# INN Hotels- Logistic Regression/Decision Tree Project

## UT Data Science & Business Analytics

May 2023

# Contents / Agenda

- Executive Summary

- Business Problem Overview and Solution Approach

- EDA Results

- Data Preprocessing

- Model Performance Summary

- Appendix

# Executive Summary

- Lead time is a major factor in customers cancelling their reservations. Hotels could place a limit on the extent of how far out a customer can book their reservation to help cut down cancellations or hotels could implement a fee for cancellations to help deter customers from cancelling.

- Staff could offer more special requests available to customers since the data shows that customers who had special requests were less likely to cancel.

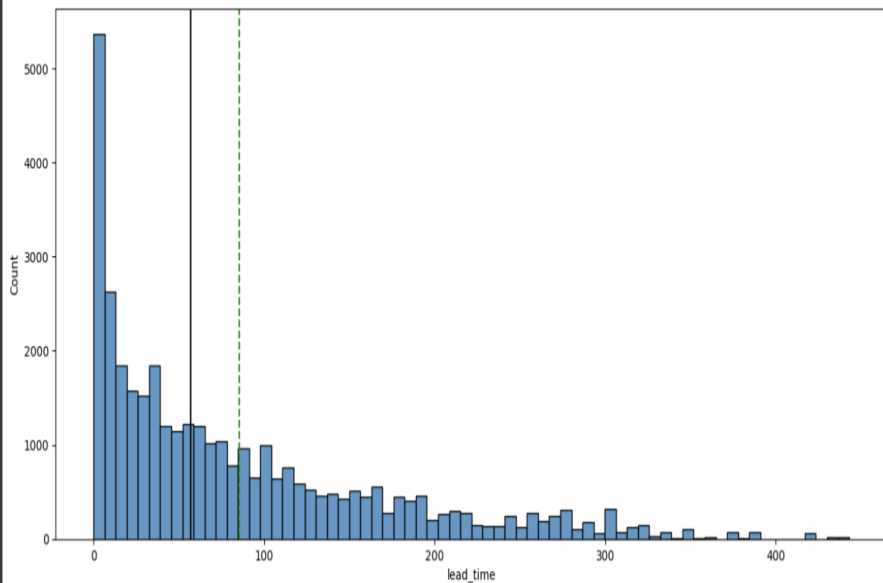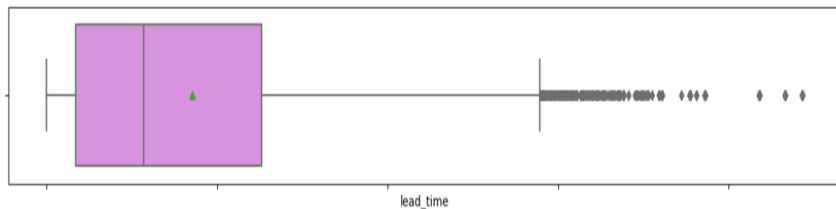# Business Problem Overview and Solution Approach

- A significant number of hotel bookings are called off due to cancellations or no-shows. The typical reasons for cancellations include change of plans, scheduling conflicts, etc. This is often made easier by the option to do so free of charge or preferably at a low cost which is beneficial to hotel guests but it is a less desirable and possibly revenue-diminishing factor for hotels to deal with. Such losses are particularly high on last-minute cancellations.

- The new technologies involving online booking channels have dramatically changed customers' booking possibilities and behavior. This adds a further dimension to the challenge of how hotels handle cancellations, which are no longer limited to traditional booking and guest characteristics.
- The cancellation of bookings impact a hotel on various fronts:
- Loss of resources (revenue) when the hotel cannot resell the room.
- Additional costs of distribution channels by increasing commissions or paying for publicity to help sell these rooms.
- Lowering prices last minute, so the hotel can resell a room, resulting in reducing the profit margin.
- Human resources to make arrangements for the guests.

- The increasing number of cancellations calls for a Machine Learning based solution that can help in predicting which booking is likely to be canceled. INN Hotels Group has a chain of hotels in Portugal, they are facing problems with the high number of booking cancellations and have reached out to your firm for data-driven solutions. You as a data scientist have to analyze the data provided to find which factors have a high influence on booking cancellations, build a predictive model that can predict which booking is going to be canceled in advance, and help in formulating profitable policies for cancellations and refunds.

# Data Overview

- There are 36,275 rows and 19 columns in the dataset.

- No duplicate values are present.

- Booking ID column has been dropped from the data frame.

# EDA- Lead Time



- Lead_time : Number of days between the date of booking and the arrival date.

- Mean lead time is ~85 days

- 25% are below 17 days, 50% are below 57 days, and 75% are below 126 days.

- The distribution is right-skewed.

- Outliers are present.

# EDA- Average Price Per Room

- The mean average price per room is ~ 103 euros

- 25% are < 80 euros

- 50% are < 99 euros

- 75% are < 120 euros

# EDA- Number of previous cancellations



- 75% of the customers had never had a previous cancellation.

- The distribution is right-skewed.

# EDA- Number of previous bookings not cancelled

- Mean ~ .15

- 75% of customers have never had
A previous booking cancelled

- The distribution is right-skewed.



no_of_previous_bookings_not_canceled

# EDA- Number of Adults



● 72% of guests have 2 adults staying

# EDA- Number of Children

- Most of the guests staying had no children.

- 33,577 of the bookings had no children, and very few had children at all.

- Travelling with children is more expensive in general, which may contribute to the number of bookings including children.

# EDA- Number of Week nights

- Most guests only book for 2 week-nights in the hotel.

- Distribution is right-skewed

# EDA- Number of Weekend Nights

- The distribution is right-skewed.

- Almost half of the guests didn't spend the weekend nights at the hotel.

# EDA- Required Car Parking Space



- Most guests require no parking space

# EDA- Type of Meal Plan

- Meal Plan 1 is the most requested meal.

# EDA- Room Type Reserved



- Over 25,000 of the rooms reserved were Type 1

- Most of the rooms were Room Type 1, with nothing else even close

# EDA- Arrival Month

- Most rooms were booked for the month of October, and the least was January

# EDA- Market Segment Type

- The online market segment had the most bookings

- Complimentary bookings were the least.

# EDA- Number of Special Requests

- Most bookings didn't have any special requests

# EDA- Booking Status

- 24, 390 reservations were not cancelled, and 11,885 were.

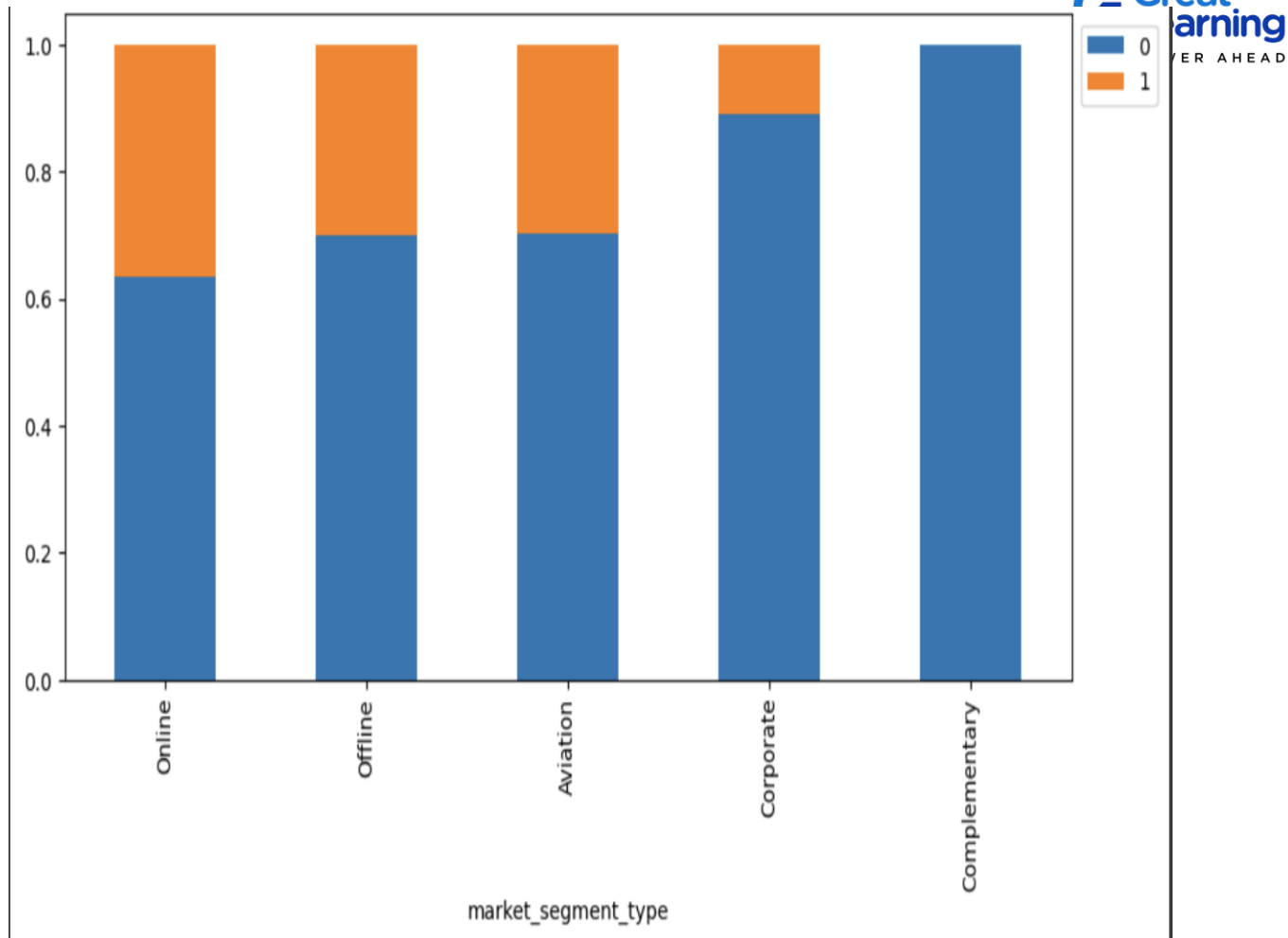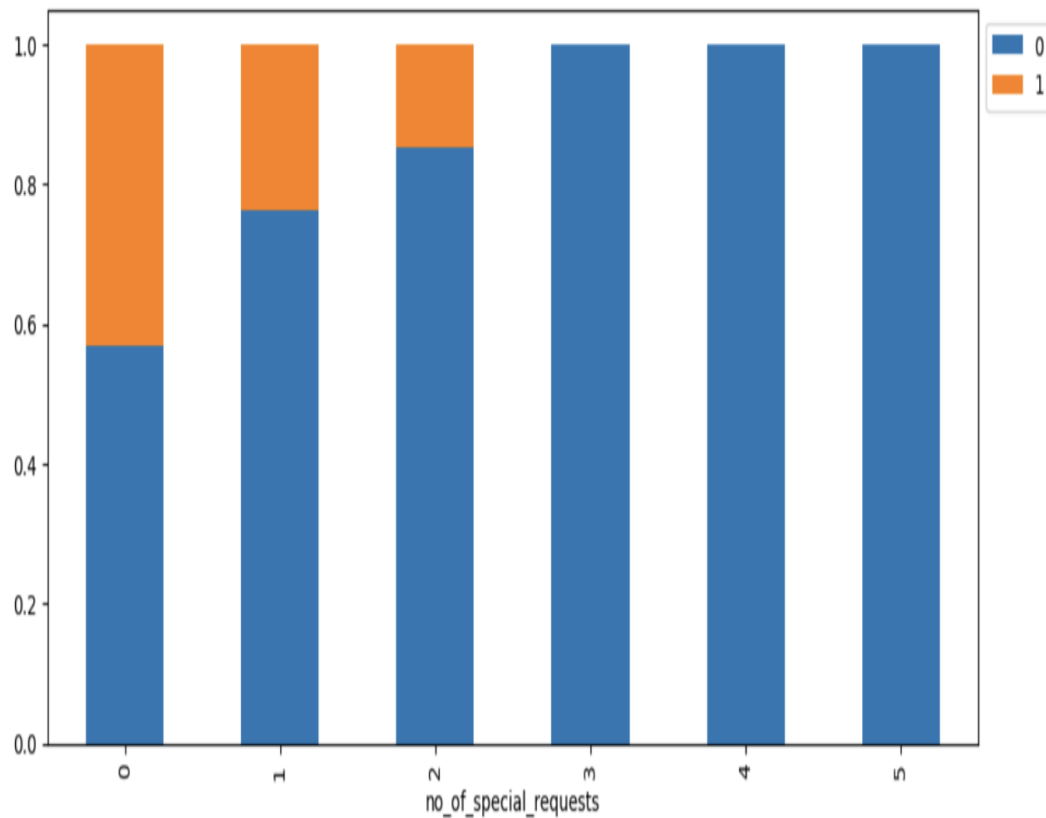- There wasn't a strong correlation between any of the variables

