# Capstone Project-1

# Hotel Booking Analysis

## Team Members

**Shivangi Mishra**

**Saksham Tripathi**

**Bindu Kovvada**

**Deepak Kumar Gautam**

**Satyajit Sahoo**

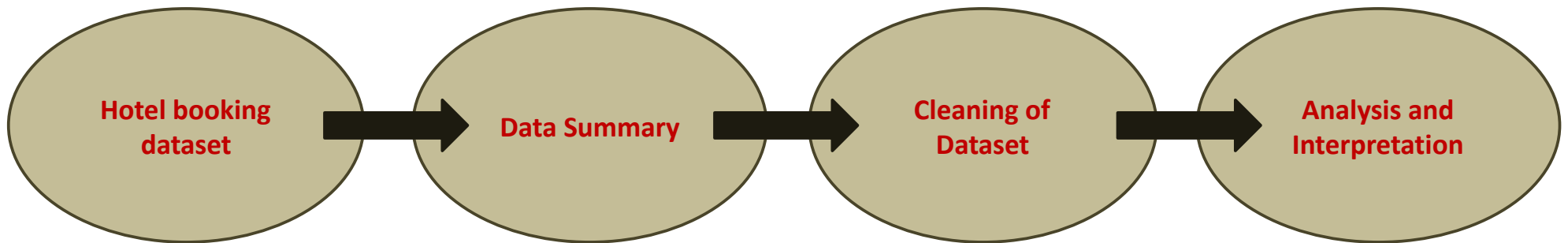# Flow of Presentation

- **Agenda**

- **Data Summary**

- **Cleaning of Dataset**

- **Data Visualization**

- **Inferences**

- **Conclusion**

# Agenda

Agenda is to discuss the given problem statement i.e., Analysis of several factors influencing hotel bookings in the given hotel booking dataset.

**APPROACH:**

# Data Summary

## Data set name:

Hotel Booking Database that contains booking information for a city hotel and a resort hotel of various countries from 2015 to 2017.

## Data shape:

Rows - 119390

Columns - 32

# Columns used:

- <u>hotel</u> - name of hotel whether City Hotel or Resort Hotel
- <u>is_canceled</u> –(0 or1) Indicates whether booking was cancelled or not.
- <u>lead_time</u> – The time between reservation and actual arrival.
- <u>arrival_date_month</u> - Month of arrival date
- <u>arrival_date_week_number</u> – Week number of year for arrival date
- <u>arrival_date_day_of_month</u> – Day of arrival date
- <u>stays_in_weekend_nights</u> – Number of weekend nights the guest stayed or booked to stay at the hotel.
- <u>Stays_in_week_nights</u> : Number of week nights the guest stayed or booked to stay at the hotel.

- <u>adults</u> : Total number of adults in hotel
- <u>children</u> : Total number of children in hotel
- <u>babies</u> : Total number of babies in hotel
- <u>meal</u> : Type of meal booked i.e. Bed & Breakfast (BB), Half Board (HB), Full Board (FB), Undefined contain no meal package
- <u>country</u> : Country of customers
- <u>market segment</u> : A group of people who share one or more common characteristics used for business
- <u>distribution channel</u> : Chain of business through which a service passes until it reaches the final buyer
- <u>previous cancellations</u> – (0 or 1)previous cancellation by customer
- <u>previous booking not canceled</u>- confirmed booked by customer

- <u>reserved_room_type</u> – Type of Room booked
- <u>assigned_room_type</u> – Type of Room assigned/alloted
- <u>booking_charges</u> – booking charges we charged
- <u>deposit_type</u> – No Deposit, Non Refund, Refundable
- <u>agent</u>- ID of travel agency
- <u>company</u>- ID of the company
- <u>days_in_waiting_list</u> – Number of days in waiting
- adr- Average daily rate
- <u>required_car_parking_space</u> – Number of car parking spaces required
- <u>total_of_special_request</u> – Number of special request requested
- <u>reservation_status</u> -  Canceled, Check out
- <u>reservation_status_date</u> – Canceled / check out
- <u>is_repeated_guest</u> –(0 or 1) contain data of repeated guest
- <u>customer_type</u> – Type of customer Contract / Group / Transient

# Cleaning of Dataset

**Dropping Duplicates**

In [ ]:
```python
# rows containing duplicate data
duplicate_rows_data = data[data.duplicated()]
print("Number of duplicate rows:", duplicate_rows_data.shape)
```

Number of duplicate rows: (31994, 32)

In [ ]:
```python
# Droping duplicate values
data.drop_duplicates(inplace=True)
```

In [ ]:
```python
data.shape
```

Out[ ]:
(87396, 32)

**Missing values**

In [ ]:
```python
# getting summary of missing values present in the dataset.
for column in data:
  if data[column].isnull().any():
    print('{0} column has {1} missing values, which are {2} % of total column'.format(column,data[column].isnull().sum(),round(data[column].isnull().s
```

children column has 4 missing values, which are 0.005 % of total column
country column has 452 missing values, which are 0.517 % of total column
agent column has 12193 missing values, which are 13.951 % of total column
company column has 82137 missing values, which are 93.983 % of total column

In [ ]:
```python
# as company column has 94% missing values so we are dropping that column
data=data.drop(['company'],axis=1)
```

# Cleaning of Dataset

```python
#filling null values
data['agent'].fillna(0, inplace = True)              # filling null values of agent with 0
data['children'].fillna(data['children'].mean(), inplace = True)   #filling null values of children with mean
data['country'].fillna('unknown', inplace = True)    # filling null values of country with 'unknown'
```

