# CAPSTONE PROJECT:

## Hotel Booking Analysis

SAMEER THETE

# Content

- Importing and loading data of Hotel booking analysis
- Data cleaning
- Data preparation
- Data visualization
- Challenges
- Conclusion
- Suggestions

# Data summary

Data set – Hotel booking analysis database includes information about the hotels booked between the year 2015 to year 2017

Shape:

Rows – 119390

Columns – 32

Important columns- lead time, arrival date, no. of persons, repeated guest, no. of kids

## Data cleaning

- we will check for duplicate values and null values.
- we will drop those columns and filter the data accordingly.

```
[ ]  duplicate_rows_df = df[df.duplicated()].shape

     print(f"the no. of duplicate rows :" , duplicate_rows_df)

     the no. of duplicate rows : (31994, 32)
```

Lets drop the duplicate values

```
[ ]  df=df.drop_duplicates()
     df.shape

     (87396, 32)
```

Data Cleaning:

Checking for Null Values and duplicate values.

Dropping the necessary columns accordingly.

Since the column named **Company and Agents** have lots of null values , we will drop these columns and
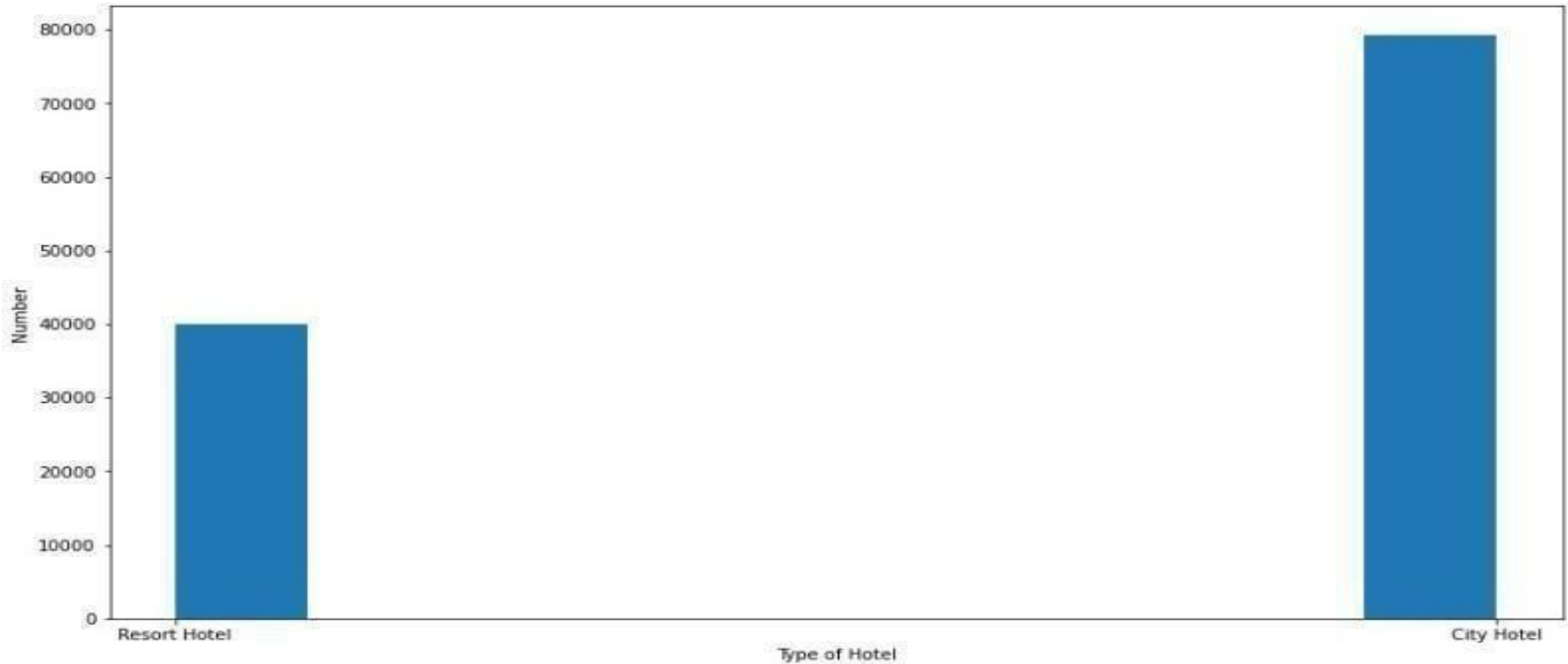
```
df = df.drop(columns=['company','agent'])
```

```
[ ] df.isnull().sum()
```

