# Star Hotels
# Business Presentation

# Contents

The presentation consists of three Four Sections:

1. Background and Business Problem Overview

2. Data Overview

3. Exploratory Data Analysis

4. Model overview and performance summary

5. Key findings and insights

6. Business recommendations

# Background of the Business and Objective

## Background

A significant number of hotel bookings with STAR Hotels are called off due to cancellations or no-shows.
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.

## Objective

Develop a Machine Learning based solution that can help in predicting which booking is likely to be canceled.

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.
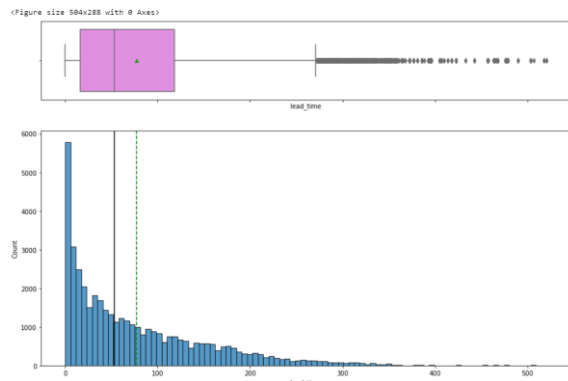
# Data Overview

| Data | Description |
|------|-------------|
| no_of_adults | Number of adults |
| no_of_children | Number of Children |
| no_of_weekend_nights | Number of weekend nights (Saturday or Sunday) the guest stayed or booked to stay at the hotel |
| no_of_week_nights | Number of week nights (Monday to Friday) the guest stayed or booked to stay at the hotel |
| type_of_meal_plan | Type of meal plan booked by the customer |
| required_car_parking_space | Does the customer require a car parking space? (0 - No, 1- Yes) |
| room_type_reserved | Type of room reserved by the customer. The values are ciphered (encoded) by Star Hotels Group |
| lead_time | Number of days between the date of booking and the arrival date |
| arrival_year | Year of arrival date |
| arrival_month | Month of arrival date |
| arrival_date | Date of the month |
| market_segment_type | Market segment designation. |
| repeated_guest | Is the customer a repeated guest? (0 - No, 1- Yes) |
| no_of_previous_cancellations | Number of previous bookings that were canceled by the customer prior to the current booking |
| no_of_previous_bookings_not_canceled | Number of previous bookings not canceled by the customer prior to the current booking |

# Data Overview

- **There are 56926 rows and 18 columns.**

- **There are 14350 duplicate values in the dataset, We removed the duplicate data before processing.**

- **Most of the data-types are either int64 or float64.**

- **type_of_meal_plan, room_type_reserved, market_segment_type and booking_status having data-types as an object, this means we need to convert these into suitable data-type before we feed our data into the model.**

- **There are no missing values in the dataset.**

- **Unique values shows that**
  **- There are 4 types of meal plans**
  **- There are 7 types of room types**
  **-There are 5 types of market segment type**
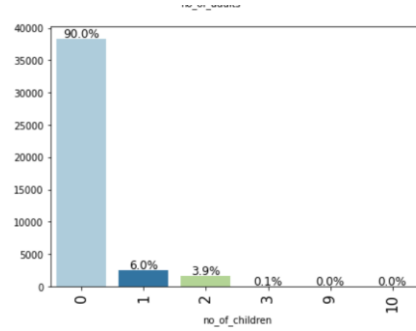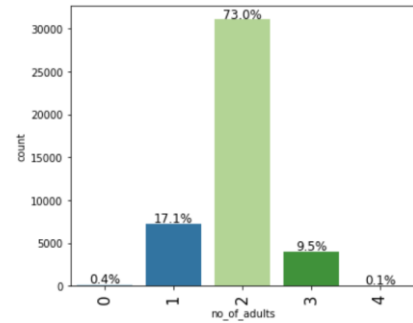
# Exploratory Data Analysis

