Jazmaine Vitta

May 16th, 2023
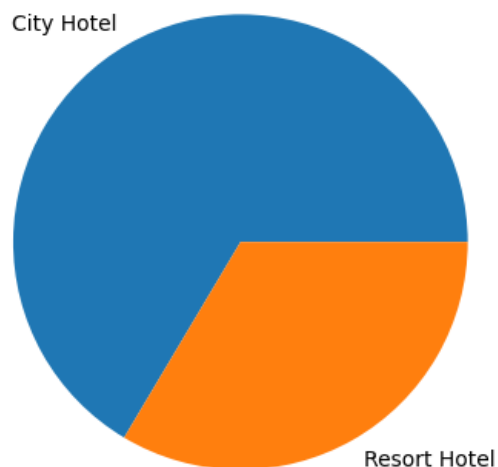
# Predicting Hotel Booking Cancellations

[Link to Kaggle Dataset](#)

The inability to anticipate whether or not a customer's hotel booking will be canceled undoubtedly causes major setbacks for the hotel business. As a result of cancellations, hotels can not realistically estimate the total number of guests each night causing them to be inadequately equipped to manage customers, resulting in a loss in revenue and countless other complications. Due to the uncertainty that comes with the possibility of cancellations, hotels are unable to accurately forecast their daily demands. In most cases when a customer does not show up on the arrival date of their stay or cancels close to the date of their booking, those rooms are either resold at a lesser rate than they were originally booked for or even worse not resold at all. These issues not only lose revenue for the business but also cause dissatisfaction with customers. To beat the negative effects of cancellations hotels attempt to overbook rooms which comes with much more issues. Cancellations are a direct cause of loss of revenue in numerous ways.
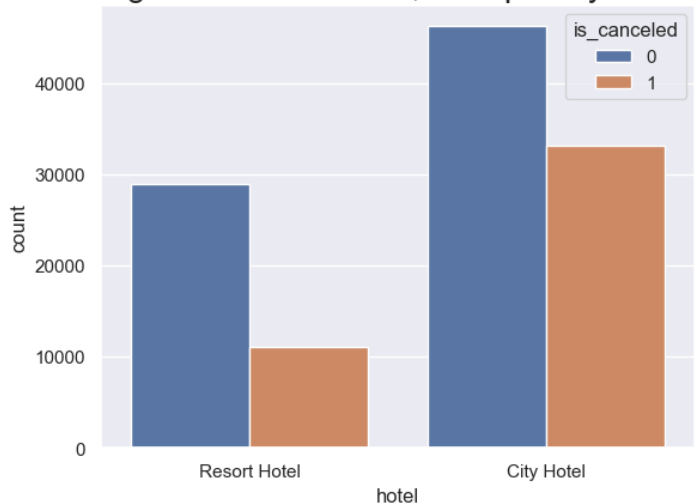
By building a classification model to predict the likelihood of a cancellation, I can assist businesses in overcoming the uncertainty and loss in revenue caused by cancellations. Using a data set from Kaggle.com the data was first cleaned and then