After Dealing with nulls ,these are some of the quick observations:
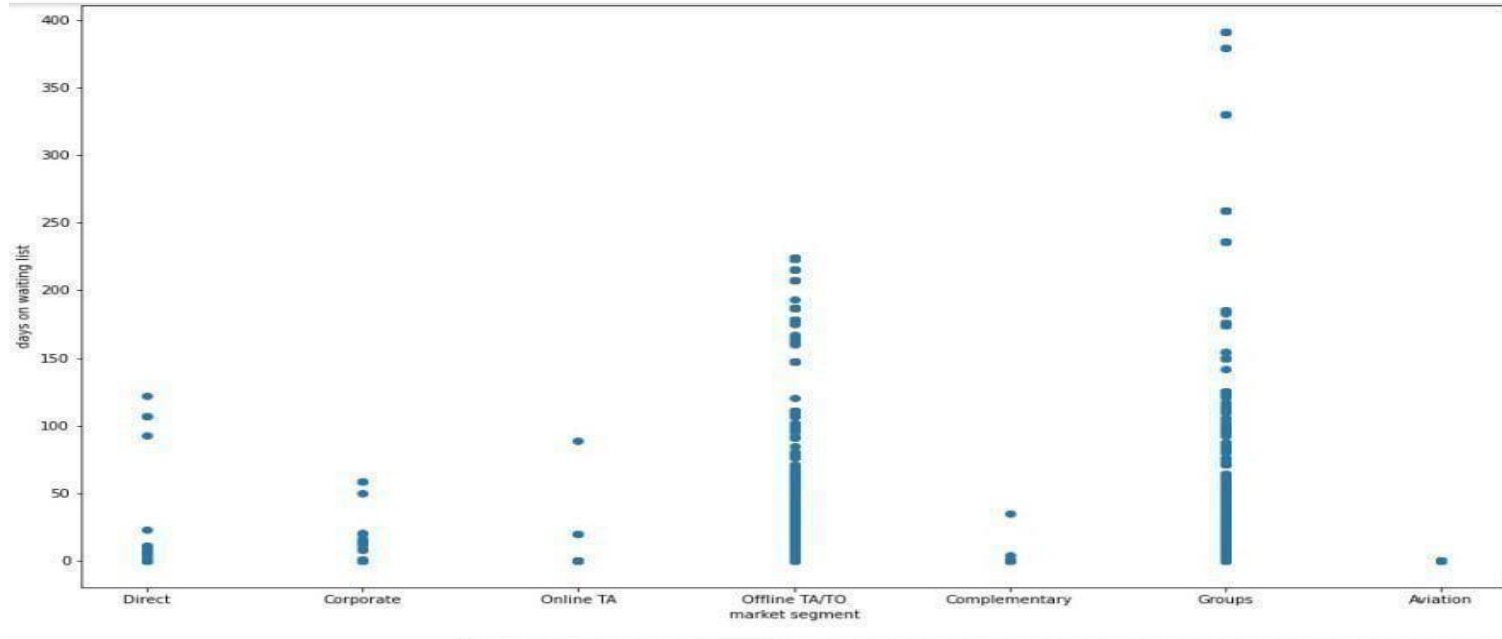
## Quick observation

