# CAPSTONE PROJECT

## HOTEL BOOKING EDA ANALYSIS

# Points to Discuss:

Agenda

Data Summary

Univariate Analysis

Bivariate Analysis

**Multivariate Analysis**

Correlation Analysis

Hypothesized Question Analysis

Conclusions and Summary

# AGENDA



This project comprises a real-world data record of hotel bookings for a city and a resort hotel from 2015 to 2017, including details such as bookings, cancellations, and guest information. The project's main goal is to comprehend and visualize data from the hotel and customer perspectives to get the proper insights from the dataset and to make proper data driven decisions.



- **We have performed various kind of analysis to manage the data in terms of data wrangling.**

- **Wrangled our raw data to proper formatted data to perform the analysis to get desired results.**
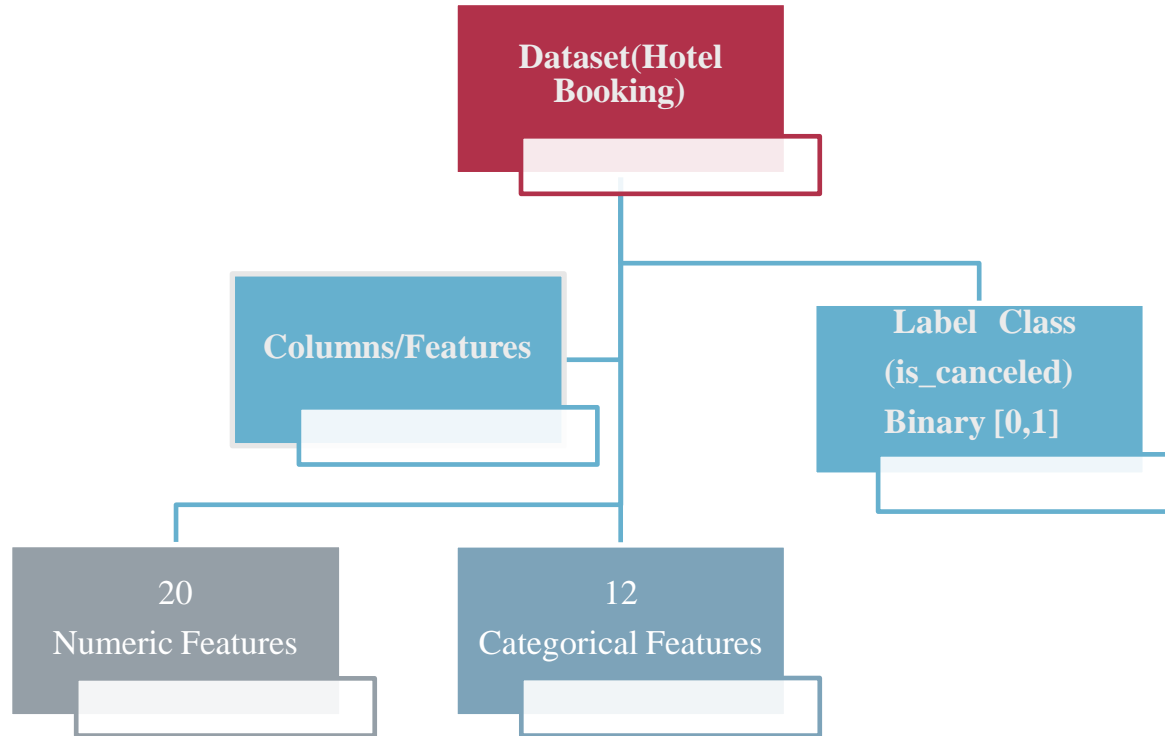
# DATA SUMMARY

- The "Hotel Booking" dataset contains very important and leading information regarding a customer prone to cancel his/her hotel booking or not.

- There are many features which have very explanatory information which helps us to get to exact flow of hotel booking occurrence.

- Some of the data fields are defined in further slides to get the proper understanding of different kind of important feature/columns:

| Features/Columns | Explanation |
|---|---|
| hotel | There are two types of hotels: resort hotels and city hotels. |
| is_cancelled | The cancellation type is indicated by the value 0 and 1 of the column.<br>1 -> Not Cancelled<br>2 -> Canceled |
| lead_time | The period between making a reservation and arriving. |
| stayed_in_weekend_nights | The number of weekend nights a reservation can stay for. |
| stayed_in_weekday_nights | Per reservation, the amount of weekday nights to stay. |
| market_segment | This column explains how reservations are made and why they are made. |
| distribution_channel | The medium of booking. |

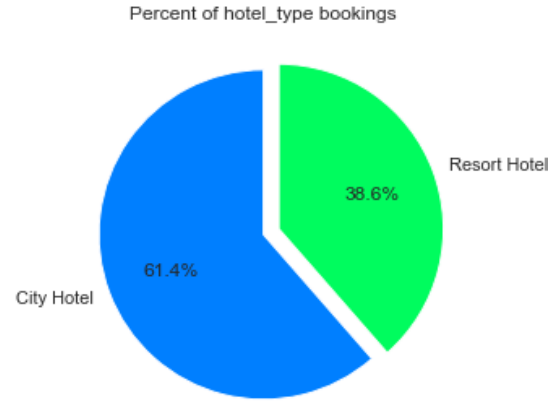| Features/Columns | Explanation |
|---|---|
| is_repeated_guest | Indicates if the guest was the one who arrived earlier or not.<br>1 -> Repeated<br>2 -> Not Repeated |
| days_in_waiting_list | The time between making a reservation and completing a transaction is measured in days. |
| customer_type | Type of customers( Transient, group, etc.) |
| reserved_room_type | The code representing the room type assigned to the reservation. Due to hotel operations, the assigned room type may differ from the reserved room type. |
| adr | column represents the average daily rate |
| assigned_room_type | Code for the type of room assigned to the booking. |

