

Analysis of Resort Hotel

Minji Song

Introduction

In this project, I used Hotel Booking Demand data which is from the course website. This dataset includes 2097 American visitor's hotel booking information. The dataset includes 18 variables. There are 6 categorical variables describing booking information of guests such as hotel, arrive date month, meal, market segment, reserved room type, and customer type. Meal indicates the type of meal booked. BB is Bed & breakfast, HB is breakfast and one other meal, FB is breakfast, lunch, and dinner and SC/Undefined indicates no meal package. Market segment has two types: Online TA, TO, direct, corporate, and groups. TA means Travel Agents and TO means Tour Operator. Customer type has four different types: Contract, group, transient, transient party. There are also 13 categorical variables including the value indicating if the booking was canceled, the number of days between booking date and arrival date, year of the arrival date, day of the arrival date, week number of the arrival date, the number of stays in weekend nights, the number of stay in weeknights, number of adults, children, and babies, average daily rate, a total of special requests.

Exploratory Data Analysis

Since I am in session 1, I analyzed corresponding to Resort Hotel which contains 479 guests' information. Thus, I split the data only for Resort Hotel through Excel and then split the data into training data and test data. My goal of this project is to determine if there is any relationship between the response variable and other predictor variables. I chose adr, average daily rate, variable as a response variable because it is important for a hotel to analyze which factors affect the most on a daily rate which is directly related to the sales. Before I do anything first, I changed some values of the is_canceled variable and meal variable. is_canceled variable is a categorical variable with 0 and 1, so I changed it to 0 to No and 1 to Yes to identify values easily. For meal variables, SC and Undefined have the same meaning but they are different values, thus I combined both values to name it as no meal. After I changed some values of categorical variables, I performed Explanatory Data Analysis by plotting each categorical variable to see the relationship with adr. As we look at Figure 1, five bar plots of categorical variables show the proportion of each variable. Meal and reserved room types are right-skewed and the market segment is left-skewed. There is no variable that has a perfect normal graph. In Figure 2, I plotted a box plot for categorical variables as well. For the arrival date month variable, we can see that August is associated with high adr, and January has low adr. Reserved room type H has a higher adr with about 180 adr. Other variables in this plot have similar trend. If we look at Figure 3, these plots show the relationship between continuous variables and adr. It seems that lead time and arrival date week numbers have distinct trend and we can assume that

they can be related to each other somehow and this is the part we need to check by using interaction term. Figure 4 shows correlation plots where we see that arrival date week number and arrival date year tend to have a negative correlation more than other variables. Residual plots indicate if any numerical variables have a nonlinear trend and we can check this in Figure 5. All plots are well distributed and there is no clear trend. Even though lead time plot looks like to be a heteroscedasticity a bit in a little plot, it is distributed well enough. From the residual plots, none of the numerical variables support the nonlinear trend.

Method

Based on explanatory data analysis, I performed linear regression using all variables and interaction terms with lead time and arrival date week numbers. The model includes 16 variables and one interaction term that I wanted to check if it is significant variable in the model. In a summary of the regression, the p-value of the interaction term is lower than 0.05, so we can conclude this is a significant variable and keep this in the model. Since we have a lot of variables, I performed the AIC model selection to reduce some variables based on the small AIC value. Through the selection, the final model has 14 variables: is canceled, lead time, arrival date year, arrival date month, arrival date week number, stay weeknights, adults, children, meal, market segment, reserved room type, customer type, a total of special requests and the interaction terms. Before we confirm the final model, I performed the VIF test to see if there is any multicollinearity with this model. It seems like there is multicollinearity in this model because some values have high VIF values. The highest VIF is arrival date week number with 81.3846, there are few more variables that have over 10 VIFs. However, those variables that have high VIF values can be ignored since those are either categorical variables with more than two levels or part of interaction terms. Other than those variables, all other variables look good. Thus, this would be my final model. With this final model and variables, I predicted how each predictor influences the response variable by using minimum, median, and maximum number of the predictors in data.

Through the prediction, it has the result that a large number of arrival date week number has a higher average daily rate than a small number of arrival date week number. When the arrival date week number is 1, 29 and 53, the predicted adr are 21.74, 103.47, and 173.52. It means the second half of the year has a higher rate. The lead time has the opposite results. The shorter lead time has a higher average daily rate by 30% than 537 days longer lead time. The shorter time between booking date and arrival date have a higher rate. This data set is for between 2015 and 2017. In 2015, the average daily rate is much lower than in 2017. In 2017, the rate has increased by about 45% compared to 2015 and it was -6.5 in 2015. The more nights staying in a week increases the average daily rate. The difference between 0 nights to 10 nights is 28.1%. More adults and children also increase the average daily rate. Also, a lot of special requests increase the daily average rate. A special request increases the rate by about 4.5%.