## Lead Time .

## Average nightly price





lead_time is right skewed.
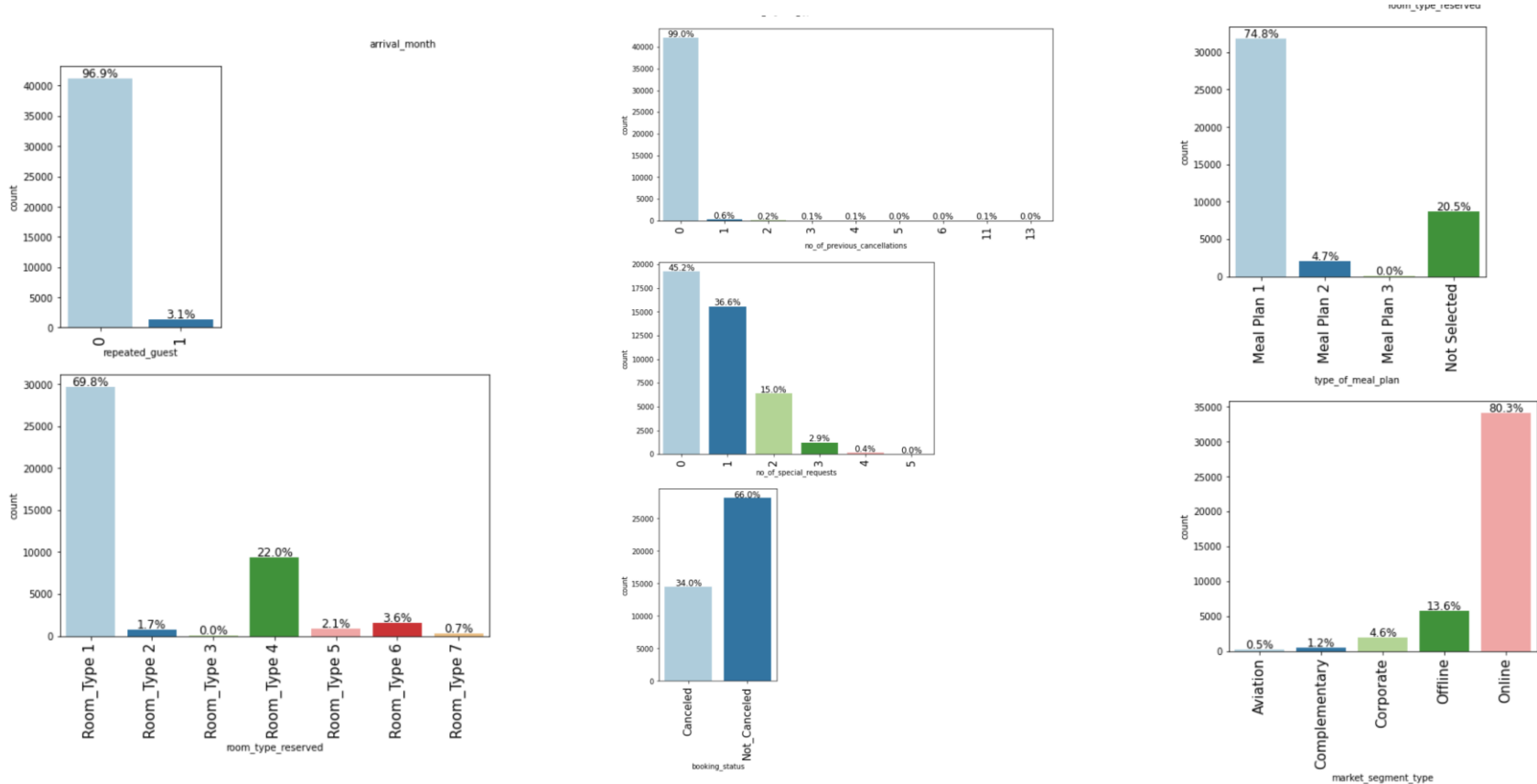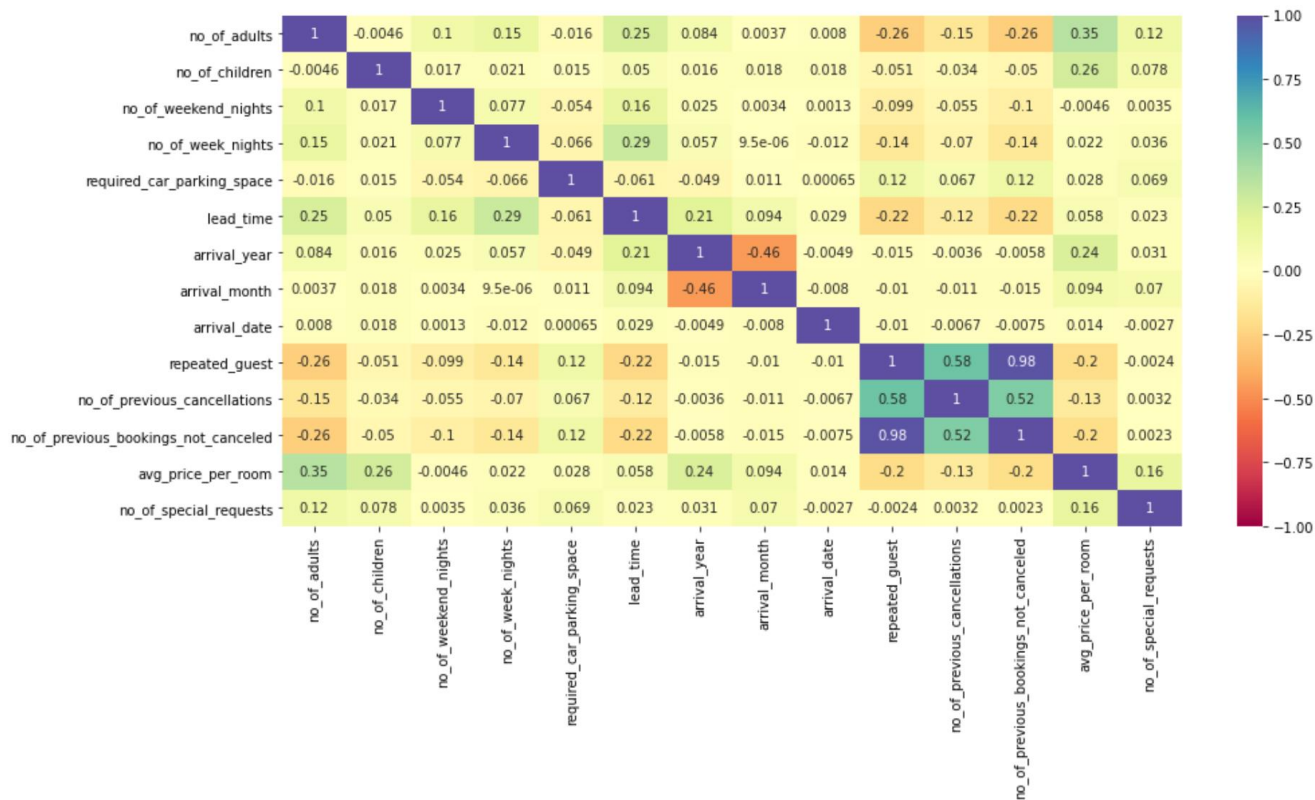There are outliers exists in avg_price_per_room.
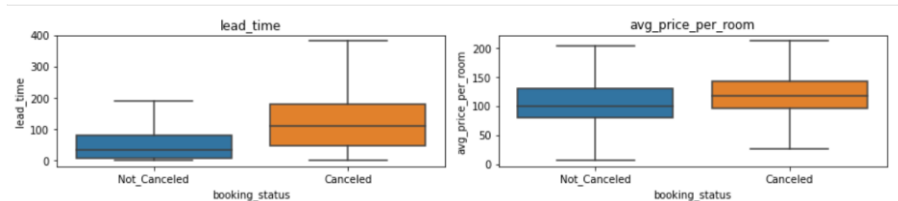
# Exploratory Data Analysis