# DATA SUMMARY

# UNIVARIATE ANALYSIS

- Because "uni" means one and variate means variable, there is only one trustworthy variable in univariate analysis. The goal of univariate analysis is to derive data, characterize and summaries it, and examine any patterns that may exist. It investigates each variable separately in a dataset. There are two types of variables that can be used: categorical and numerical.

- We have performed a lot of univariate analysis on various features:

1. **Observation from hotel type that which type of hotel has how many bookings**

2. **The count of bookings is_canceled**

3. **Observation on Market Segment wise bookings**

4. **Observation on Distribution Channel wise bookings**

# Hotel Type Booking Count



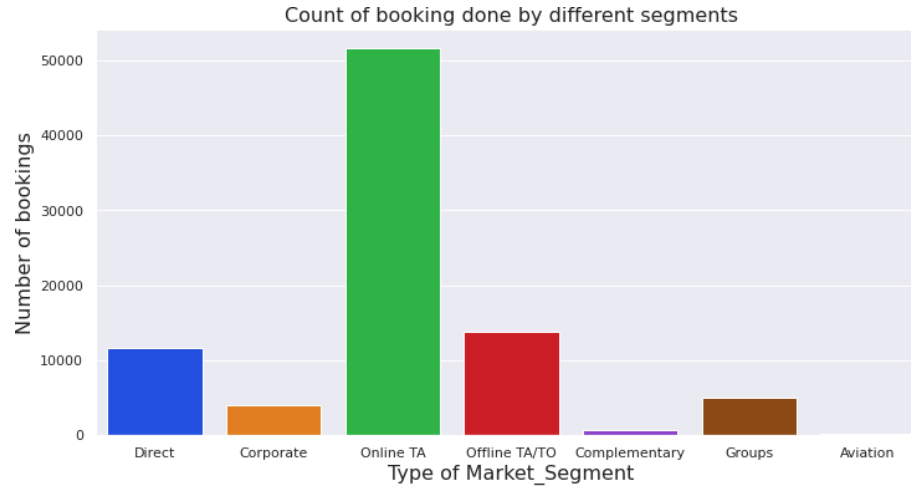Percent of hotel_type bookings

Percent of hotel_type bookings

- City hotel has bookings above 50000 which is around 61.4% and Resort hotel has less than 35000 which is around 38.6% bookings

- Which implies that City Hotel has more number of bookings than Resort Hotel
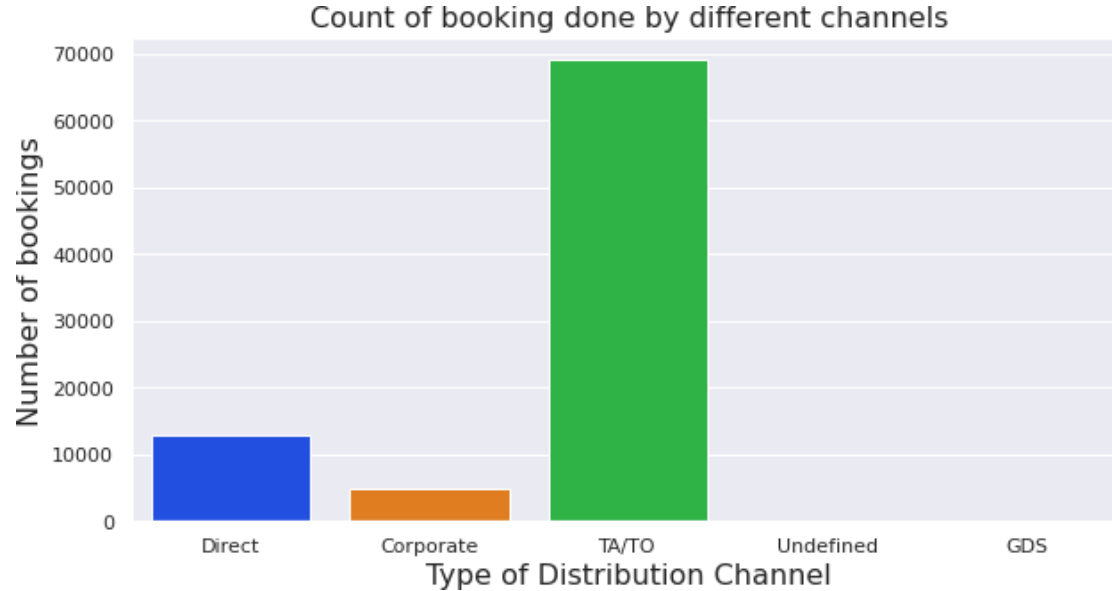
# Counting ratio of cancelled and not cancelled bookings