The next regression I performed is regression trees using rpart function. This regression split the model into several steps and the first node of the variable we split on is arrival date month with 181 observations. The second branch is also an arrival date month with 202 observations. It is

telling that the most important variable in this regression is the arrival date month. This tree plot can be seen in Figure 6. As we look at the size of the tree plot with 11 terminal nodes in Figure 7, error reduction diminish as the tree grows deeper. We can fit a bigger tree by pruning the tree using $cp=0.001$ instead of $cp=0.01$. To determine the optimal tree size, we can pick either optimal tree size with the smallest Cross-Validation error or smallest tree size with CV error within 1 standard error of the smallest CV error. In this model, I chose the second option and the red line of the plot in Figure 8 indicates the CV error within 1 standard error of the smallest CV error. Figure 9 is the smallest tree plot based on the CV error that looks simpler than the first tree plot. The optimal cp value has to be below red line in plot and the final plot for the model can be seen in Figure 10. Using the same predictors as linear regression in a tree model, I fitted the random forests to see how training error and cross-validation test errors are different. The training error is 101966.9 and the test error is 381974.4. Thus, we can see that training error is much smaller than the test error, and the training dataset is more accurate.

Result

This hotel booking demand dataset was analyzed through linear regression analysis and regression tree analysis to interpret the relationship between the average daily rate and other predictors. In linear regression analysis, 14 continuous and categorical variables are selected by AIC model selection and this became the final model. To be more specific, it turns out that a large number of arrival date week increase the average daily rate, which means the second half of the year have higher average daily rate than the first half of the year. For the lead time, the number of days between the date of customer book and the arrival date, shorter lead time have a higher average daily rate than longer lead time. The spontaneous customer would pay a higher rate than a planned customer. The more nights customers stay in a week, the higher the average daily rate they have to pay. Also, more adults or more babies have to pay a higher average daily rate. Increasing the total of special requests increase the average daily rate as well.

Through regression tree analysis with the same predictors in a previous regression model, the most important variable in determining the average daily rate is the arrival date month. The customer who visit the hotel in August, July and June tend to pay a higher rate than other periods because those three months are the time people usually go on a vacation. Reserved room type is the second important variable in determining the average daily rate. If the room type is not A or D, it has a higher average daily rate. However, we do not have any details about the room type, so we assume that relatively cheap rooms which is types A and D might have less room quality than other room types. Two analyses I performed have different results, however, this would give the hotel ideas what factors play an important role in determining an average daily rate that is related to increase the sales.

