# Capstone Project
## Hotel Booking Analysis

By Ayush Kumar

# Points for Discussion

- Data Summary
- Data Preparation & Cleaning
- Hotel Types
- Cancellations
- Cancellations by Deposit Type
- Lead Time & Cancellations
- Bookings & Cancellations by Months
- Avg. Daily Rate over Months
- Bookings by Country
- Bookings & Cancellations by Market Segments
- Correlation of features

- Distribution of ADR
- Car Parking Space
- Room Types
- Top Agents
- Special Request by Country
- Lead Time
- Bookings by Distribution Channel
- Meals
- Previous cancellations & Cancellations
- Repeated Guests
- Conclusions

**AI**

# Data Summary

The **Hotel Booking Dataset** contains booking information for city and resort type hotels mainly in the years 2015, 2016 and 2017. It includes information like when the booking was made, length of stay, the number of adults, children and babies, and the number of available parking spaces, meals etc. All personally identifying information has been removed from the data. This dataset contains more than 119,000 rows of data.

# Data Preparation & Cleaning

- There were 31994 duplicate data, so deleted those rows.

- 4 columns (*company*, *agent*, *country* & *children*) had missing values.

- There may be some cases when customer didn't booked hotel via any agent or via any company. Hence we will replace null values by 0 in these columns.

```
company              82137
agent                12193
country                452
children                 4
reserved_room_type       0
assigned_room_type       0
booking_changes          0
```

- children column is numeric and skewed, hence choose median for imputing missing values.

- country column is categorical column. We use mode of *'country'* column, But it can lead to bias towards a specific country that occurs most frequently in the data. Hence I created a new category *'others'* for missing values.

# Data Preparation & Cleaning

- There were 166 rows with 0 adults, children and babies which seems unlikely hence dropped them.

- Created new column for *total_guest = (adults + children)*. Ignored babies because generally they are not charged.

- Some values in 'adr' was negative, which must be an error so imputed them with median since it was a skewed data.

- Created new column for *total_stays_night = (weekend_nights + week_nights)* to analyze average length of stay.

- Assigned appropriate data types for some columns.

# Hotel Types



Total bookings by hotel types

There are two types of hotels:

- City Hotel

- Resort Hotel

There are around 52k bookings in City Hotels which is 61% of total bookings in the dataset, similarly around 34k bookings (39%) in Resort Hotels.
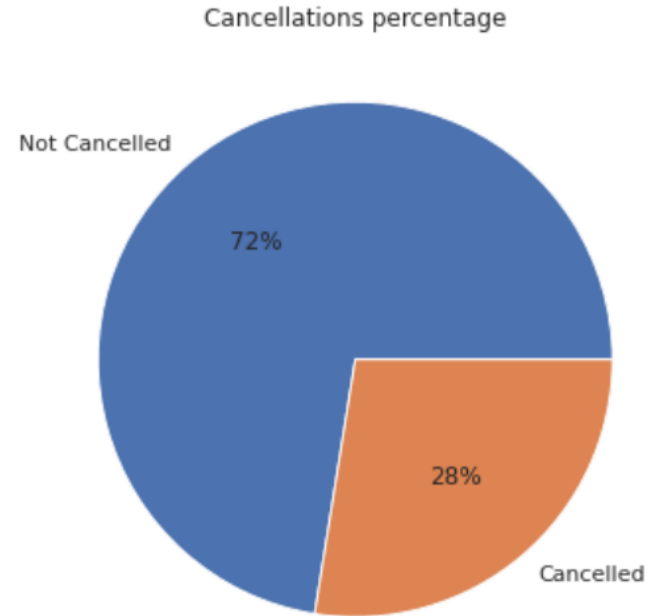


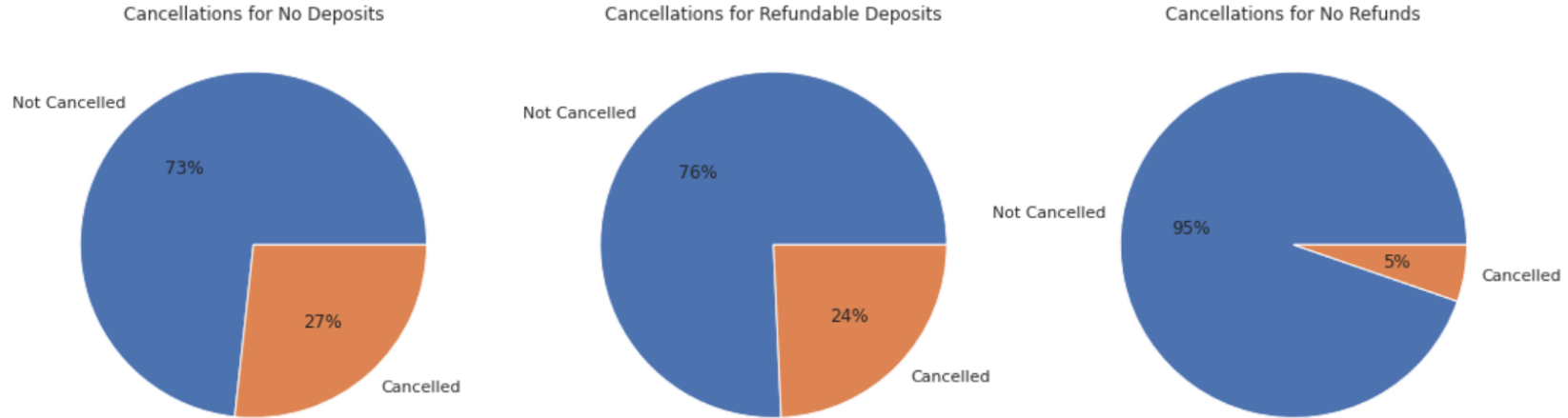Percentage shares of bookings by hotel types