trained to find any correlations or causations to the change in booking status. The data set used to build this model contains booking information during the years of 2015 to 2017. The data is made up of a total of 119,389 booking observations with 32 columns, of which 20 are categorical and 12 are numerical. Throughout my process of preparing the data I was able to uncover the answers to many interesting questions. Based on the information the target feature is the 'is_canceled' column.
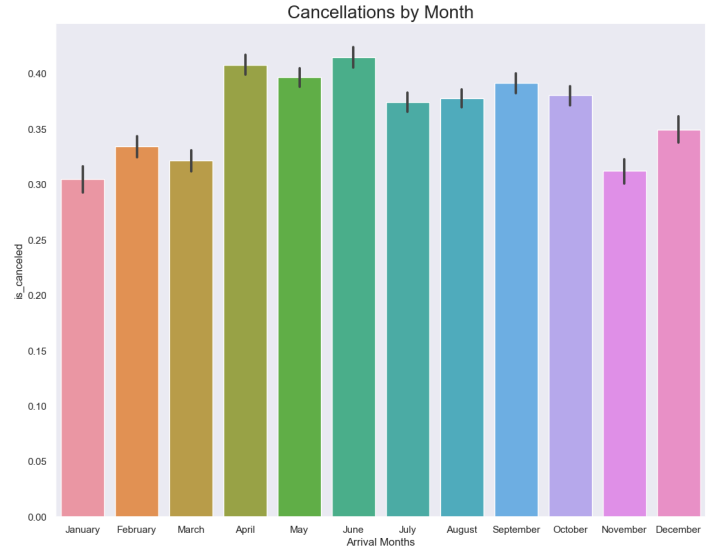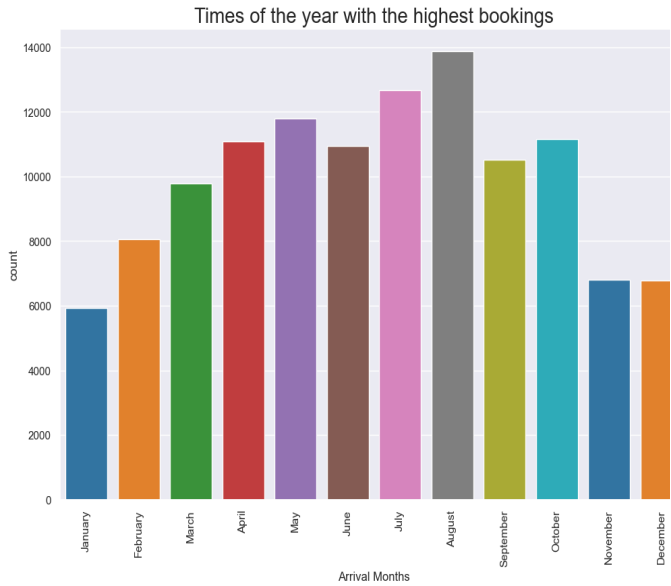
The data contains more than double the amount of city hotels compared to the amount of resort hotels. Between the two, City hotels experience almost as many cancellations as they do bookings, which is a big contrast compared to resort hotels who receive around a third of their bookings canceled.
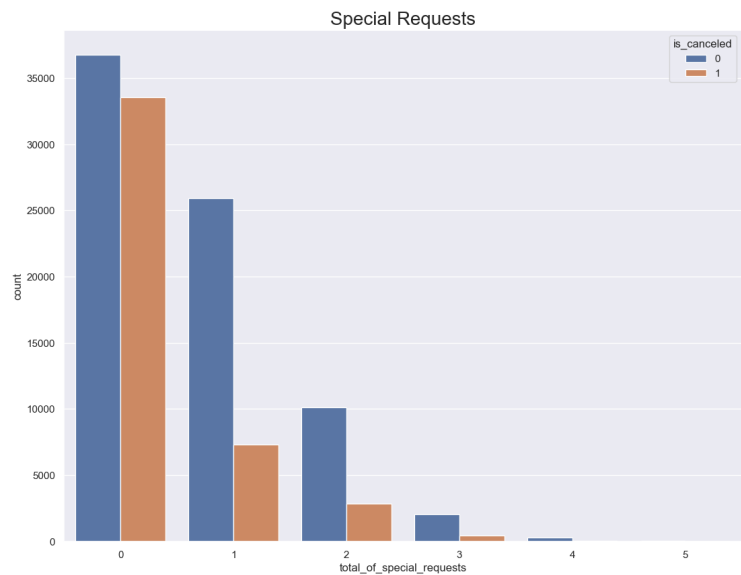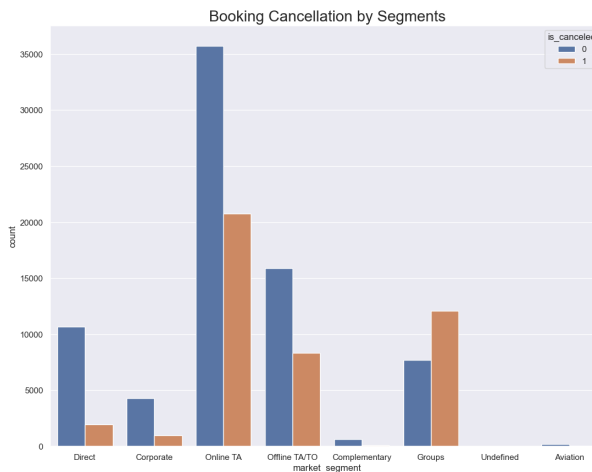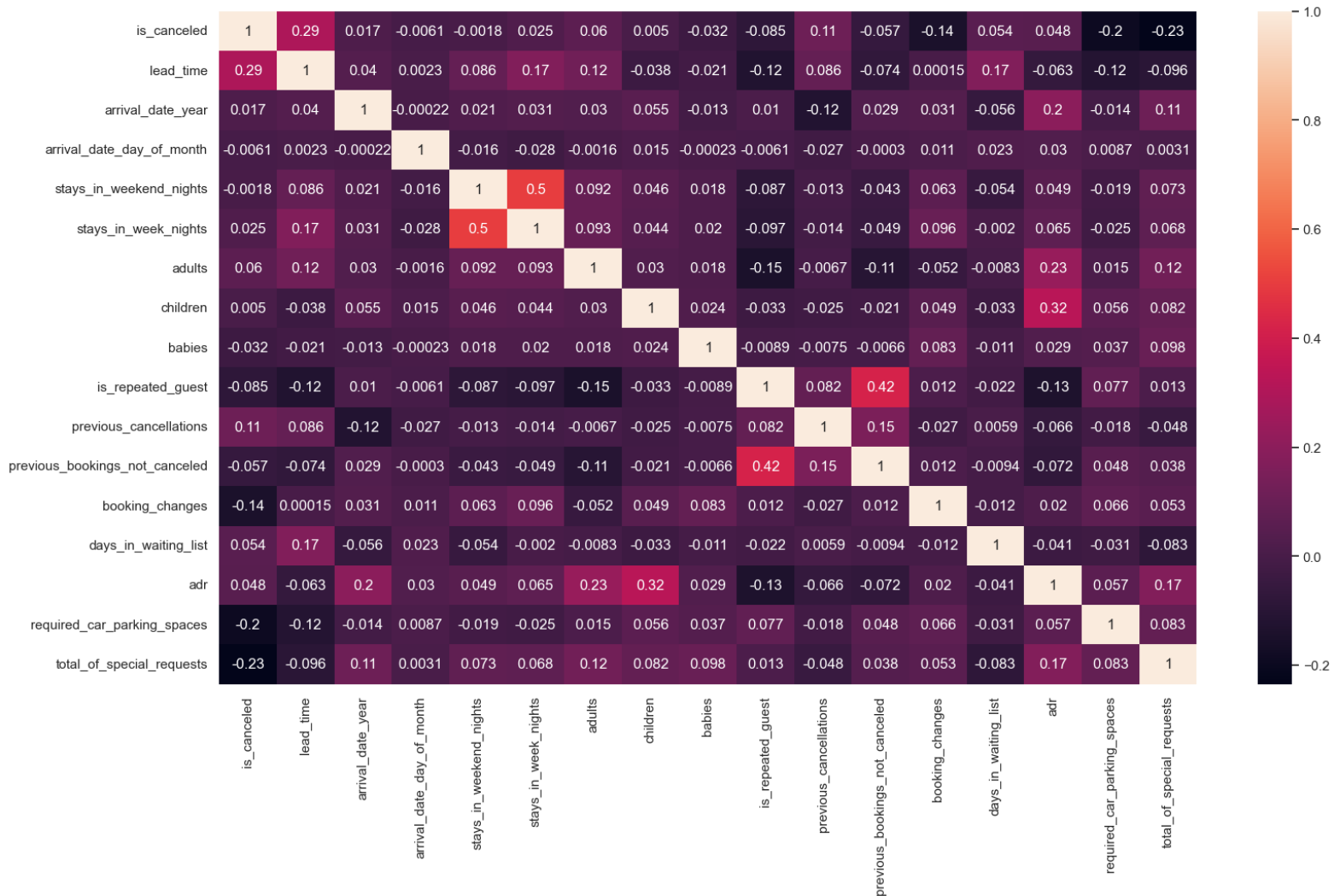