Appendix

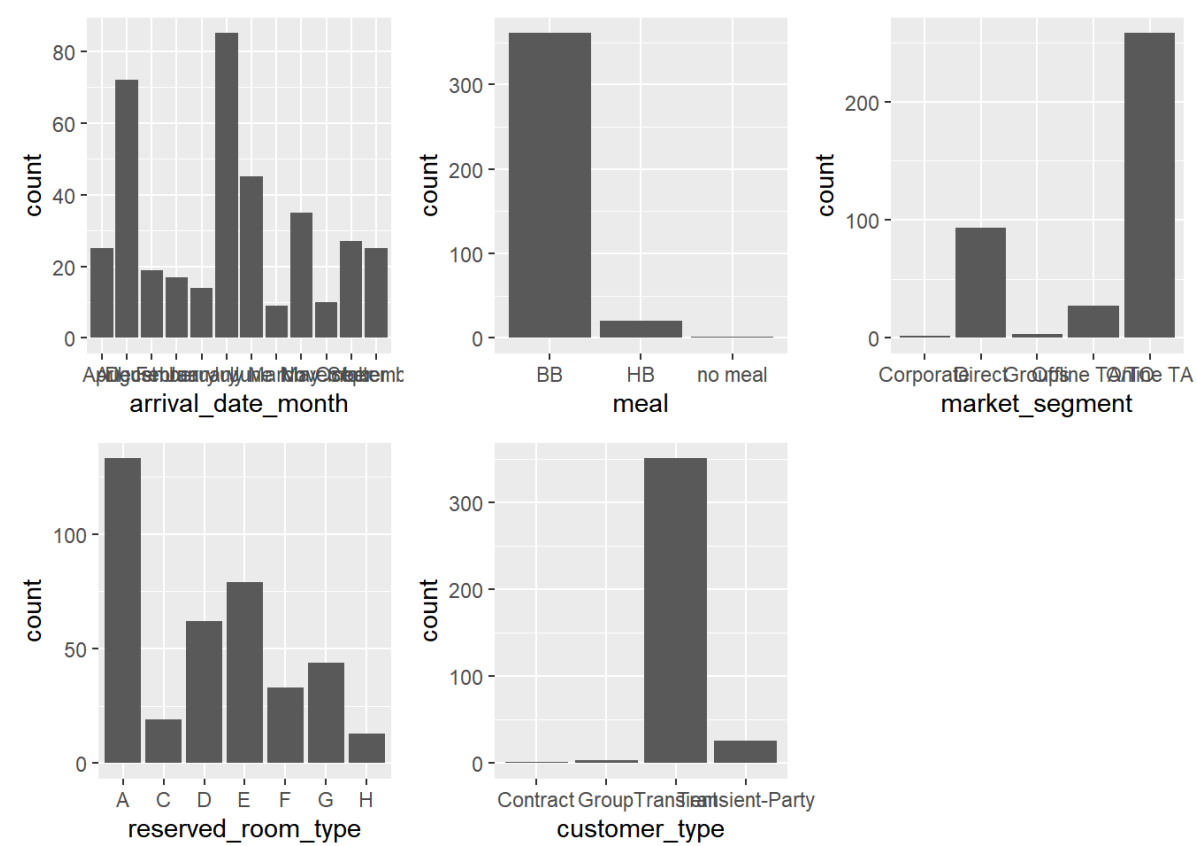


Figure 1

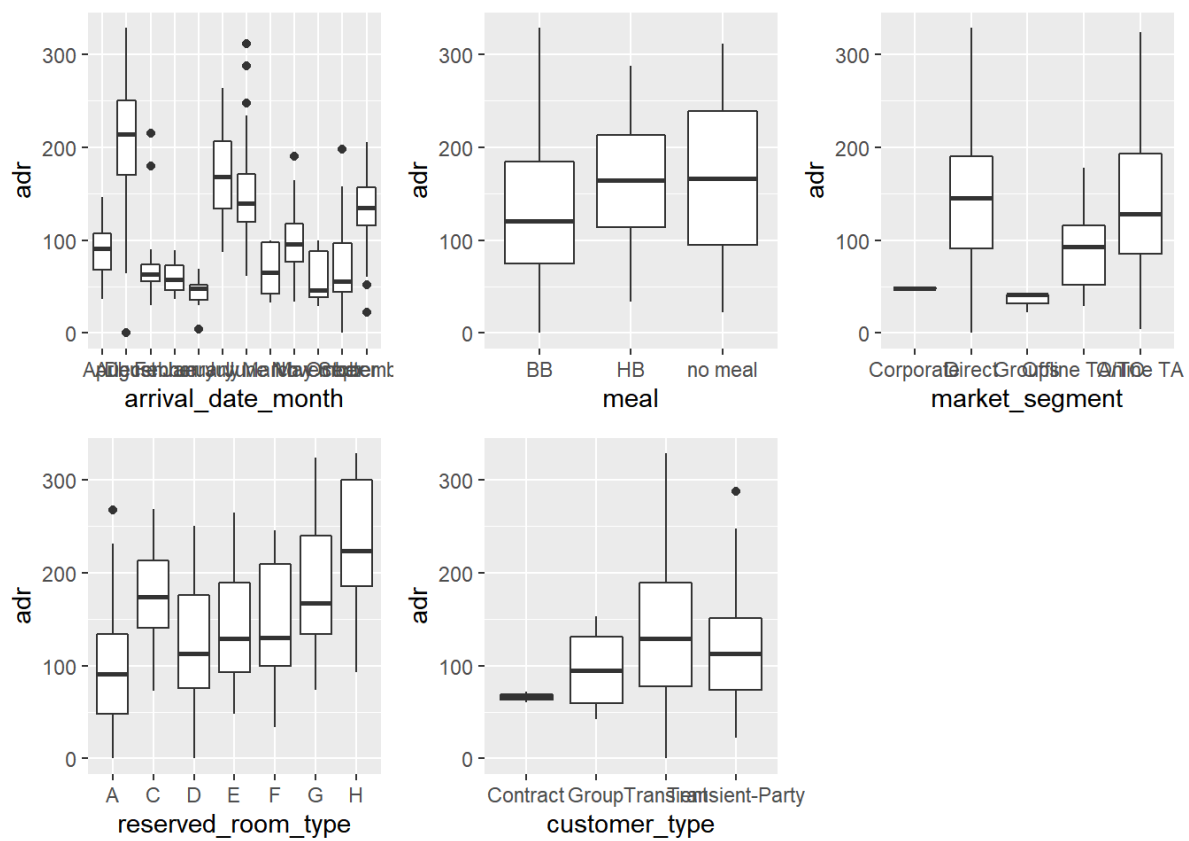