Count of hotel booking cancellation

- 0 ---> Non Cancelled Booking | 1 ---> Cancelled Booking

- Hotels are less cancelled as compare to cancellation bar

# Market Segment Booking Count
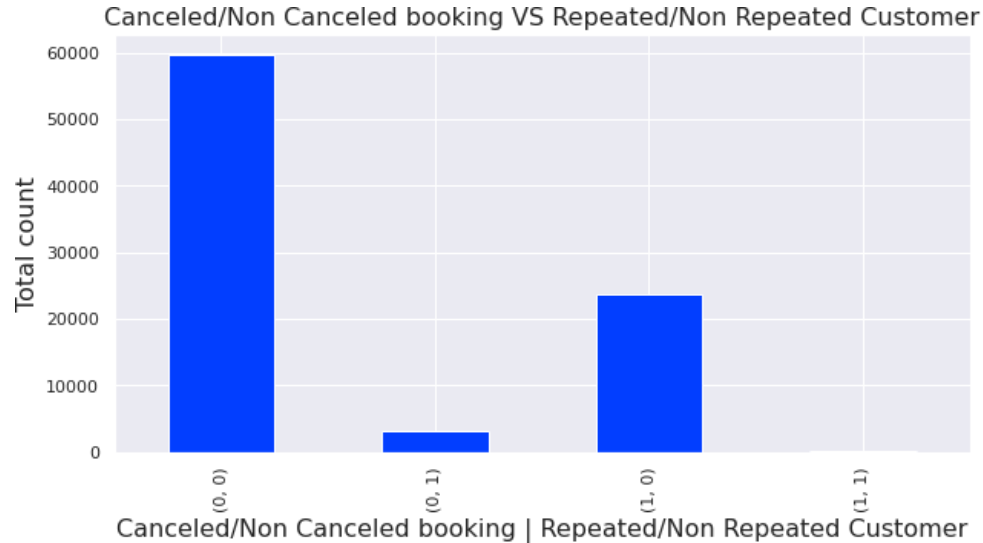
Count of booking done by different segments



- Online TA --> Online Travel Agent
- Offline TA/TO--> Offline Travel Agent and Travel Operator

- Online TA and Offline TA/TO has made to maximum number of hotel bookings

# Distribution Channel Booking Count

Count of booking done by different channels



- Offline TA/TO has made to maximum number of hotel bookings

# Repeated Guest count who cancelled the booking



Canceled/Non Canceled booking VS Repeated/Non Repeated Customer

- The lowest cancellation occurs when the customer is repeated
- That means the existing customers are very less likely to cancel the booking
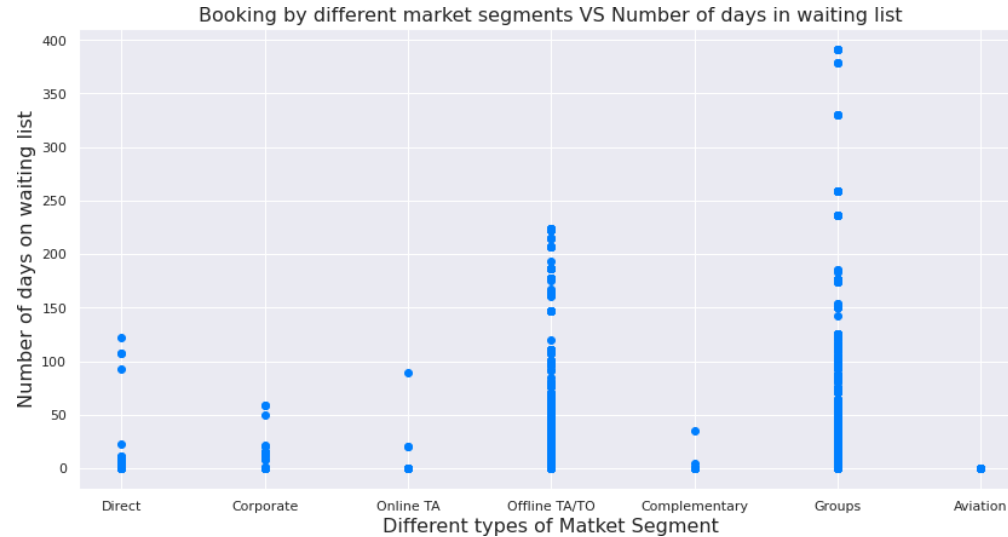
# BIVARIATE ANALYSIS

There are two variables since bi means two and variate implies variable. The investigation focuses on the root of the problem and the relationship between the two factors. Bivariate analysis can be divided into three categories.

Analysed different features together and we got interesting insights about the data.

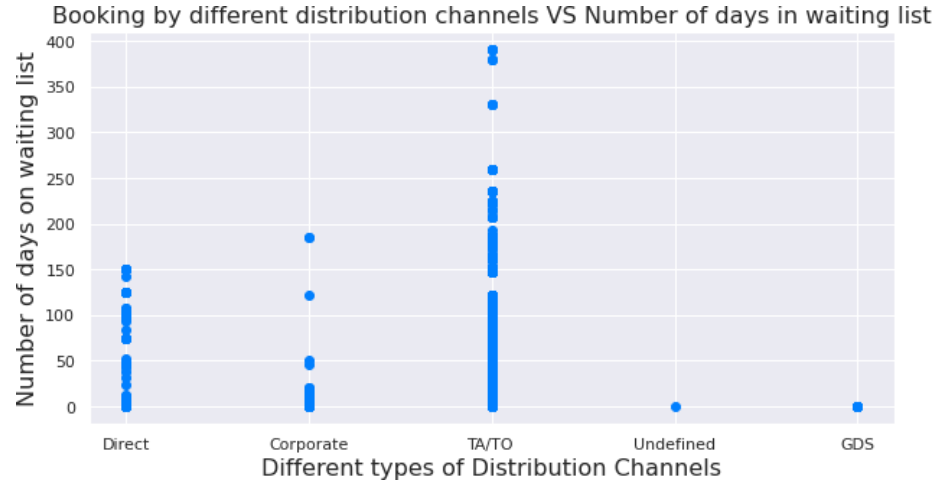In this data we have done the following steps mentioned below