# Exploratory Data Analysis - Univariate analysis
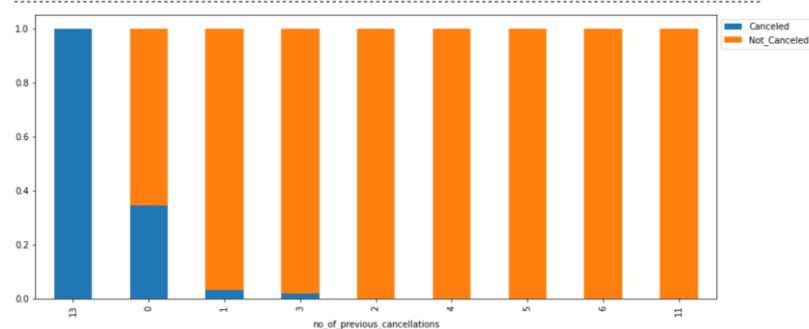
# Exploratory Data Analysis - Bivariate Analysis

# Exploratory Data Analysis - Bivariate Analysis graphs

# Exploratory Data Analysis - Bivariate Analysis graphs

# Exploratory Data Analysis - Observations

Meal Plan 1 is mostly selected, Also 10072 bookings not selected for any meal plans.

Room_Type 1 is highest slected among bookings 42807.

Room_Type 3 is least selected among bookings.

Highest market_segment type is through online booking.

35378 bookings are not cancelled vs 21548 Cancelled.

There are outliers exists in the lead time

lead_time is right skewed.

avg_price_per_room is also right skewed.

average nightly price is 112.37

73% of the booking is done where occupants are 2 adults in the room

90% room booked with by the occupants with no kids.

41% of the time week ends nights are not booked.

weekends nights occupency for 1 night and 2 nights both are 28.8% and 28.2%

highest occupancy are for 2 weeknights which is 27.6% followed by 1 night which is 25.6%

2018 and 2019 were among the busiest year with 51.9% and 38.9% arrivals.

July and August are busiest months with 11.1% and 12.5% arrivals.

96.9% guest were not a repeated guest only 3.1% guest were repeated guest.

Out of 7 catagories room types, Room catagory type 4 has highest booking.

# Exploratory Data Analysis - Observations

74.8% customers selected meal plan 1 and 20.5% customers did not selected any meal plan.

80.3 % bookins is done through online.

36% guest has some special request

66% customers not cancelled the booking v/s 34% customers cancelled.

no_of_repeated_guests is highly co-related with no_of_previous_bookings_not cancelled

max number of booking cancelled by customers which are booked for 2 adults stay in the room.

max number of booking cancelled by customers is in August month and also max booking took place in August.

max number of booking cancelled by customers who were not repeated guests, Repeated guests hardly cancelled any bookings.

max number of booking cancelled by customers who booked Room_Type 1 which is also had highest room type booked.

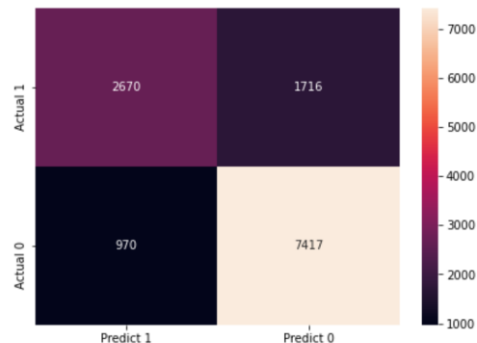max number of booking cancelled by customers who opted for meal plan 1

max number of booking cancelled by customers who booked online

max number of booking cancelled by customers which got the rooms booked for an average_price_per_room of $115

# Logistic regression overview and performance summary

- We used regular Logistic regression using sklearn and statsmodels library to build our machine learning model.
- We replaced the categorical data into the continuous format to fit into the models using the onehot columns method.
- We also split the data into the training and testing set in a ration of 70-30%.
- We build confusion matrix for both the models on and also checked the performance on train and test data.