# EDA- Market Segment Type v. Average Price/Room



- Mean Average price per room is highest for online reservations

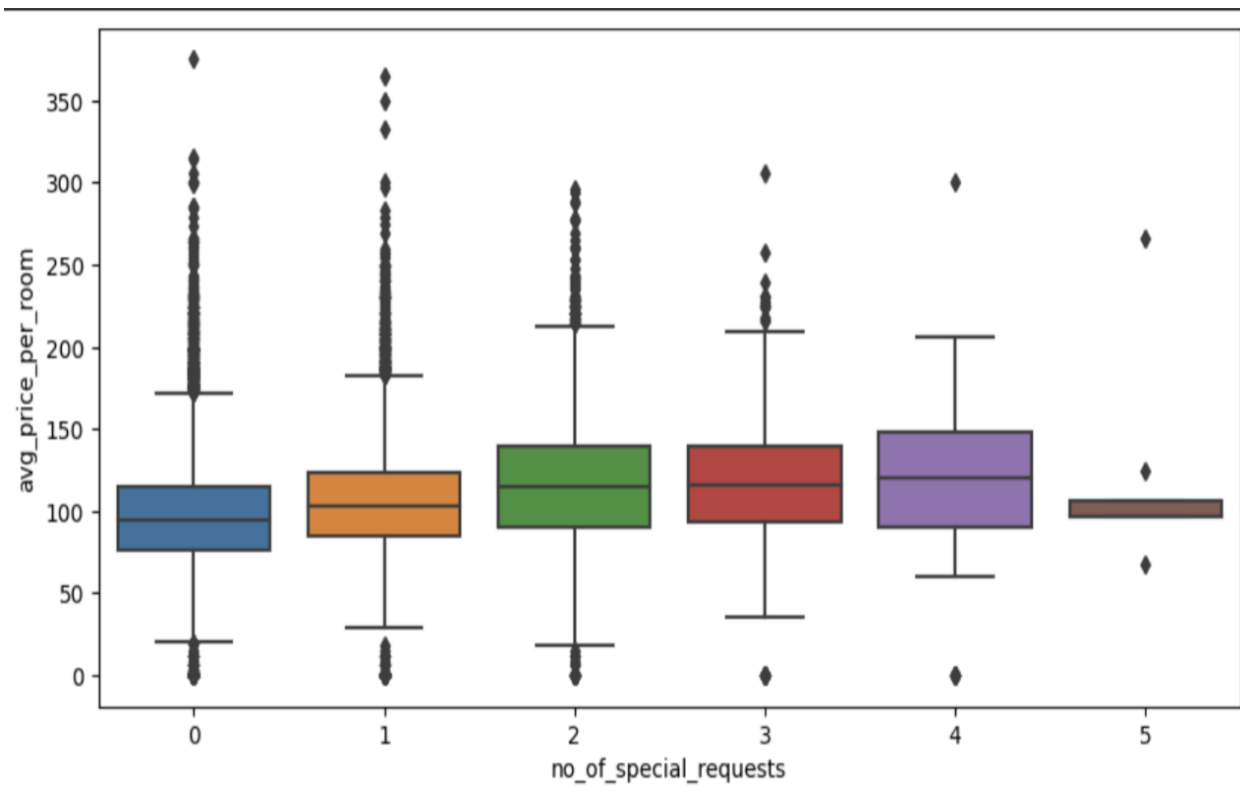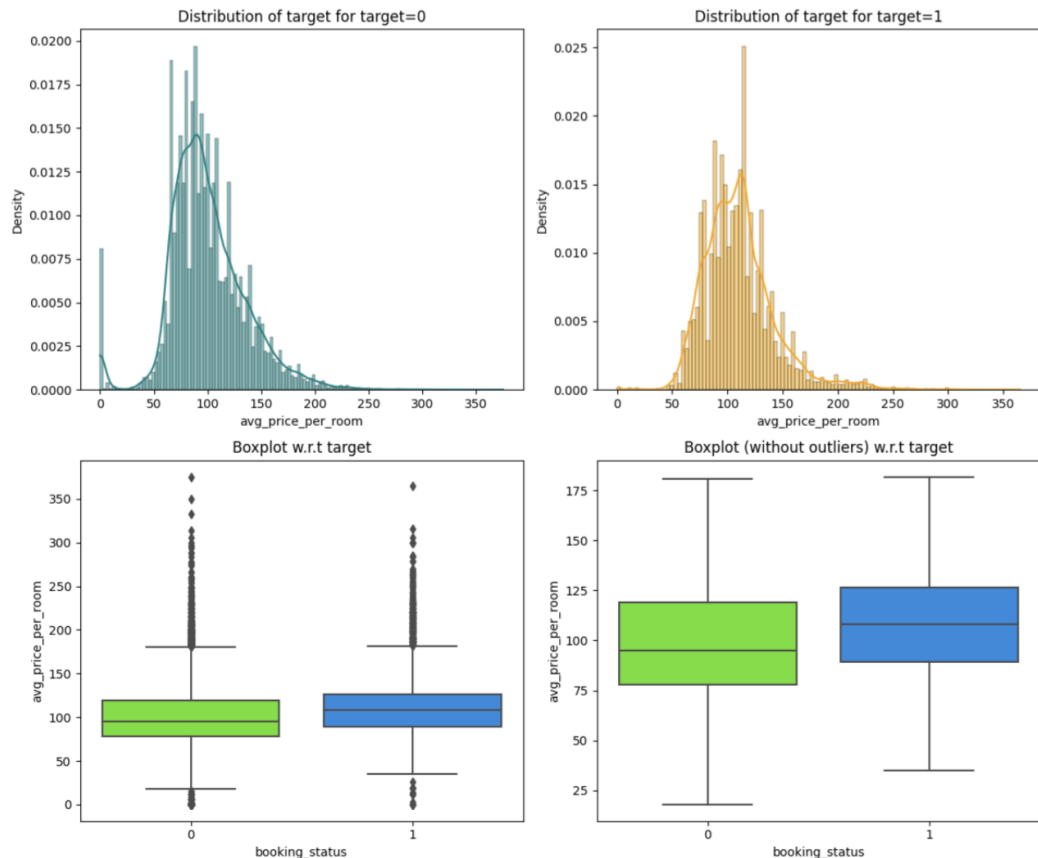- There are many outliers with the online reservation segment.

# EDA-

.

- The more specific a guest is with their special requests, the less likely they are to cancel.

# EDA- Number of Special Requests v. Average Price/Room

- There is nothing distinct about this comparison.

- Most reservations h ad no special requests.

