

## Kaggle report

### 1. Introduction

This report summarizes the analysis process of Airbnb rentals' price prediction in NYC and the reflection after the project. Appendix contains the R code details of the model with lowest Root Mean Square Error (RMSE) in test data prediction. Some early stage coding is not included in the appendix, but the logic behind is stated clearly in the analysis process.

### 2. Data analysis process

#### 1. Exploratory analysis

- 1) Use '*glimpse*' function from dplyr package to get a glimpse of the data, and diagnose which variables are more or less predictive. For example, host ID is generated randomly when hosts registered their account, so it is not predictive in price.
- 2) Use '*summary(is.na())*' and '*describe*' function to explore the variables with missing data. If a variable contains more than 1/3 missing values, it is risky to make predictions, thus, the variable will be considered to be dropped. For example, 'host\_response\_time'.
- 3) Test collinearity by '*VIF*' function, and use 10 as benchmark. In the dataset, the VIF between availability\_60 and availability\_90 is 22, thus the variable with less significance/importance level is dropped.

Correlation plot is also used when comparing correlation between variables. However, since the number of variables is more than 40, it is difficult to read the graph. By framing the 1600 correlation data (40\*40) into a data frame and using '*filter*' function to rank from the highest to lowest correlation, it is easier to find out the pairs of variables with highest correlation.

However, due to the nature of housing variables, many variables are correlated, such as 'bedroom' and 'bathroom', 'bedroom' and 'guest included'.

*Note: 'cor' function can only be applied on numeric variables*

#### Things to improve in this future:

- 1) To understand the overall distribution of variables and find the unusual outliers, it is efficient to use boxplot to view data visually. In this way, the price below 0 can be discovered and replaced.

- 2) For collinearity, it is difficult to determine the criteria to decide whether or not to keep the correlated variables with high importance. In other words, there is a tradeoff between importance level and correlation. I assume this requires more advanced knowledge to solve.

## 2. Variable selection

After exploration process and knowing more about the dataset and variables. Variables that are unable to contribute to prediction, such as host ID, are removed. Next step is to further select variables and clean data. Variable selection is conducted before cleaning data, because it is not necessary to clean the data that will not be used later.

Attempts: At the first attempt, common sense is used to select variables and it turns out to be wrong, because common sense creates bias. At the second attempt, feature selection is used to select variables. However, due to the collinearity issue in the variables, feature selection method could not work. Finally, I realized that because the dataset and relationship among all these variables are very complex, more flexible models such as random forest should be used, rather than linear regression model. And because the nature of linear regression and random forest are very different, variables selected by feature selection may not work for random forest. So, I decided to fit these variables into random forest model firstly, and then use the ‘importance’ function to select variables.

So far, the exploratory process gave me a basic idea about the variables I need to drop, clean, and further investigate: (1) Drop 31 variables (2) Clean 5 variables (3) Further explore 43 other variables.

### Things to improve in the future:

Deciding to drop the 31 columns by assessing its NA values/ VIF /common sense still involves bias. For example, ‘zipcode’ is not included in selected variables, because I assumed neighborhood is already included to measure location. This is an invalid assumption.

Later I realized, because location is a critical indicator of pricing in real estate industry and ‘zipcode’ contains so much specific information, it is accurate to use zip code measure location.

Overall, in future, common sense should be used less, and more scientific methods need to be used to select variables.

### 3. Text analysis

Many valuable information is contained in text, such as description, access, and amenities. Extracting text data from this dataset (csv file) do not need advanced techniques such as NPL, instead, a simple function in R can solve this problem, the *'grepl'* function.

- 1) Amenities: extract the top 10 amenities, such as Airconditioning, Kitchen, and Elevator etc.
- 2) Neighborhood: extract neighborhoods that contains more than 1000 housings, such as Chelsea, Bushwick, and Harlem etc.
- 3) Description, Access, Transit, and Host\_verification are also valuable information. However, due to the complexity, they are measured by wordcount. This is to assess how detail the information is provided by hosts, and it is initially assumed that more wordcount would lead to higher price.