# Cancellations

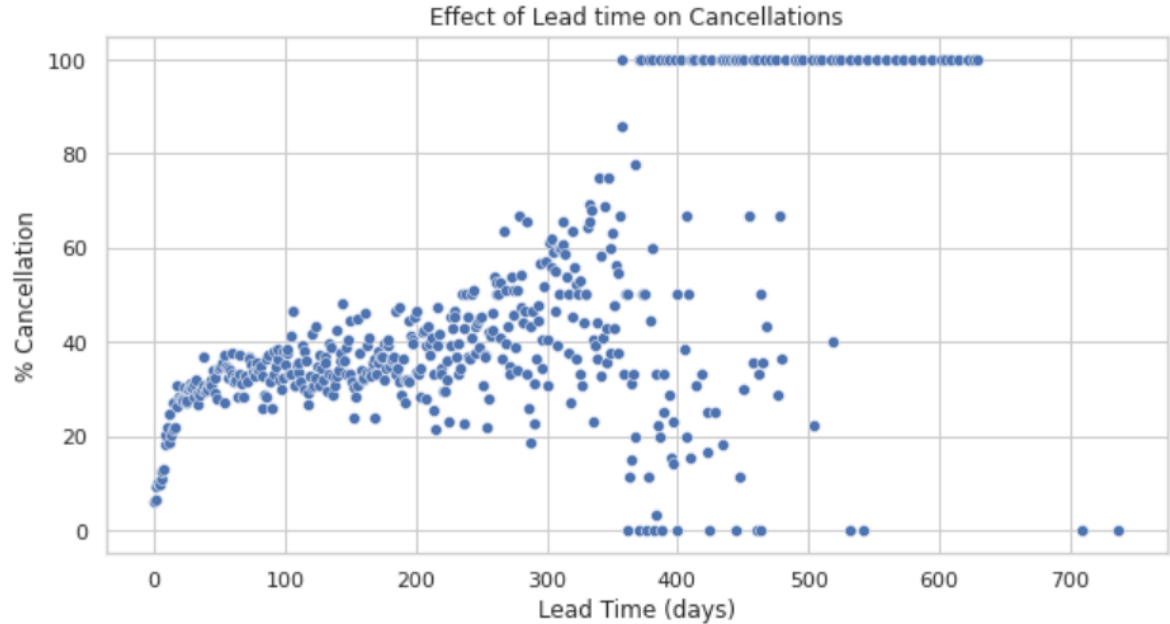Out of all the bookings, 28% of them was cancelled and 72% was not cancelled.



Cancellations percentage

# Cancellations by Deposit Type



Cancellations for No Deposits

Not Cancelled 73%

Cancelled 27%

Cancellations for Refundable Deposits

Not Cancelled 76%

Cancelled 24%

Cancellations for No Refunds

Not Cancelled 95%

Cancelled 5%

Total cancellation percentage : 28%

Around 1/4ᵗʰ of bookings were cancelled where customer will get the money back. Only 5% cancels their bookings when they can't get their money back.

# Lead Time vs Cancellations

Bookings having less lead time are less often cancelled whereas high lead time bookings are more often cancelled.
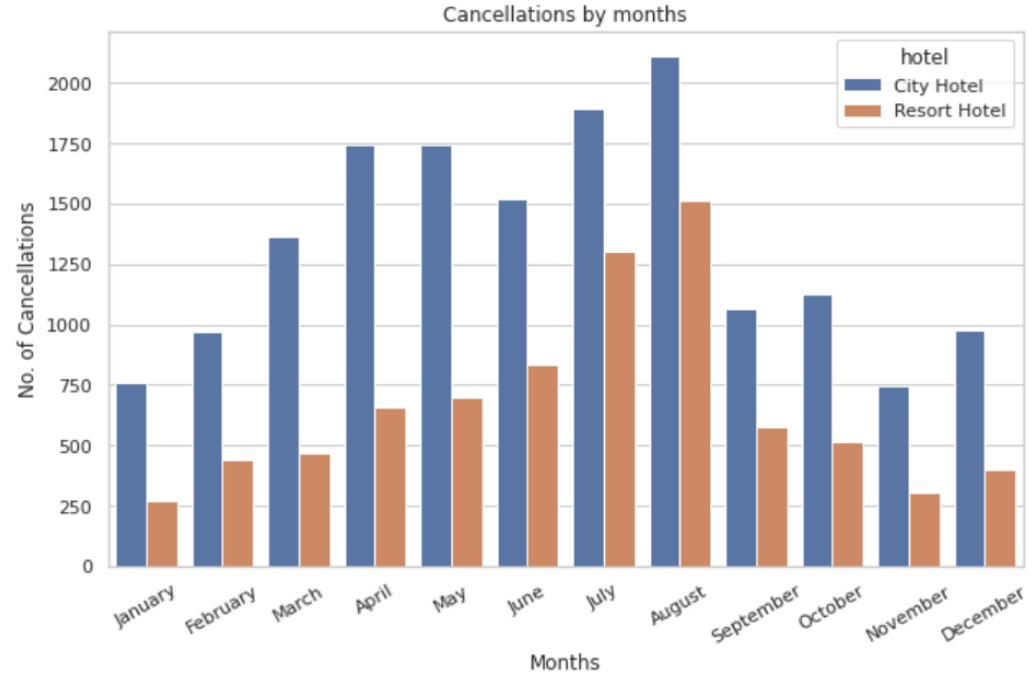


Effect of Lead time on Cancellations

# Monthly Bookings

For both hotel types, winter months have less bookings compared to mid-year (or summer) months.
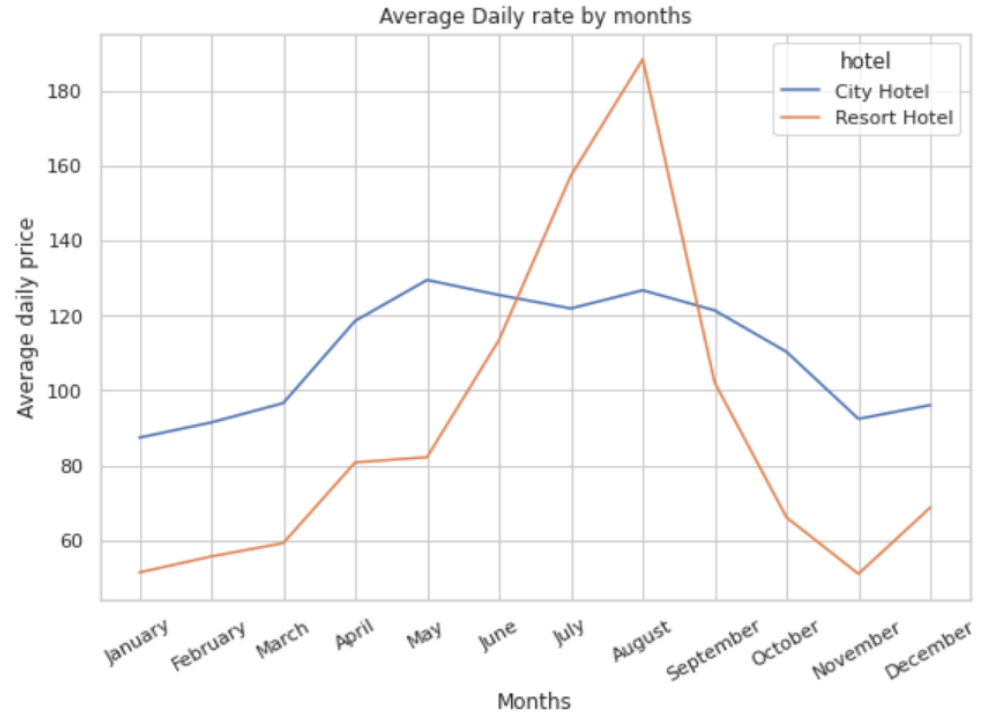