Figure 2

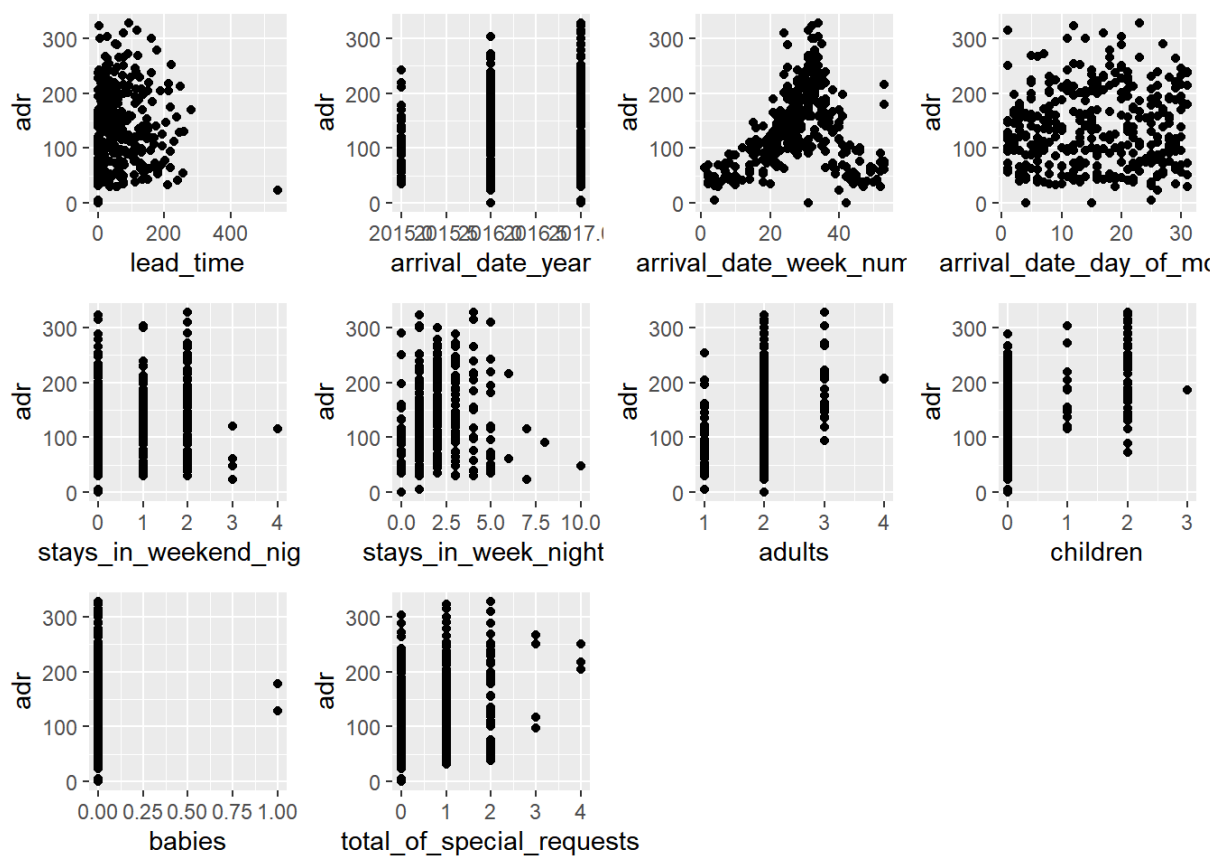


Figure 3