- 27.4 % of the people have cancelled their booking as per the dataset.
- Avg. lead time is 80 days.
- Only 4% of the guests are repeated.
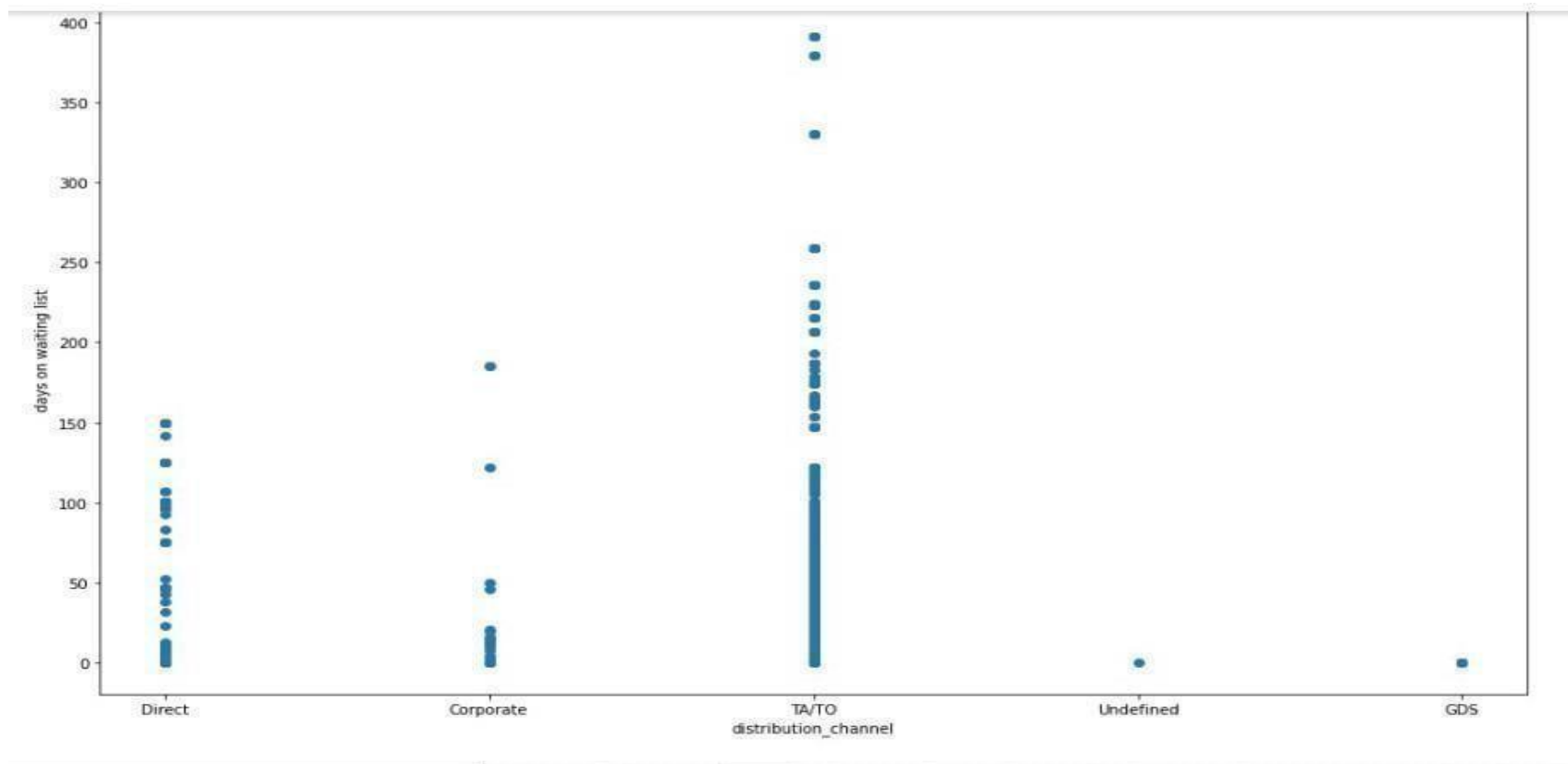- Each booking has on an average 1.8 adults and 0.13 children.