Bookings by months

# Monthly Cancellations

Graphs for total bookings and total cancellations for both hotel types are similar which means with more number of bookings, more are the cancellations.
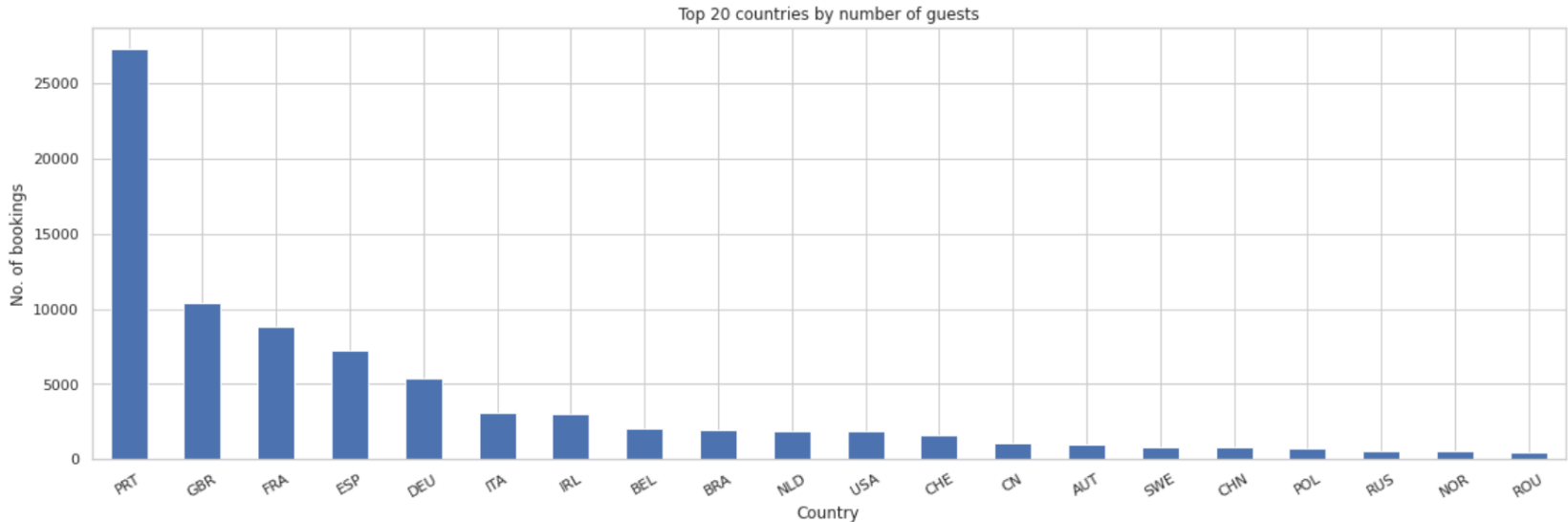


Cancellations by months

# Daily Rates by Months

adr is higher in summer months and lower in winter months, and we saw the similar pattern in total bookings over months. Simple market principle i.e., more demand more price and less demand less price.
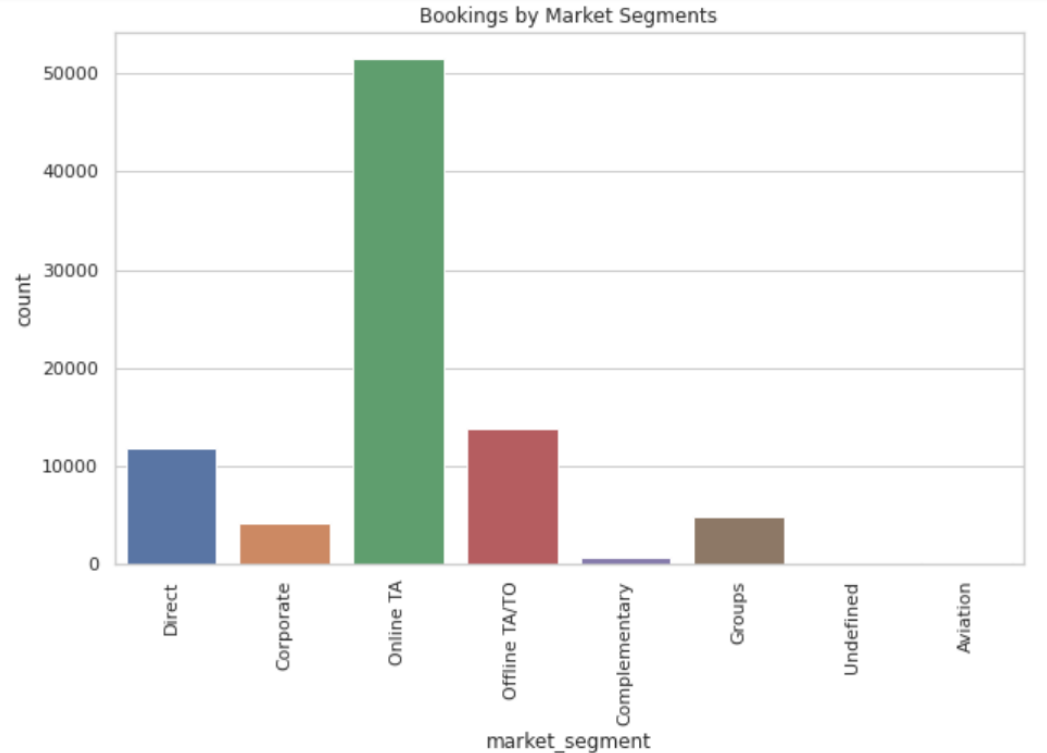


Average Daily rate by months

# Bookings by Country

Portugal has significantly higher number of bookings compared to other countries. Next major countries are GBR, FRA, ESP etc. Most of the top countries are European countries.



