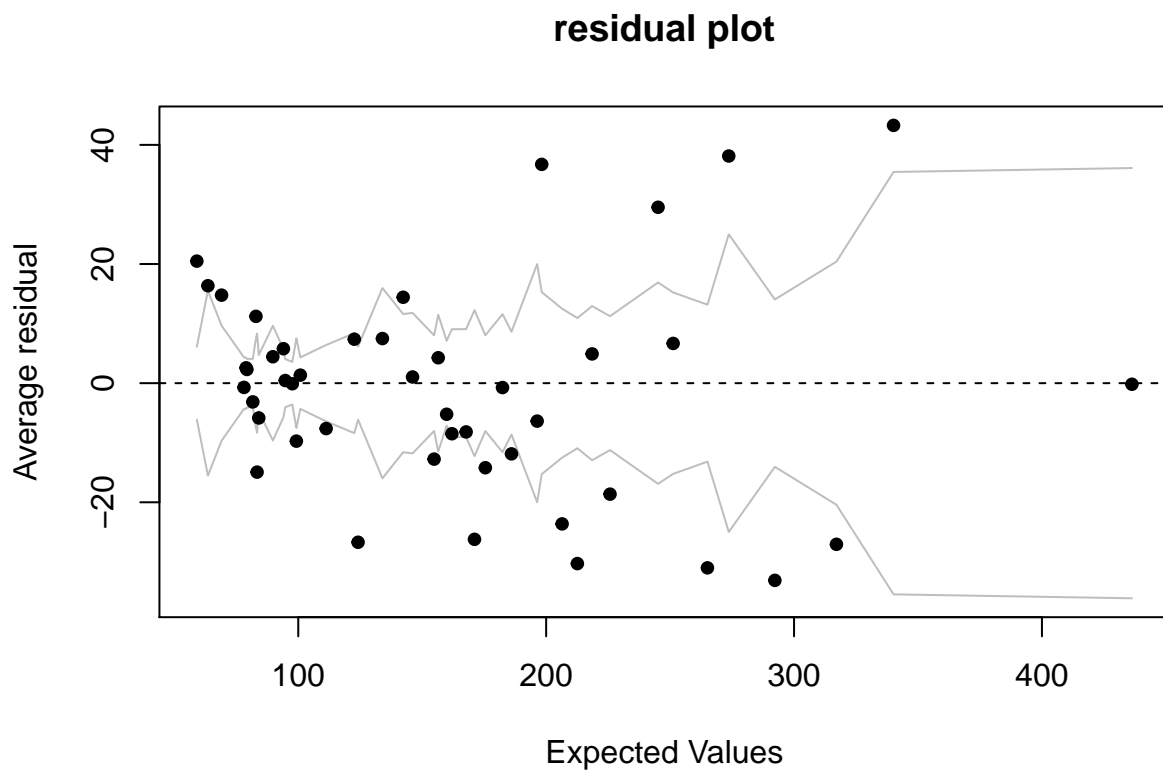
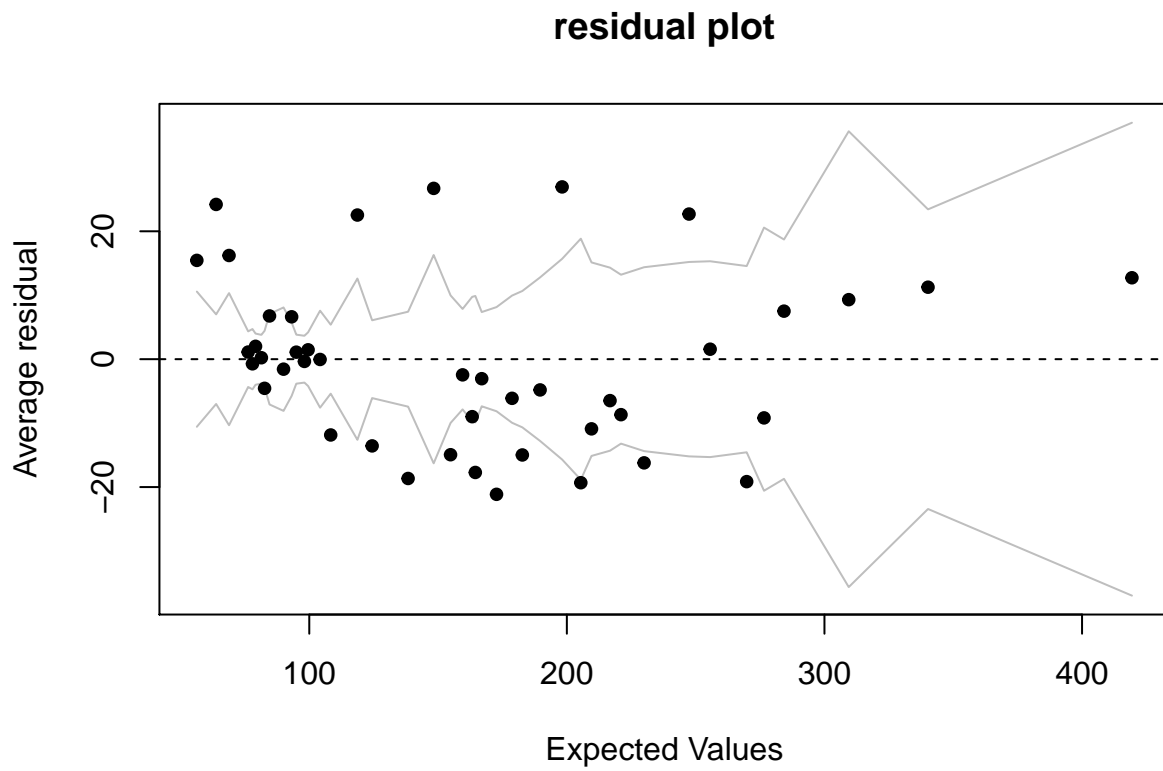
```
## lmer(formula = price ~ neighborhood + bedrooms + room_type +
##       (1 + bedrooms | accommodates), data = CAMB_main, REML = FALSE)
##
##               coef.est coef.se
## (Intercept)      136.22     5.48
## neighborhoodArea 2/MIT      27.99    13.03
## neighborhoodArea Four       5.79     5.59
## neighborhoodCambridge Highlands -26.37    13.14
## neighborhoodCambridgeport     9.34     5.21
## neighborhoodEast Cambridge    35.91     5.29
## neighborhoodMid-Cambridge     11.66     5.07
## neighborhoodNeighborhood Nine  11.02     5.78
## neighborhoodNorth Cambridge   -5.28     5.83
## neighborhoodRiverside        6.44     5.42
## neighborhoodStrawberry Hill  -26.94     9.11
## neighborhoodWellington-Harrington -8.95     5.80
