# Kaggle Project Report
## —— Predicting Airbnb Rental Price

Name: Zeying Liu

UNI: zl3121

Email: zl3121@columbia.edu

Professor: Vishal Lala

Associate: Davin Kaing, Lei Yu

Date: 2021. Dec. 8th

## Table of Contents

# 1. Introduction

## 1.1. Goal of the project

Use the models learned in the course to build a price prediction model for the Airbnb rental data given by the Kaggle platform. The model is repeatedly adjusted to minimize the RMSE of the final prediction results.
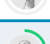
## 1.2. Result

I ranked 29/581, with 58.614 RMSE.

| | | | | | | |
|---|---|---|---|---|---|---|
| 26 | ▼ 17 | Yutao Yan | | 58.40633 | 22 | 18d |
| 27 | ▲ 80 | Kai Tung Chuang | | 58.43853 | 6 | 19d |
| 28 | ▲ 16 | brmiao | | 58.55902 | 9 | 17d |
| 29 | ▲ 23 | Zeying Liu | | 58.61444 | 19 | 18d |
| 30 | ▼ 2 | ElaineShi | | 58.62545 | 18 | 18d |
| 31 | ▲ 25 | Biyu Wang | | 58.62696 | 10 | 18d |
| 32 | ▼ 15 | Joey05 | | 58.65953 | 48 | 18d |

*Figure 1 Ranking*

## 1.3. Final model

The final selected variables for achieving the above results is:

```
1.  skim(data)
2.  #input variables without too much NA (keep scoring data same with training data, except change 'price'
    to 'id'), Variable number:  42/91
3.  #After cleaning and selection, Variable number: 35/91
4.  colnames(data)
```

```
[1] "price"                                    "neighbourhood_group_cleansed"
 [3] "neighbourhood_cleansed"                   "zipcode"
 [5] "room_type"                                "property_type"
 [7] "accommodates"                             "guests_included"
 [9] "extra_people"                             "bathrooms"
[11] "bedrooms"                                 "beds"
[13] "bed_type"                                 "host_response_rate"
[15] "host_is_superhost"                        "host_has_profile_pic"
[17] "host_identity_verified"                   "security_deposit"
[19] "cleaning_fee"                             "number_of_reviews_ltm"
[21] "instant_bookable"                         "cancellation_policy"
[23] "review_scores_rating"                     "review_scores_accuracy"
[25] "review_scores_cleanliness"                "review_scores_location"
[27] "maximum_nights_avg_ntm"                   "availability_60"
[29] "calculated_host_listings_count_entire_homes" "calculated_host_listings_count_private_rooms"
[31] "calculated_host_listings_count_shared_rooms"  "is_location_exact"
[33] "require_guest_profile_picture"            "require_guest_phone_verification"
[35] "host_listings_count"
```

The final xgboost model for achieving the above results is:

```
1.  set.seed(1234)
2.  xgb_nrounds = xgb.cv(data=as.matrix(train_input),
3.                       label = train$price,
4.                       nrounds = 100,
5.                       nfold = 5,
6.                       eta = 0.22,
7.                       gamma = 0,
8.                       verbose = 0)
9.
10. which.min(xgb_nrounds$evaluation_log$test_rmse_mean)
11.
12. xgboost2 = xgboost(data=as.matrix(train_input),
13.                    label = train$price,
14.                    nrounds= 99,
15.                    eta = 0.22,
16.                    gamma = 0,
17.                    verbose = 0)
```

To get the final results, I went through the process of variable screening, data cleaning, feature selection, changing the model type, and model tuning. The details will be expanded in the following parts.