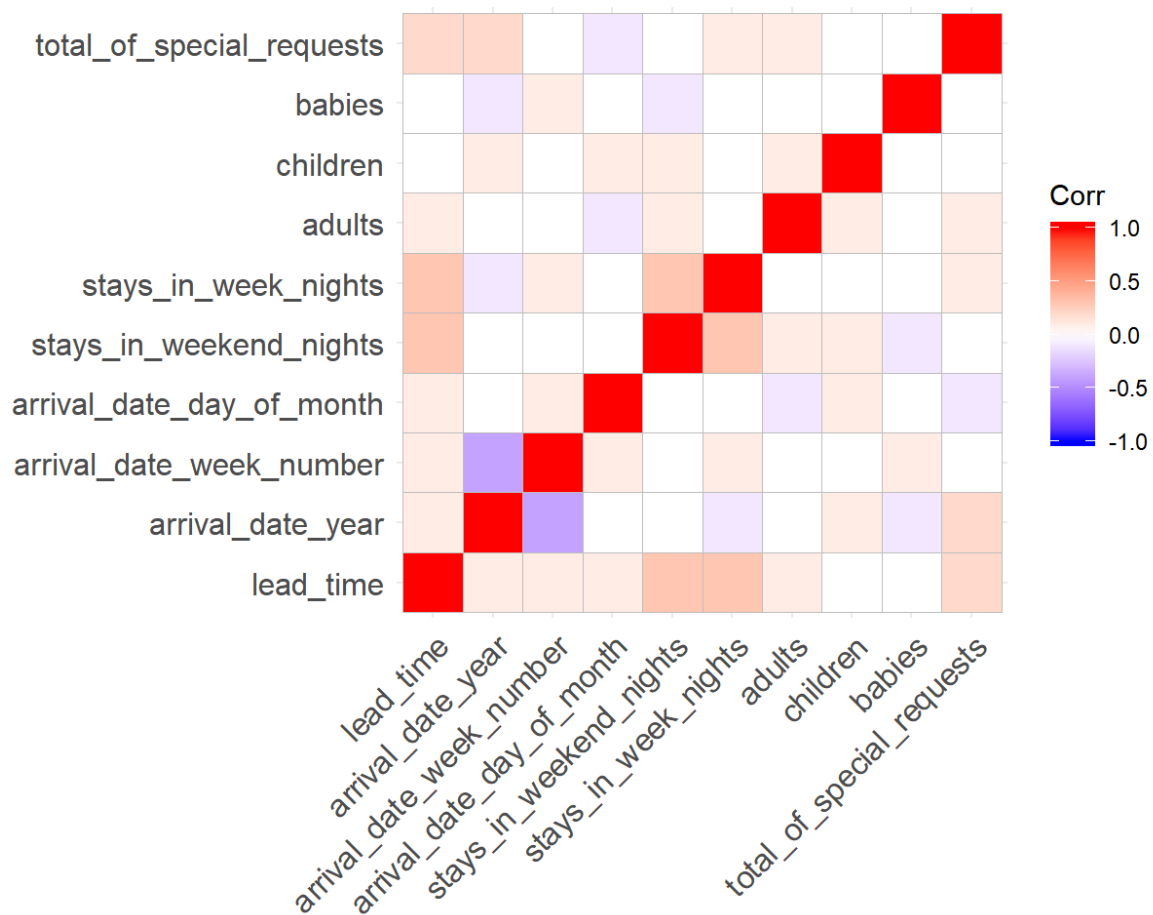


Figure 4

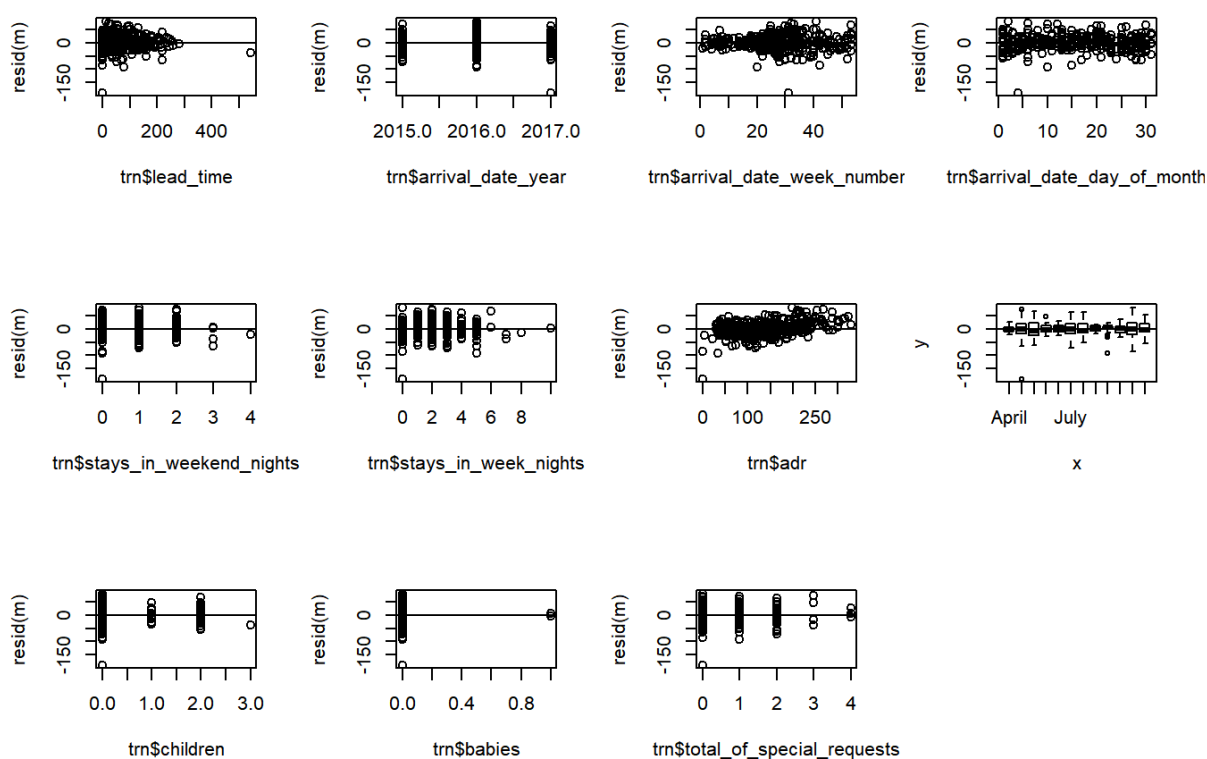