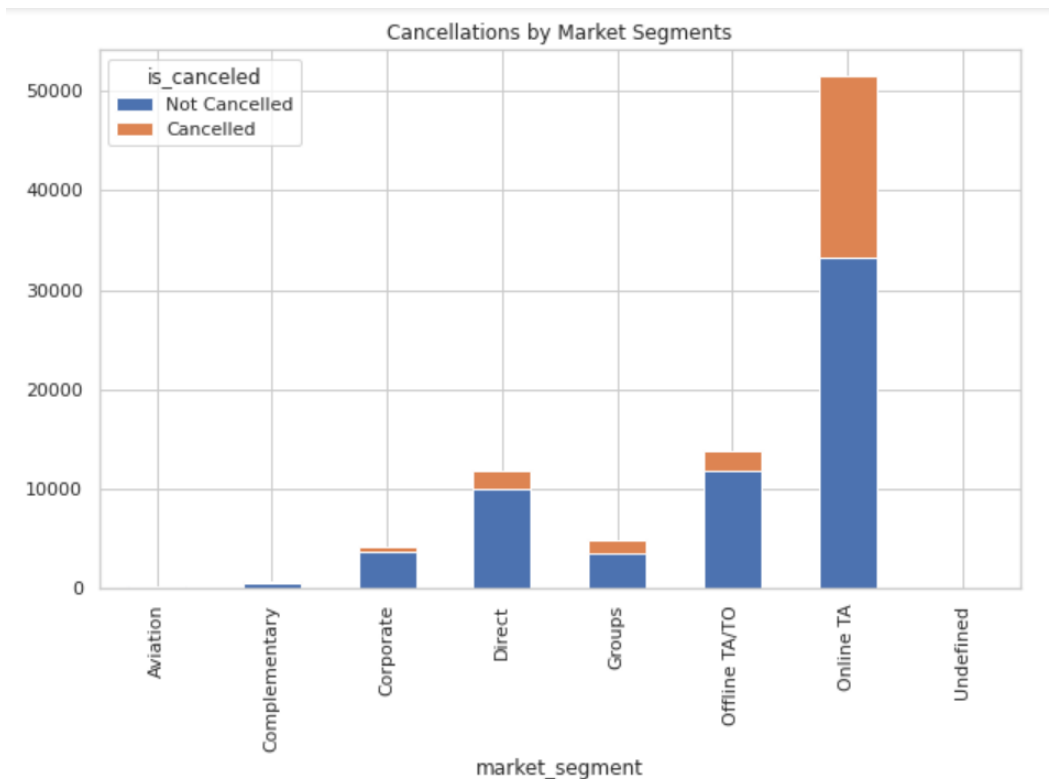Top 20 countries by number of guests

# Bookings by Market Segments

Indirect bookings through travel agents (specially Online) are higher compared to other bookings like direct, corporate, group etc.



Bookings by Market Segments

# Cancellations by Market Segments

We see that Online TA has highest proportion of cancellations (around 35%) but at the same time number of non cancelled bookings are also the highest for Online TA. Cancellation proportions for other segments are less.
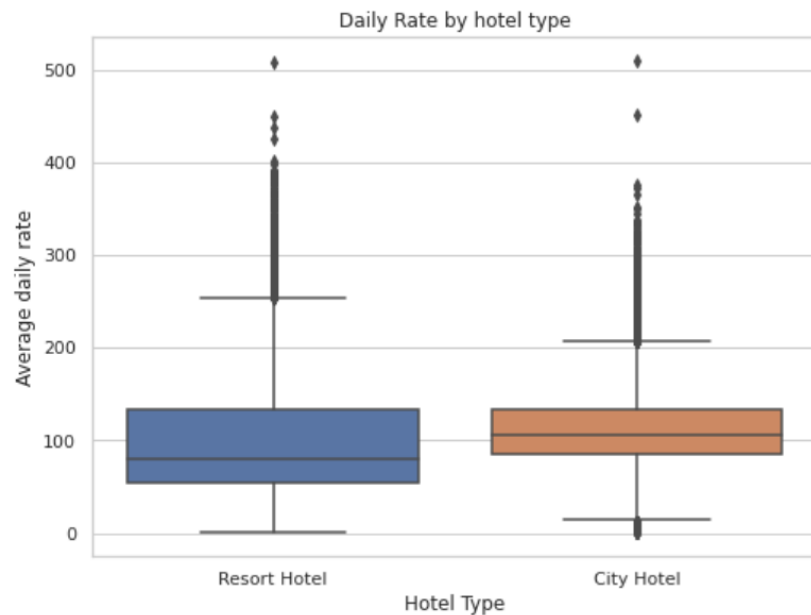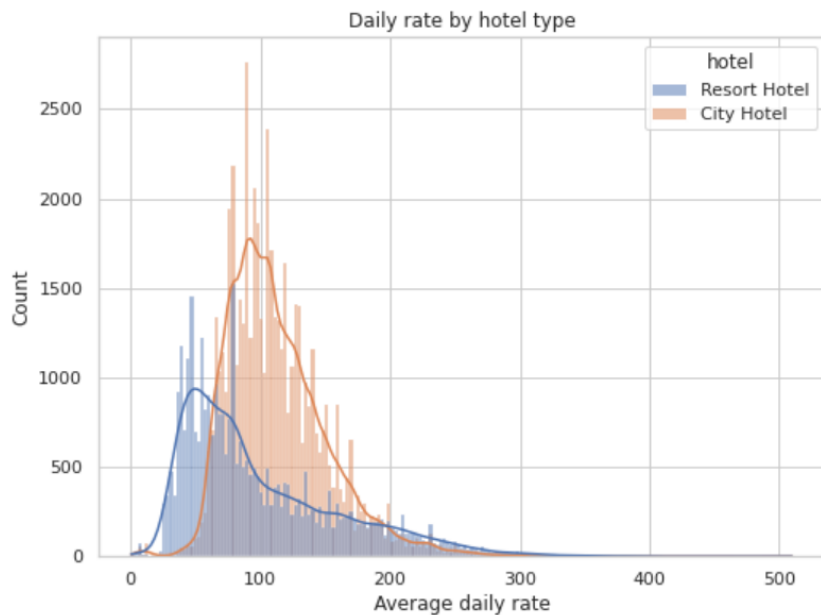


Cancellations by Market Segments

# Correlation of features

Lead time and total stays night have little correlation. This may indicates that people who wants to stay longer generally plans early hence longer lead time.

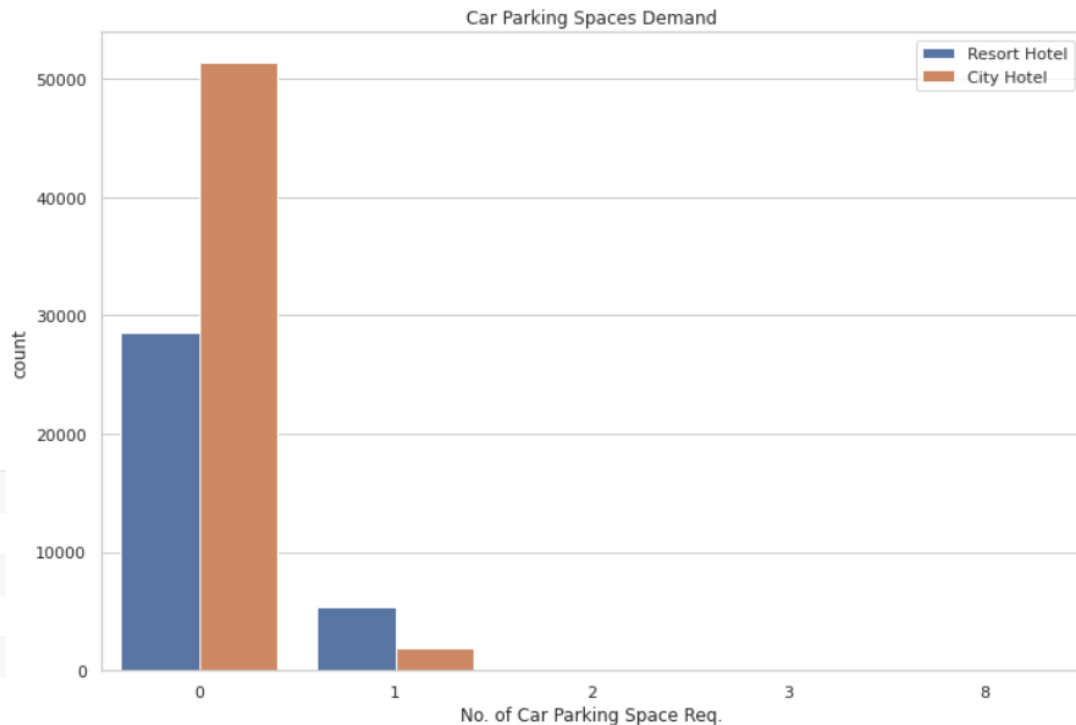adr and total guests are also slightly correlated, which makes sense as more number of people means more cost/charges.