# 2. Data Preparation

## 2.1. Variables Filtering

First, looking at the provided data, there are a large number of null values, which are not meaningful for the training model. Meanwhile, text information is difficult to be processed directly into recognizable numeric or logical values and also needs to be removed. Moreover, some numerical variables have potential multicollinearity problems are needed to further exam, such as the maximum or minimum number of days customers stay, the availability of places, and the listing counts. I finally chose the following 42 variables:

```
1.  > sum(which(is.na(data)))
2.  [1] 592195197166
```

```
1.  #checking correlation between alike variables
2.  #1. stay nights
3.  round(cor(data[,c(55:62)]),2)
4.  #keep 'minimum_nights_avg_ntm','maximum_nights_avg_ntm' which is adjusted, no multicollinearity, and
    cover other variables
5.
6.  #2. availability in days
7.  round(cor(data[,c(65:68)]),2)
8.  #keep availability_60
9.
10. #3. listing counts
11. round(cor(data[,c(88:91)]),2)
12. #keep 89:91
```

```
1.  colnames(data)
2.   [1] "price"                                      "host_since"
3.   [3] "neighbourhood_group_cleansed"               "neighbourhood_cleansed"
4.   [5] "zipcode"                                    "room_type"
5.   [7] "property_type"                              "accommodates"
6.   [9] "guests_included"                            "extra_people"
7.  [11] "bathrooms"                                  "bedrooms"
8.  [13] "beds"                                       "bed_type"
9.  [15] "host_response_rate"                         "host_is_superhost"
10. [17] "host_has_profile_pic"                       "host_identity_verified"
11. [19] "security_deposit"                           "cleaning_fee"
12. [21] "number_of_reviews_ltm"                      "instant_bookable"
13. [23] "cancellation_policy"                        "review_scores_rating"
14. [25] "review_scores_accuracy"                     "review_scores_cleanliness"
15. [27] "review_scores_checkin"                      "review_scores_location"
16. [29] "review_scores_communication"               "review_scores_value"
17. [31] "minimum_nights_avg_ntm"                     "maximum_nights_avg_ntm"
18. [33] "availability_60"                            "calculated_host_listings_count_entire_homes"
19. [35] "calculated_host_listings_count_private_rooms" "calculated_host_listings_count_shared_rooms"
20. [37] "last_review"                                "first_review"
21. [39] "is_location_exact"                          "require_guest_profile_picture"
22. [41] "require_guest_phone_verification"           "host_listings_count"
```

The processing of the scoring data is identical to the analysis data. In order to be more streamlined, the code of processing scoring data will not be shown in this paper.

## 2.2. Date Conversion

For date-type data, the format first needs to be converted to standard format. And in order to show the time from the corresponding date to the present day, I then ran them separately with the date at that time. The larger the value, the more experienced the hosts are.

```
1.  #calculate days
2.  #1. host days
3.  data = data %>%
4.    mutate(host_since = mdy(host_since)) %>%
5.    mutate(host_long = today()-host_since)
6.  data = data[,-2]
7.  data$host_long[is.na(data$host_long)] = 0
8.  data$host_long = as.numeric(data$host_long)
9.  # table(data$host_long)
10.
11. #2. first review
12. data = data %>%
13.   mutate(first_review = mdy(first_review)) %>%
14.   mutate(first_review_long = today()-first_review)
15. data = data[,-37]
16. data$first_review_long[is.na(data$first_review_long)] = 0
17. data$first_review_long = as.numeric(data$first_review_long)
18.
19. #3. last review
20. data = data %>%
21.   mutate(last_review = mdy(last_review)) %>%
22.   mutate(last_review_long = today()-last_review)
23. data = data[,-36]
24. data$last_review_long[is.na(data$last_review_long)] = 0
25. data$last_review_long = as.numeric(data$last_review_long)
```

In addition, the above variables also have multicollinearity. By following the same method, keep the host days and last review variables.

```
1.  round(cor(data[,c(40:42)]),2)
2.  #keep36,38
3.  data = data[,-41]
4.  scoringData = scoringData[,-41]
```

## 2.3. Grouping low-frequency values

Firstly, I <u>filled</u> the zip code variable in Excel using the VLOOKUP function based on the neighborhood for the null values. In addition, for the zip codes that were not found, the corresponding zip codes were <u>searched</u> on the internet using other location information. After processing, due to the low number of occurrences of some zip codes, some of the them in the test set did not appear in the training set when the data was split. Therefore, I <u>modified</u> the zip codes with less than 3 occurrences to "other". A similar problem occurs with communities and properties, which are implemented as follows:

```
1.  # zipcode
2.  data = data %>%
3.    group_by(zipcode) %>%
4.    mutate(zipcode =  replace(zipcode, n()<=2, 'other')) %>%
5.    ungroup()
6.  # table(data$zipcode)
7.
8.  # neighborhood
9.  # table(data$neighbourhood_cleansed)
10. data = data %>%
11.   group_by(neighbourhood_cleansed) %>%
12.   mutate(neighbourhood_cleansed =  replace(neighbourhood_cleansed, n()<=2, 'other')) %>%
13.   ungroup()
14. # table(data$neighbourhood_cleansed)
15.
16. # property
17. # table(data$property_type)
18. data = data %>%
19.   group_by(property_type) %>%
20.   mutate(property_type =  replace(property_type, n()<=2, 'other')) %>%
21.   ungroup()
22. # table(data$property_type)
```