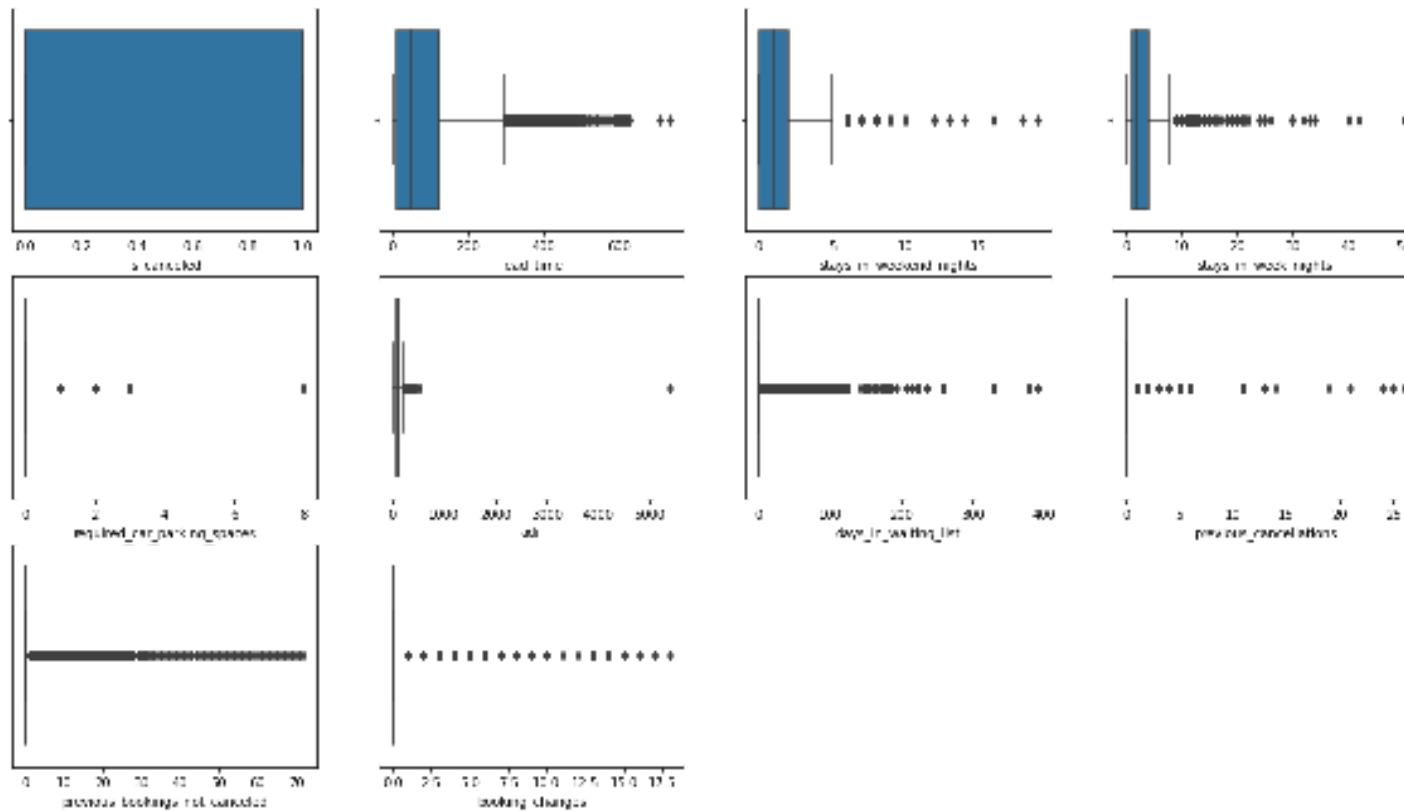**Converting columns to appropriate datatypes.**

```python
#now let's convert datatype of columns 'children' and 'agent' from float to int.
data[['children','agent']]=data[['children','agent']].astype(int)
```

```python
# changing datatype of column 'reservation_status_date' to date-type.
data['reservation_status_date'] = pd.to_datetime(data['reservation_status_date'], format = '%Y-%m-%d')
```

**Adding important columns.**

```python
#let's group some of the columns which can be useful in analysis as grouped element
data['total_stay']=data['stays_in_weekend_nights']+data['stays_in_week_nights']
data['total_guests']=data['children']+data['babies']+data['adults']
```

# Cleaning of Dataset



✔ **As we can see that this dataset has many outliers and thus we can infer that this dataset is not very reliable.**

# Data Visualization



✔ Lead time and total stay have slight correlation. This means when people want to stay little longer they plan little before than actual arrival.

✔ Total guests and average daily rate has some correlation. This means the when the number of guests increases, adr will also increase.

# Data Visualization

Following insights are pulled out from this analysis:

1. What is the percentage of bookings in each hotel type?

2. From which country most guests come?

3. Which is the busiest month for hotels?

4. Which room type is in most demand ?

5. Which room type generate highest adr?

6. Which meal type is most preferred meal of customers?

7. How many booking were cancelled?

8. Which type of customers are most repeated?

9.  Booking cancellation and Repeated guest

# Data Visualization

10. Checking whether not getting allotted the same room type as demanded is the cause of cancellation of bookings?

11. Does not getting same room affects the adr.?

12. Which is the most common channel for booking hotels?

13. Which Distribution channel has highest no. of days in waiting list?

14. Which distribution channel has highest cancellation percentage?

15. Which Market segment has highest no. of days in waiting list?

16. Which is preferred stay length in each hotel?

17. Does lead time have effect on cancellation?

# Data Visualization

❏ **Percentage of bookings in each hotel type**



✔ Around 60% bookings are for City Hotel and 40% are for Resort Hotel.

# Data Visualization

❑ **From which country most guests come?**



✔ **Portugal** is the country from where most guests come. Around 31.36 % of guests come from Portugal, followed by Great Britain and France.
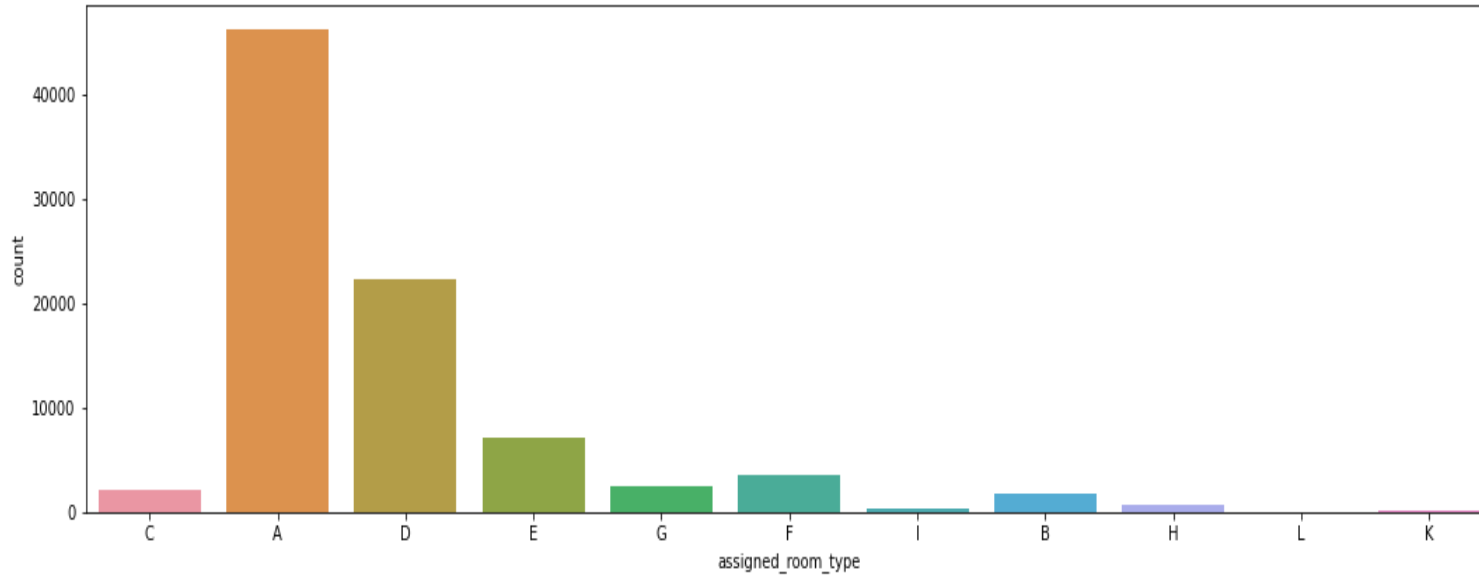
# Data Visualization