# Daily Rate Distribution

# Car Parking Space

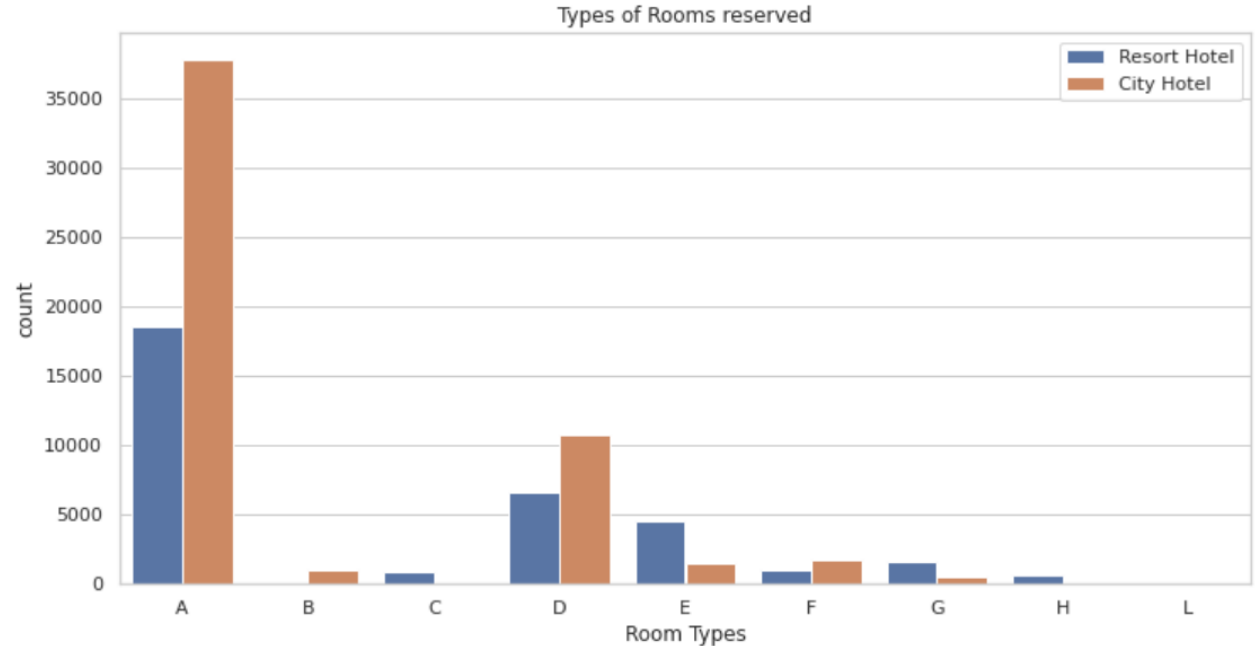More than 91% guests don't need parking space.

Approx. 8% guests need parking space for 1 car.

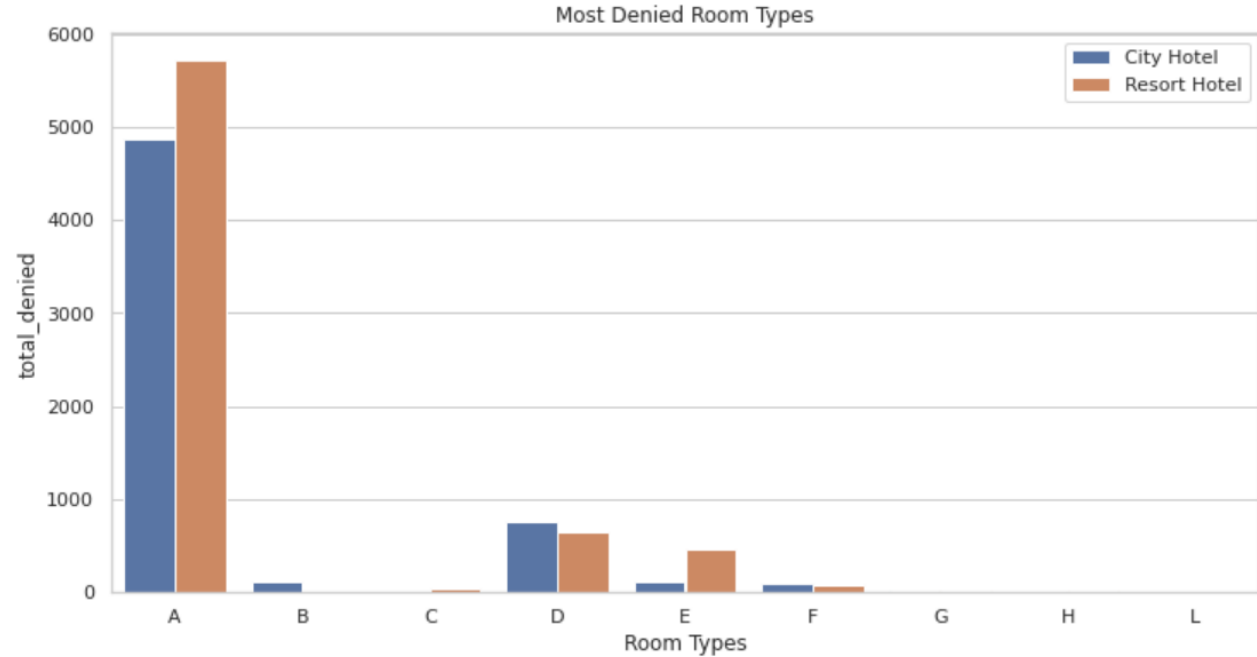| | required_car_parking_spaces | total_demand | demand_percentage |
|---|---|---|---|
| 0 | 0 | 79924 | 91.624441 |
| 1 | 1 | 7273 | 8.337728 |
| 2 | 2 | 28 | 0.032099 |
| 3 | 3 | 3 | 0.003439 |
| 4 | 8 | 2 | 0.002293 |