## 2.4. Assigning values to null

It is simpler to assign values to the numeric variables than to zip codes. Since the case of too many null values has been removed, even if the assignment does not exactly match the reality it will not have much effect on the prediction results. However, in order to minimize the error, we choose to use the value that accounts for the largest proportion of each variable to assign the value.

```
1.  # get number of response rate
2.  data$host_response_rate[is.na(data$host_response_rate)] = 0
3.  data = data %>%
4.    mutate(host_response_rate = parse_number(host_response_rate))
```

```
1.  # give values to NA
2.  #1. super host?
3.  # table(data$host_is_superhost)# f is majority
4.  table(data$host_is_superhost)
5.  data$host_is_superhost[is.na(data$host_is_superhost)] = 'f'
6.  table(data$host_is_superhost)
7.
8.  #2.profile?
9.  # table(data$host_has_profile_pic)# t is majority
10. data$host_has_profile_pic[is.na(data$host_has_profile_pic)] = 't'
11. # data = data %>%
12. #   mutate(host_has_profile_pic = ifelse(host_has_profile_pic == 'f',0,1))
13. table(data$host_has_profile_pic)
14.
15. #3. identity
16. # table(data$host_identity_verified) # f is majority
17. table(data$host_identity_verified)
18. data$host_identity_verified[is.na(data$host_identity_verified)] = 'f'
19. # data = data %>%
20. #   mutate(host_identity_verified = ifelse(host_identity_verified == 'f',0,1))
21.
22. #4. beds
23. data$beds[is.na(data$beds)] = 1
24.
25. #5. listing counts
26. data$host_listings_count[is.na(data$host_listings_count)] = 0
27. scoringData$host_listings_count[is.na(scoringData$host_listings_count)] = 0
```

After completing the above operation, the processing of null values is finished.

```
1.  > sum(which(is.na(data)))
2.  [1] 0
```

The following part of the article will further select the above variables and pick the appropriate prediction model.

However, it is also worth noting that the aforementioned variables were not all selected at the beginning. Initially, I used only 31 of them and used a linear regression model to predict prices. The RMSE of the first submission was 79. In the later model optimization, I chose to add and correct variables to make the prediction more accurate in the first place. However, since the most original code was not kept, I am able to show the modification process here.

```
1.  #Linear regression
2.  model1 = lm(price~., train)
3.  pred_test = predict(model1, newdata = test)
4.  rmse_test = sqrt(mean((pred_test - test$price)^2));rmse_test
```

# 3.  Exploratory

## 3.1.Reduced variables: Correlation

As mentioned above, the process of increasing variables was not documented. However, the process of data deletion was kept. Since too many variables may cause overfitting problems, I, therefore, chose

to look at the correlation between the numeric variables and AIC of the model for data deletion, and excluded the irrelevant variables. Code is shown below:

```
1.  #Selecting numeric data
2.  data_cor = data[,c(1, 7:12,14, 18:20,23:35, 39:41)]
3.  #Visualizing
4.  ggcorrplot(cor(data_cor),
5.             method = 'square',
6.             type = 'lower',
7.             show.diag = F,
8.             colors = c('#e9a3c9', '#f7f7f7', '#a1d76a'))
9.  res = cor(data_cor)
10. a = round(res, 3)#check the correlations
```

## 3.2. Reduced variables: Feature Selection

```
1.  #select features
2.  library(leaps)
3.  start_mod = lm(price~1,data=data_cor)
4.  empty_mod = lm(price~1,data=data_cor)
5.  full_mod = lm(price~.,data=data_cor)
6.  backwardStepwise = step(start_mod,
7.                          scope=list(upper=full_mod,lower=empty_mod),
8.                          direction='both')
```

By combining the results of the above two methods the final decision was made to delete the following variables:

```
11.  #remove: review_scores_value(-0.004)
12. which( colnames(data)=="review_scores_value" )
13. data = data[,-29]
14. #remove: host_long, last_review_long
15. data = data[,-c(39:40)]
16. #remove:review_scores_checkin, minimum_nights_avg_ntm,review_scores_communication
17. which( colnames(data)=="review_scores_checkin" )
18. data = data[,-26]
19. which( colnames(data)=="minimum_nights_avg_ntm" )
20. data = data[,-28]
21. which( colnames(data)=="review_scores_communication" )
22. data = data[,-27]
```