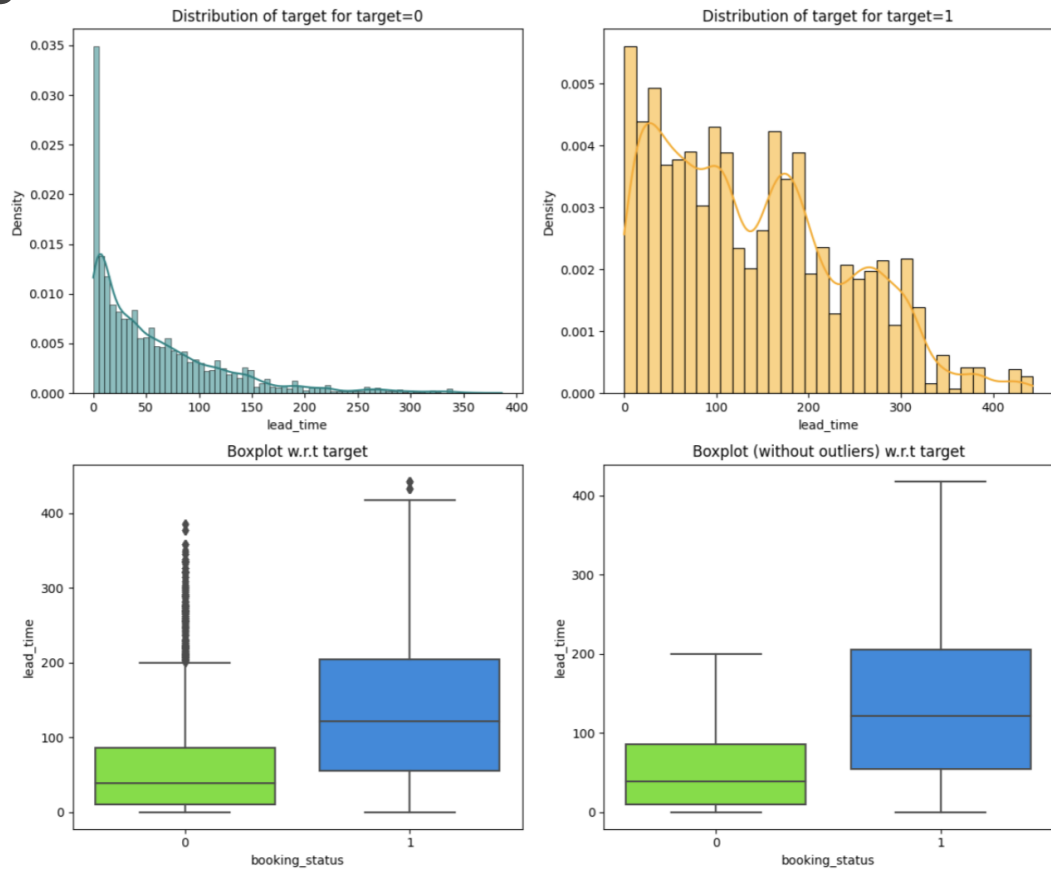- There is no correlation between average room price and number of special requests.

# EDA- Avg. Price/Room v. Booking Status



- Booking status and average price per room is slightly right-skewed. The rooms that weren't cancelled are cheaper than the average price for cancelled rooms.
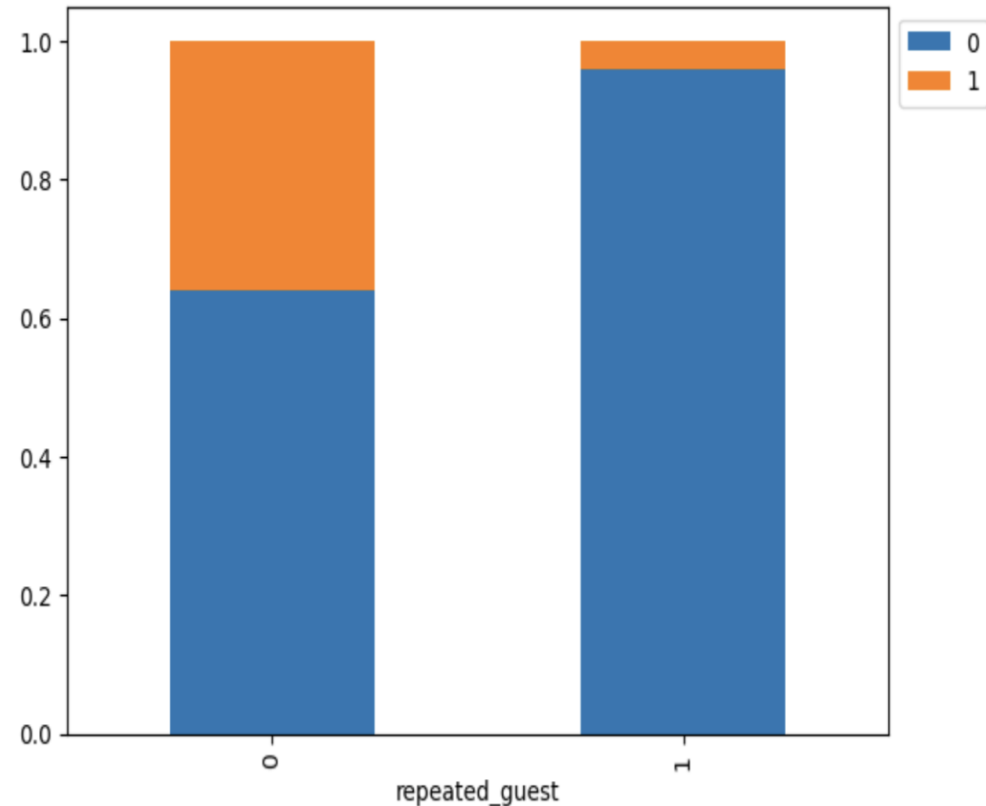
# EDA- Lead Time v. Booking Status

- The lead time distribution is right-skewed.

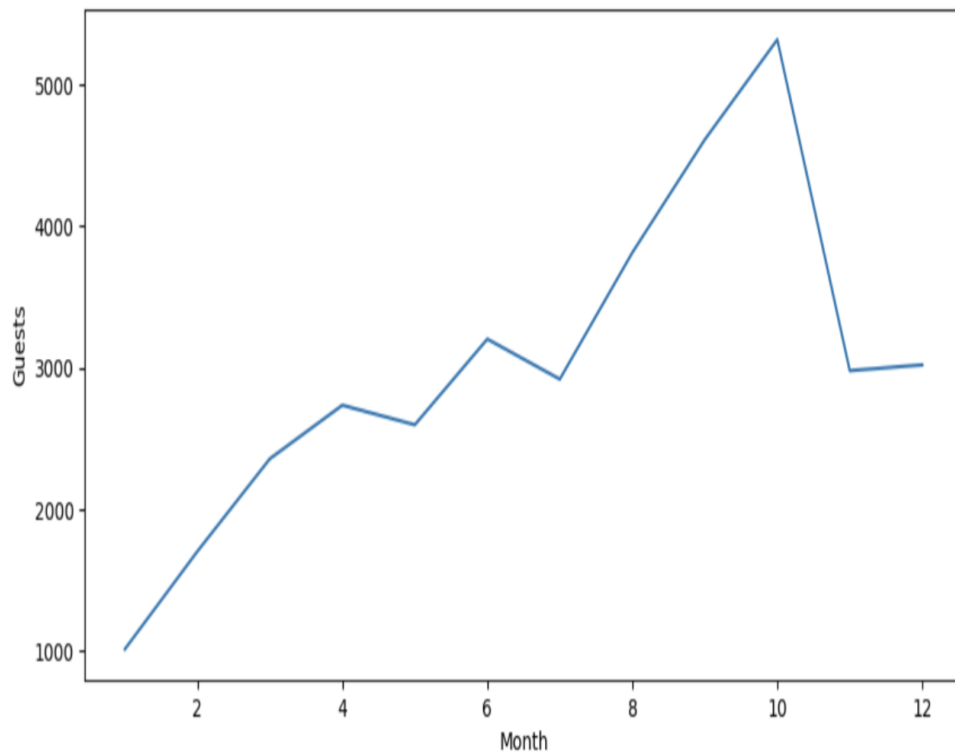- The greater the lead time, the greater likelihood of cancellation

.

# EDA- Repeated Guest

- Repeat guests have a very low number of cancellations.

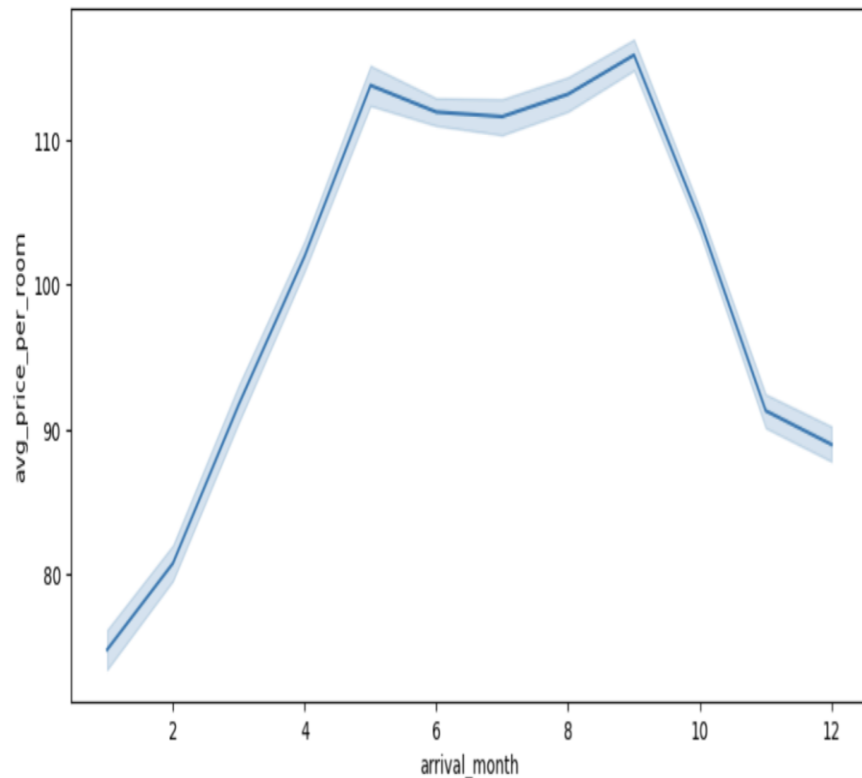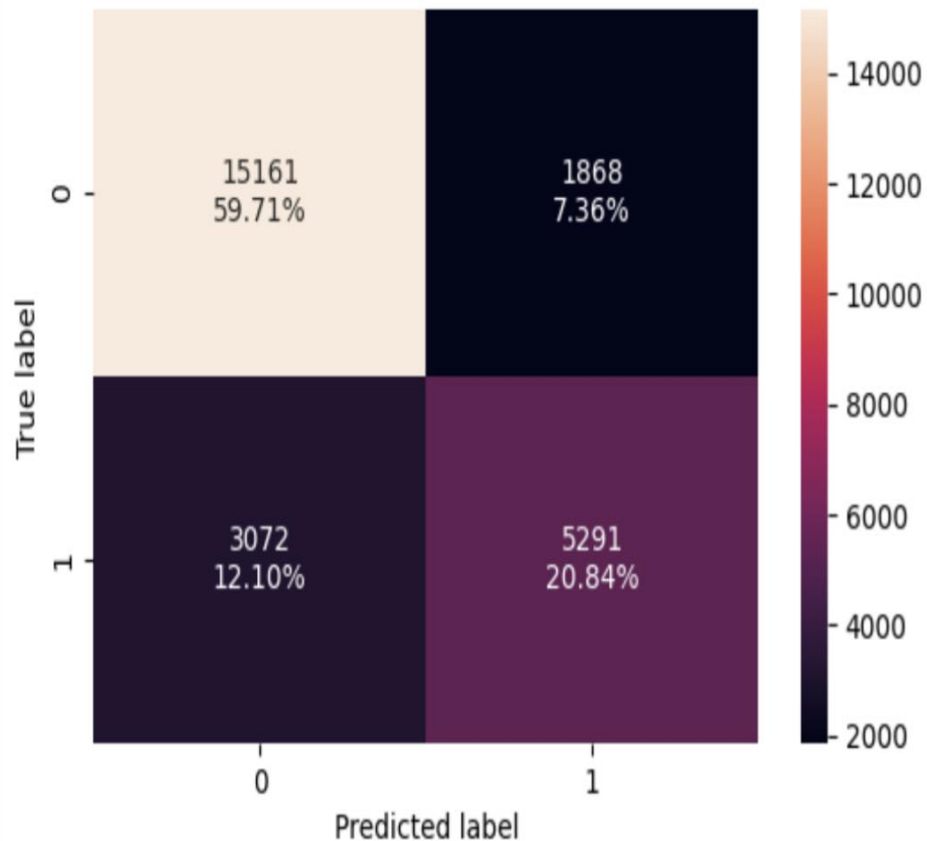- Most cancellations come from new customers.