- Market segment and days on waiting list for each of them
- Comparing Distribution Channels and days on waiting list for each of them
- Per month arrivals hotel bookings
- Market Segment wise bookings for each Hotel Type
- Distribution Channel wise bookings for each Hotel Type
- Months of the year with the lead time
- Arrival of customers per day of the months to the hotels
- Demand of parking space by customers in each type of hotels

# MARKET SEGMENT AND DAYS ON WAITING LIST



Booking by different market segments VS Number of days in waiting list

- The aviation industry has the minimum number of days in waiting list.

- It might due to when a flight landed they have to provide immediate accommodation for there working staff like pilots and air hostess and they do not go for that hotels which put them into higher days on waiting list.

# DISTRIBUTION CHANNELS AND DAYS ON WAITING LIST



Booking by different distribution channels VS Number of days in waiting list

Here we can also see that the TA -> Travel Agent and TO -> Travel Operator have the higher number on days on waiting list
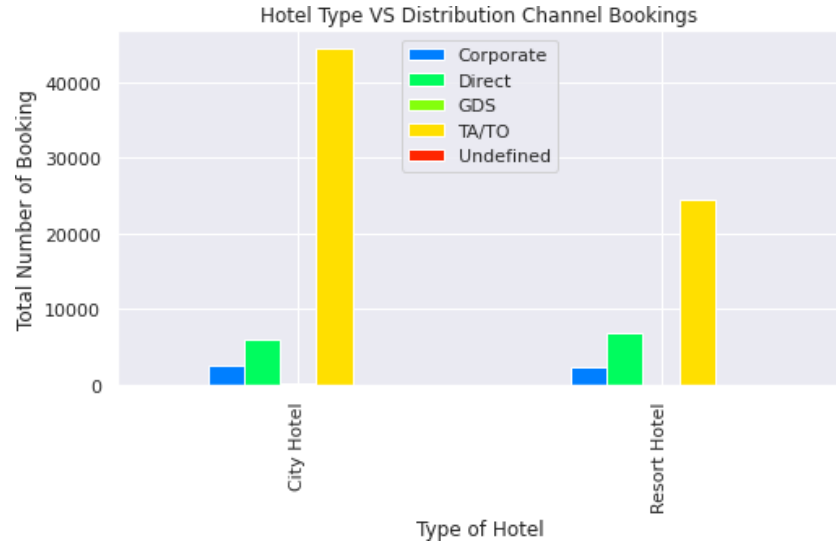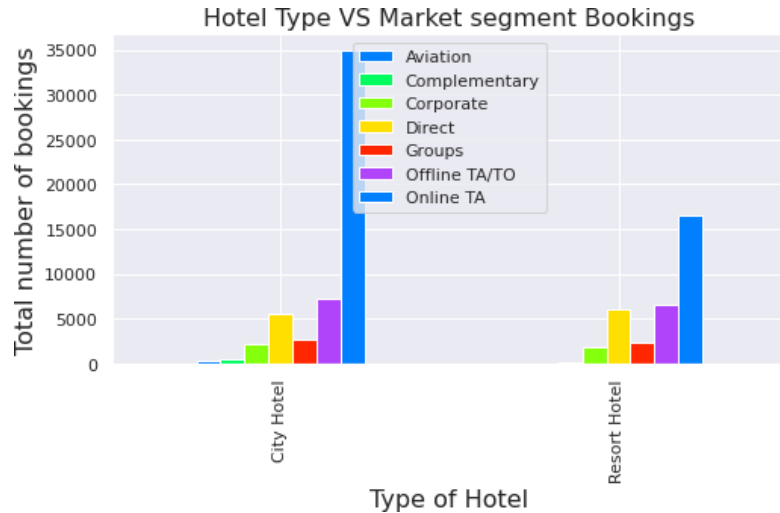
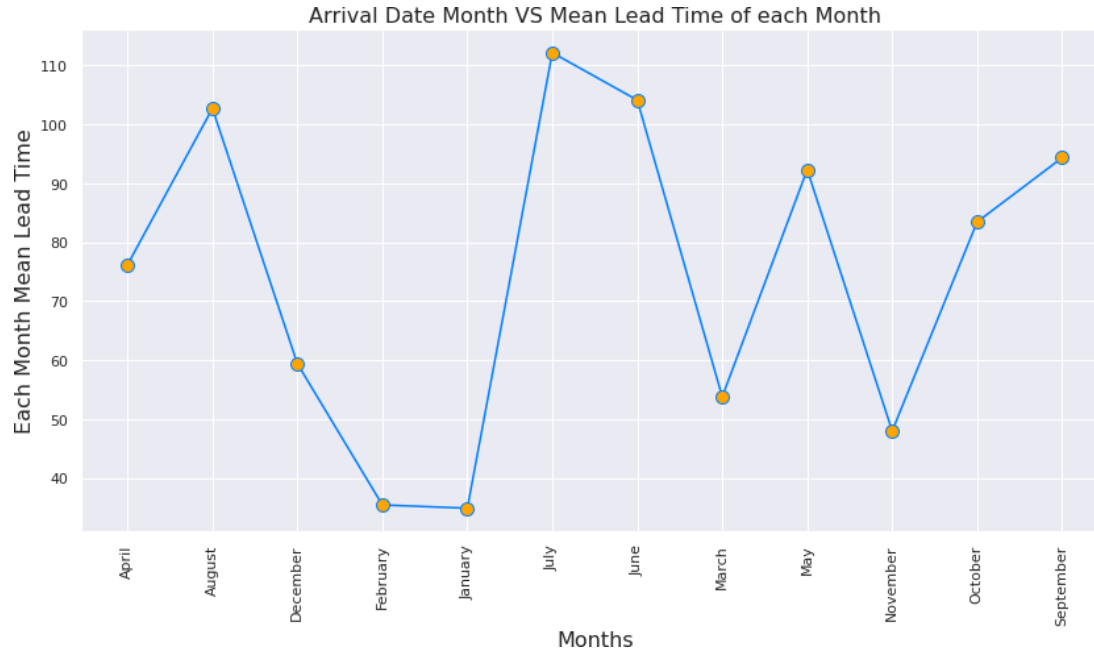number of booking done VS type of hotels