❑ **Which is the busiest month for hotels?**


Booking Trend (Monthly)

✔ As we can see most bookings were made from **July to August**. And the least bookings were made at the start and end of the year.

# Data Visualization

❏ **Which room type is in most demand ?**



✔ **Most demanded room** type is 'A'. Hotel should increase room type A to increase the revenue.

# Data Visualization

❑ **Which room type generate highest adr?**



ADR according to room type

✔ Although room type A was oh high demand and most booked, Highest revenue was produced by room type H followed by G. We can say that room type **H and G are premium rooms.**
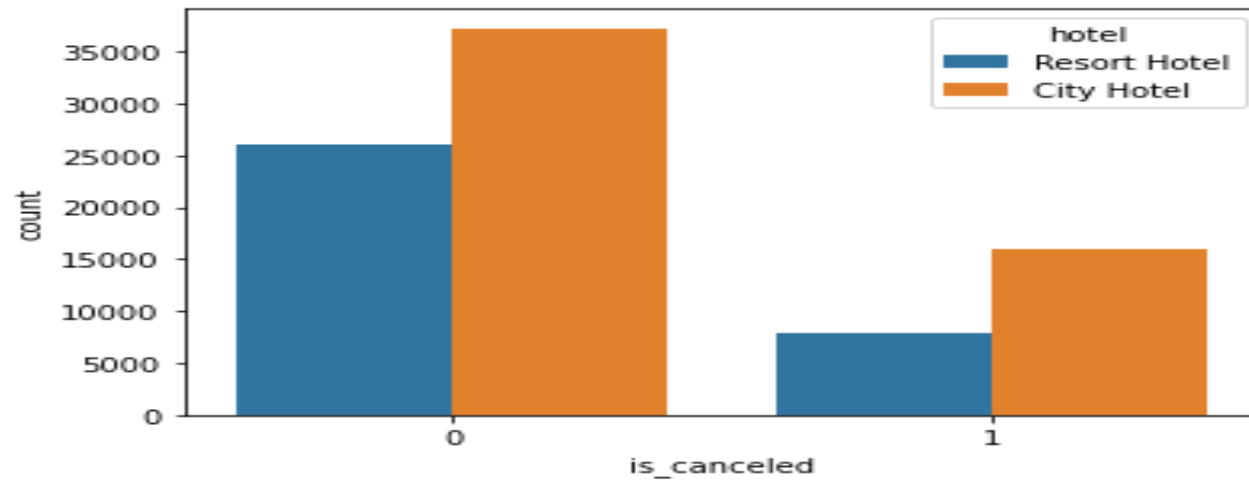
# Data Visualization

❑ **Which meal type is most preferred meal of customers?**



✔ Most preferred meal type is BB i.e., **Bed and breakfast**
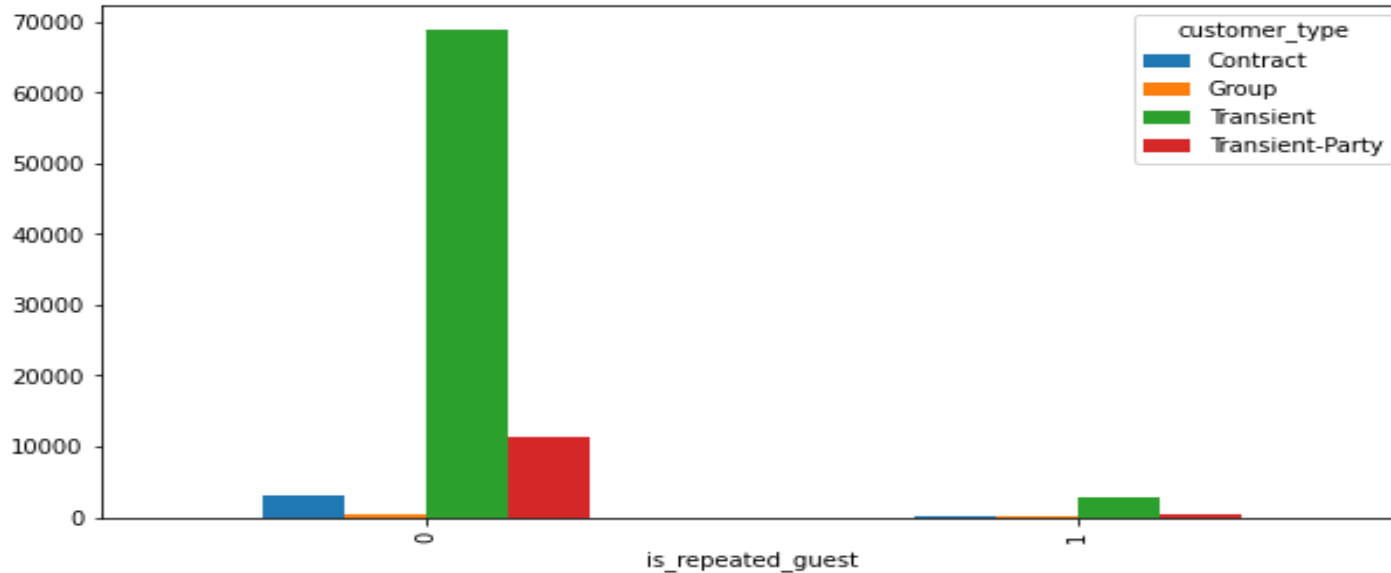
# Data Visualization