# EDA- Month v. Guests



- Most bookings occur between August

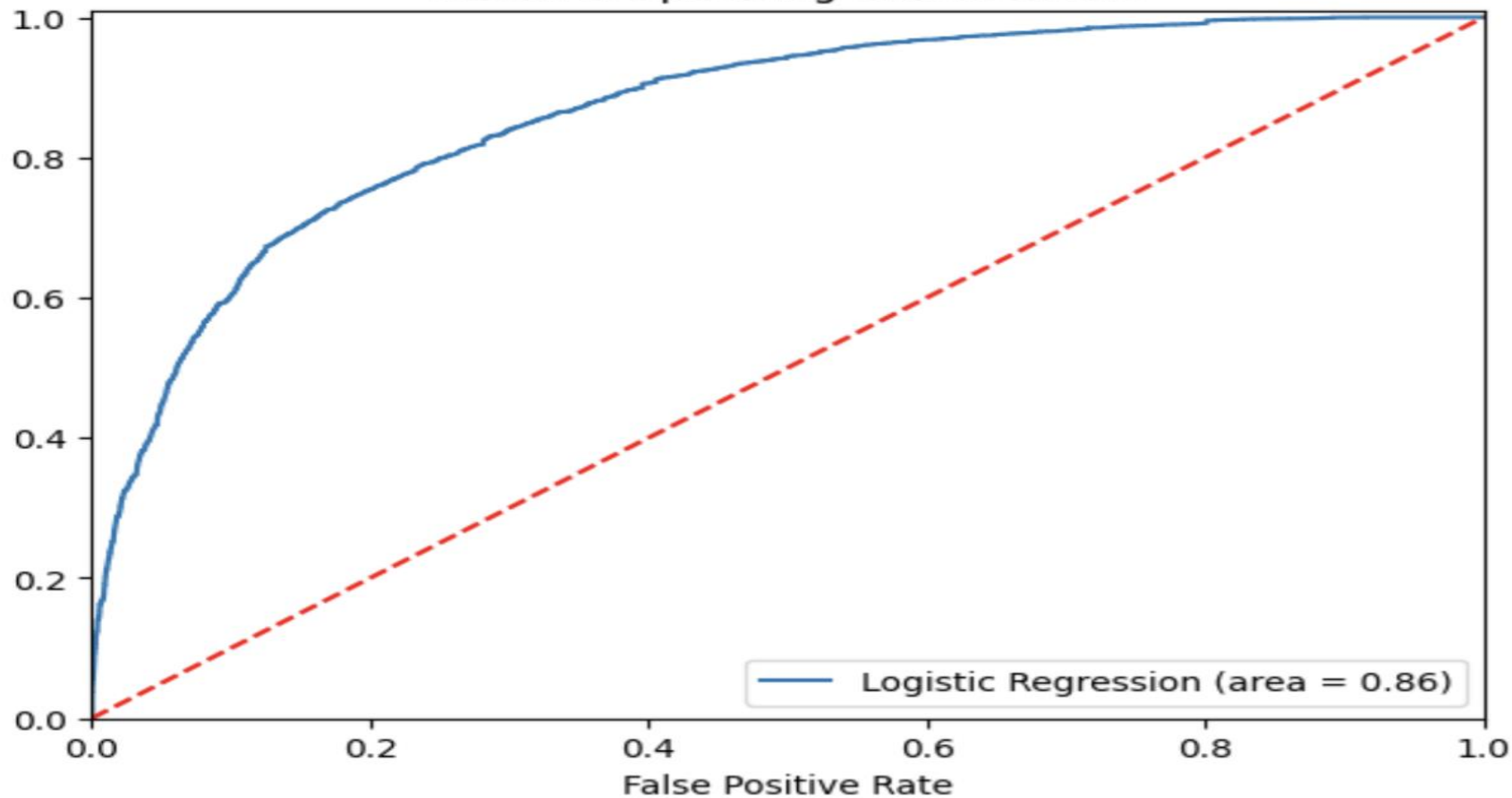- And October.
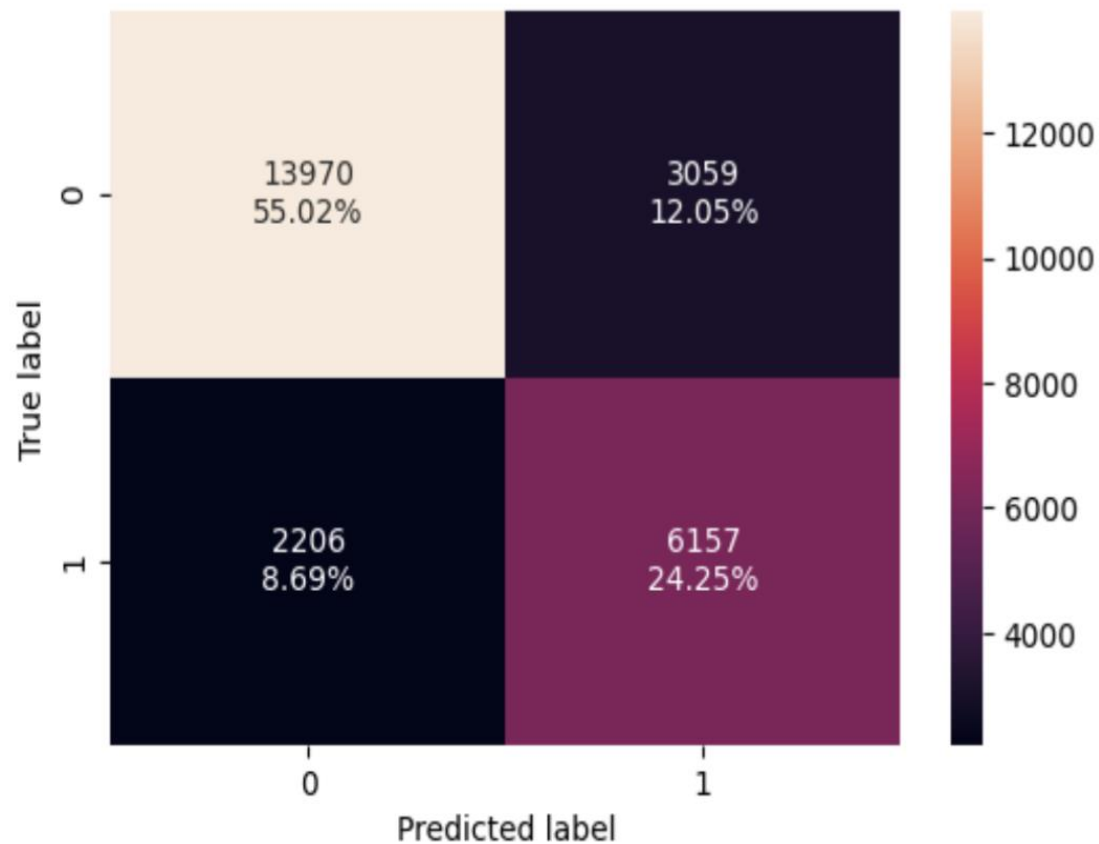
# EDA- Arrival Month v. Avg. Price/Room



- Most families travel in the summer for summer vacations. Many hotels increase prices during this time to maximize profits.

