# Capstone Project
## EDA on Hotel Booking Analysis

**Nehal S Jambhulkar**

# Problem Statement

✓ For this project we will be analyzing Hotel Booking data. This data set contains booking information for a city hotel and a resort hotel, and includes information such as when the booking was made, length of stay, the number of adults, children, and/or babies, and the number of available parking spaces.

✓ Hotel industry is a very volatile industry and the bookings depends on above factors and many more.

✓ The main objective behind this project is to explore and analyze data to discover important factors that govern the bookings and give insights to hotel management ,which can perform various campaigns to boost the business and performance.

➢ So we will divide our work flow into following 3 steps

| Data Collection and Understanding | Data Cleaning and Manipulation | Exploratory Data Analysis(EDA) |
| --- | --- | --- |

**EDA will be divided into following 3 analysis.**

1. **Univariate analysis:** Univariate analysis is the simplest of the three analyses where the data you are analysing is only one variable.
2. **Bivariate analysis:** Bivariate analysis is where you are comparing two variables to study their relationships.
3. **Multivariate anlysis:** Multivariate analysis is similar to Bivariate analysis but you are comparing more than two variables.

# Data Collection and Understanding:

❖ After collecting data it's very important to understand your data. So we had hotel Booking analysis data. Which had 119390 rows and 32 columns. So let's understand this 32 columns.

## Data Description:

**hotel :** Resort Hotel or City Hotel

**is_canceled :** Value indicating if the booking was canceled (1) or not (0)

**lead_time :** Number of days that elapsed between the entering date of the booking and the arrival date

**arrival_date_year :** Year of arrival date

**arrival_date_month :** Month of arrival date

**arrival_date_week_number :** Week number of year for arrival date

**arrival_date_day_of_month :** Day of arrival date

**stays_in_weekend_nights :** Number of weekend nights

**stays_in_week_nights :** Number of week nights.

**adults :** Number of adults

**children :** Number of children

**babies :** Number of babies

**meal :** Type of meal booked.

**country :** Country of origin.

**market_segment :** Market segment designation. (TA/TO)

**distribution_channel :** Booking distribution channel.(T/A/TO)

**is_repeated_guest :** is a repeated guest (1) or not (0)

**previous_cancellations :** Number of previous bookings that were cancelled by the customer prior to the current booking

**previous_bookings_not_canceled :** Number of previous bookings not cancelled by the customer prior to the current booking

**reserved_room_type :** Code of room type reserved.

**assigned_room_type :** Code for the type of room assigned to the booking.

**booking_changes :** Number of changes made to the booking from the moment the booking was entered on the PMS until the moment of check-in or cancellation

**deposit_type :** No Deposit, Non Refund , Refundable.

**agent :** ID of the travel agency that made the

**booking company :** ID of the company/entity that made the booking .

**days_in_waiting_list :** Number of days the booking was in the waiting list before it was confirmed to the

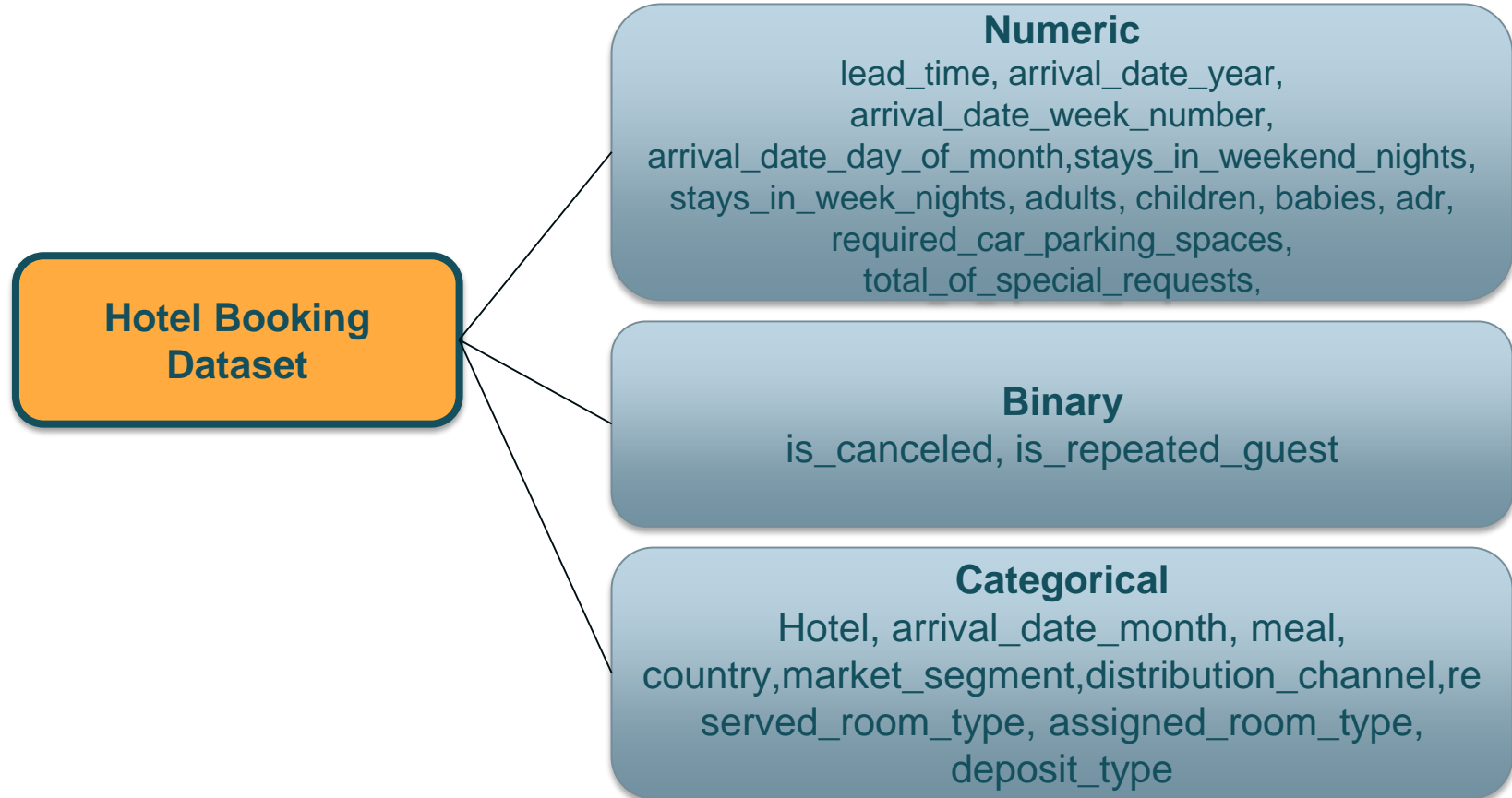**customer customer_type :** type of customer,Contract,Group,transient,Transient party.