❑ **How many booking were cancelled?**



✔ Around 8000 **Resort Hotels** and 16000 **City Hotels** got cancelled.
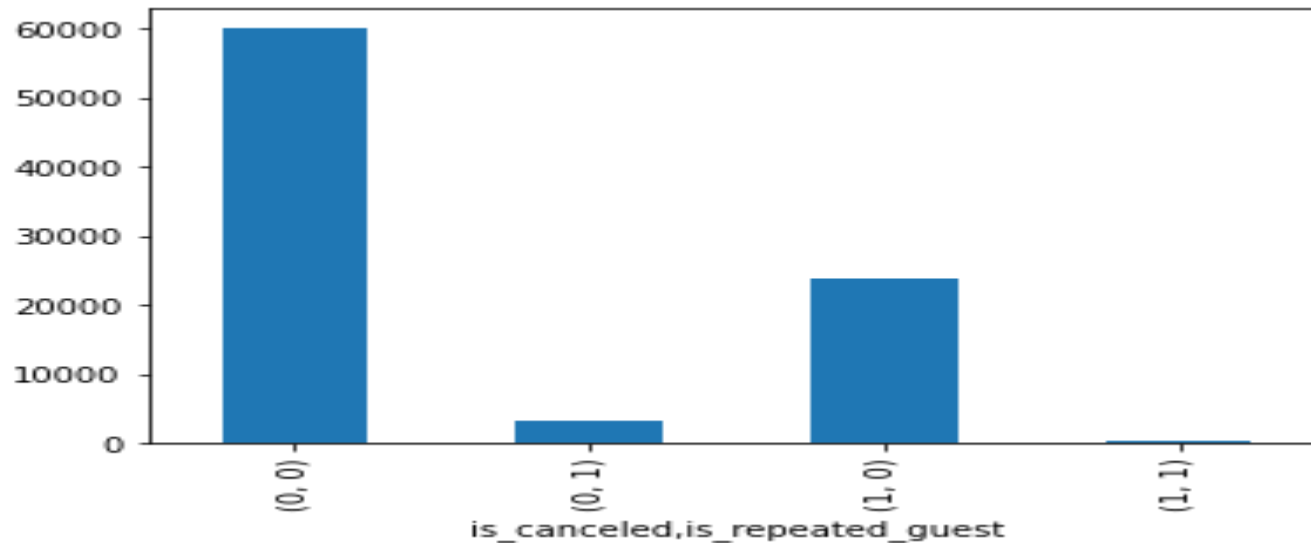
# Data Visualization

❑ **Which type of customers are most repeated?**



✔ Here, we can see that the maximum number of repeated guests are "**Transient type**" i.e., the "**Short-time customers**"
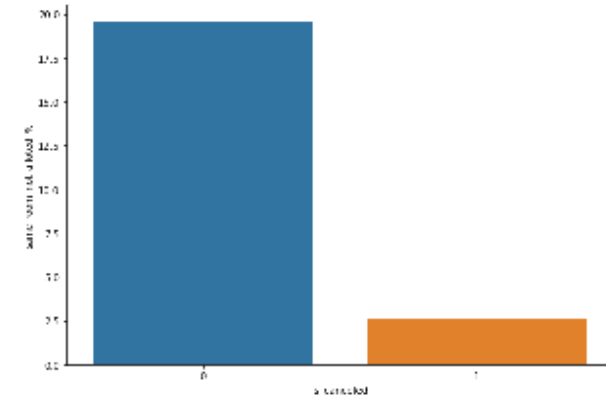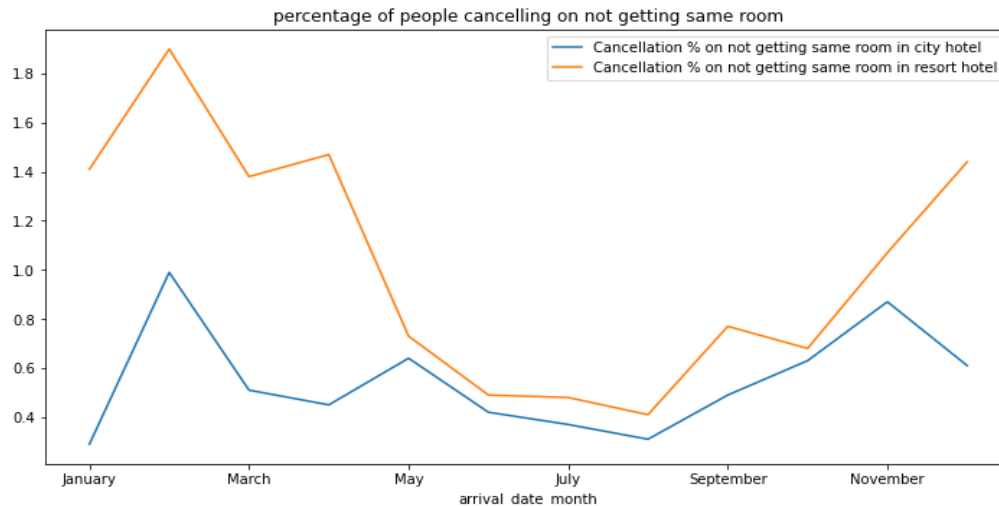
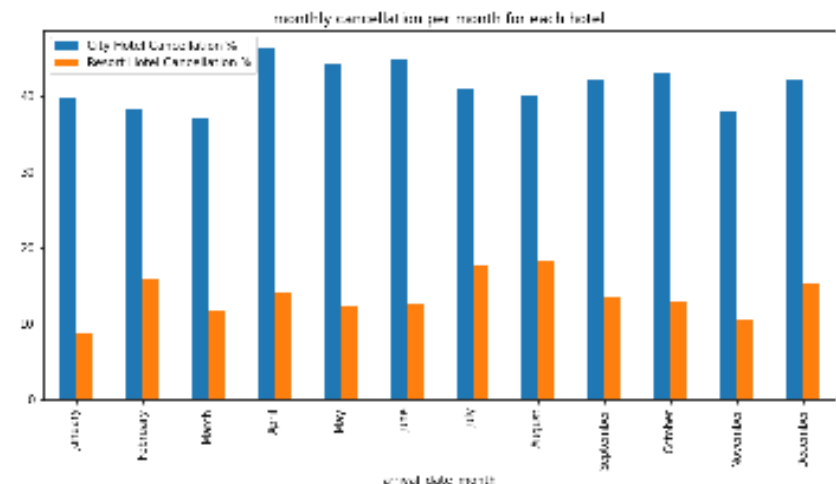# Data Visualization

❑ **Booking cancellation and Repeated guest**



✔ We see that when a hotel booking is cancelled and the customer is a repeated guest, the entries are almost '0', which means that **repeated guest is very less likely to cancel his booking with the hotel**.