Figure 5

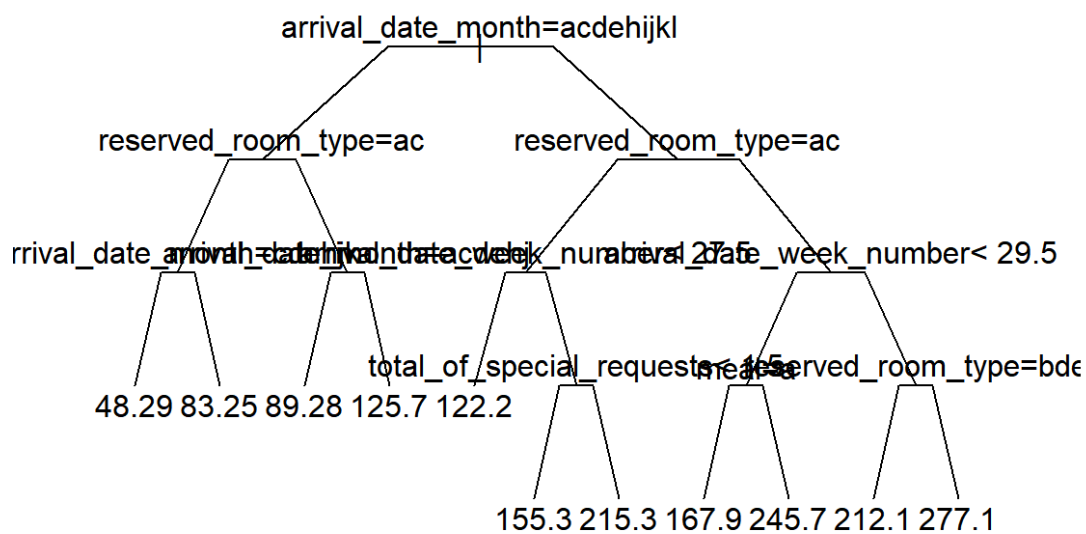


Figure 6

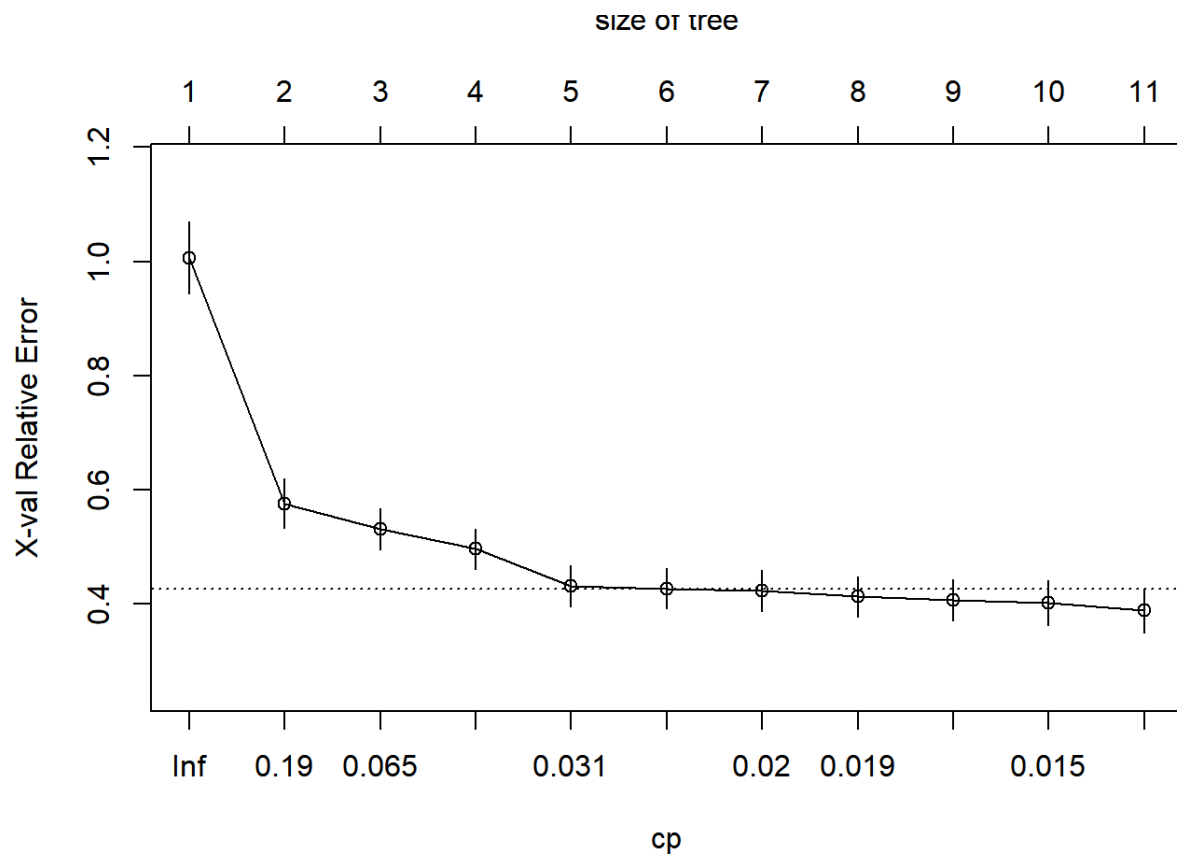