**adr :** Average Daily Rate as defined by dividing the sum of all lodging transactions by the total number of staying nights

**required_car_parking_spaces :** Number of car parking spaces required by the customer

**total_of_special_requests :** Number of special requests made by the customer (e.g. twin bed or high floor)

**reservation_status :** Reservation last status.
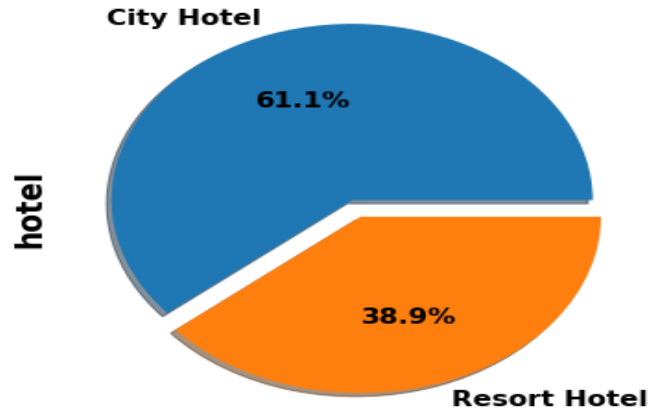
# Dataset Input data summary

**AI**

**Hotel Booking Dataset**

**Numeric**
lead_time, arrival_date_year,
arrival_date_week_number,
arrival_date_day_of_month,stays_in_weekend_nights,
stays_in_week_nights, adults, children, babies, adr,
required_car_parking_spaces,
total_of_special_requests,

**Binary**
is_canceled, is_repeated_guest

**Categorical**
Hotel, arrival_date_month, meal,
country,market_segment,distribution_channel,re
served_room_type, assigned_room_type,
deposit_type

# Data Cleaning and Manipulation:

1. Company, agent, country and children columns with missing values.I replaced  missing  values as per requirement.

2. Dropping company column because more then 90% data is missing

3. Data had 31994 duplicates values. So I dropped it from the data.

4. I created 2 new columns

   A)'total_People' = from the Children, adults, babies.

   B) 'total_stay ' = From weekend nights and weekdays night.

# Exploratory Data Analysis (EDA) :

**AI**

## Univariate Analysis



**Pie Chart for Most Preffered Hotel**

City Hotel — 61.1%
Resort Hotel — 38.9%

**Cancellation and non Cancellation**

is_canceled
0 — 72.5%
1 — 27.5%

**Conclusions:**

1) City hotels is the most preferred hotel type by the guests. We can say City hotel is the busiest hotel.

2) 27.49 % bookings were got cancelled out of all the bookings

**Percentgae (%) of repeated guests**

is_repeated_guest

0 — 96.1%
1 — 3.9%

**% Distribution of Customer Type**

customer_type

- Transient
- Transient-Party
- Contract
- Group

82.4%
13.4%
3.6%
0.6%

**Conclusion:**
➢ Only 3.9 % people were revisited the hotels. Rest 96.1 % were new guests. Thus retention rate is low.
➢ Most of the customers/guests were Transient type(82.4%). And transient party were 13.4% and 0.6 belongs to group. Remaining guests belongs to Contract type.

