# CAPSTONE PROJECT- I

## EDA on Hotel Booking Analysis

BY

Monika Ransing

AImaBetter

# PROBLEM STATEMENT:

- In this project I am going to analyze Hotel Booking dataset. This dataset contains information of city hotel and resort hotel, and includes information of booking time, length of stay, number of adults, children and/or babies, also have information of available parking space, among other thing.

- The objective of this project is explore and analyze the data to discover important factors that govern the booking.

# WORK FLOW:

➡ I divide my workflow into 3 steps.
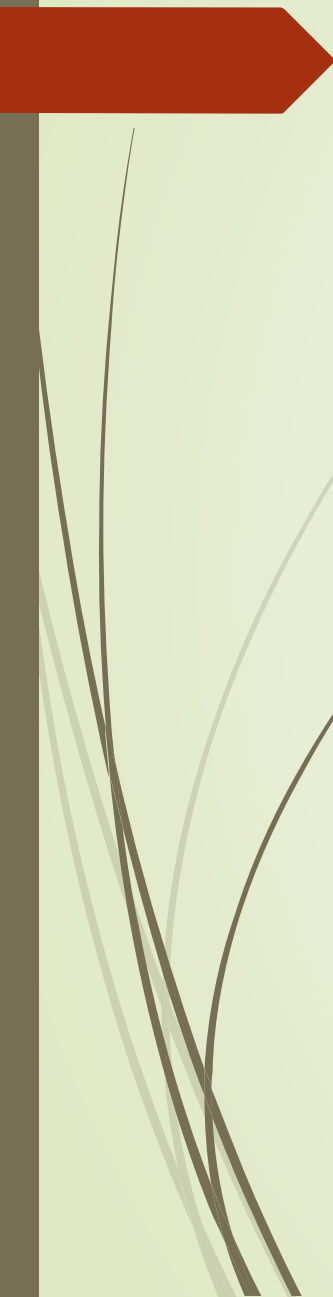
| Data Collection and Understanding | Data cleaning and Manipulation | Exploratory Data Analysis(EDA) |

➢ **EDA is divided into 3 analysis:**

 1) **Univariate Analysis:** Here I am analysing only one variable.

 2) **Bivariate Analysis:** Here I am analysing two variables and their relationship.

 3) **Multivariate Analysis:** Here I am analysing more than two variables.

# DATA COLLECTION AND UNDERSTANDING:

➡ Data collection and understanding are very important. So I have Hotel Booking data. This data contains **119390 rows** and **32 columns**. Let's understand the columns.

➡ **Dataset Description:**

o **Hotel :** Type of hotel(City or Resort)

o **is_canceled :** If the booking was canceled (1) or not (0)

o **lead_time:** Number of days before the actual arrival of the guests

o **arrival_date_year :** Year of arrival date

o **arrival_date_month :** Month of arrival date

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

o **arrival_date_day_of_month :** Day of arrival date

o **stays_in_weekend_nights :** Number of weekend nights (Saturday or Sunday) spent at the hotel by the guests.

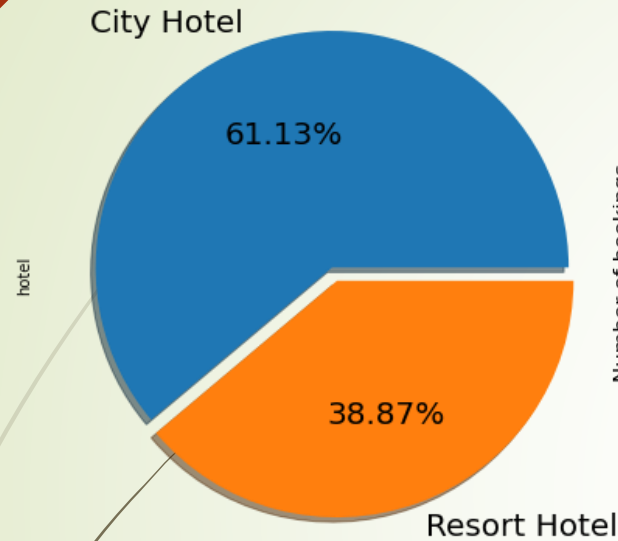o **stays_in_week_nights :** Number of weeknights (Monday to Friday) spent at the hotel by the guests.