# Number of Bookings for various types of hotels (Histogram)

# Scatter plot between Type of market segment and waiting list for the booking

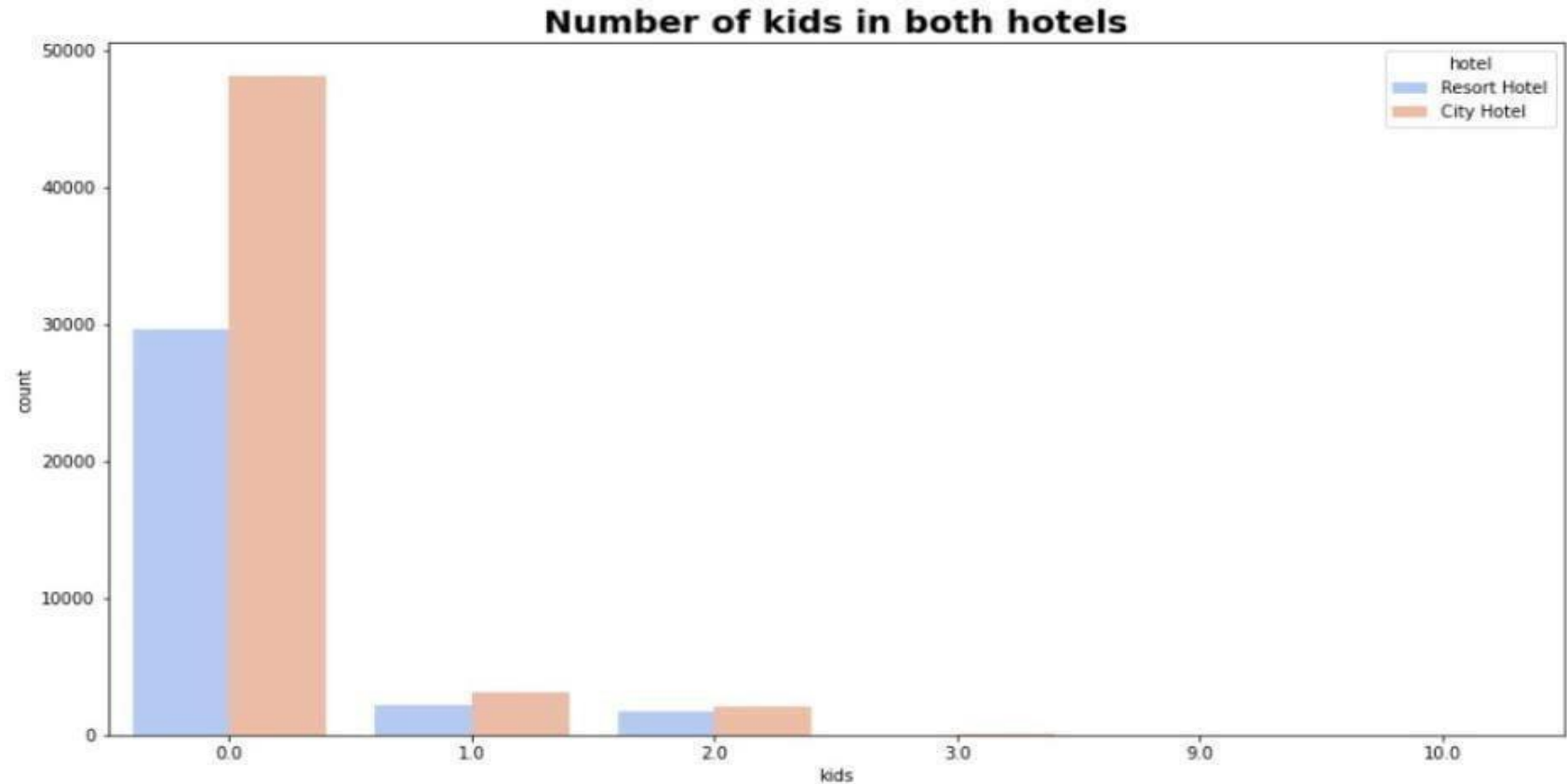# Scatter Plot between Distributing Channel and Days on the waiting