# Data Visualization

❑ **Checking whether not getting allotted the same room type as demanded is the cause of cancellation of booking?**



percentage of people cancelling on not getting same room
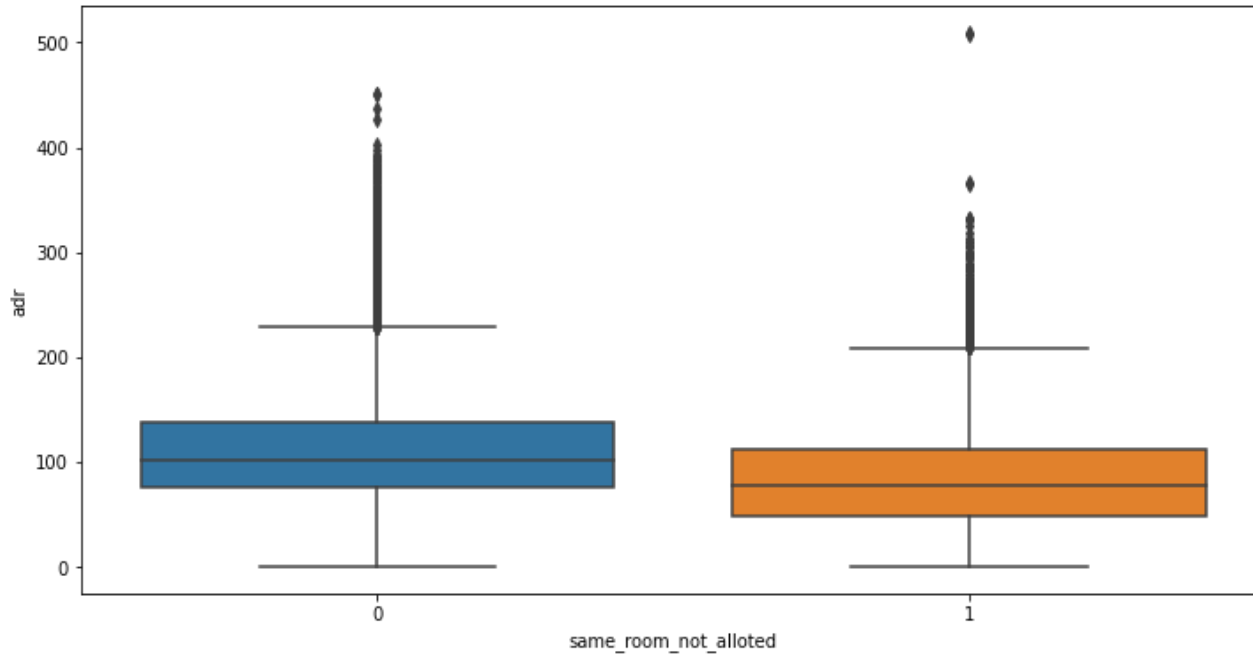


monthly cancellation per month for each hotel

✔ Here we see that not getting same room as demanded is not a case of room cancellation. A significant percentage of bookings are not cancelled even after getting different room as demanded.

✔ we see that the cancellation are very less for resort and city hotels, which is quite significant. **Guests are willing to keep their reservation even if they do not get the chosen accommodation.**
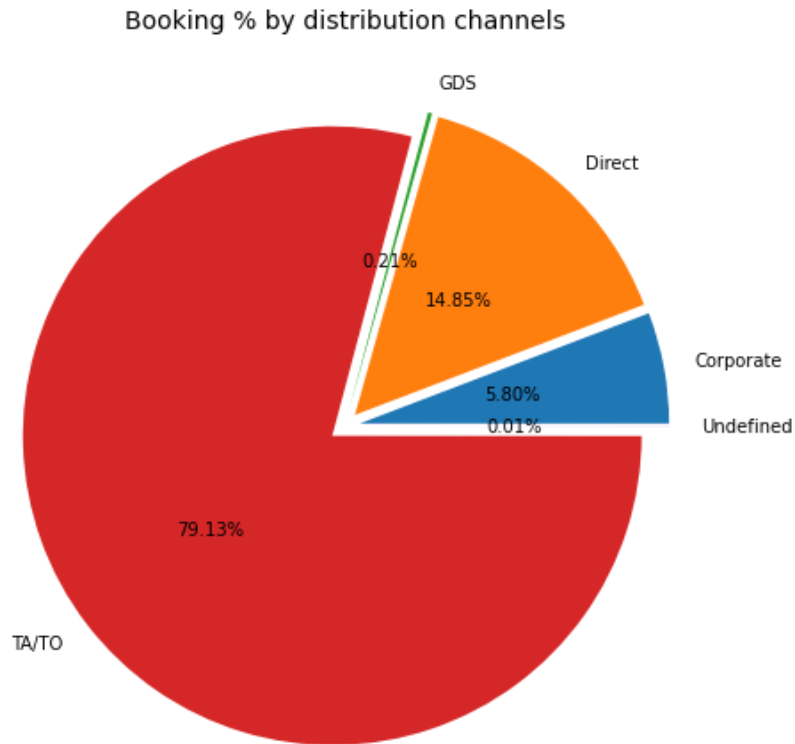
# Data Visualization

❑ **Does not getting same room affects the adr?**



✔ Not getting same room **do affects the adr**, if a guest does not get the desired room, he will pay a bit less.
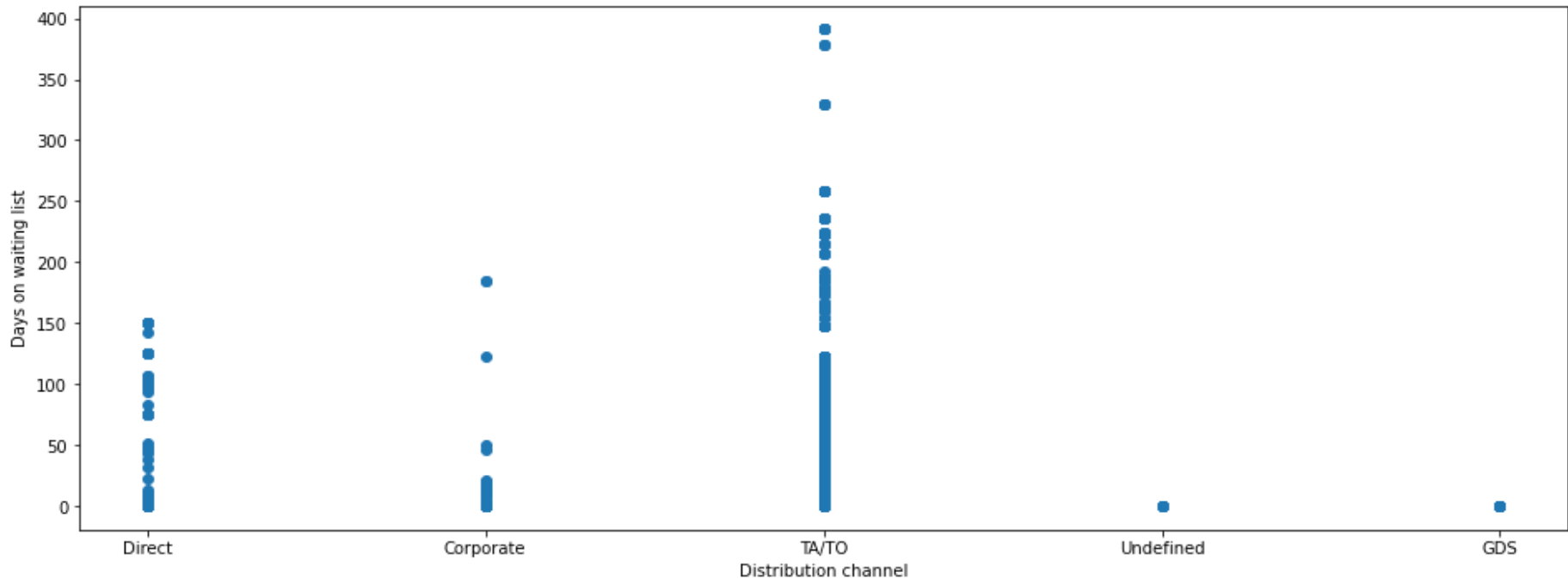
# Data Visualization