## 3.3. Data split

Stratified 80% of data into training set, and restore other records into test set.

```
1.  set.seed(1031)
2.  split = createDataPartition(y=data$price,p = 0.8,list = F,groups = 100)
3.  train = data[split,]
4.  test = data[-split,]
```

## 3.4.Xgboost

When I first used Xgboost directly I found that the logical type could not be applied to the model, so the conversion was finished by the following code:

```
1.  trt_data = designTreatmentsZ(dframe = data,
2.                               varlist = names(data)[c(2:35)])
3.  newvars.data = trt_data$scoreFrame[trt_data$scoreFrame$code%in% c('clean','lev'),'data.varName']
4.
5.  data_input = prepare(treatmentplan = trt_data,
6.                       dframe = data,
7.                       varRestriction = newvars.data)
8.
9.  train_input = prepare(treatmentplan = trt_data,
10.                       dframe = train,
11.                       varRestriction = newvars.data)
12.
13. test_input =  prepare(treatmentplan = trt_data,
14.                        dframe = test,
15.                        varRestriction = newvars.data)
16.
17. scoringData$price = as.integer(0)
18. scoring_input = prepare(treatmentplan = trt_data,
19.                         dframe = scoringData,
20.                         varRestriction = newvars.data)
```

Finally, apply the transformed data to the Xgboost model. Using the training set and test set to adjust and tune model, and find out the lowest RMSE.

```
1.  xgb_nrounds = xgb.cv(data=as.matrix(train_input),
2.                       label = train$price,
3.                       nrounds = 100,
4.                       nfold = 5,
5.                       eta = 0.22,
6.                       gamma = 0,
7.                       verbose = 0)
8.
9.  which.min(xgb_nrounds$evaluation_log$test_rmse_mean)
10.
11. xgboost2 = xgboost(data=as.matrix(train_input),
12.                    label = train$price,
13.                    nrounds= 99,
14.                    eta = 0.22,
15.                    gamma = 0,
16.                    verbose = 0)
17.
18. pred_data = predict(xgboost2, as.matrix(test_input))
19. rmse_data = sqrt(mean((pred_data - test$price)^2)); rmse_data
```

Then, apply the data set to the corresponding model, and eventually output the predicted document.

```
1.  xgb_nrounds_final = xgb.cv(data=as.matrix(data_input),
2.                       label = data$price,
3.                       nrounds = 250,
4.                       nfold = 5,
5.                       eta = 0.22,
6.                       gamma = 0,
7.                       min_child_weight = 4,
8.                       verbose = 0)
9.
10. which.min(xgb_nrounds_final$evaluation_log$test_rmse_mean)
11.
12. xgboost3 = xgboost(data=as.matrix(data_input),
13.                       label = data$price,
14.                       nrounds= 132,
15.                       eta = 0.22,
16.                       gamma = 0,
17.                       min_child_weight = 4,
18.                       verbose = 0)
19. pred_data2 = predict(xgboost3, as.matrix(data_input))
20. rmse_data2 = sqrt(mean((pred_data2 - data$price)^2)); rmse_data2
```

# 4. Conclusion

In this Kaggle project, two models are applied to predict the Airbnb Rental Price. The best model so far is Xgboost with final RMSE of 58.614. I gained a lot from this competition, both in terms of knowledge and in terms of life.

First of all, this project helped me to apply what I learned in the course to a practical problem, from data cleaning and variable filtering, to model building and model tuning. At the same time, I also felt the importance of data cleaning for analysis. It is the most time-consuming step, but also the most essential one. In addition, computing with models is a process that requires a lot of arithmetic power. Before finally using the model, be sure to confirm the input data and the parameters of the model. Otherwise, you will spend a lot of time in waiting.

Frankly, my first submission was earlier than most of classmates, and at that time my second submission with linear regression as the prediction model already had me ranked at about 15. This made me slacken off for a few days. In the following several days, I was just kept adding some variables a day and submitted the results directly. But as time went by, other students worked harder and harder. The ranking before going to bed was so different with the ranking after waking up. Later, I not only changed to a more complex Xgboost model, but also kept tuning it in an attempt to improve the accuracy of prediction even just a little. Looking back now, I miss and enjoy that time of hard work and persistence.

In a nutshell, I would like to thank my professor Vishal for the very practical knowledge he taught and for preparing us for this interesting competition. Thanks to Lei and Davin for answering all my doubts this semester. Wish you all the best in the work and life.