Car Parking Spaces Demand

# Room Types

Most of guests reserves rooms of type 'A', also significant number of guests also reserves room type 'D' & 'E'. Demand for rest types are very minimal.
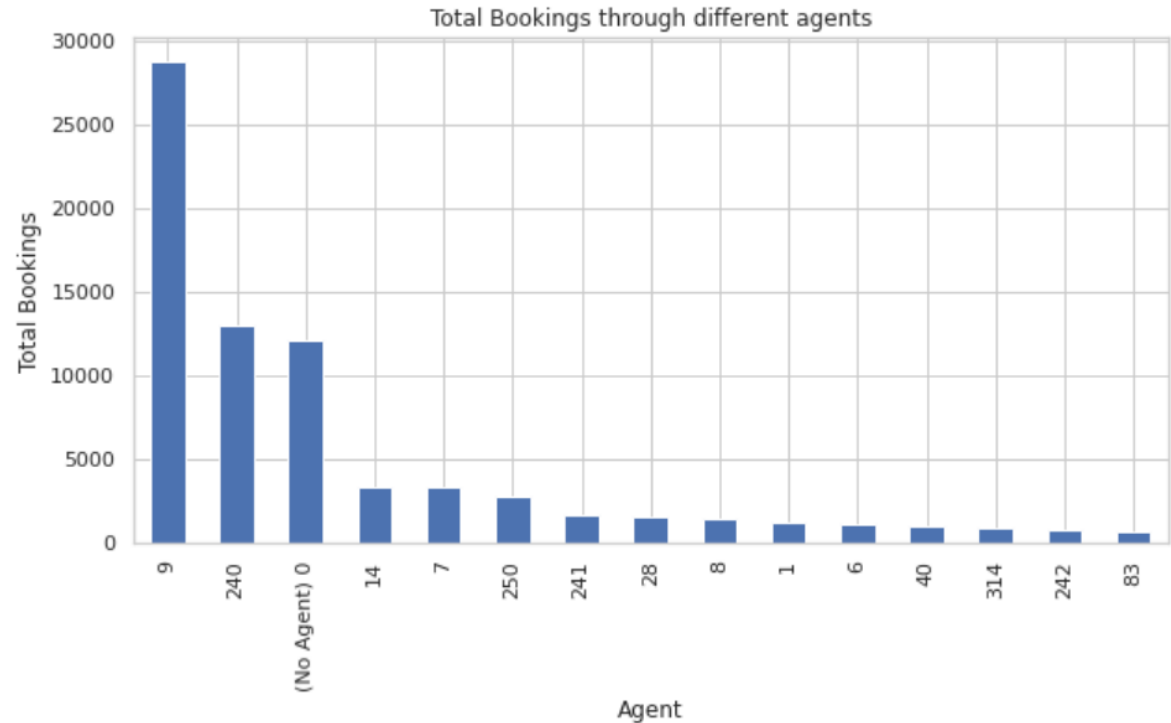


Types of Rooms reserved

# Room Types

Similar to number of bookings, room type 'A' has most number of denials. Followed by 'D' and 'E' etc.