**Most Booking Made By the agent**

**Conclusions:**

➤ Agent Id no -9 made the highest bookings which is more than 28721.

**% of Booking change**

**Conclusions:**

➤ The percentage of 0 changes made in the booking was more than 82 %.

● Percentage of Single changes made was about 10%.

% Distribution of required car parking spaces

**Conclusions:**

➢ Most of the customers(91.6%) do not require car parking spaces. Only 8.3 % people required only 1 car parking space.

**Mostly Used Distribution Channel for Hotel Bookings**

- TA/TO, 79.1%
- Direct, 14.9%
- Corporate, 5.8%
- GDS, 0.2%
- Undefined, 0.0%

**Most preferred room Type**

## Conclusions:

➢ 79.1 % bookings were made through TA/TO (travel agents/Tour operators). Second most channel is direct.

➢ Room type 'A' is most preferred by the guests second most preferred is 'D'.

**Conclusions:**

➤ Almost 98.7% of the guests prefer 'No deposit' type of criterion while booking hotels.

PRT- Portugal
GBR- United Kingdom
FRA- France
ESP- Spain
DEU - Germany
ITA –Italy
IRL - Ireland
BEL -Belgium
BRA -Brazil
NLD-Netherlands

**Number of guest from diffrent countries**

## Conclusions:

➢ Most of the guests are coming from Portugal i.e. more 25000 guests are from Portugal

**Conclusion:**

➢ Most of the bookings for City hotels and Resort hotel were happened in 2016. As we can see Most of the bookings were for City hotels
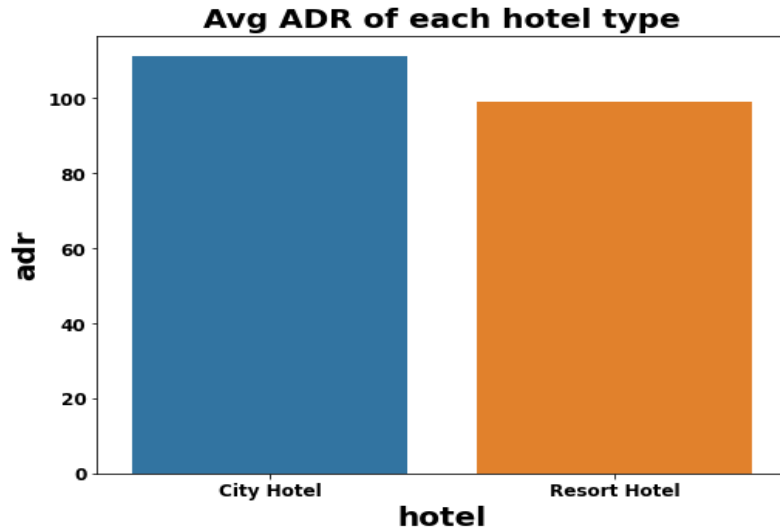
**Conclusions:**

➢ BB( Bed & Breakfast) is the most preferred type of meal by the guests.

➢ Full Board i.e. FB is least preferred.

➢ HB (Half Board) and SC(Self Catering) are equally preferred.



Preferred Meal Type

➢ As we can see in the line chart,from June to September most of the bookings happened.
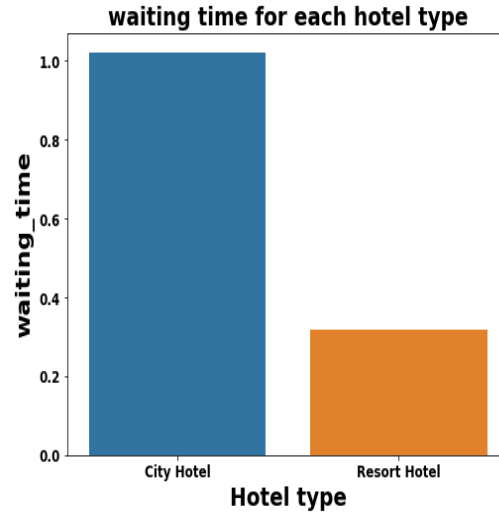It's Summer time.
After September bookings Starts declining.