### 3. Prepare the data

During the attempt to run random forest, many error messages appear. It's time to further prepare the data.

#### 1. Missing values

The first reason fitting model does not work is due to the 'NA' values. If not consider the advanced machine learning methods, the NA can be replaced by mode or median, depending on how many values are missing and the type of variable. For example, 'cleaning fee' has high importance level in random forest model, but contains more than 1500 missing values. I replaced the missing values with 0, since most housings do not have cleaning fee. Other variables such as 'if host a super host' has one missing value, and the type of the variable is logical, then it is also appropriate to replace it by 'F', the mode.

## 2. Transform format

Random forest model has restrictions on the type of predictors. In this case, 'character' variables are transformed to 'factors' by '*as.factor()*' function. The newly added dummy variables are transformed to numeric data type by '*as.numeric()*' function

## 3. Good to go: Able to run random forest model

After getting the model, '*importance()*' is used to further select variables. The variables with least importance level need to be dropped.

Note: (1) documentary is necessary for future review. (2) '*rm(list = ls())*' is necessary if need to start over the programming.

## 4. Prepare test data

Prepare test data to be ready in same order of preparing the train data, (1) create new variables from text data, (2) replace NA values, (3) transform data types to numeric or factor.

### New levels in test data:

During the prediction process, the error message indicates new levels in test data that train data does not have. By '*glimpse*' function, the variable 'property type' is found to have 3 new levels. By replacing them with the similar and existed levels in train data, the problem is solved. For example, replace 'Castle' by 'Resort'.

## **4. Conclusion**

After 31 submissions and attempts, my final score is 58.77. To further improve, it is suggested to

- (1) Work on turning parameters, and try boosting model with cross validation
- (2) Try more variable selection methods, such as Lasso
- (3) Use less common sense, more specific data information, machine learning methods
- (4) Use more advanced skills to balance the tradeoff between lower correlation between predictors and their importance level