o **adults :** Number of adults among guests

o **children :** Number of children among guests

o **babies :** Number of babies among guests

o **meal :** Type of meal booked

o **country :** Country of guests

o **market_segment :** Designation of the market segment

o **distribution_channel :** Name of booking distribution channel

o **is_repeated_guest :** If the booking was from a repeated guest (1) or not (0)

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

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

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

o **assigned_room_type :** Code of room type assigned

o **booking_changes :** Number of changes/amendments made to the booking

o **deposit_type :** Type of the deposit made by the guest

o **agent :** ID of the travel agent who made the booking

o **company :** ID of the company that made the booking

o **days_in_waiting_list :** Number of days the booking was on the waiting list

o **customer_type :** Type of customer, assuming one of four categories

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

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

o **total_of_special_requests :** Number of special requests made by the customer

o **reservation_status :** Reservation status (Canceled, Check-Out or No-Show)

o **reservation_status_date :** Date at which the last reservation status was updated
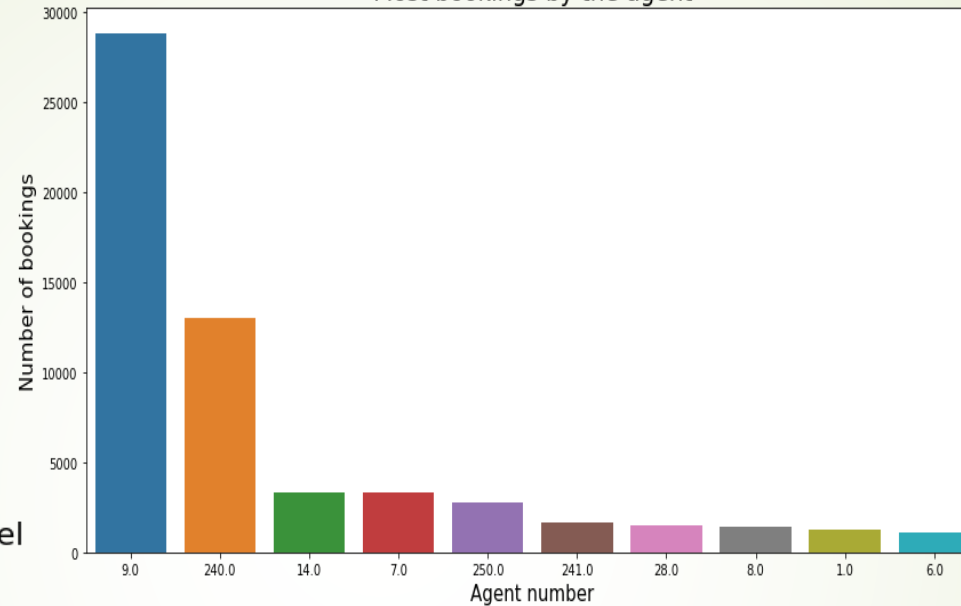
# DATA CLEANING AND MANIPULATION:

- Given data has **31994** duplicate values. So the duplicate values are dropped by using **drop_duplicates**().

- Data has 4 columns that have missing values and those columns are **company(82137)**, **agent(12193)**, **country(452)** and **children(4)**. So the missing values of columns company, agent and children are replaced by **0** and the missing values of column country are replaced by others by using **.fillna**().

- There are some rows in data that have the total number of adults, children or babies equal to zero this means there is no booking made. So I remove such rows.

- In the dataset I add 2 important columns and those columns are **"Total stay"** and **"Total People"**. For **total stay** column I add **'stays_in_weekend_night'** and **'stays_in_week_night'** columns and for **total people** column I add **'adults', 'children', 'babies'** columns.