Although the summer months are the most popular time of the year for bookings, the spring months seem to receive the most cancellations.

Times of the year with the highest bookings



Cancellations by Month

Customers who book directly with the hotel are much less likely to cancel as opposed to customers who book with a travel agent or online via a third party website. Bookings made with at least one special request are drastically less likely to cancel compared to those made with zero requests.



Booking Cancellation by Segments



Special Requests

As shown in the heatmap below, the features most correlated to the target feature 'is_canceled' are 'lead_time', 'total_of_special_requests', 'booking_changes' and 'required_car_parking_spaces'.



With the knowledge of columns that have the most relation to the target feature I differentiated which features are categorical and which are numerical. This was useful in determining the type of visualization that I later performed. The data in the categorical features were One Hot End Coded to be able to make proper estimates.The data was then split into training data and testing data. Using the now split data a was able to build a Logistic Regression model which  predicted booking cancellations with an accuracy of only

81%. With all the data I know the accuracy for predicting could be at least 90%. A Random Forest Model was built and we can now predict whether or not a new booking will be canceled with 93% accuracy. The Random Forest Classifier is no doubt the better model. This model has a much lower cross-validation of 0.926. It also exhibits less variability with a mean absolute error of 0.0700226.

Predicting cancellations is a real problem for the hotel industry! Having a good understanding of this problem and the features that closely relate with cancellations will be very useful to your business. The model I created will no doubt assist in decreasing the possibility of your business losing revenue by being under/over prepared. My Random Forest Model has a high accuracy of 93.09%, to predict hotel booking cancellation. Gathering more observations for resort hotels to help balance the data would also assist in strengthening the predictability of the model. In turn, this recommendation would help the model become even more of an asset to your business.