- Number of bookings are highest for city hotels in the month of August, July and May.
- And the Resort hotel has also number booking higher in the month of July and August.

# MARKET SEGMENT/DISTRIBUTION CHANNEL WISE BOOKINGS FOR EACH HOTEL TYPE
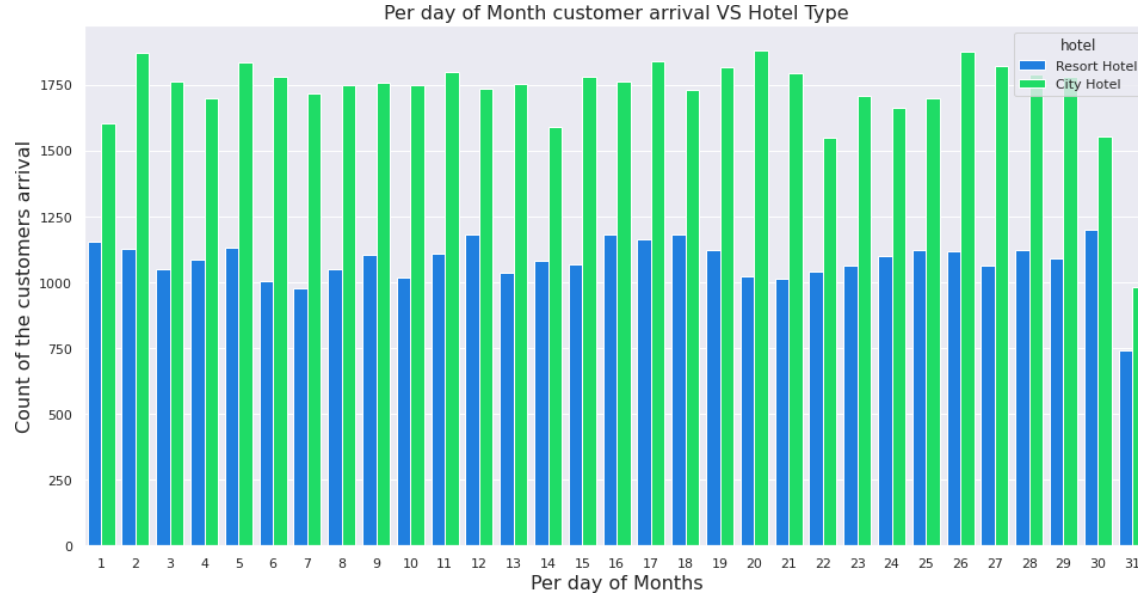


- Maximum number of booking is done by the Online TA(Travel Agent) for the City Hotels as well as for the Resort Hotels.
- The most number of booking is done by the TA/TO Distribution Channels for both City and Resort Hotel.
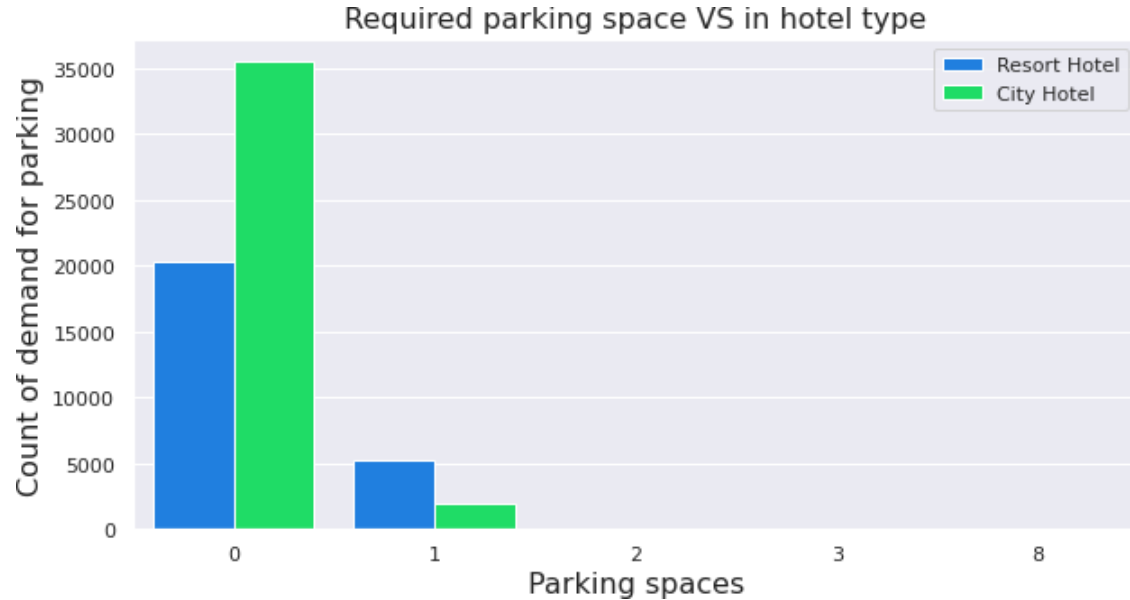
Arrival Date Month VS Mean Lead Time of each Month