# Logistic regression model sklearn and performance summary



**True Positives (TP): we correctly predicted that customer cancelled for 2670 bookings.**

**True Negatives (TN): we correctly predicted that customer will not cancel 7417**

**False Positives (FP): we incorrectly predicted that customer would cancel (a "Type I error") 970 Falsely predicted Type I error.**

**False Negatives (FN): we incorrectly predicted that customer will not cancel (a "Type II error") 1718 Falsely predicted negative Type II error.**

**The training and testing recall are 61.3% and 60.98% respectively.**

**The training and testing f1_scores are 0.66**

**Recall on the train and test sets are comparable. f1_score on the train and test sets are comparable.**

# Logistic regression statsmodel performance summary



**None of the variables exhibit high multicollinearity, so the values in the summary are reliable.**

**We removed the features which had p-value > 0.05, both the results are attached.**

**None of the variables exhibit high multicollinearity, so the values in the summary are reliable.**

**Coefficient of no_of_children,**

**no_of_week_nights,lead_time, no_of_previous_cancellations, avg_price_per_room, market_segment_type_1, market_segment_type_2, market_segment_type_4 are positive an increase in these will lead to increase in chances of a person cancelling the bookings.**

**Coefficient of required_car_parking_space, arrival_year,arrival_month, no_of_special_requests, type_of_meal_plan, room_type_reserved and repeated_guest are negative and an increase in these will lead to decrease in chances of a person cancelling booking**

Training performance comparison:

| | Logistic Regression sklearn | Logistic Regression-0.30 Threshold | Logistic Regression-0.41 Threshold |
|---|---|---|---|
| **Accuracy** | 0.793309 | 0.773110 | 0.789551 |
| **Recall** | 0.613009 | 0.806059 | 0.705376 |
| **Precision** | 0.733389 | 0.628969 | 0.683716 |
| **F1** | 0.667817 | 0.706587 | 0.694377 |

# Model overview and performance summary

## Conclusion

We kept improving on our model and reached to the model which provided a generalized performance on training and test set.
The highest recall was 80% on the training set.

Using the model with default threshold the model gave low recall but good precision scores - This model will achieve the hotel save resources but lose on potential customers.

Using the model with 0.30 threshold the model gave a high recall but low precision scores - This model will help the hotel identify potential cancellations effectively, but the cost of resources will be high.

Using the model with 0.41 threshold the model will give a balance recall and precision score - This model will help the hotel to maintain a balance in identifying potential customer and the cost of resources.

# Decision Tree Model overview and performance summary

We also evaluated the data with the decision tree model and figured out that with the default depth the model was overfitted.

We then pruned the decision tree with maximum depth 3 and 5 respectively and strived to get the best fitted model with ideal recall score on train and test data.

Recall on training set was reduced from 0.99 to 0.75 so we reached to the conclusion that this model is not overfitting and we have a generalized model.

Finally, we got the most important feature with the feature of importance which can main effective features in hotel booking cancellation.

# Decision Tree Model overview and performance summary

We also evaluated the data with the decision tree model and figured out that with the default depth the model was overfitted.

We then pruned the decision tree with maximum depth 3 and 5 respectively and strived to get the best fitted model with ideal recall score on train and test data.

Recall on training set was reduced from 0.99 to 0.75 so we reached to the conclusion that this model is not overfitting and we have a generalized model.