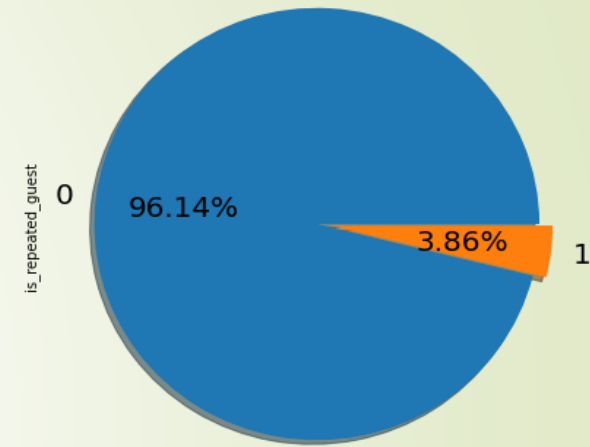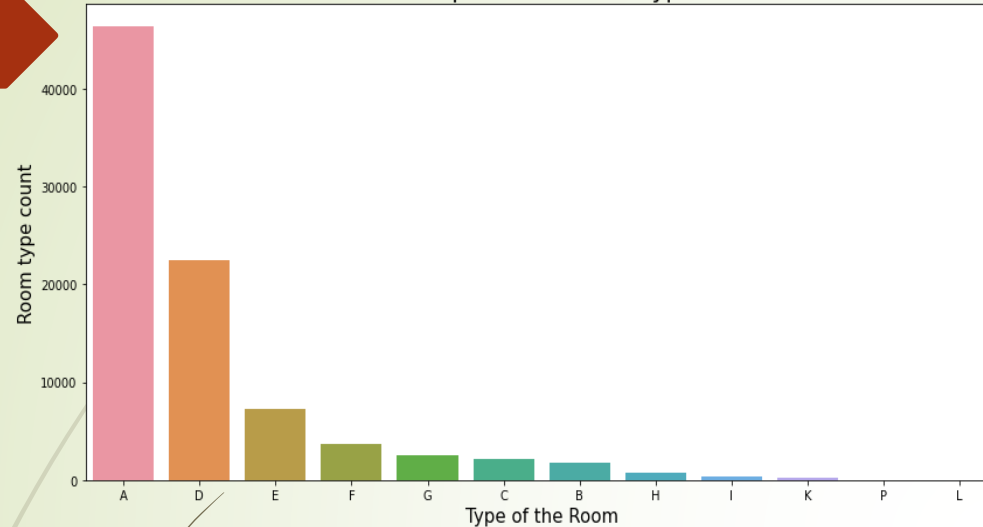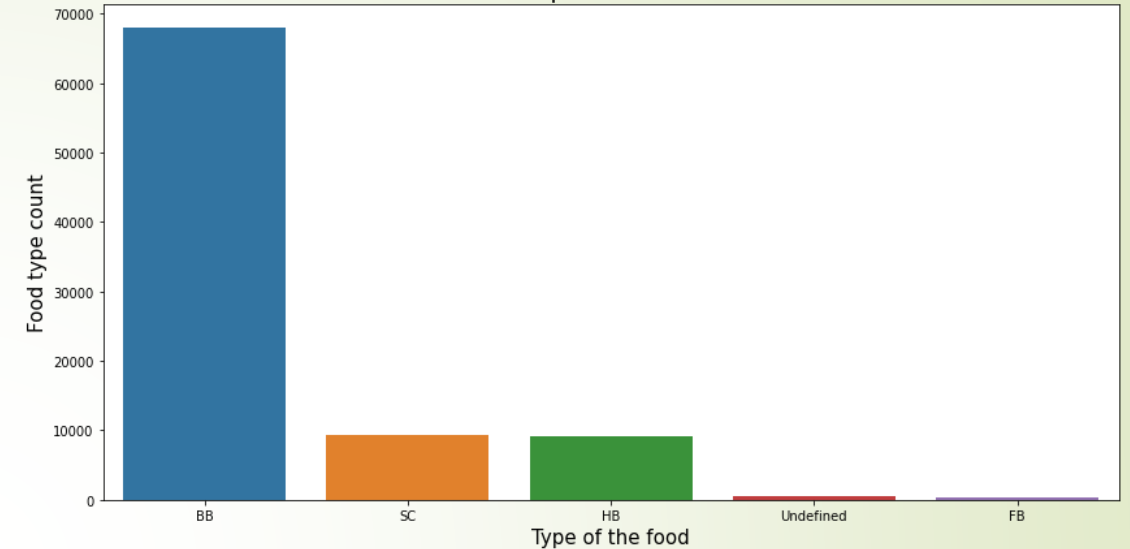# EDA: UNIVARIATE ANALYSIS