- The lead time for the month of July is very high and for January and February Month is the least.

Per day of Month customer arrival VS Hotel Type

- The least number of bookings are occurred in the month ends.
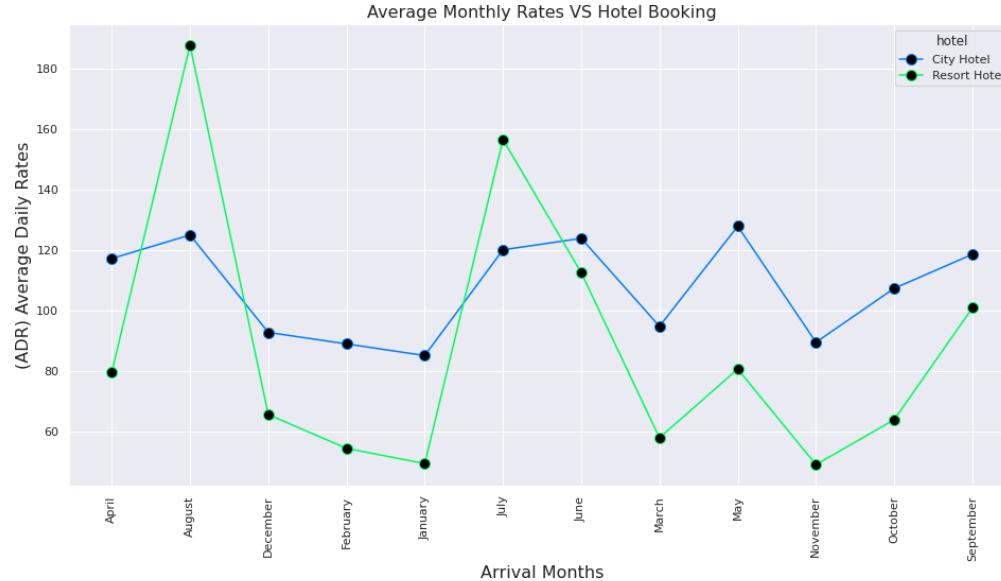
Required parking space VS in hotel type

- This chart tells us that most number of the guest has no demand for parking space in both the hotel type.

# MULTIVARIATE ANALYSIS

- Multivariate indicates that numerous dependent variables are combined to produce a single result. This explains why the vast majority of real-world situations are multivariate.

- We have done Multivariate analysis to check hotel bookings occurred according to the average daily rates per month.

- By performing this multivariate analysis we came to know about how different hotel type has impacted in each month in terms of average daily rates.

- It explains us that how customers are prone to book hotels which month they prefer and when not.