Number of bookings across each month

AI

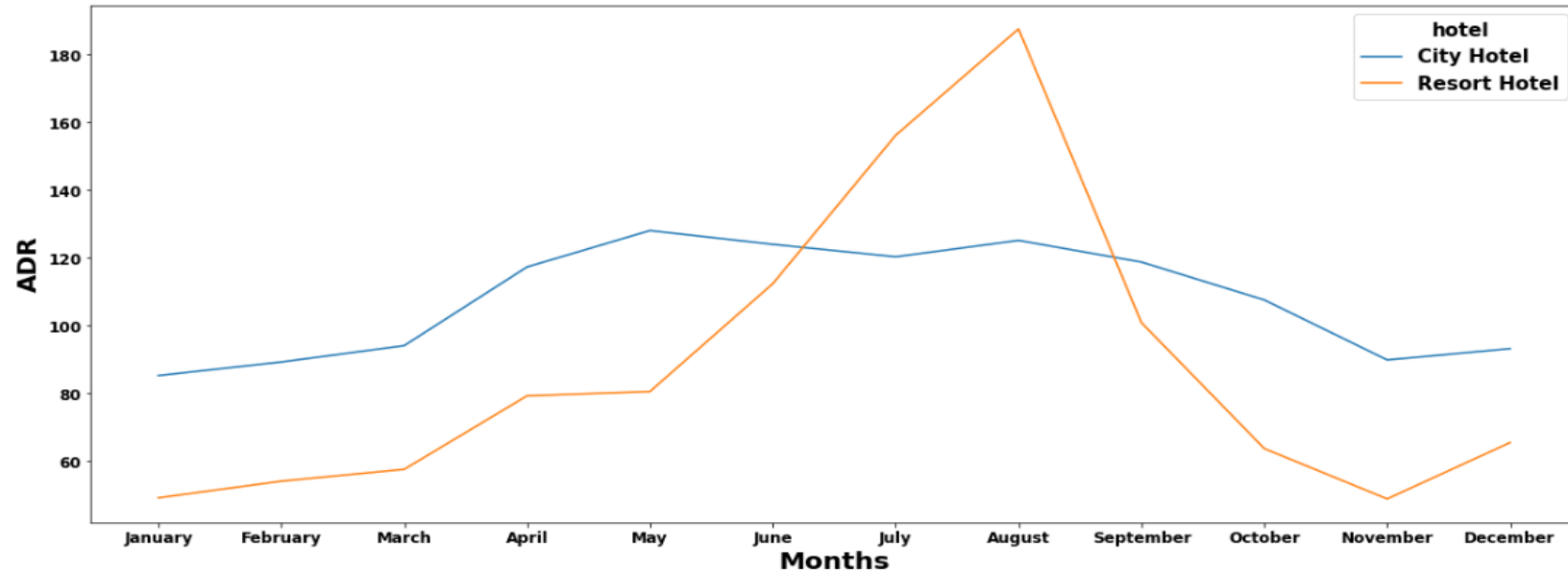# Bivariate and Multivariate Analysis



**Conclusions:**

➢ Average ADR for city hotel is high as compared to resort hotels. These City hotels are generating more revenue than the resort hotels.

➢ Average lead time for resort hotel is high. It means people plan their trip too early. Usually people prefer resort hotels for longer stays. That's why people plan early

Percentage of booking cancellation — waiting time for each hotel type — Most repeated guest for each hotel

**Conclusions:**

➢ Booking cancellation rate is high for City hotels which almost 30 %.

➢ Waiting time period for City hotel is high as compared to resort hotels. That means city hotels are much busier than Resort hotels.

➢ Resort hotels has the most repeated guests. In order to get increase the count of repeated guests hotel management need to take the valuable feedbacks from the guests and try to give good service.
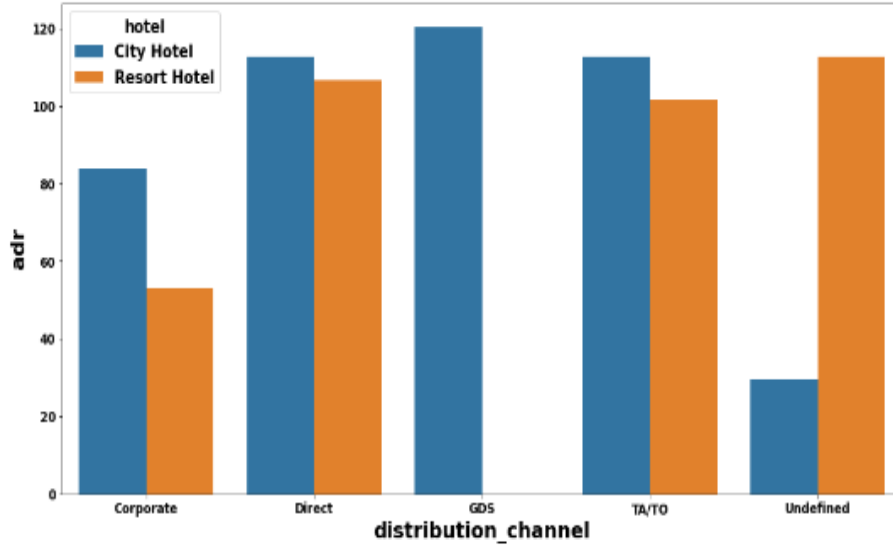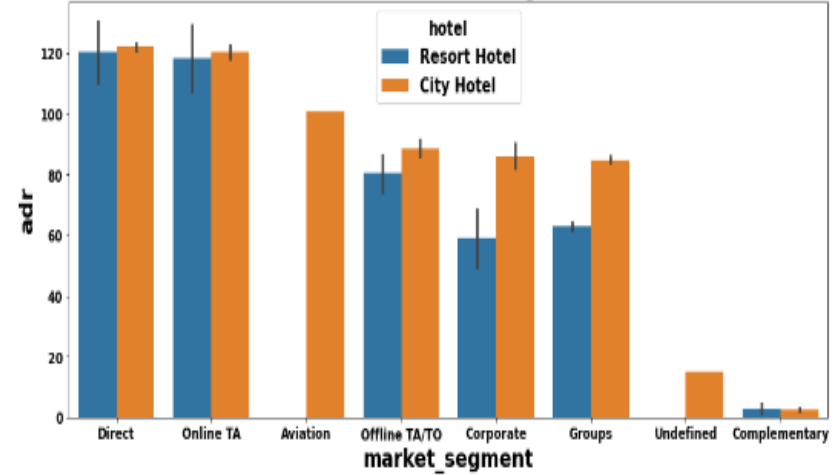
**ADR across each Month**