## neighborhoodWest Cambridge   -8.99     6.34
## bedrooms          59.35     6.93
## room_typePrivate room   -61.54     2.73
## room_typeShared room   -77.15     9.00
##
## Error terms:
## Groups      Name      Std.Dev. Corr
## accommodates (Intercept) 2.44
##               bedrooms    26.76    1.00
## Residual                80.64
## ---
## number of obs: 7112, groups: accommodates, 16
## AIC = 82733.6, DIC = 82693.6
## deviance = 82693.6
```

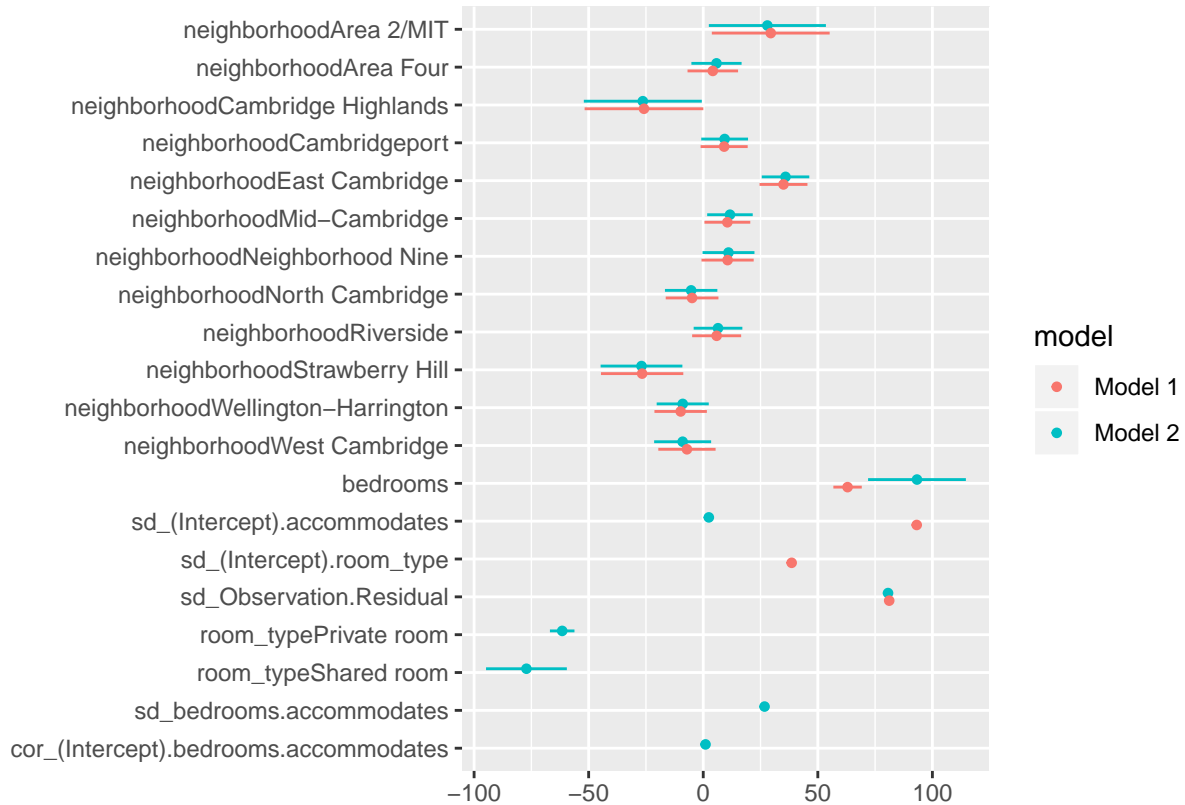
6. Interpretation

For Model 2, if I would like to rent an airbnb MIT, I would expect the price to be $136.22 + 27.99 + 59.35 \text{ Bedrooms} - 61.54 \text{ Private room} - 77.15 \text{ shared room}$

By increasing one bedroom, I would expect the price to rise by 59.35, if it is a private room, the price would drop by 61.54, and if it is a shared room, the price would drop by 77.15. It makes sense because the shared room is normally cheaper than a private one.

7. Model check





From the plots we can see that the second regression model have a better and more reliable result.

8. ANOVA

```
anova(mulre1, mulre2)
```

```
## Data: CAMB_main
## Models:
## mulre1: price ~ neighborhood + bedrooms + (1 | room_type) + (1 | accommodates)
## mulre2: price ~ neighborhood + bedrooms + room_type + (1 + bedrooms |
## mulre2: accommodates)
##      Df   AIC   BIC logLik deviance  Chisq Chi Df Pr(>Chisq)
## mulre1 17 82834 82951 -41400    82800
## mulre2 20 82734 82871 -41347    82694 106.85      3 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Model two has a DF of 3 and a smaller p-value which lead to the conclusion that Model 2 has a better fit.