Average Monthly Rates VS Hotel Booking

- The average daily rate is more expensive during June, July, August and September for Resort Hotels
- The average daily rate is more expensive during June, July and May for City Hotels
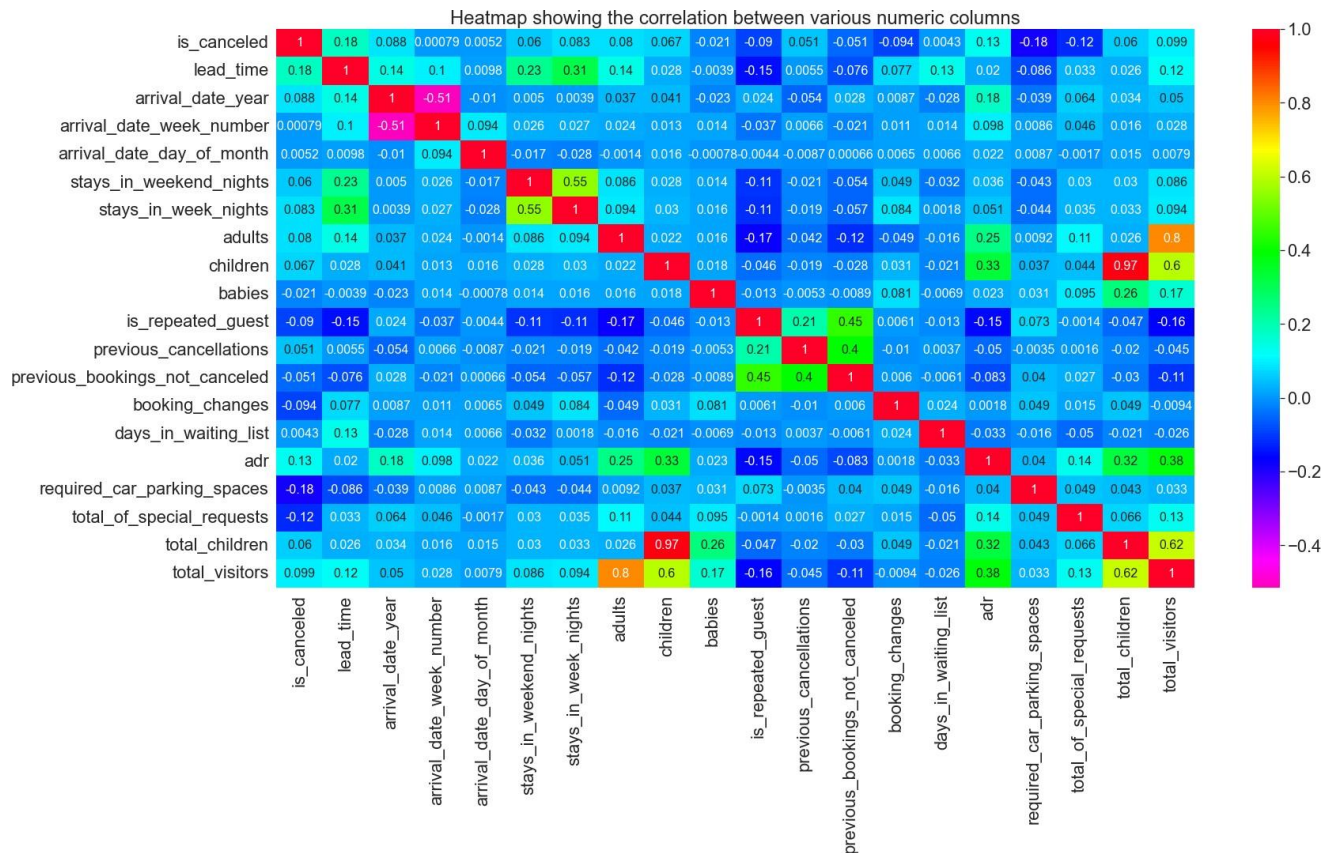
Correlation is a statistical term that expresses how closely two variables are related in a linear fashion (meaning they change together at a constant rate). It's a typical way of describing simple relationships without stating a cause and effect relationship.

- The negative correlation coefficient between the data pair LEAD TIME versus ADR indicates that as the time between room booking and check-in rises, the ADR for that hotel room will fall by a modest amount..

- is_canceled is not influenced but little bit by lead_time feature.

# Correlation Matrix:



Heatmap showing the correlation between various numeric columns

# HYPOTHESIZED QUESTIONS TO CROSS-CHECK SOME ANOTHER ASPECTS OF OUR ANALYSIS:

1. Do customers have any preference with or without children's.

2. How daily average rate is impacting the reserved room type in hotels.

3. Does total stay per month impacted by customers with and without children's.

4. Do reservation status impacted by type of customers.

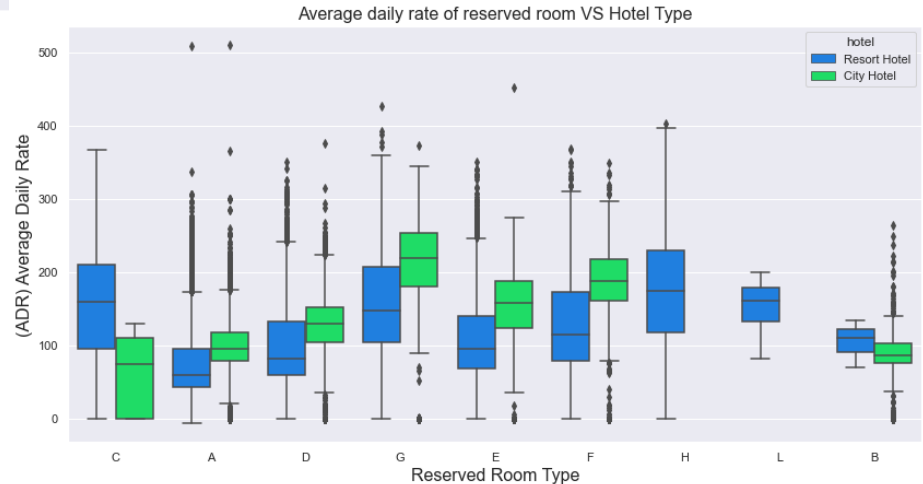5. Do customers cancel their booking if they are allotted with different room type.

Preference for Hotels VS Customers with childrens

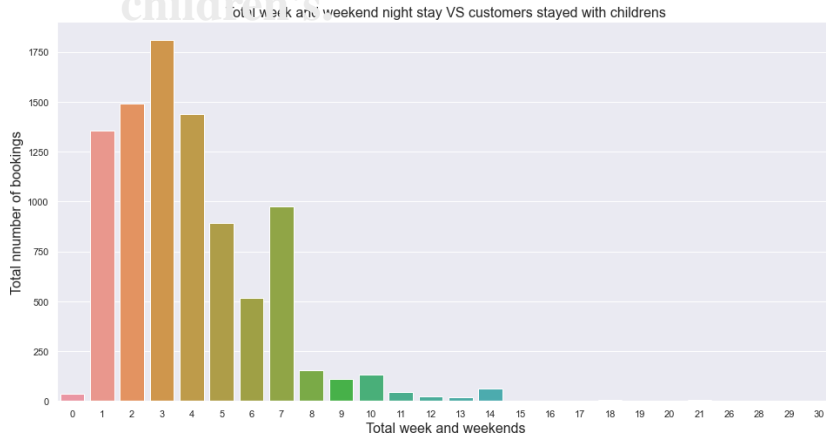- Customers with children's doesn't have that much preference for the hotel type.