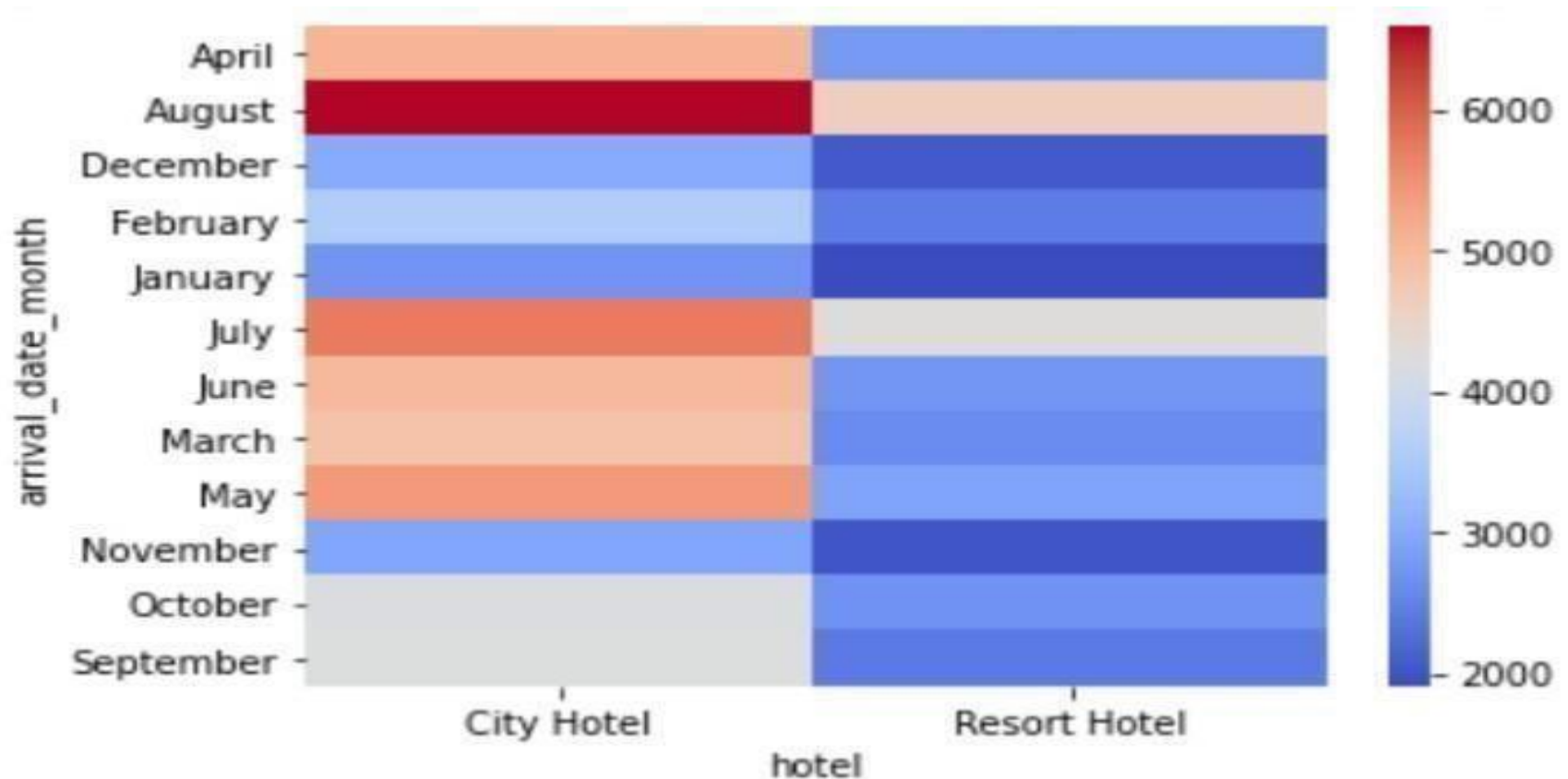For Displaying the number of kids in both hotels the new data frame is created :