Finally, we got the most important feature with the feature of importance which can main effective features in hotel booking cancellation.

Precision Recall curve at  threshold around 0.41 shows that we got a balanced recall and precision.

# Key Findings and Insights

We analyzed the "StarHotels Group" booking data using different techniques and used Decision Tree Classifier to build a predictive model for the same.

The model built can be used to predict if a customer is going to cancel the booking or not.

We visualized different trees and their confusion matrix to get a better understanding of the model. Easy interpretation is one of the key benefits of Decision Trees.

We observed the fact that much lesser data preparation is needed for Decision Trees and such a simple model gave good results even with no outliers treatment which shows the robustness of Decision Trees.

Lead Time, Market Segment Type, Avg Price per room, and No_Of_Special Requests are the most important variable in predicting the customers cancellation behavior

We established the importance of hyper-parameters/ pruning to reduce overfitting.

# Key Findings and Insights

We can make observations such as from decision tree that gives us an indication of kind of bookings we need to observe:

Lead Time > 9.5 and lead_time <= 150.50
no_of_special_requests <= 0.50 ,
market_segment_type_1 > 0.50  then the booking is more likely to cancel

If lead_time > 150.50 but avg_price_per_room <= 100.04, customer has a higher chance of cancelling.
if lead_time > 150.50 and even if avg_price_per_room > 100.04, The customer has lesser special requests than 2.5, then booking is likely to cancel.

If lead_time > 150.50 but avg_price_per_room <= 100.04, customer has a higher chance of cancelling
if lead_time > 150.50 and even if avg_price_per_room > 100.04, if the customer has lesser special requests than 2.5, she is very likely to cancel

# Recommendations

Customer booking Lead Time is one of the most important feature to determine the cancellation odds i.e.if the customer will be cancelling the booking

Market_Segment_Type is another important field for determining the customer cancellation trend

Avg_Price Per room is indicated by more than one tree models as an important field for determining the customer cancellation trend

No_Of_Special Request is a factor in determining the customer cancellation trend.

Lead Time is a key indicator . Stat Hotel can run marketing schemes to incentivize customers to book in advance along with special offers and requests to lock the booking and make them aware through channels to ensure that they are more engaged and less likely to cancel

It is observed that 36% guests had special requests and were less likely to cancel the bookings. This depicts clear trend across Customers who do adequate planning with room details are less likely to cancel.

Star Hotel should engage with the customers and increase this data collection resulting in deeper engagement.

Repeat guests % is lower (3.1%) but they were least likely to cancel. Star Hotel should run special discounts for customers with schemes for Repeat guests as it may improve the non-cancellation trends.

max number of booking cancelled by customers is in August month and also max booking took place in August so to compensate this loss Hotel can come up with a Hotel Credit which can be used for future bookings.

# Recommendations

Repeated guests hardly cancelled any bookings, but repeated guest bookings are also very less, Hotel should focus on this and they should improve the count on repeated guest by providing them more intensives by offering them memberships.

Online Bookings are the maximum counts which are getting cancelled. Hotel can use engagement techniques such as booking confirmation calls, follow-ups etc. to ensure that there is deeper engagement with the customer.

Average nightly rates is an important factor influencing cancellation. We observe that Room Type 1 (potentially the lowest segment type) is cancelled maximum times. Hotel can deploy an algorithm which ensure that cancellation factor is marginally included in the pricing algorithm.

Most booking is cancelled by customers couples ( 2 adults) with zero children. Star Hotels pricing algorithm should take this into account while pricing for couple vs families. That way it can recover some of the costs. The model can follow grades for cancellation charges with lowest refunds on last minute cancellations, 20-50% refunds within a week of booking vs 70% or more for cancellations done prior to 3 weeks of booking start date

Since there is 41% vacancy over weekends, the Hotel can offer weekend bookings instead of cancellation refunds/charges. That way, customer would be locked