- **Conclusions:**
- ➢ Resort hotels had the highest adr in June ,July and August than the City hotels. But in other months adr of Resort hotel was less than the City hotels.
- ➢ Thus we can say that, the January, February, March, April ,November and December are the good months for customers to get good adr

ADR across Distribution channel / Adr across market segment
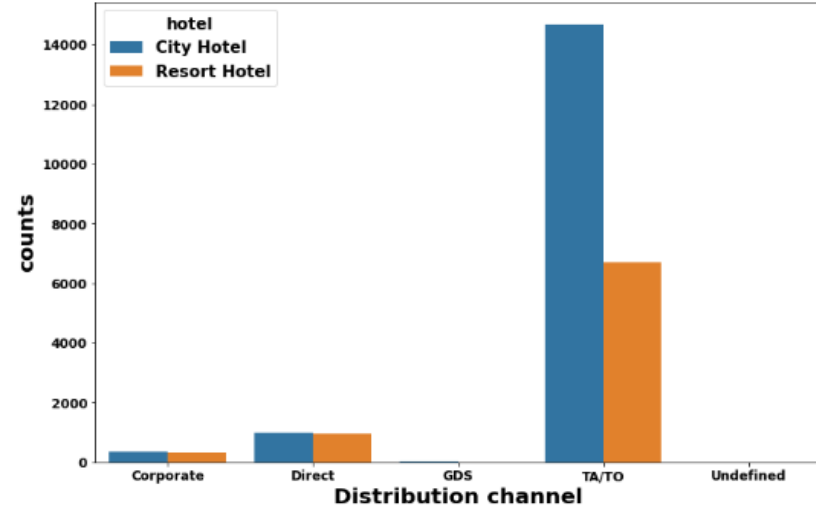
## Conclusions:

**Distribution channel:**

➢ 'Direct' and 'TA/TO' has almost equal adr in both type of hotels which is high among other channels.

➢ GDS has high adr in 'City Hotel' type. GDS needs to increase Resort Hotel bookings. From this we can say that "Direct" and 'TA/TO' are generating more revenue than the other channels.
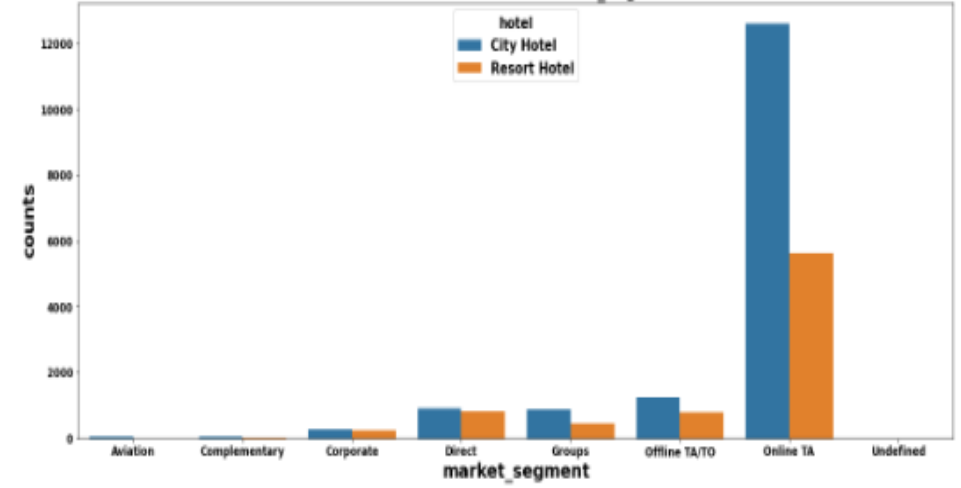
**Market Segment:**

➢ Here "Direct" and 'Online Travel Agency' has high adr for both hotel types. Aviation segment needs to increase Resort hotel bookings.