## Conclusion:

1) **City hotel** is the most preferred hotel and the percentage is **61.13%** means city hotel is the busier hotel type.

2) **Agent no. 9** made the most bookings and the count is **28759** means agent no 9 is the most preferred agent for booking.

3) Percentage of repeated guests is very less which is **3.86%** means repeated guests do not prefer the same hotel for their stay.
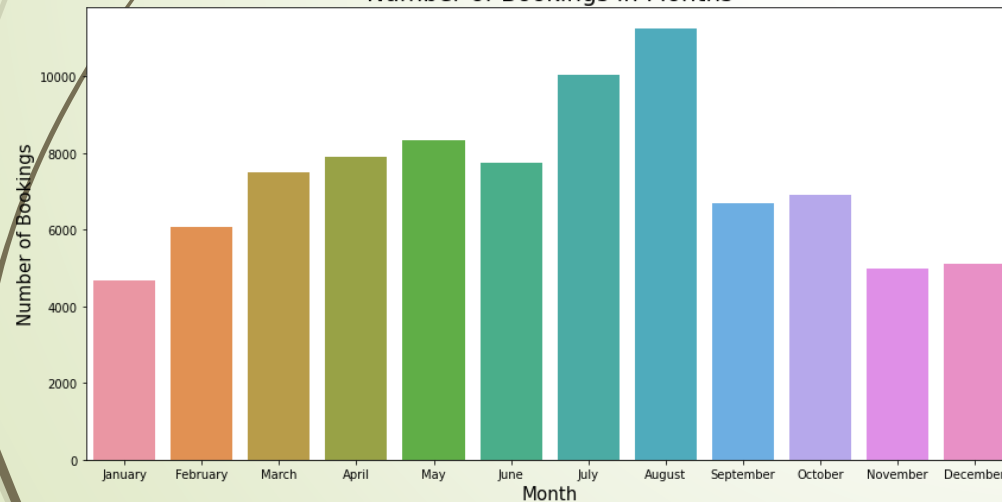
**AImaBetter**
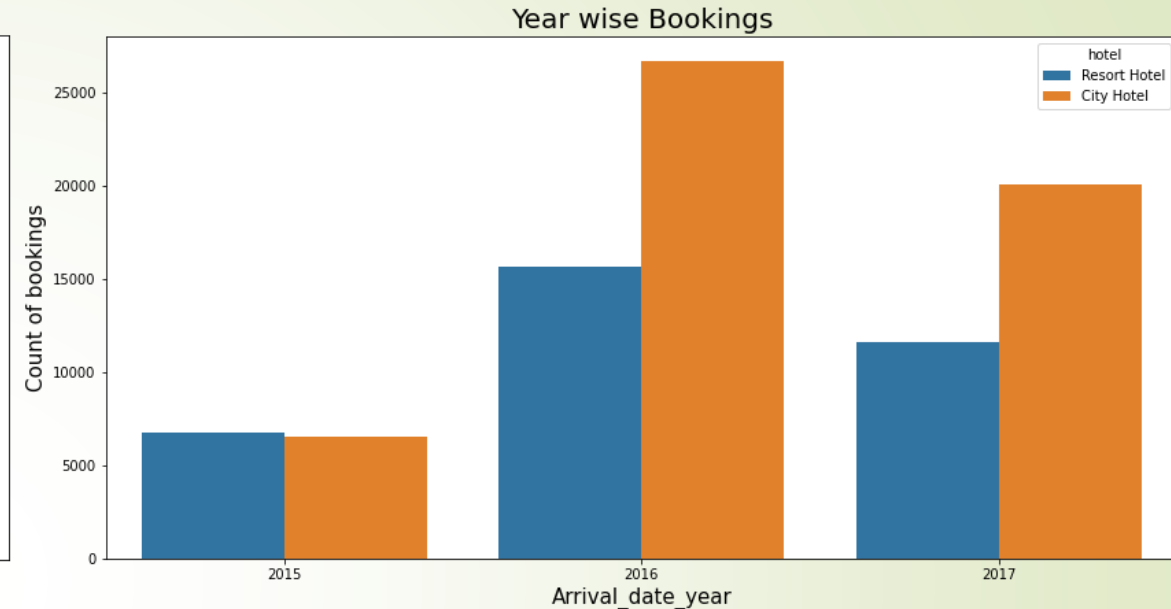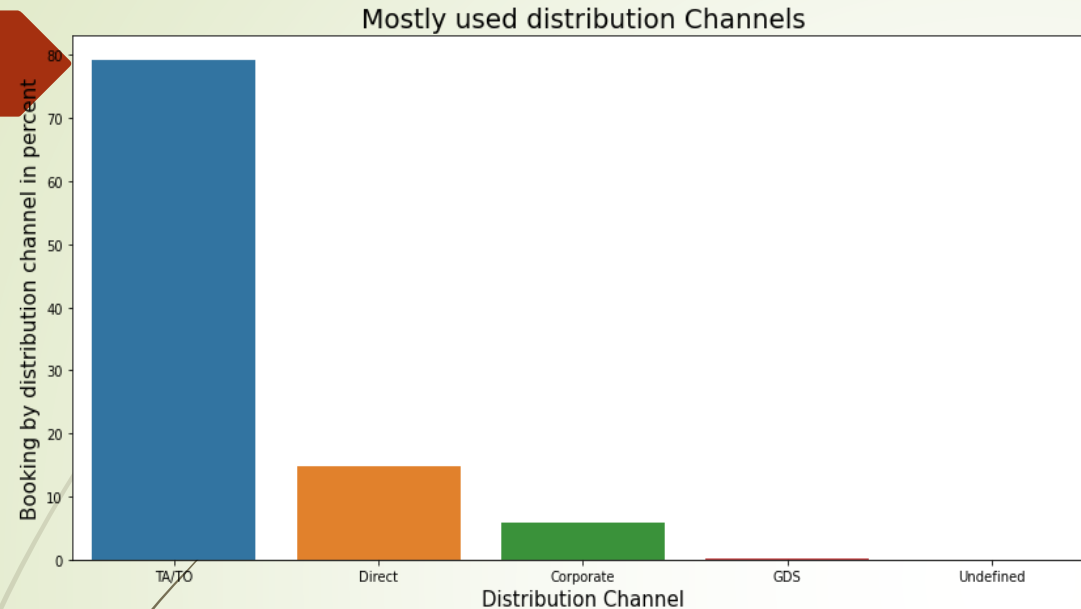