- The training data predicts a true positive ~ 60%
- Training data predicts True negative ~21%
- Training data predicts false positive at ~7%
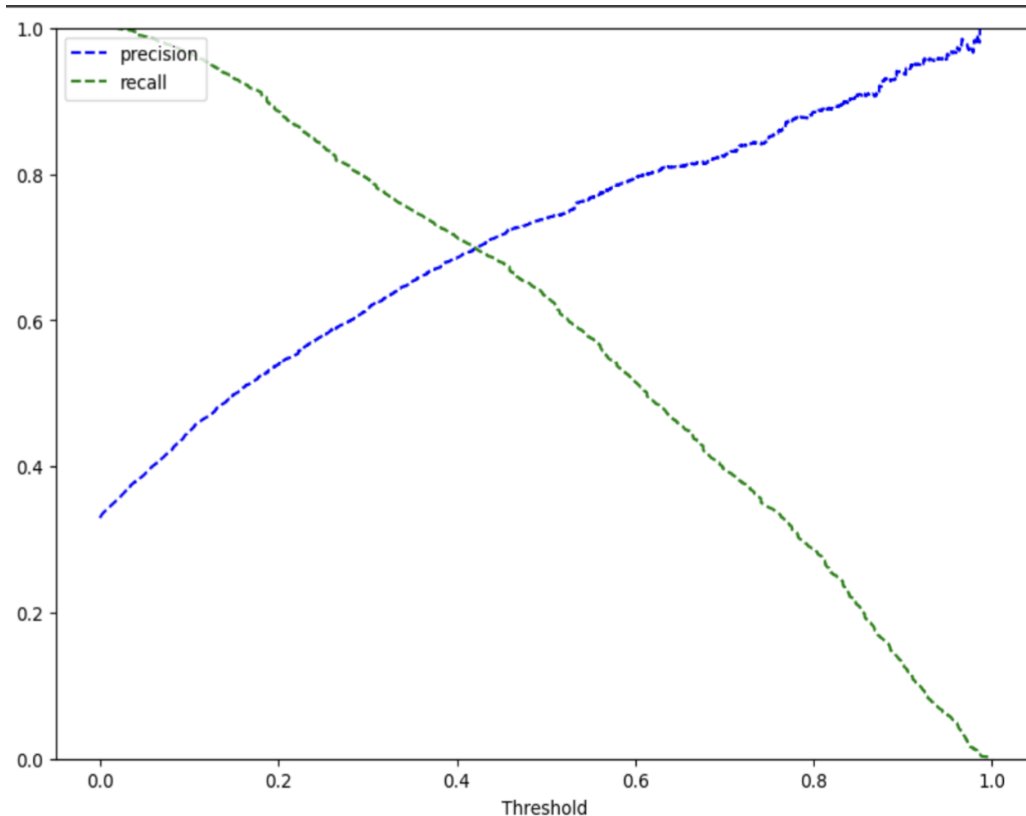- Training data predicts false negative at 12%
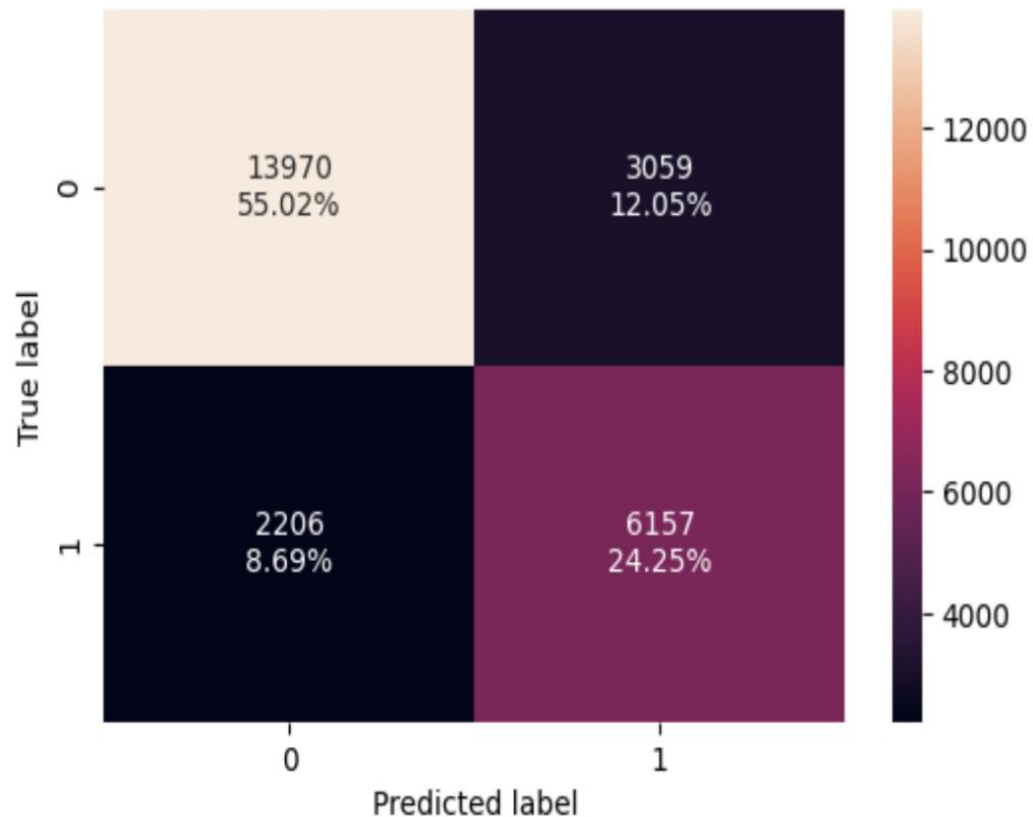
Receiver operating characteristic

Logistic Regression (area = 0.86)

False Positive Rate

- Confusion matrix on the test model shows:
- True Positive: ~55%
- False Positive: ~12%
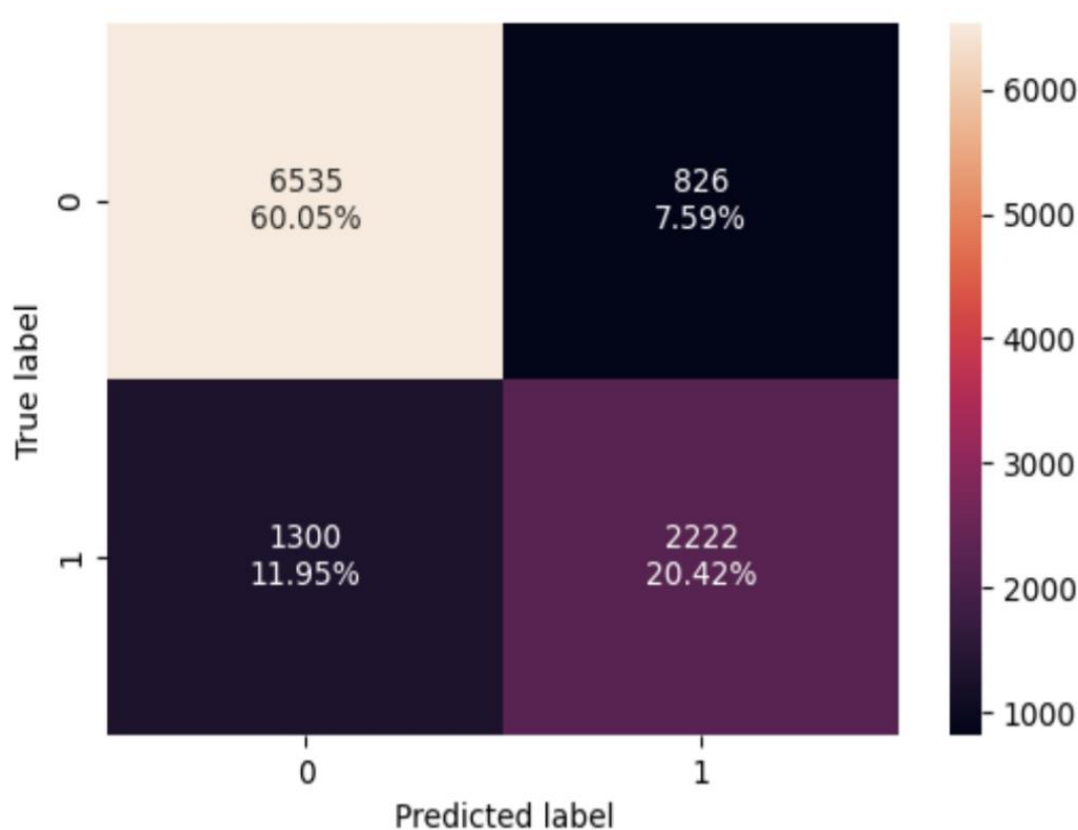- False Negative: ~ 24%
- True Negative: ~9%

- Treshold of 0.42 gives equal precision and recall.