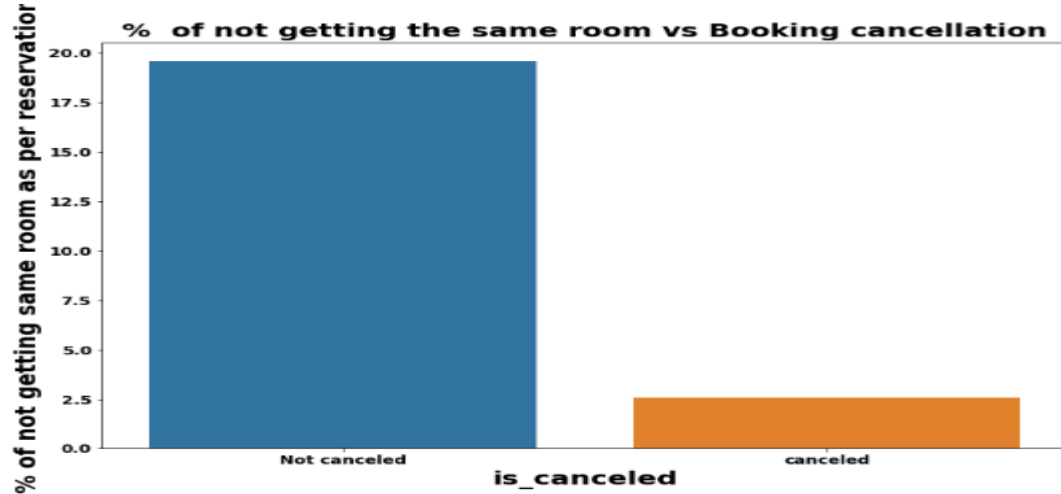**Cancellation Rate Vs Distribution channel** / **cancellation Rate Vs market_segment**

## Conclusions:

**Distribution channel:**

➢ 'TA/TO' distribution channel has highest cancellations for city hotels and more than 6000 cancellations for resort hotels. In order to reduce the cancellations they should improve their cancellation policies and deposit policies.
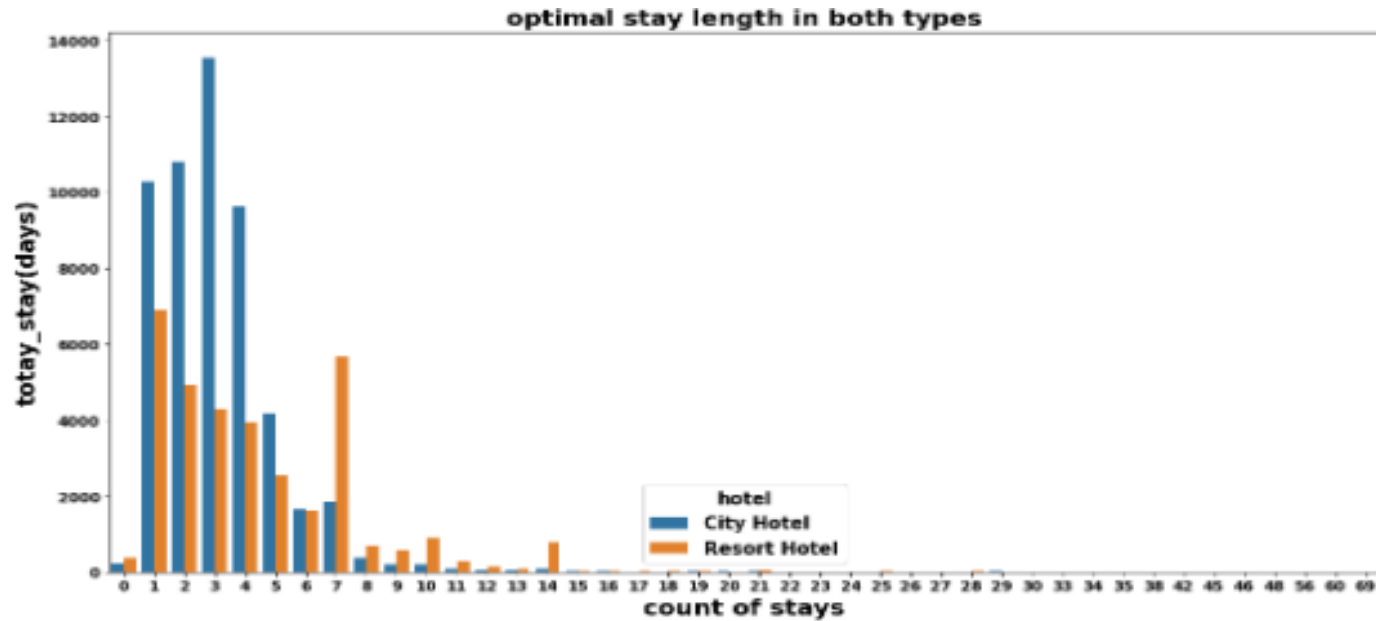
**Market Segment:**

➢ 'Online TA/TO' market segment has highest cancellations for city hotels.