Most preferred Room type
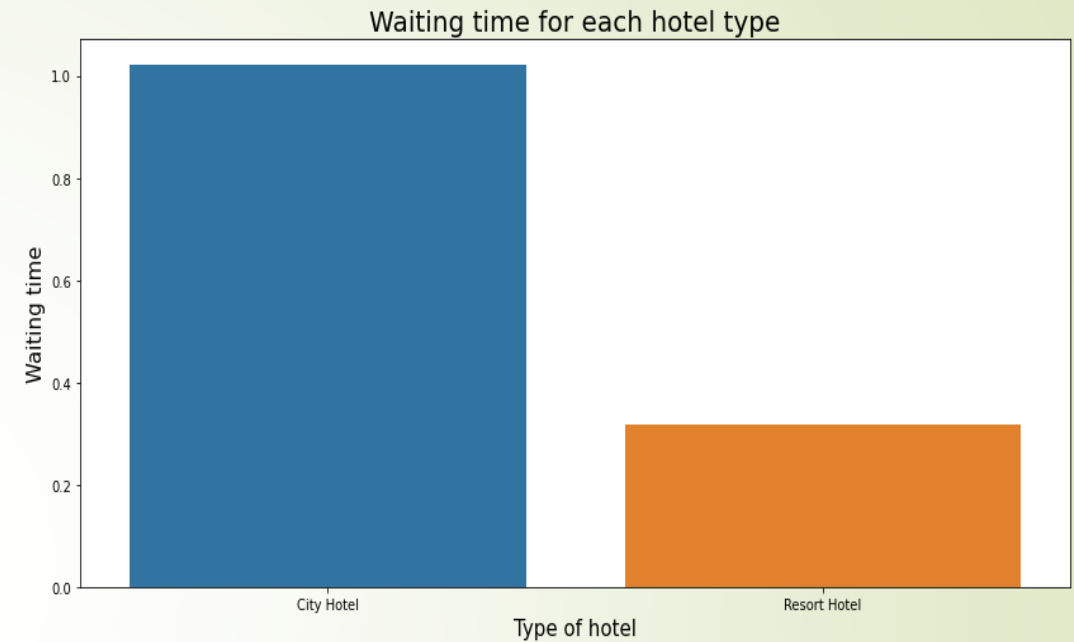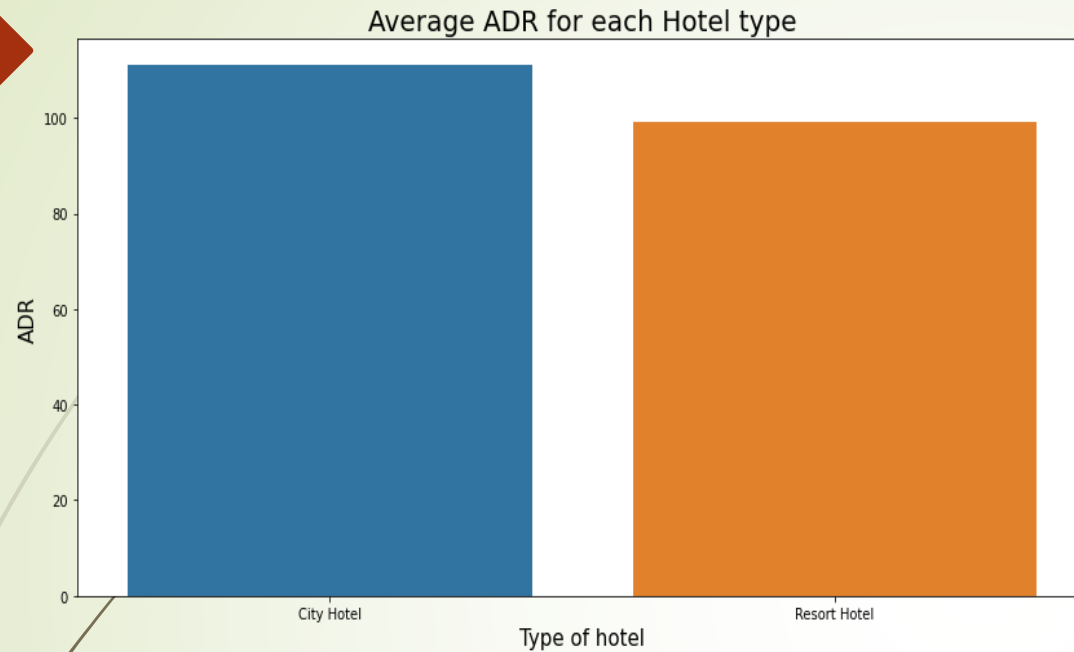