Figure 7

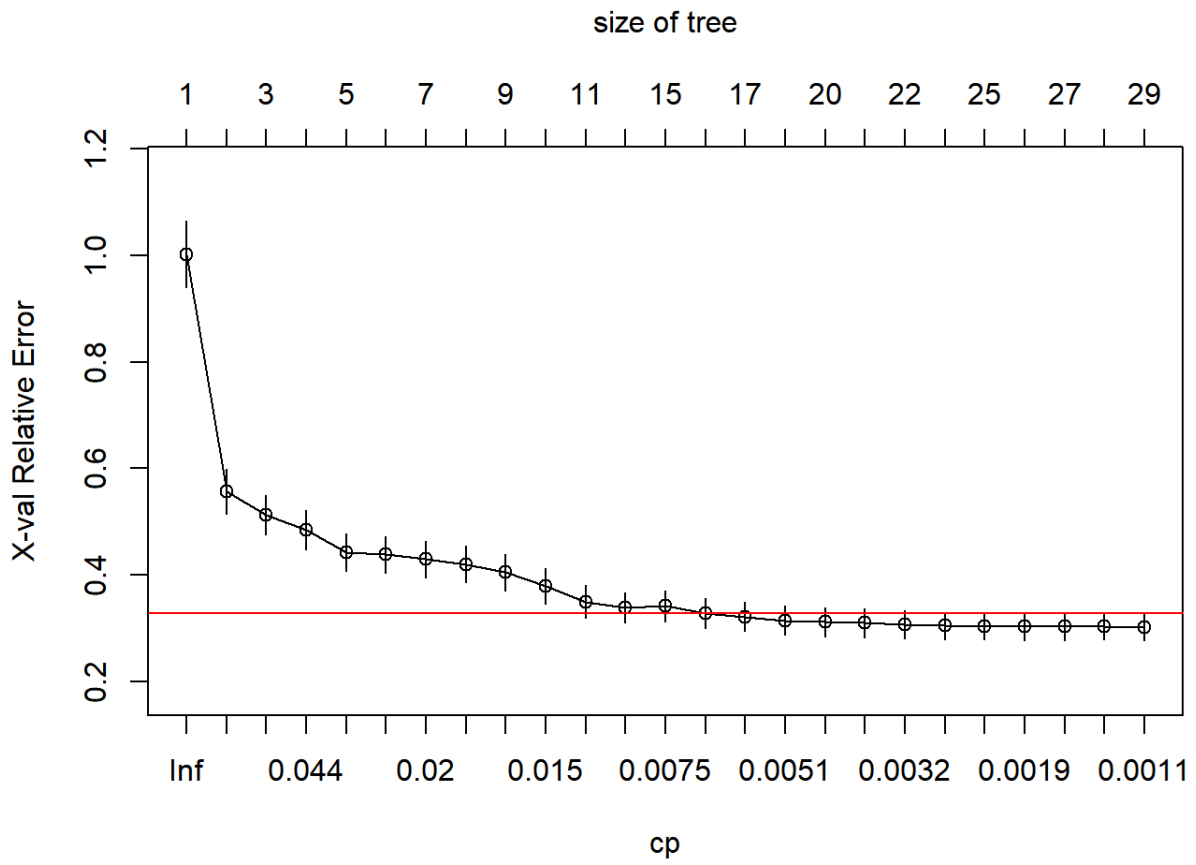


Figure 8