❑ **Which is the most common channel for booking hotels?**

Booking % by distribution channels



✔ Guests use different channels for making bookings out of which most preferred way is **TA/TO** (i.e., Travel Agents and Tour Operators)
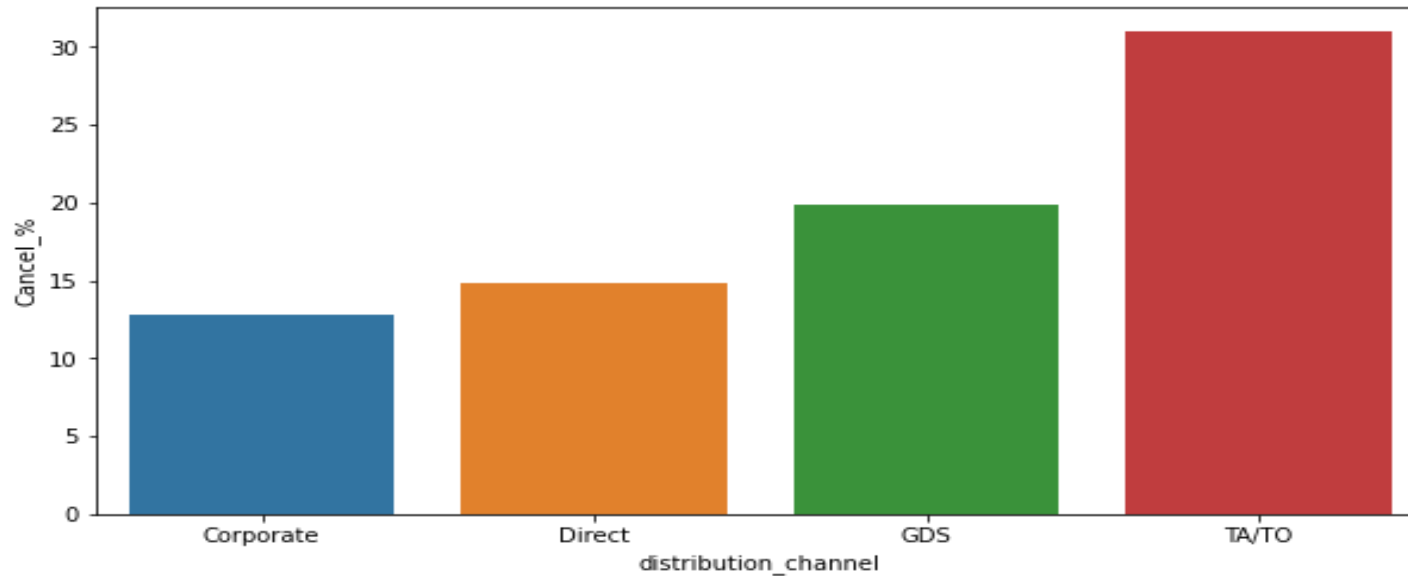
# Data Visualization

❑ **Which Distribution channel has highest no. of days in waiting list?**



✔ We see that the **'Travel Agent' and 'Tour Operators'** are the distribution channels for which the highest number of days are there on the waiting list.

# Data Visualization

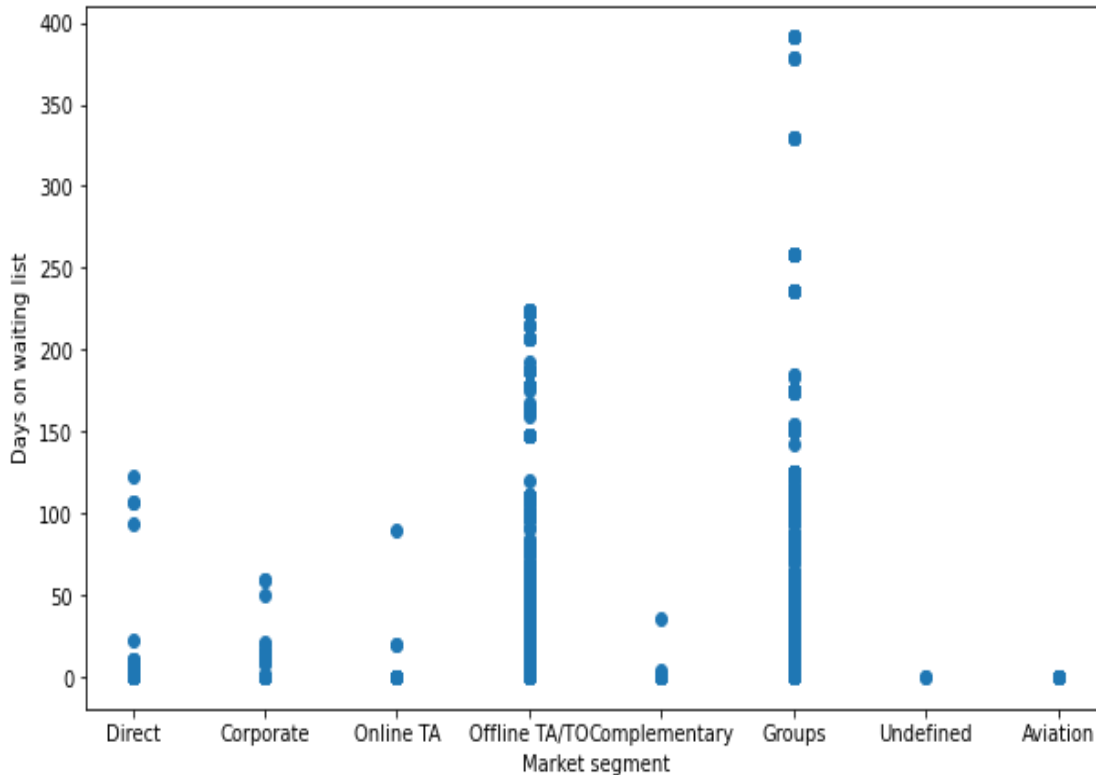❑ **Which distribution channel has highest cancellation percentage?**