- Room type 'G' has the highest and Room type 'C' has the lowest Average Daily Rates (ADR) in City Hotels.

- Room type 'C' and 'H' has the highest and Room type 'A' has the lowest Average Daily Rates (ADR) in Resort Hotels.



Average daily rate of reserved room VS Hotel Type

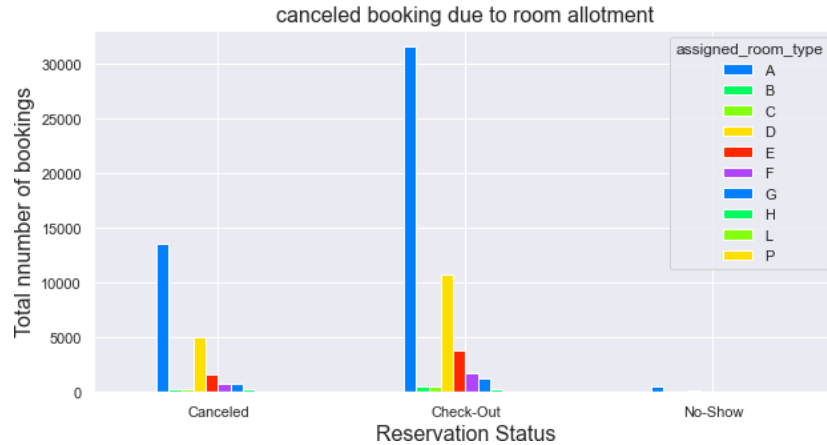Total week and weekend night stay VS customers stayed with childrens

- There is very much similarity in both cases customers stayed with and without children's about 1-7 days have the count for stays with and without children.

- If we go for the average stay then there would be 3-4 days of stay in which people have there children with them otherwise there is only adults staying
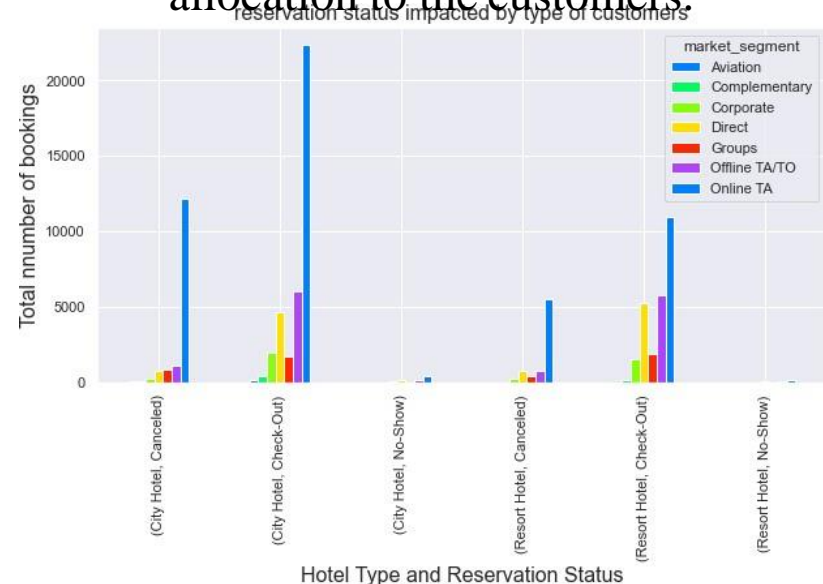


Total week and weekend night stay VS customers stayed without childrens

canceled booking due to room allotment

- Hence there is not that much impact of room type allocation to the customers.

reservation status impacted by type of customers

- Customer group from Online TA has the highest number of booking for city hotels.
- Customer group that booked directly is highest at resort hotels.
- But We can Clearly see that most of the bookings are provided by Online

# CONCLUSIONS AND SUMMARY

- The vast majority of reservations are for hotels in cities. Resort hotels have fewer cancellations than city hotels. The aviation industry has the shortest wait time.

- The months of August, July, and May saw the most hotel bookings in the city. The lowest cancellation rate occurs when a consumer is repeated. The lead time for July is quite long, whereas the lead time for January and February is extremely short. In July, August, and September, the average daily rate for Resort Hotels is higher.

- In June, July, and May, the average daily rate for City Hotels is higher. Customers travelling with children have little preference for the type of hotel they stay in.

# Conclusion and Summary(contd..)

- In both cases, the number of stays with and without children is rather equal. Customers who stayed for 1-7 days with or without children had the same count.

- If we take the average stay, three to four days will be spent with youngsters, while the rest of the time will be spent with adults solely.

- Online TA customer group has the highest number of hotel reservations in the city and at resort hotels. As a result, the impact of room type allocation on clients is minimal. The majority of appointments are made by consumers of the Online TA group who book directly through the website.

# THANK YOU