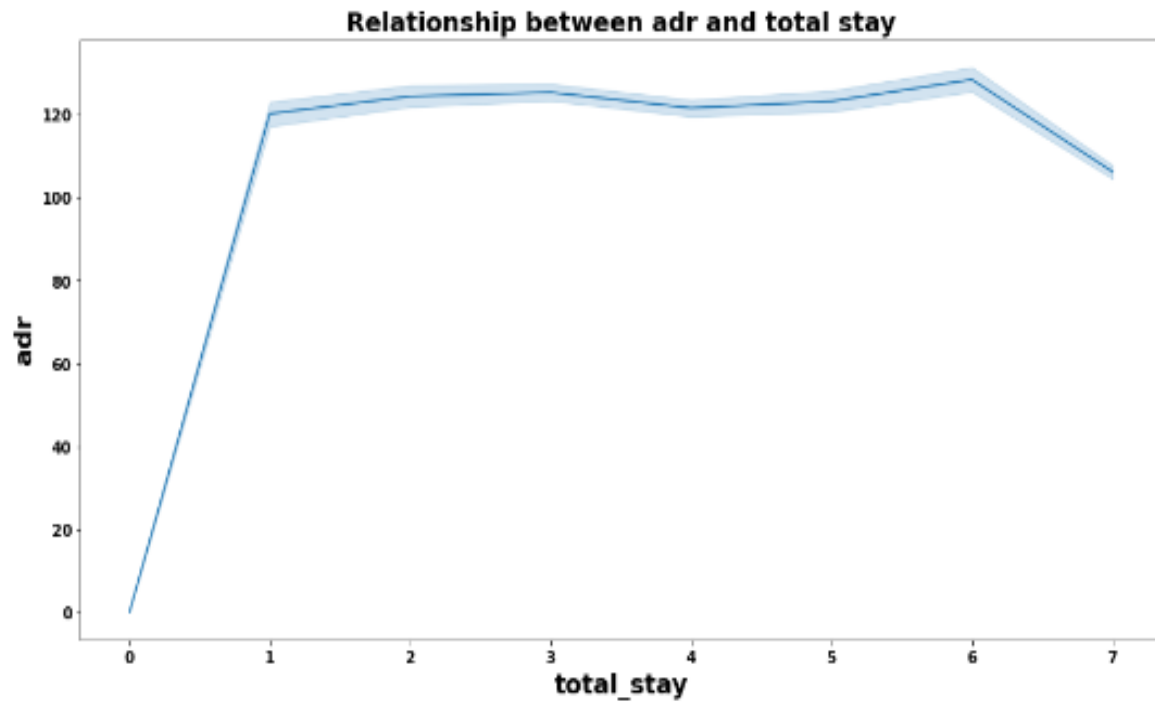% of not getting the same room vs Booking cancellation

**Conclusions:**

➤ Almost 19 % people did not canceled their bookings even after not getting the same room which they reserved while booking hotel. Only 2.5 % people cancelled the booking.

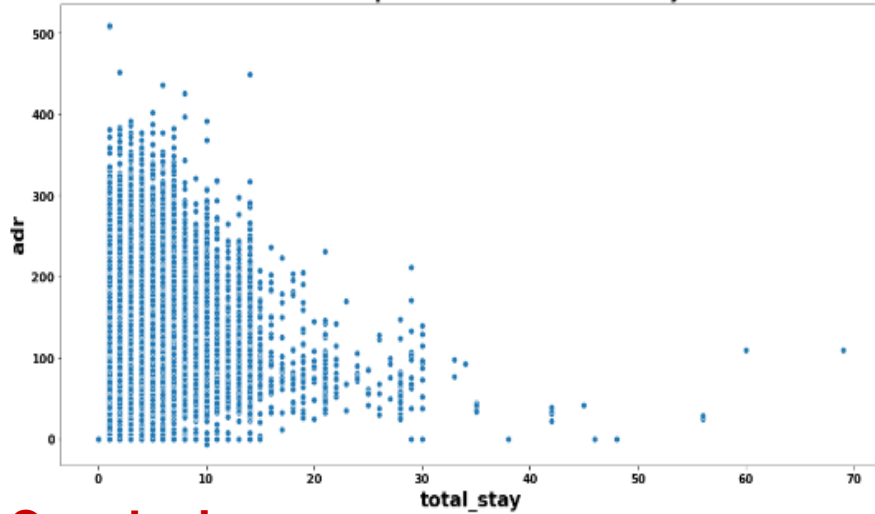➤ Thus not getting the same room as per reserved room is not the reason for booking cancellations.