Most Denied Room Types
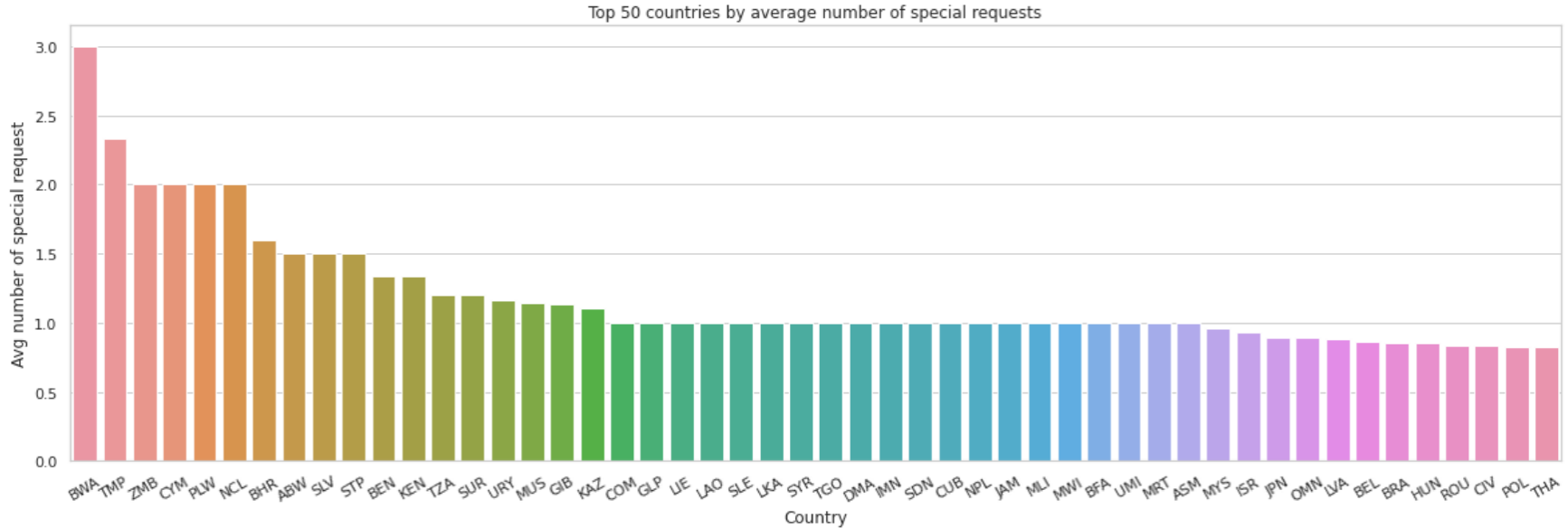
# Top Agents

Agent no. 9 has made most no. of bookings followed by 240, 14 etc.
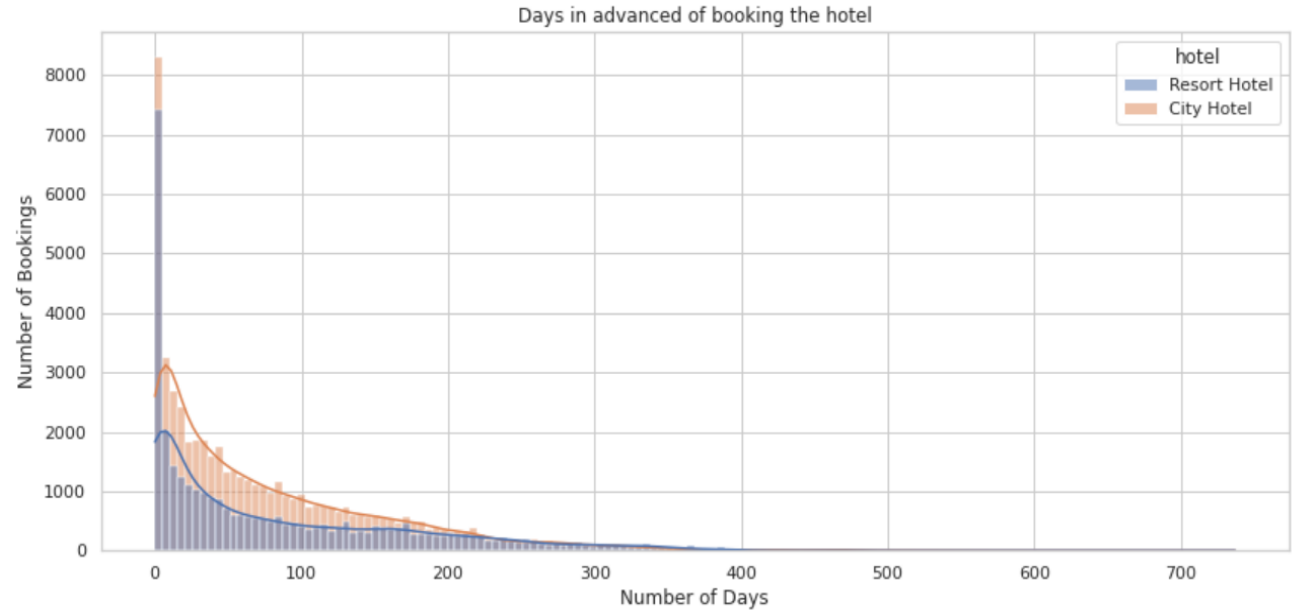


Total Bookings through different agents

# Special Requests by Country

Botswana, East Timor, Zambia are some countries with highest number of special requests.



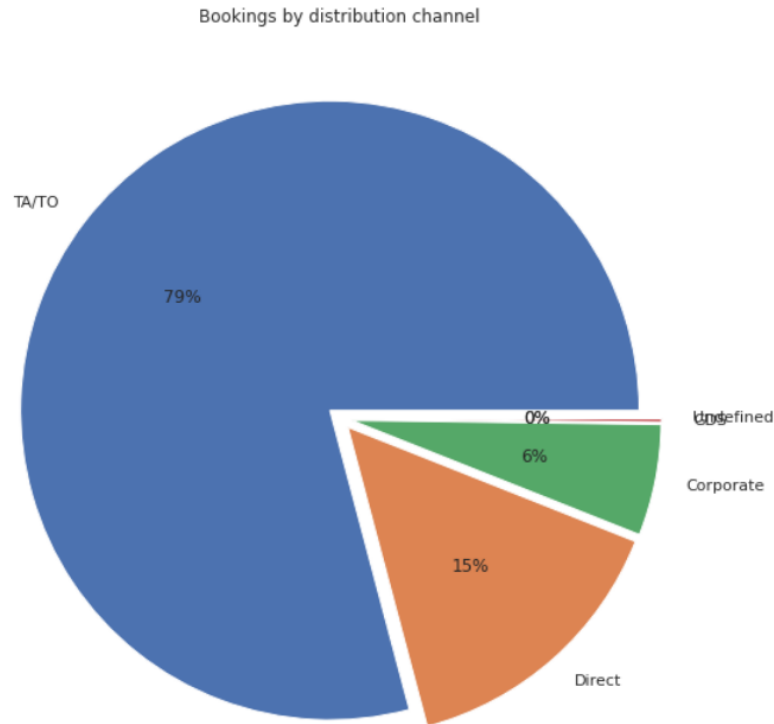Top 50 countries by average number of special requests

# Lead Time

Majority of the bookings are done within 100 days of check-in date to the hotel.



Days in advanced of booking the hotel

# Bookings by Distribution Channel

Around 79% of bookings are done through TA/TO distribution channel.

Other major distribution channels are Direct and Corporate.
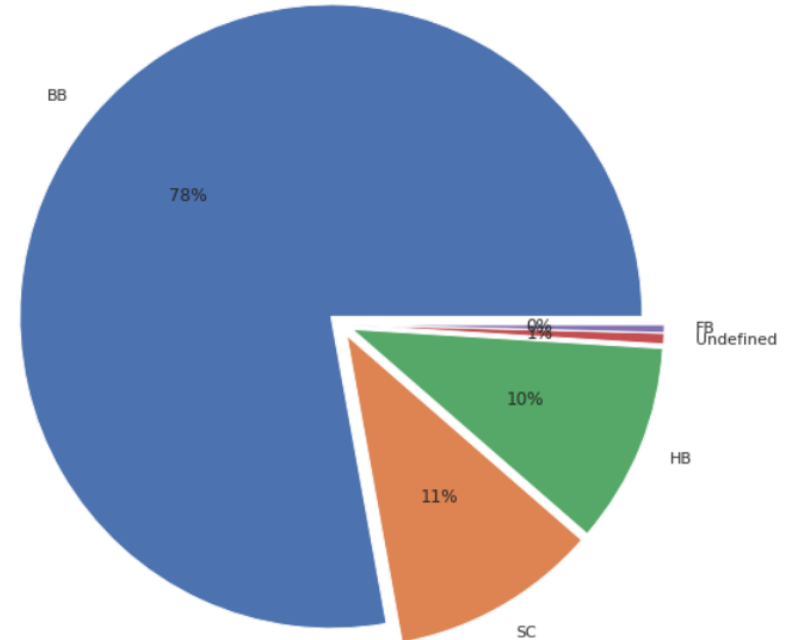
Bookings by distribution channel

# Meals

Almost 90% of total guests books meal in their bookings, among which majority books meal of type BB (78%) followed by SC (11%), HB (10%) etc.
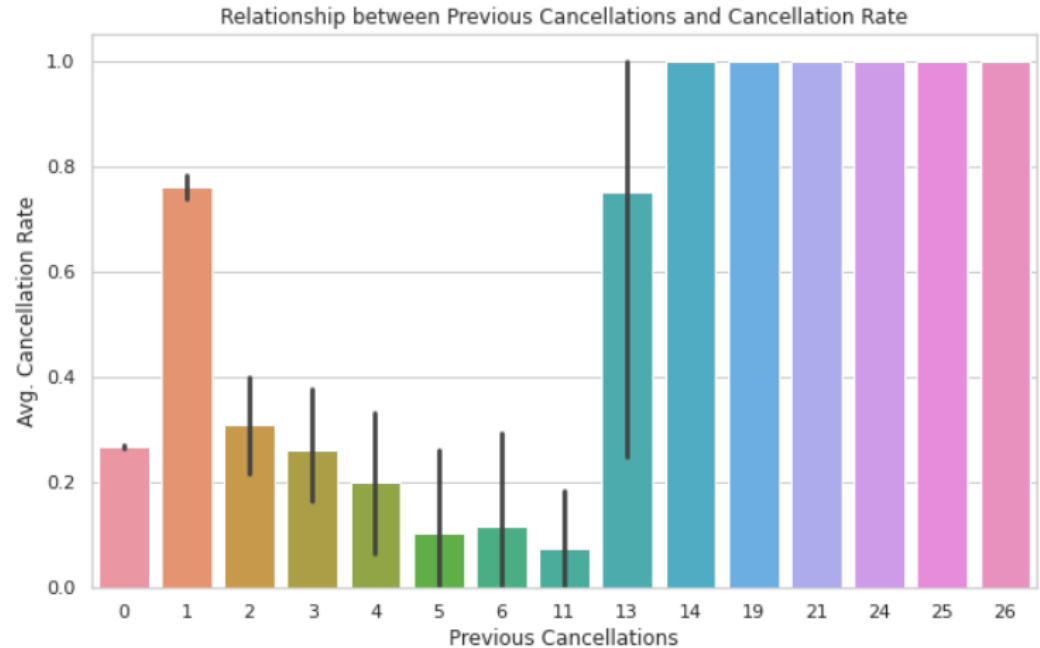
BB: Bed and breakfast.
FB: Full board (breakfast, lunch and dinner).
HB: Half board (breakfast and lunch/dinner).
SC: Room only with no meals included.
Undefined: No meal or undefined.