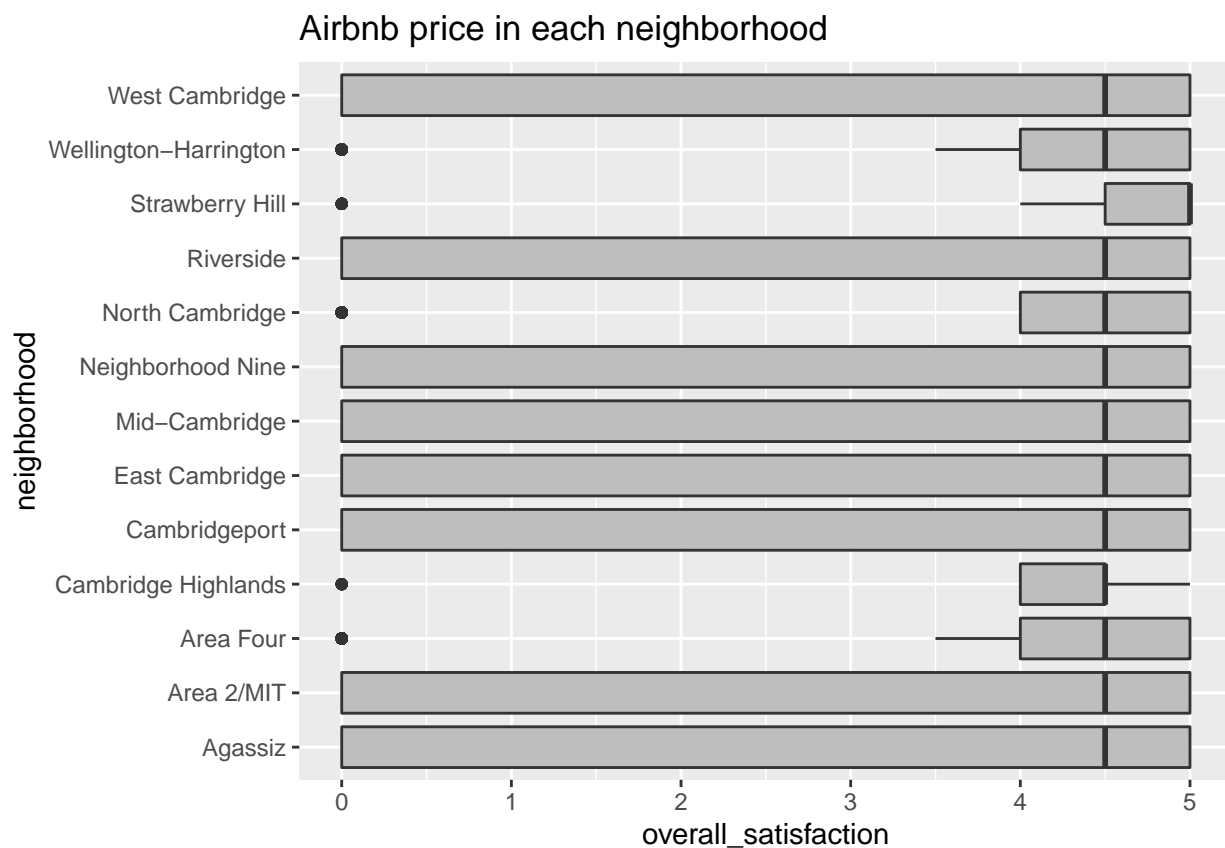
9. Discussion

#9.1 From further research on google map, airbnb near Harvard and MIT tend to have a higher price comparing to other areas in Cambridge, which makes sense considering that those areas are more popular for visitors

#9.2 The result of the fit shows some residuals outside the model which indicate some limitation of the model

#9.3 My initial thought is to analyze the relationship between neighborhood and the rating, however, as I looked into the dataset and tried to do some plot, I found out that the rating is not sufficient

enough for me to keep my work (see the graph below), also the result does not make any sense.



10. Future Direction

I am still very interested in the relationship between rating and other factors. I would like to find some more reliable dataset if possible, and keep up the work.