## **Appendix: R code of the prediction with lowest rmse**

```
1. # load the data
2.
3. install.packages('dplyr')
4. library(dplyr)
5.
6. # exploratory data analysis
7. glimpse(analysisData)
8. str(analysisData)
9. sum(is.na(analysisData))
10. summary(is.na(analysisData))
11.
12. # explore correlation
13. data_numeric = analysisData %>%
14.   select(price,
15.     beds,
16.     cleaning_fee,
17.     reviews_per_month,
18.     accommodates,
19.     bathrooms,
20.     bedrooms,
21.     guests_included,
22.     extra_people,
23.     minimum_nights,
24.     maximum_nights,
25.     availability_30,
26.     availability_365,
27.     number_of_reviews,
28.     number_of_reviews_ltm,
29.     review_scores_rating,
30.     review_scores_accuracy,
31.     review_scores_cleanliness,
32.     review_scores_checkin,
33.     review_scores_communication,
34.     review_scores_location,
35.     review_scores_value,
36.     calculated_host_listings_count,
37.     calculated_host_listings_count_entire_homes,
38.     calculated_host_listings_count_private_rooms,
39.     calculated_host_listings_count_shared_rooms)
40. describe(data_numeric)
41. # correlation plot
42. library(corrplot)
43. corrplot(cor(data_numeric[, -25]), method = 'square', type = 'lower', diag = F)
44.
45. # calculate correlation
46. cor(data_numeric[, -25])
```

```
47.
48. # clean and select data in cleandata1
49. cleandata1 = analysisData %>%
50.   select(price,
51.     host_is_superhost,
52.     neighbourhood_group_cleansed,
53.     cleaning_fee,
54.     reviews_per_month,
55.     is_location_exact,
56.     room_type,
57.     accommodates,
58.     bathrooms,
59.     bedrooms,
60.     minimum_nights,
61.     maximum_nights,
62.     availability_365,
63.     availability_30,
64.     number_of_reviews,
65.     property_type,
66.     review_scores_communication,
67.     review_scores_value,
68.     review_scores_rating,
69.     instant_bookable,
70.     cancellation_policy,
71.     calculated_host_listings_count,
72.     host_listings_count,
73.     host_has_profile_pic,
74.     host_identity_verified,
75.     guests_included,
76.     extra_people)
77.
78. # find the median of missing values
79. median(analysisData$host_listings_count, na.rm = T) # median
80. summary(analysisData$host_has_profile_pic, na.rm = T)
81. summary(analysisData$host_identity_verified, na.rm = T)
82.
83. # replace missing values by median
84. cleandata1[is.na(cleandata1$host_listings_count), ]$host_listings_count <- 1 # median
85. cleandata1[is.na(cleandata1$host_has_profile_pic), ]$host_has_profile_pic <- TRUE
86. cleandata1[is.na(cleandata1$host_identity_verified), ]$host_identity_verified <- TRUE
87. cleandata1$host_has_profile_pic <- as.factor(cleandata1$host_has_profile_pic)
88. cleandata1$host_identity_verified <- as.factor(cleandata1$host_identity_verified)
89.
90.
91. # add new variable (1): wordcount
92. install.packages('ngram')
93. library(ngram)
94. cleandata1$wordcount_description <- sapply(strsplit(analysisData$description, " "), length)
```

```

95. cleandata1$wordcount_transit<- sapply(strsplit(analysisData$transit, " "), length)
96. cleandata1$wordcount_interaction<- sapply(strsplit(analysisData$interaction, " "), length)
97. cleandata1$wordcount_access<- sapply(strsplit(analysisData$access, " "), length)
98. cleandata1$wordcount_host_about<- sapply(strsplit(analysisData$host_about, " "), length)
99. cleandata1$wordcount_host_verifications<- sapply(strsplit(analysisData$host_verifications, ","),
length)
100.
101. # add new variable (2): neighborhood
102. cleandata1$Bedford_Stuyvesant <- ifelse(grepl('Bedford-Stuyvesant',
analysisData$neighbourhood_cleansed), "1","0")
103. cleandata1$Bushwick <- ifelse(grepl('Bushwick', analysisData$neighbourhood_cleansed),
"1","0")
104. cleandata1$Chelsea <- ifelse(grepl('Chelsea', analysisData$neighbourhood_cleansed),
"1","0")
105. cleandata1$Crown_Heights <- ifelse(grepl('Crown Heights',
analysisData$neighbourhood_cleansed), "1","0")
106. cleandata1$East_Village <- ifelse(grepl('East Village',
analysisData$neighbourhood_cleansed), "1","0")
107. cleandata1$Harlem <- ifelse(grepl('Harlem', analysisData$neighbourhood_cleansed),
"1","0")
108. cleandata1$West_Village <- ifelse(grepl('West Village',
analysisData$neighbourhood_cleansed), "1","0")
109.
110. # add new variable (3): amenities
111. cleandata1$Airconditioning <- ifelse(grepl('Air conditioning', analysisData$amenities),
"1","0")
112. cleandata1$Kitchen <- ifelse(grepl('Kitchen', analysisData$amenities), "1","0")
113. cleandata1$Essentials <- ifelse(grepl('Essentials', analysisData$amenities), "1","0")
114. cleandata1$Elevator <- ifelse(grepl('Elevator', analysisData$amenities), "1","0")
115. cleandata1$Doorman <- ifelse(grepl('Doorman', analysisData$amenities), "1","0")
116. cleandata1$Gym <- ifelse(grepl('Gym', analysisData$amenities), "1","0")
117. cleandata1$Internet <- ifelse(grepl('Internet', analysisData$amenities), "1","0")
118. cleandata1$Smoking <- ifelse(grepl('Smoking allowed', analysisData$amenities), "1","0")
119. cleandata1$parking <- ifelse(grepl('Free street parking', analysisData$amenities), "1","0")
120. cleandata1$Heating <- ifelse(grepl('Heating', analysisData$amenities), "1","0")
121.
122. sum(is.na(cleandata1))
123.
124.
125. # continue to replace NA by median in cleandata1
126. cleandata1[is.na(cleandata1$cleaning_fee), ]$cleaning_fee <- 0 # 20
127. cleandata1[is.na(cleandata1$host_is_superhost), ]$host_is_superhost <- FALSE # after help
128. cleandata1[is.na(cleandata1$reviews_per_month), ]$reviews_per_month <- 0.96
129. median(analysisData$reviews_per_month, na.rm = T) #0.96
130.
131. # convert data to factors (in order to run random forest)
132. cleandata1$neighbourhood_group_cleansed <-
as.factor(cleandata1$neighbourhood_group_cleansed)