Most preferred Food


Number of Bookings in Months

# Conclusion:

1) **Room type A** is the most preferred room type and bookings are **46283**.
2) Most preferred food type is **BB type** and **67907** guests preferred BB type food.
3) **August** month has a maximum number of bookings and the count is **11242** means the august month is the busiest month in the year.
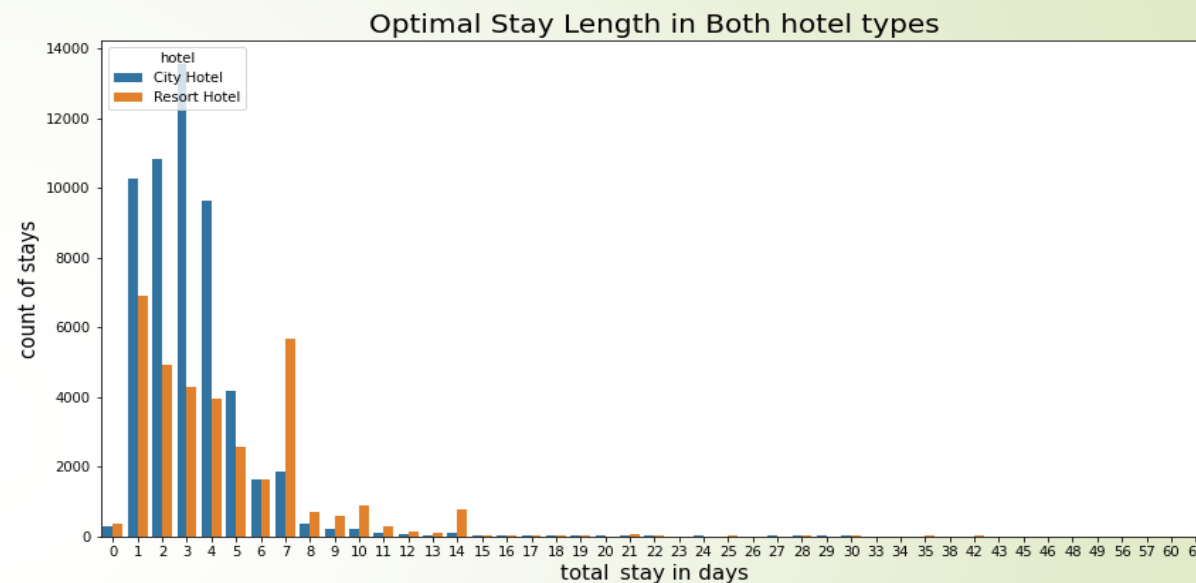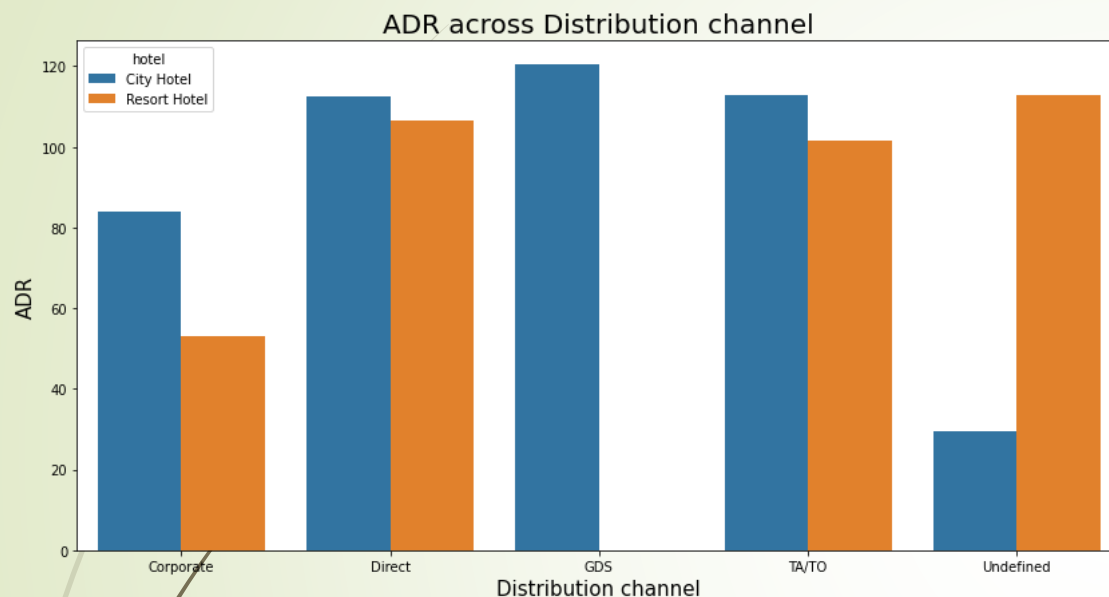
# Conclusion:

1) Mostly preferred booking channel is **TA/TO** which has a percentage of booking is **79.13%** and the second preferred channel is the **Direct** channel which made **14.85%** of bookings

2) The **Year 2016** had the highest booking in the city as well as in resort hotel type and the bookings are **42313** and the year **2015** had **fewer bookings**.

## Conclusion:

1) **City hotel** has the highest ADR and the average ADR for city hotels is **110.98**. High ADR means high revenue so the city hotel has high revenue.

2) **City hotel** has a long waiting time so the city hotel is the busier hotel.

## Conclusion:

1) **GDS distribution channel** contributed more to ADR in a city hotel and the **Direct & TA/TO** distribution channel has nearly equal contribution to ADR in both hotel types.

2) Optimal stay length in both hotel types is less than **7 days**.

**AImaBetter**



## Conclusion:

1) Repeated guests do not cancel their previous bookings but non-repeated guests cancel their bookings.

2) If the number of people increases ADR is increasing and due to this revenue also increases.

**AImaBetter**



Correlation of the columns

## Conclusion:

1) arrival_date_year and arrival_date week_number columns has na egative correlation which is -0.51.

2) Stays_in_week_nights and total_stay has a positive correlation which is 0.95.

# BUSINESS OBJECTIVE:

1) To increase hotel business some factors are important like high revenue, generation, customer satisfaction, facilities provided by the hotel, etc.

2) I am able to achieve the same things by showing the client which hotel is most preferred, the percentage of repeated guests, mostly preferred food by guests, then which hotel has the highest ADR, etc.

3) Most preferred room type is achieved by counterplot so the client can be well prepared in advance and this insight helps the client for further enhancement of their hospitality.

4) I am able to show which food type is mostly preferred so the client can offer the most preferred food to the guests.

5) Most preferred months are shown by barplot so the client can be well prepared in advance so that minimum grievances would be faced by the client.

6) Using barplot I am able to show which hotel type has high adr so the client can analyze which hotel has a high income.

7) I am able to show which hotel is the busiest hotel so the client can do relatable changes in facilities in less busy hotel types.

8) I am able to show the relationship between repeated guests and previous bookings not canceled so the client can prefer repeated guests.

9) Using barplot relationship between ADR and the total number of people is shown so the client can prefer the maximum number of people.

![AImaBetter]

# CONCLUSION:

1) City hotel is mostly preferred hotel by guests.

2) Agent no. 9 made the most bookings.

3) Percentage of repeated guests is less which is 3.86%.

4) Room type A is mostly the preferred room type.

5) Mostly preferred food type is BB type food.

6) August month has the most bookings and after august july has the most bookings.

7) TA/TO distribution channel is mostly used and the percentage age is 79.13%.

8) City hotel has the highest ADR. The highest ADR means more revenue.

9) 2016 year had the highest number of bookings and bookings were 42313.

10) City hotel has higher waiting time means city hotel is the busier hotel.

11) GDS distribution channel contributed most in ADR in city hotels but no contribution in the resort hotel.

12) Optimal stay length in both hotel types is less than 7 days.

13) Repeated guests do not cancel their bookings but not repeated guests cancel.

14) If the number of people is more then ADR also increases means revenue increases.

15) arrival_date_year and arrival_date_week_number columns have a negative correlation which is -0.51.

16) stays_in_weel_nights and total_stay columns have a positive correlation which is 0.95.