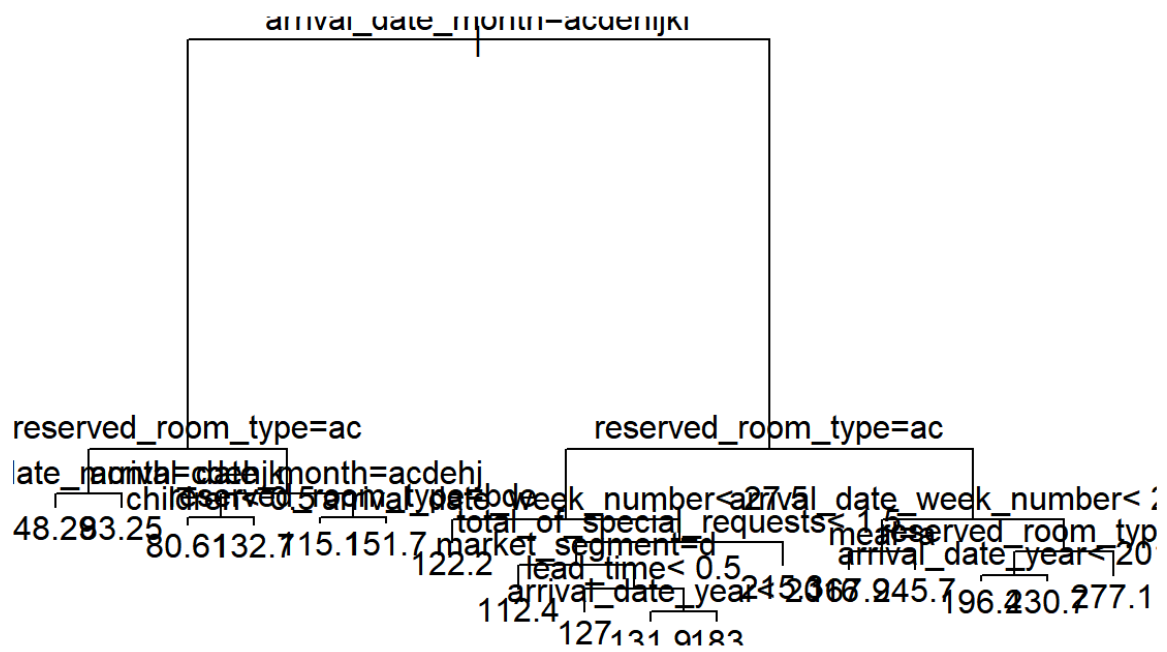


Figure 9

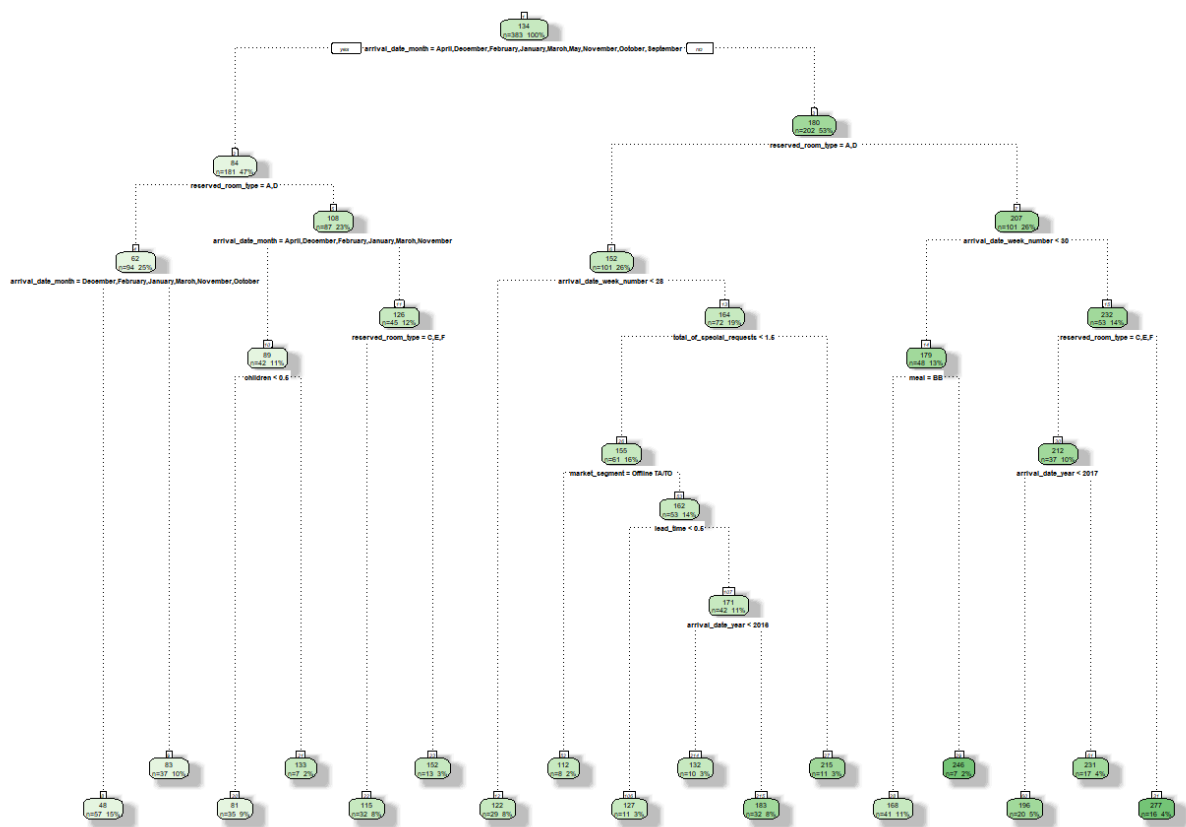


Figure 10