```

```
133. cleandata1$room_type <- as.factor(cleandata1$room_type)
134. cleandata1$cancellation_policy <- as.factor(cleandata1$cancellation_policy)
135. cleandata1$property_type <- as.factor(cleandata1$property_type)
136. cleandata1$host_is_superhost <- as.factor(cleandata1$host_is_superhost)
137. cleandata1$is_location_exact <- as.factor(cleandata1$is_location_exact)
138. cleandata1$instant_bookable <- as.factor(cleandata1$instant_bookable)
139. cleandata1$Airconditioning <- as.numeric(cleandata1$Airconditioning)
140. cleandata1$Kitchen <- as.numeric(cleandata1$Kitchen)
141. cleandata1$Essentials <- as.numeric(cleandata1$Essentials)
142. cleandata1$Elevator <- as.numeric(cleandata1$Elevator)
143. cleandata1$Doorman <- as.numeric(cleandata1$Doorman)
144. cleandata1$Gym <- as.numeric(cleandata1$Gym)
145. cleandata1$Internet <- as.numeric(cleandata1$Internet)
146. cleandata1$Smoking <- as.numeric(cleandata1$Smoking)
147. cleandata1$parking <- as.numeric(cleandata1$parking)
148. cleandata1$Heating <- as.numeric(cleandata1$Heating)
149.
150. # explore clean dataset
151. str(cleandata1)
152. sum(is.na(cleandata1))
153. glimpse(cleandata1)
154.
155. # run random forest model
156. install.packages('randomForest')
157. library(randomForest)
158. set.seed(66)
159. forestmodel = randomForest(price~., data = cleandata1, ntree = 300)
160. plot(forest)
161. varImpPlot(forest)
162. importantlist3 = importance(forestmodel)
163. varImpPlot(forest); importance(forest) ## see variable importance
164.
165.
166.
167. # clean data in test (scroing data)
168. cleandata1_test = scoringData %>%
169.   select(host_is_superhost,
170.          neighbourhood_group_cleansed,
171.          cleaning_fee,
172.          reviews_per_month,
173.          is_location_exact,
174.          room_type,
175.          accommodates,
176.          bathrooms,
177.          bedrooms,
178.          minimum_nights,
179.          maximum_nights,
180.          availability_365,
```



```

181.     availability_30,
182.     number_of_reviews,
183.     property_type,
184.     review_scores_communication,
185.     review_scores_value,
186.     review_scores_rating,
187.     instant_bookable,
188.     cancellation_policy,
189.     calculated_host_listings_count,
190.     host_listings_count,
191.     host_has_profile_pic,
192.     host_identity_verified,
193.     guests_included,
194.     extra_people)
195.
196. # replace missing values by median
197. cleandata1_test[is.na(cleandata1_test$host_listings_count), ]$host_listings_count <- 1 #
    median
198. cleandata1_test[is.na(cleandata1_test$host_has_profile_pic), ]$host_has_profile_pic <-
    TRUE
199. cleandata1_test[is.na(cleandata1_test$host_identity_verified), ]$host_identity_verified <-
    TRUE
200. cleandata1_test$host_has_profile_pic <- as.factor(cleandata1_test$host_has_profile_pic)
201. cleandata1_test$host_identity_verified <- as.factor(cleandata1_test$host_identity_verified)
202.
203. median(cleandata1_test$host_listings_count, na.rm = T) # median
204. summary(scoringData$host_has_profile_pic, na.rm = T)
205. summary(scoringData$host_identity_verified, na.rm = T)
206.
207. # (1) add new variable: wordcount
208. library(ngram)
209. cleandata1_test$wordcount_description<- sapply(strsplit(scoringData$description, " "),
    length)
210. cleandata1_test$wordcount_transit<- sapply(strsplit(scoringData$transit, " "), length)
211. cleandata1_test$wordcount_interaction<- sapply(strsplit(scoringData$interaction, " "),
    length)
212. cleandata1_test$wordcount_access<- sapply(strsplit(scoringData$access, " "), length)
213. cleandata1_test$wordcount_host_about<- sapply(strsplit(scoringData$host_about, " "),
    length)
214. cleandata1_test$wordcount_host_verifications<-
    sapply(strsplit(scoringData$host_verifications, ","), length)
215.
216. # (2) add new variable: neighborhood
217. cleandata1_test$Bedford_Stuyvesant <- ifelse(grepl('Bedford-Stuyvesant',
    scoringData$neighbourhood_cleansed), "1", "0")
218. cleandata1_test$Bushwick <- ifelse(grepl('Bushwick',
    scoringData$neighbourhood_cleansed), "1", "0")