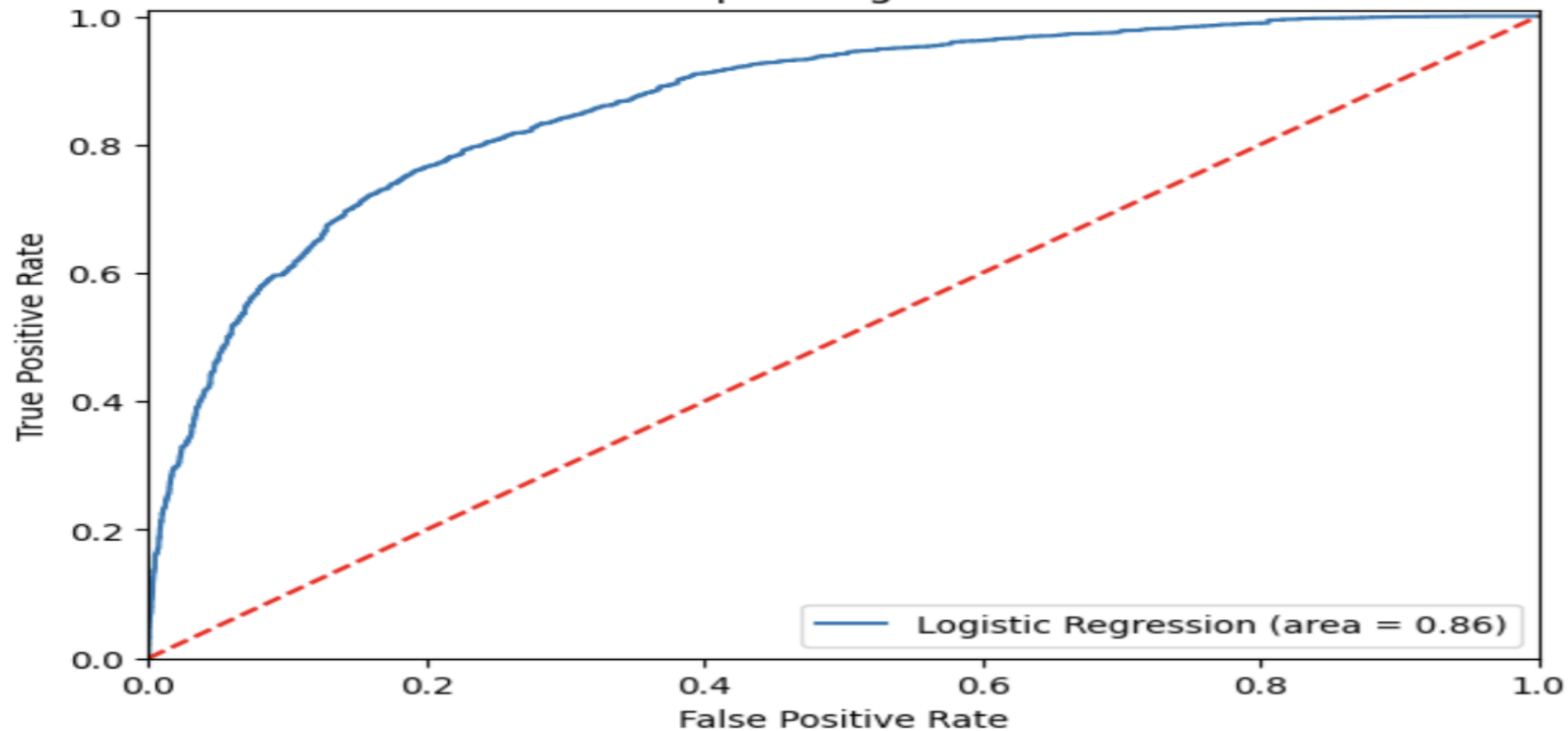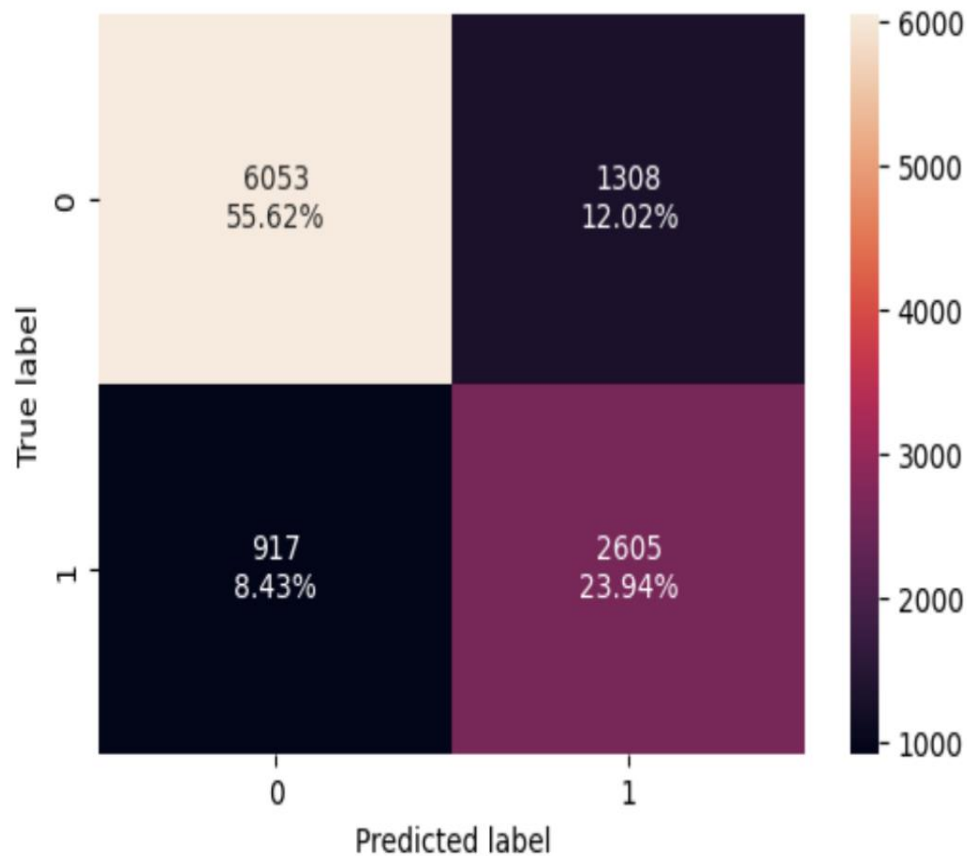- The treshold will be set at 0.42

- Confusion Matrix shows:

- True Positive: 55%

- True Negative: ~ 24%

- False Positive: 12%

- False Negative: ~9%

- Confusion Matrix shows:

- True Positive: 60%

- True Negative: 20.42%

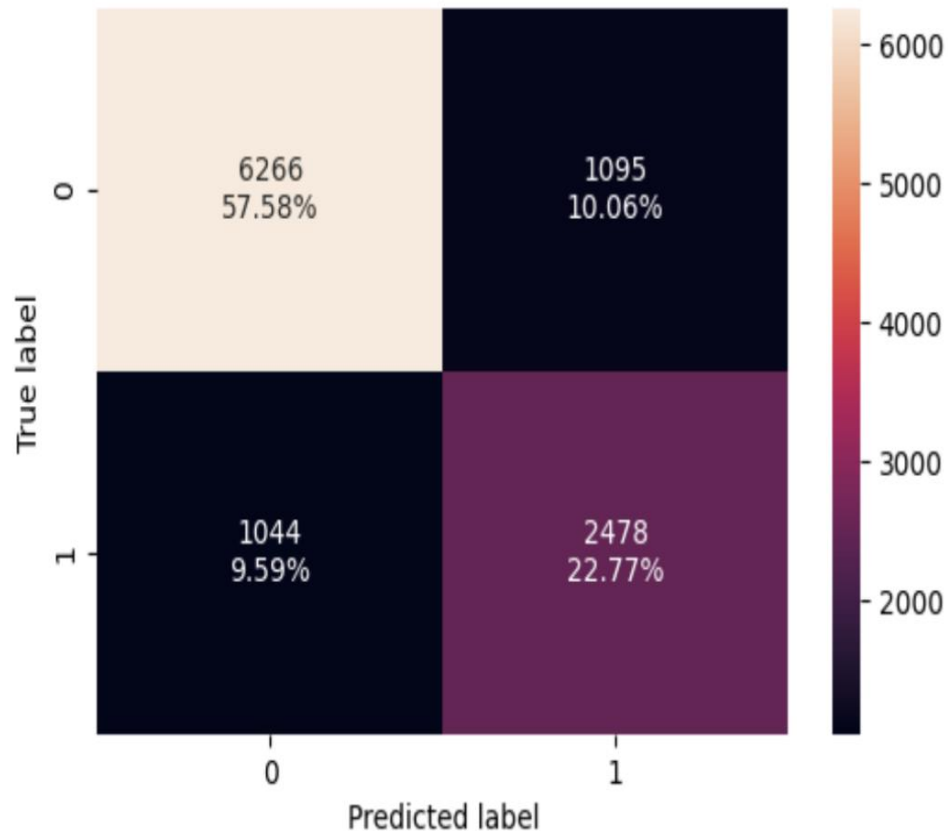- False Positive: ~8%

- False Negative: ~12%
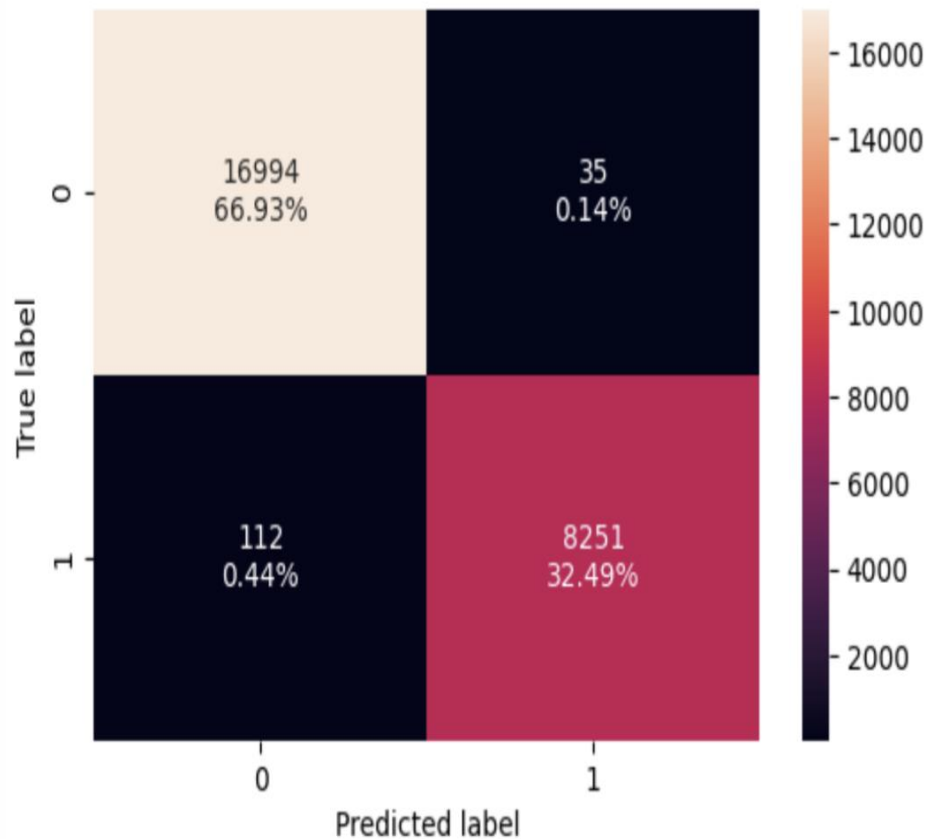
Receiver operating characteristic

- Confusion Matrix shows:

- True Positive: ~56%

- True negative: ~24%

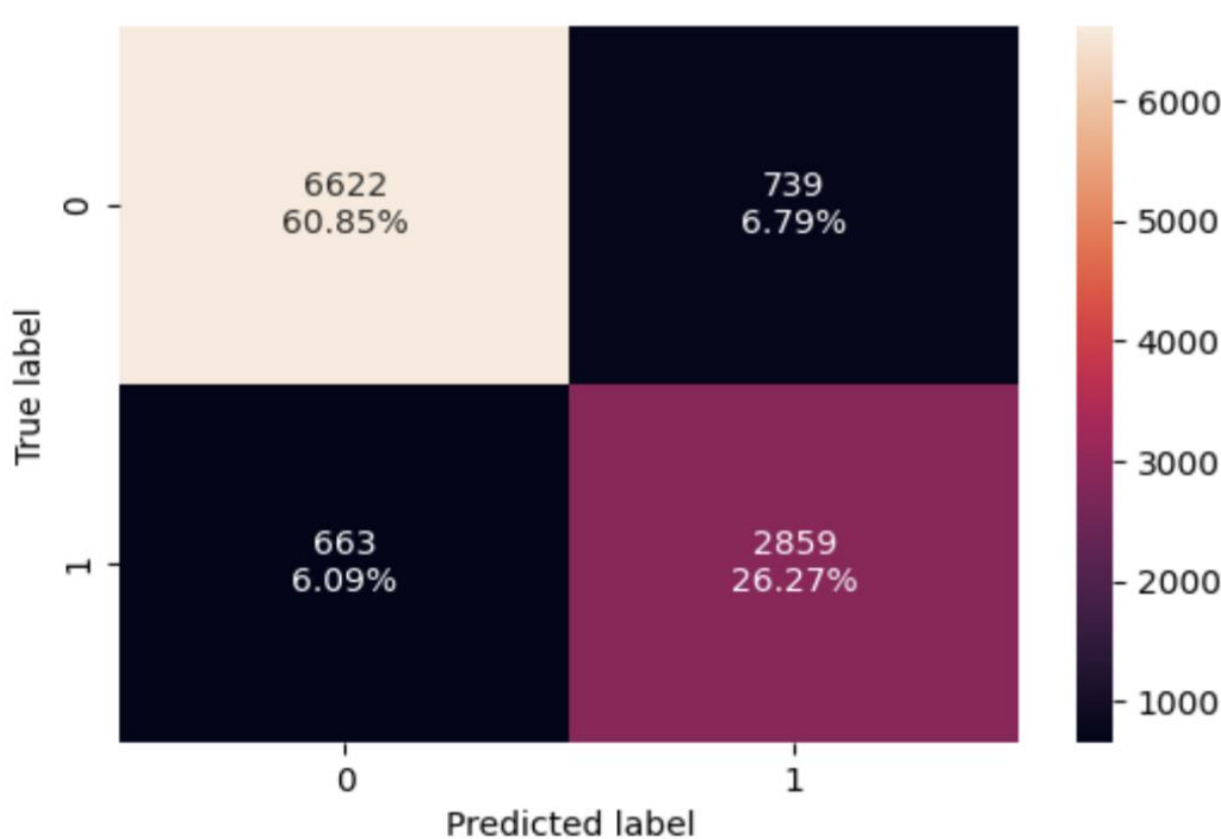- False Positive: 12%

- False Negative: ~8%

- Confusion Matrix shows:

- True Positive: ~57%

- True negative: ~23%

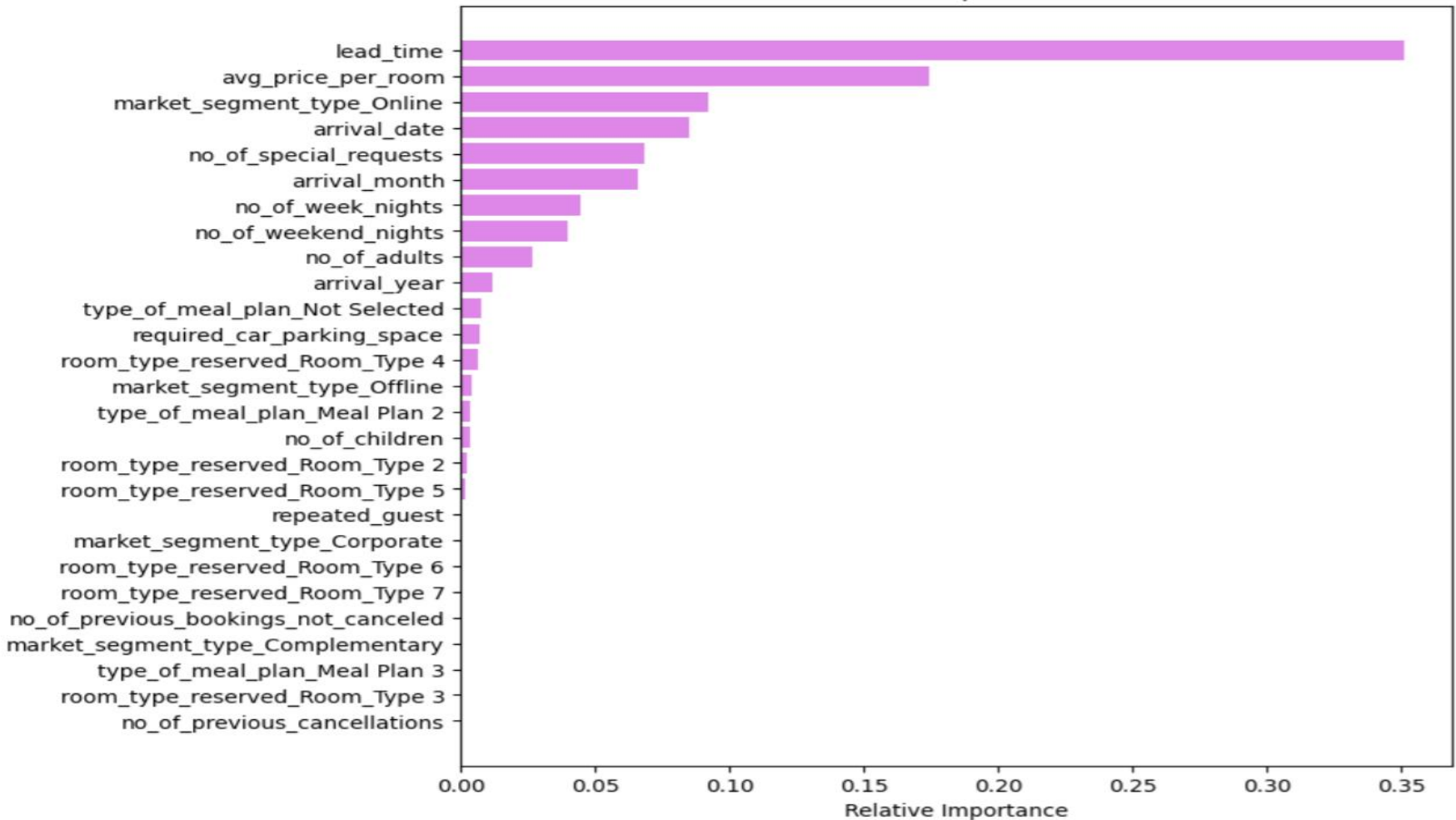- False Positive: 10%

- False Negative: ~10%

- Confusion Matrix shows:

- True Positive: ~67%

- True negative: ~33%

- False Positive: 0.14%
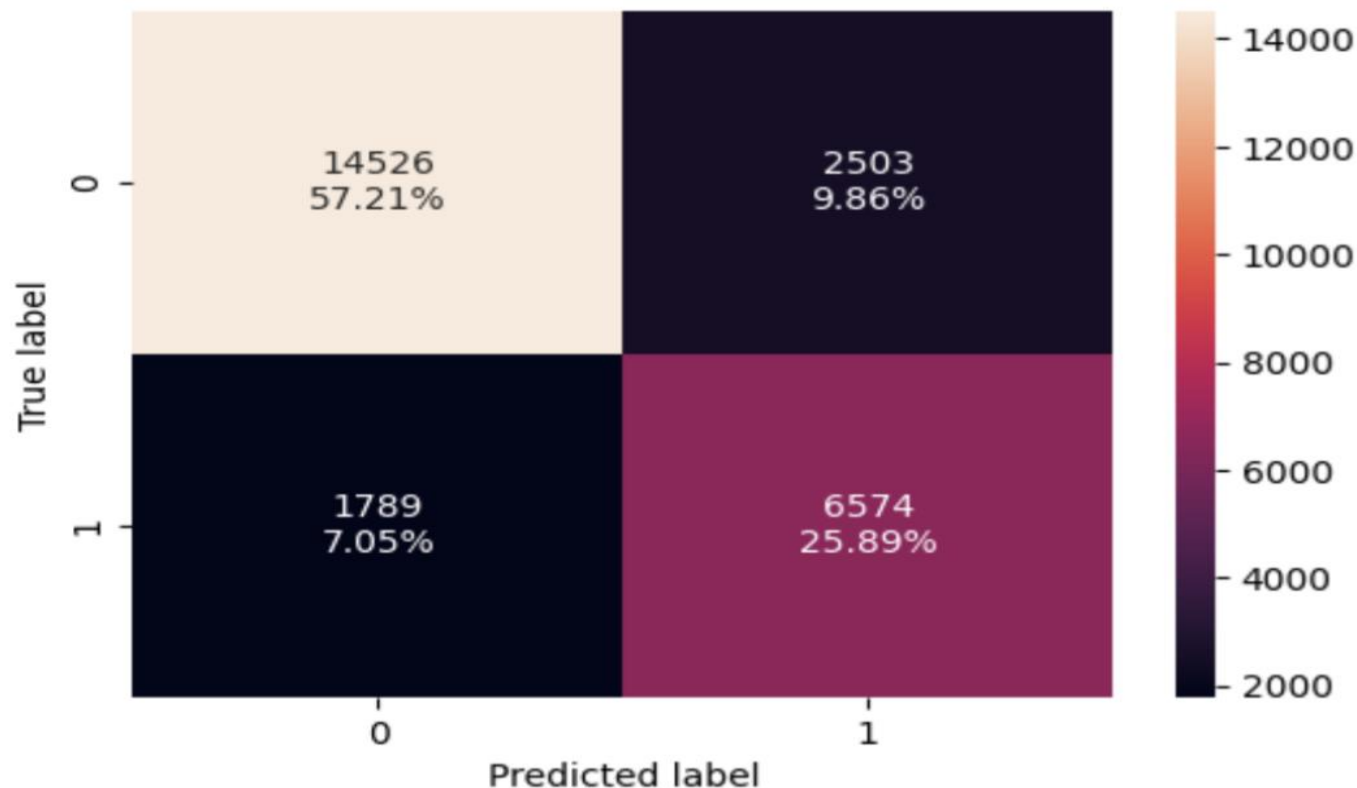
- False Negative: ~0.44%

- Confusion Matrix shows:

- True Positive: ~61%

- True negative: ~26%

- False Positive: ~7%

- False Negative: ~6%

Confusion matrix values:

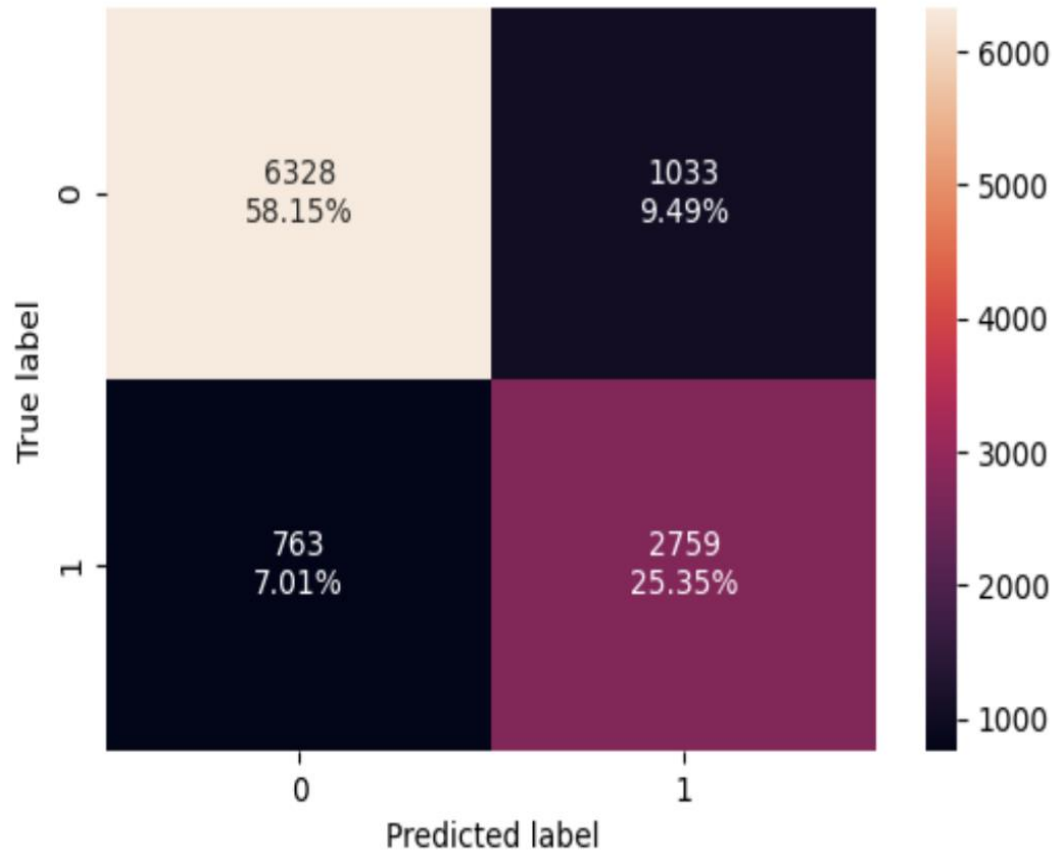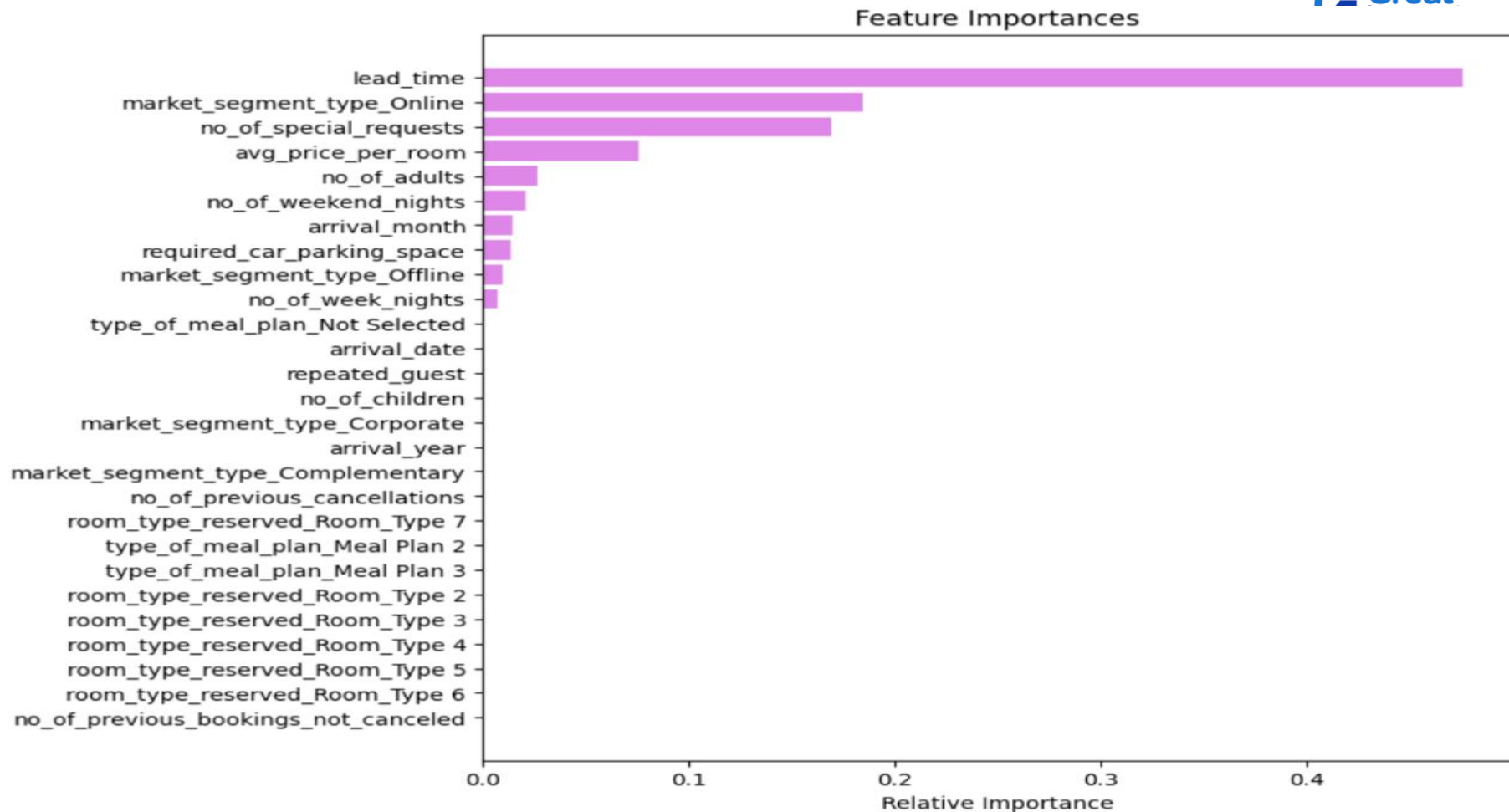| | Predicted 0 | Predicted 1 |
|---|---|---|
| True 0 | 6622 (60.85%) | 739 (6.79%) |
| True 1 | 663 (6.09%) | 2859 (26.27%) |

## Feature Importances

- Confusion Matrix shows:

- True Positive: ~57%

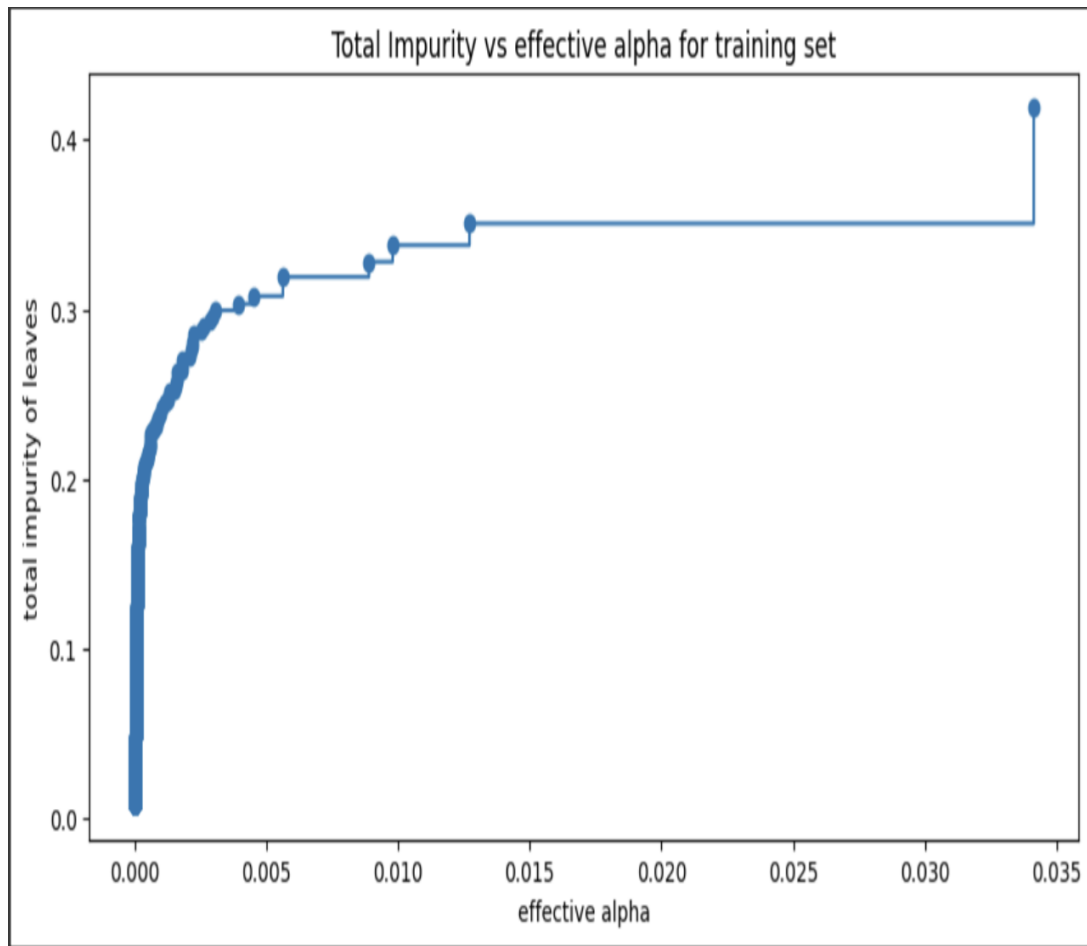- True negative: ~26%

- False Positive: ~10%

- False Negative: ~7%

- Confusion Matrix shows:

- True Positive: ~58%

- True negative: ~25%

- False Positive: ~10%

- False Negative: ~7%

## Feature Importances

Total Impurity vs effective alpha for training set

**Happy Learning !**