optimal stay length in both types

**Conclusions:**

➢ Optimal stay in both the type hotel is less than 7 days.
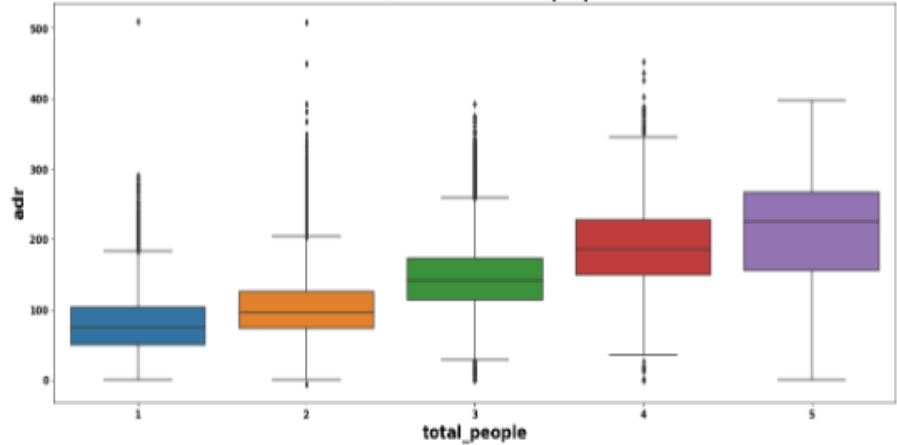
Relationship between adr and total stay

**Conclusions:**

➢ As the total stay increases the adr also increases.
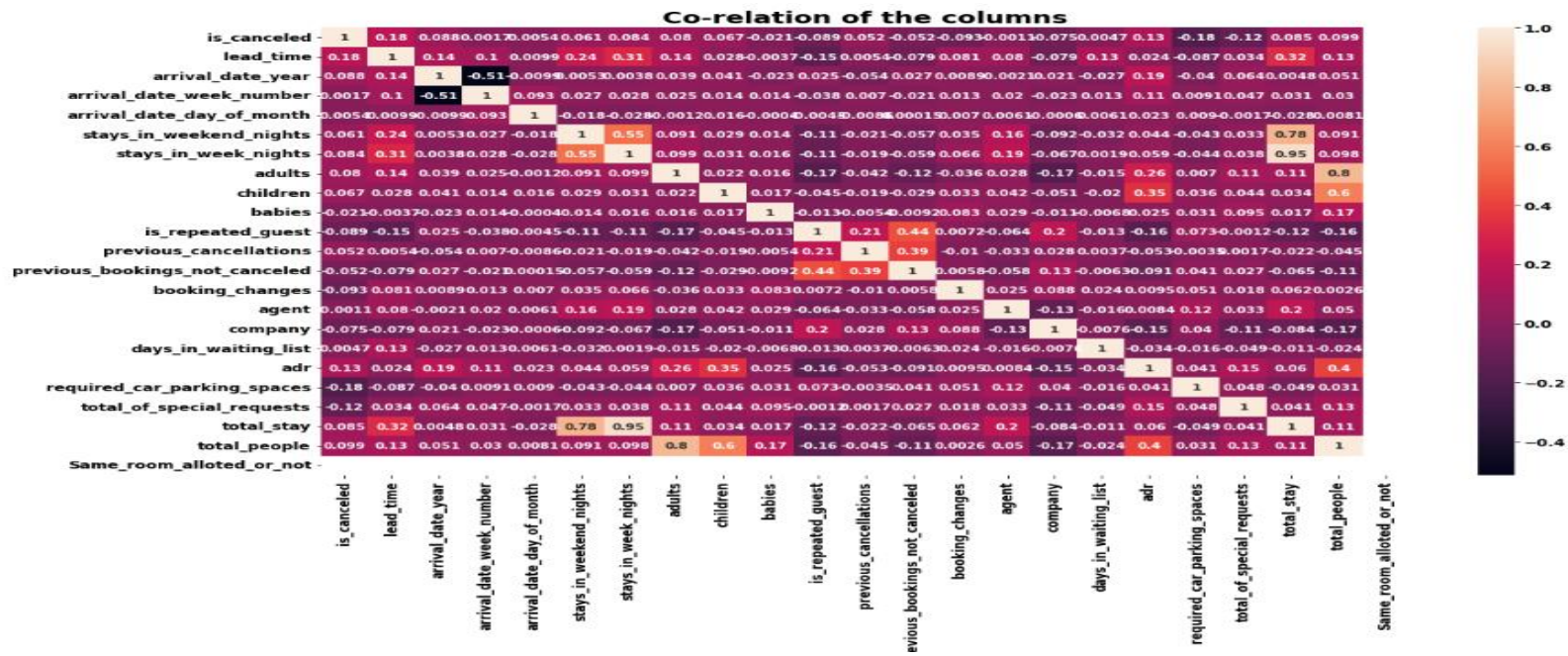
Relationship between adr and total stay



ADR v/s Total Number of people

**Conclusions:**

➤ As the total number of people increases adr also increases.

➤ Thus adr and total people are directly proportional to each other.

➤ From above scatter we can say that as the stay increases adr is decreasing. Thus for longer stays customer can get good adr.

Co-relation of the columns

## Conclusions:

➢ is_canceled and same_room_alloted_or_not are negatively corelated. That means customer is unlikely to cancel his bookings if he don't get the same room as per reserved room. We have visualized it above.

➢ lead_time and total_stay is positively correlated. That means more is the stay of customer more will be the lead time.

➢ adults, childrens and babies are corelated to each other. That means more the people more will be adr.

➢ is_repeated guest and previous bookings not canceled has strong correlation. may be repeated guests are not more likely to cancel their bookings.