```

```

219. cleandata1_test$Chelsea <- ifelse(grepl('Chelsea', scoringData$neighbourhood_cleansed),
    "1","0")
220. cleandata1_test$Crown_Heights <- ifelse(grepl('Crown Heights',
    scoringData$neighbourhood_cleansed), "1","0")
221. cleandata1_test$East_Village <- ifelse(grepl('East Village',
    scoringData$neighbourhood_cleansed), "1","0")
222. cleandata1_test$Harlem <- ifelse(grepl('Harlem', scoringData$neighbourhood_cleansed),
    "1","0")
223. cleandata1_test$West_Village <- ifelse(grepl('West Village',
    scoringData$neighbourhood_cleansed), "1","0")
224.
225. # (3) add new variable: amenities
226. cleandata1_test$Airconditioning <- ifelse(grepl('Air conditioning', scoringData$amenities),
    "1","0")
227. cleandata1_test$Kitchen <- ifelse(grepl('Kitchen', scoringData$amenities), "1","0")
228. cleandata1_test$Essentials <- ifelse(grepl('Essentials', scoringData$amenities), "1","0")
229. cleandata1_test$Elevator <- ifelse(grepl('Elevator', scoringData$amenities), "1","0")
230. cleandata1_test$Doorman <- ifelse(grepl('Doorman', scoringData$amenities), "1","0")
231. cleandata1_test$Gym <- ifelse(grepl('Gym', scoringData$amenities), "1","0")
232. cleandata1_test$Internet <- ifelse(grepl('Internet', scoringData$amenities), "1","0")
233. cleandata1_test$Smoking <- ifelse(grepl('Smoking allowed', scoringData$amenities), "1","0")
234. cleandata1_test$parking <- ifelse(grepl('Free street parking', scoringData$amenities),
    "1","0")
235. cleandata1_test$Heating <- ifelse(grepl('Heating', scoringData$amenities), "1","0")
236.
237. cleandata1_test$Airconditioning <- as.numeric(cleandata1_test$Airconditioning)
238. cleandata1_test$Kitchen <- as.numeric(cleandata1_test$Kitchen)
239. cleandata1_test$Essentials <- as.numeric(cleandata1_test$Essentials)
240. cleandata1_test$Elevator <- as.numeric(cleandata1_test$Elevator)
241. cleandata1_test$Doorman <- as.numeric(cleandata1_test$Doorman)
242. cleandata1_test$Gym <- as.numeric(cleandata1_test$Gym)
243. cleandata1_test$Internet <- as.numeric(cleandata1_test$Internet)
244. cleandata1_test$Smoking <- as.numeric(cleandata1_test$Smoking)
245. cleandata1_test$parking <- as.numeric(cleandata1_test$parking)
246. cleandata1_test$Heating <- as.numeric(cleandata1_test$Heating)
247. sum(is.na(cleandata1_test))
248.
249. # (4) replace NA value
250. cleandata1_test[is.na(cleandata1_test$cleaning_fee), ]$cleaning_fee <- 0 # after help
251. cleandata1_test[is.na(cleandata1_test$host_is_superhost), ]$host_is_superhost <- FALSE #
    after help
252. cleandata1_test[is.na(cleandata1_test$reviews_per_month),]$reviews_per_month <- 0.95
253. median(cleandata1_test$reviews_per_month, na.rm = T) #0.95
254.
255. # (5) convert to factors
256. cleandata1_test$neighbourhood_group_cleansed <-
    as.factor(cleandata1_test$neighbourhood_group_cleansed)
257. cleandata1_test$room_type <- as.factor(cleandata1_test$room_type)