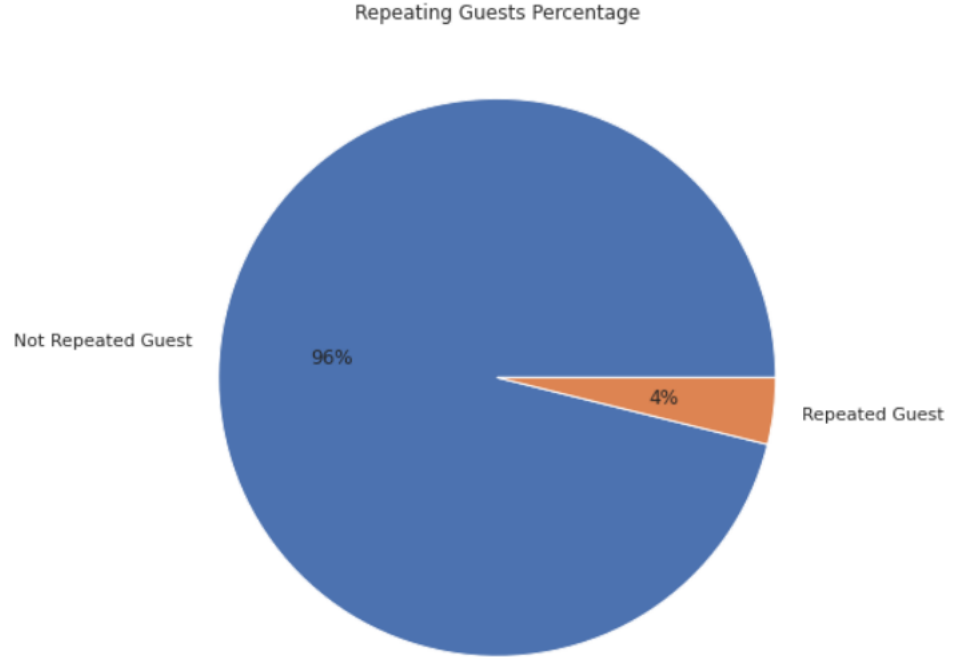Portion of bookings with meals and its type

# Previous and Current cancellations

Those who have cancelled previously once have high chance of cancellation compared to 0 previous cancellations. For higher values of previous cancellations, numbers are very less and can be ignored.



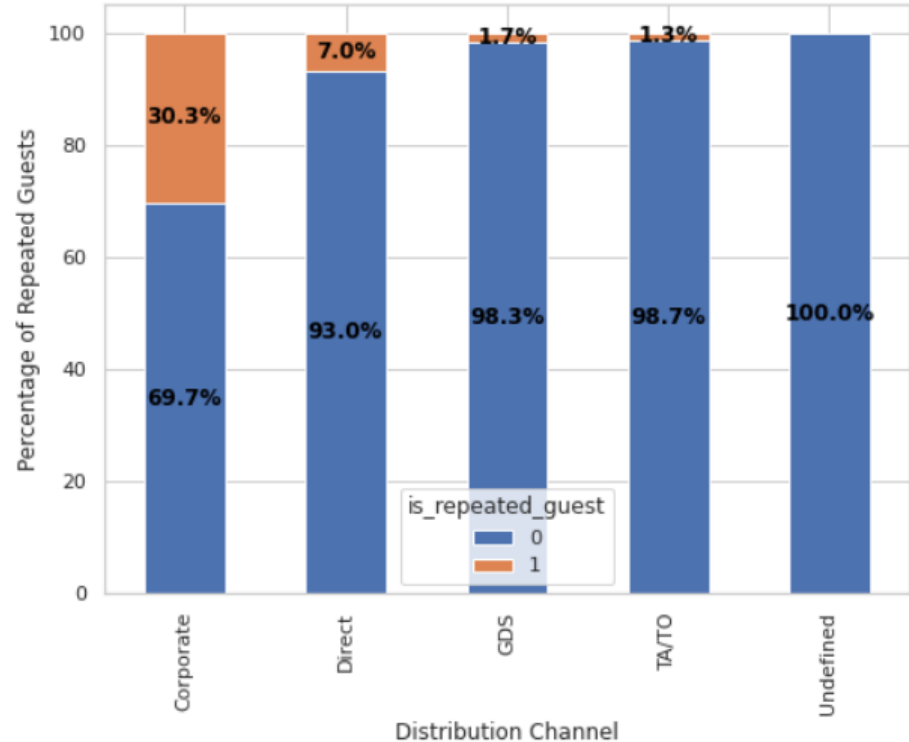Relationship between Previous Cancellations and Cancellation Rate

# Repeated Guests

Around 4% of total bookings are made by repeated guests. So majority bookings are from new customers.

Repeating Guests Percentage

Not Repeated Guest 96%

4% Repeated Guest

# Repeated Guests by Distribution Channel

Corporate has most number of repeated guests followed by Direct, GDS etc.

# Conclusions

- City Hotels have more bookings (61%) compared to Resort Hotels (39%).

- Non-refundable deposits decreases significantly the chances of cancellation.

- Longer lead time increases the chances of cancellations.

- Summer months (April to August) have more bookings compared to winter months.

- Portugal has the highest number of customers followed by Britain, France etc.

- More bookings comes through TA/TO (Online/Offline) compared to Direct, Corporates etc.

- Most customers don't need parking space around 91%.

- Room Type A is the most preferred room type among customers.

- Agent no. 9 has made most no. of bookings followed by 240, 14 etc.

- Almost 90% of total guests books meal in their bookings, among which majority books meal of type BB (78%) followed by SC (11%), HB (10%) etc.

- Corporate has the most percentage of repeated guests while TA/TO has the least whereas in the case of cancelled bookings TA/TO has the most percentage while Corporate has the least.

# Thank You!