✔ TA/TO has highest booking cancellation %. As a result, **booking made through TA/TO are 30% more likely to be cancelled.**
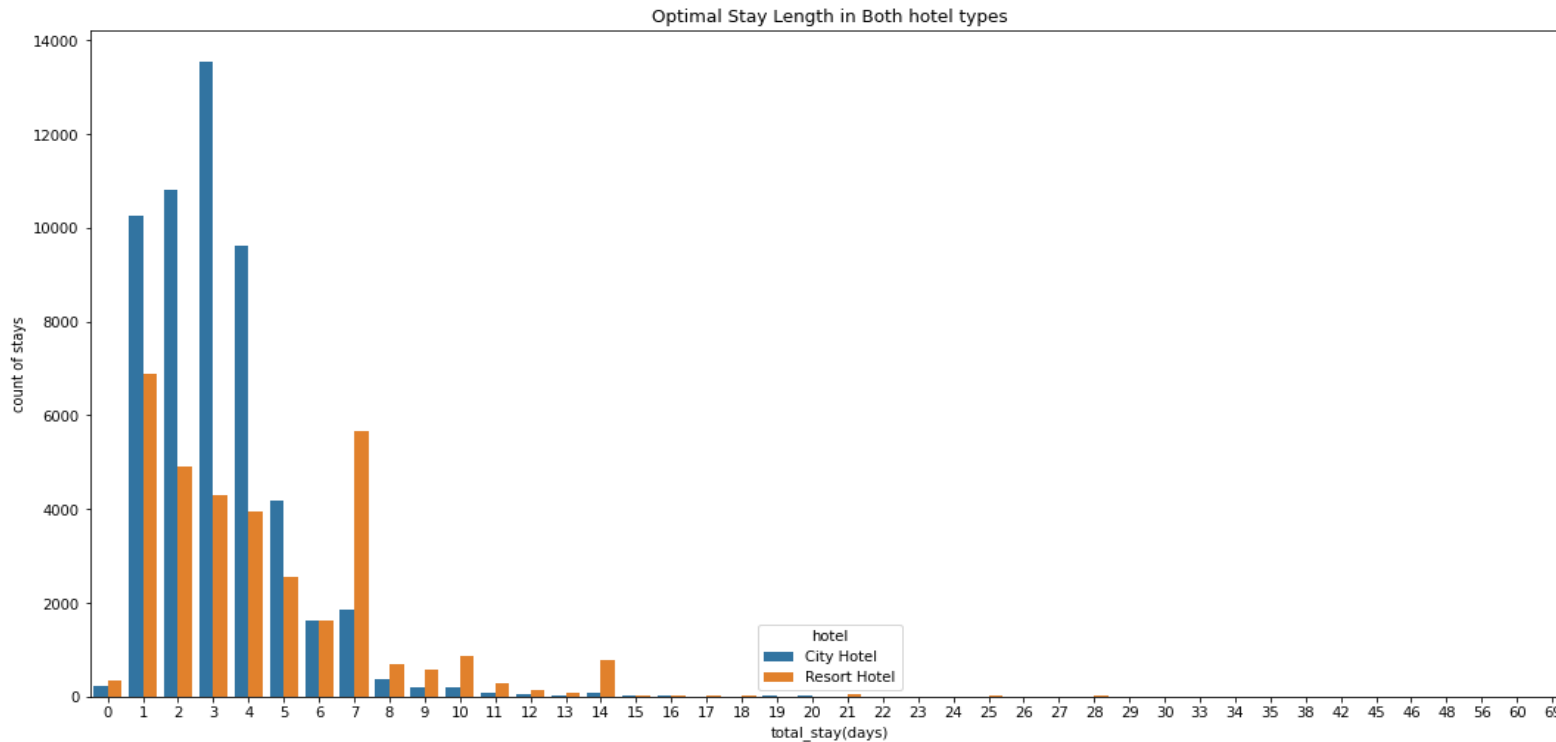
# Data Visualization

❑ **Which Market segment has least no. of days in waiting list?**



✔ Here, we see that **Aviation industry has the minimum number of days on the waiting list.** The reason for this could be that when a flight has to land at the location, it has to provide immediate accommodation to all of its working staff such as pilots and air hostages therefore they do not entertain hotels that have a long waiting list. So, in general, the hotel management sees to it that their needs are satisfied immediately and that they have almost no days on the waiting list.