```

```

258. cleandata1_test$ cancellation_policy <- as.factor(cleandata1_test$ cancellation_policy)
259. cleandata1_test$ bed_type <- as.factor(cleandata1_test$ bed_type)
260. cleandata1_test$ property_type <- as.factor(cleandata1_test$ property_type)
261.
262. # (6) Deal with new levels in Property type
263. # replace with levels that are similar
264. cleandata1_test[cleandata1_test$ property_type == 'Castle',]$ property_type <- 'Resort'
265. cleandata1_test[cleandata1_test$ property_type == 'Farm stay',]$ property_type <- 'House'
266. cleandata1_test[cleandata1_test$ property_type == 'Casa particular (Cuba)',]$ property_type
    <- 'House'
267. dada = cleandata1_test[cleandata1_test$ property_type == 'Casa particular (Cuba)',]
268. levels(cleandata1$ property_type)
269. levels(cleandata1_test$ property_type)
270. str(cleandata1)
271. str(cleandata1_test)
272. ?droplevels.factor
273. library(gdata)
274. drop.levels(cleandata1_test$ property_type)
275.
276. x <- x[x!="Casa particular (Cuba)"]
277. df <- data.frame(a=x,b=x,c=x)
278. df <- drop.levels(df)
279. factor
280. cleandata1_test$ property_type <- droplevels(cleandata1_test$ property_type )
281. cleandata3_test = select(cleandata1_test , -c('minutes','train','subway'))
282.
283.
284. # convert test data to types of train data
285. cleandatatestaaaa = cleandata1_test
286. cleandatatestaaaa <- rbind(cleandata1[1, -1], cleandatatestaaaa)
287. cleandatatestaaaa <- cleandatatestaaaa[-1,]
288. str(cleandata3_test)
289. str(cleandata1)
290. sum(is.na(cleandata1_test))
291.
292. # predict testing data
293. install.packages('caret')
294. install.packages('ggplot2')
295. library(randomForest)
296. library(caret)
297.
298. library(ggplot2)
299. library(caret)
300. predforest2 = predict(forestmodel,newdata=cleandatatestaaaa)
301. predforest2
302. # score = 58.77108
303.
304. # write file in csv

```

Huizhe ZHU (hz2657)

```
305. submission14 = data.frame(id=scoringData$id, price = predforest2)
306. write.csv(submission14,'14submission.csv',row.names = F)
```