# Data Visualization

❑ **Which is preferred stay length in each hotel?**



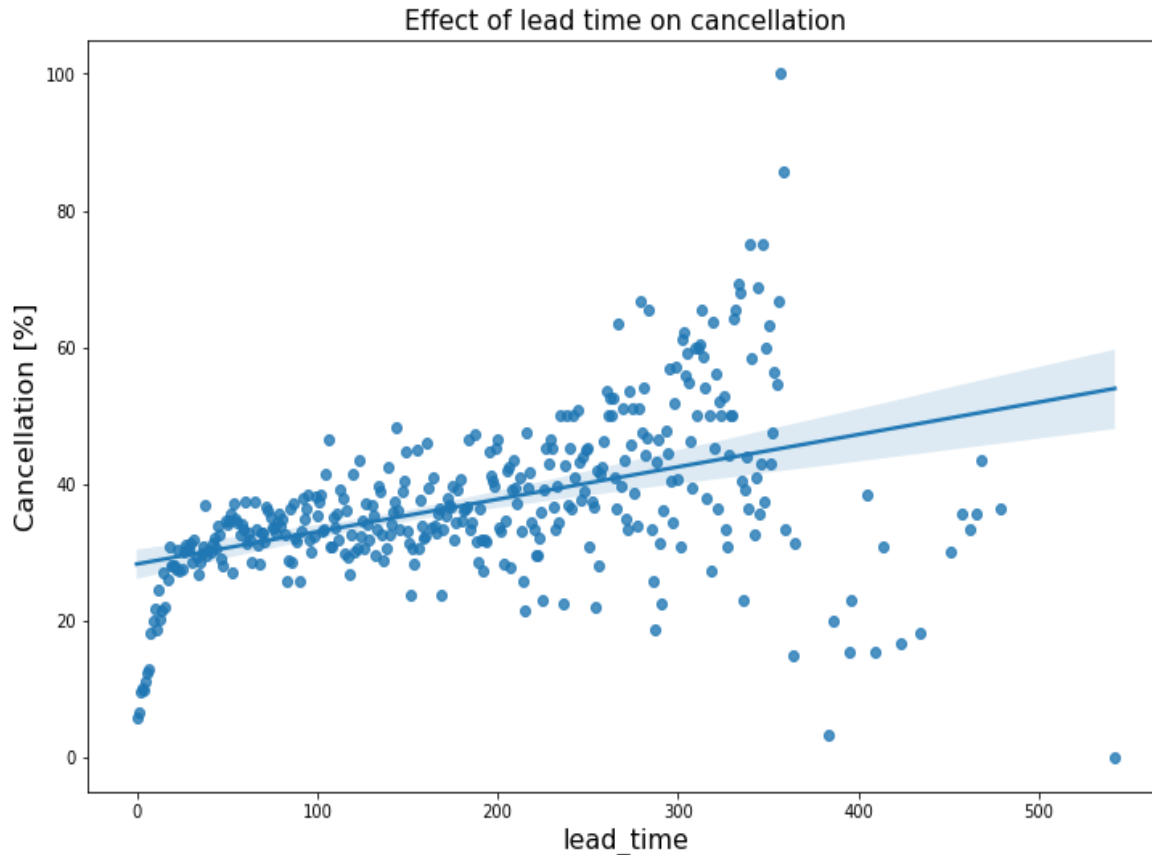Optimal Stay Length in Both hotel types

✔ Most people prefer to stay in hotels **for at least 5 days in both types of hotel**. However, people tend to stay longer at resort hotels.
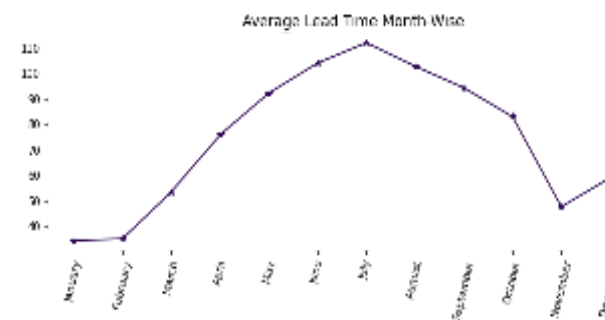
# Data Visualization
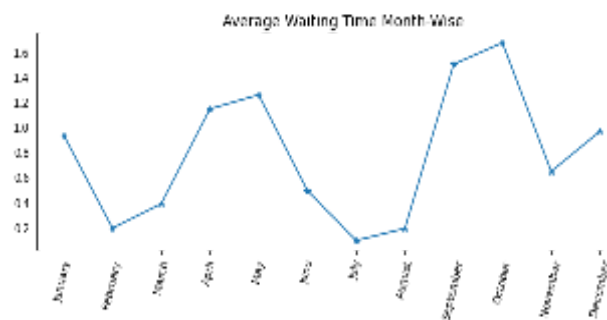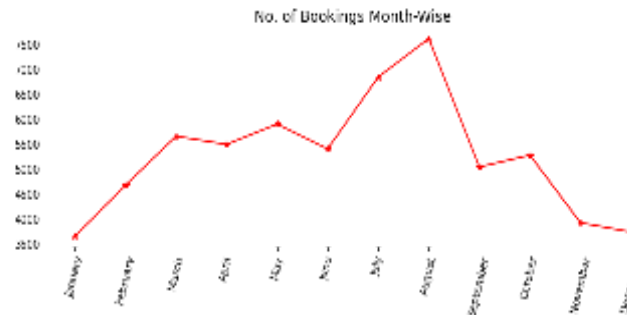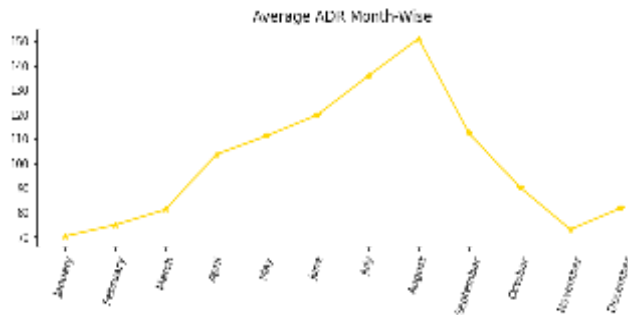
☐ **Does lead time have effect on cancellation?**



Effect of lead time on cancellation

✔ We see that the Effect of Lead time on Cancelation is **POSITIVE** which means that as the lead time increases, so does the number of cancelation.

# Data Visualization

❑ **Month-wise Analysis**


Average ADR Month-Wise


No. of Bookings Month-Wise


Average Waiting Time Month-Wise


Average Lead Time Month-Wise

When compared to previous months, we can see that the price and number of bookings are reduced from October to February. During this time, the average lead time is also very low.

The only concern is the waiting time, which is slightly longer from September to October, whereas the waiting time is the lowest and the price is the highest in July and August.

So the best time to book a hotel would be between October to February to get cheaper price and some privacy.

# Inferences

1. The most popular hotel type is **'City hotel'** which is preferred by **61.07% of total visitors.**

2. **96.15% of guests were not repeated guests**.

3. **27.08%** customers have **cancelled their booking**.

4. Market was widely captured by distribution channels such as TA or TO, especially **"Online TA" which captured 59% percent of the market.**

5. Around **53%** customers were assigned A room type and **65%** has reserved it. This also indicates the fact that **room type A is in higher demand.**

6. Almost **82.4% of customers are transient,** with very few booking as a group.

# Conclusion

1. **Higher lead time** has higher chance of **cancellation.**

2. **July-August are the busiest** and most profitable months for both the hotels.

3. The best time to book a hotel would be between October to February to get for cheaper price.

4. Most people prefer to **stay in the hotels for at least five days** in both types of hotel. However, people tend to stay longer at resort hotels.

5. Maximum number of **repeated guests are "Transient type"** i.e., the "Short-time customers"

6. City hotels receive around 60% of bookings, while Resort hotels receive 40%, hence city hotels are busier than Resort hotels.

7. **City hotel's total adr is slightly higher than Resort hotel's.**

8. Most of the guests were from European countries, with the **most coming from Portugal** followed by Great Britain and France.

9. Guests use different channels for making bookings out of which most preferred way is **TA/TO.**

10. Almost 30% of bookings via TA/TO are cancelled.

11. Bookings are not affected if they do not receive the same room as reserved. Although different room allotment does lessen the adr.