Display the number of kids in both hotels

```
# Create a new dataframe to display hotel, adults, children, and babies only.
df1 = df[['hotel', 'adults', 'children', 'babies']]
df1['kids'] = df1['children'] + df1['babies']
df1
```
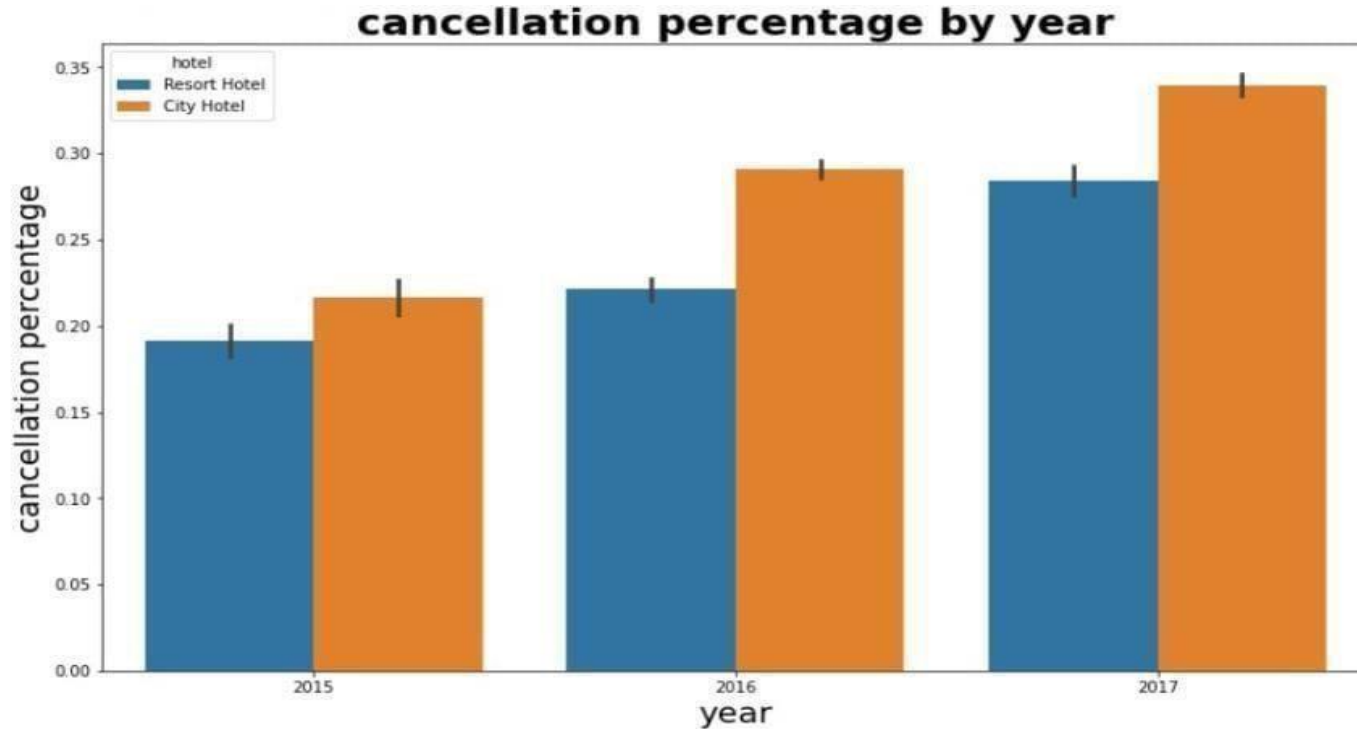
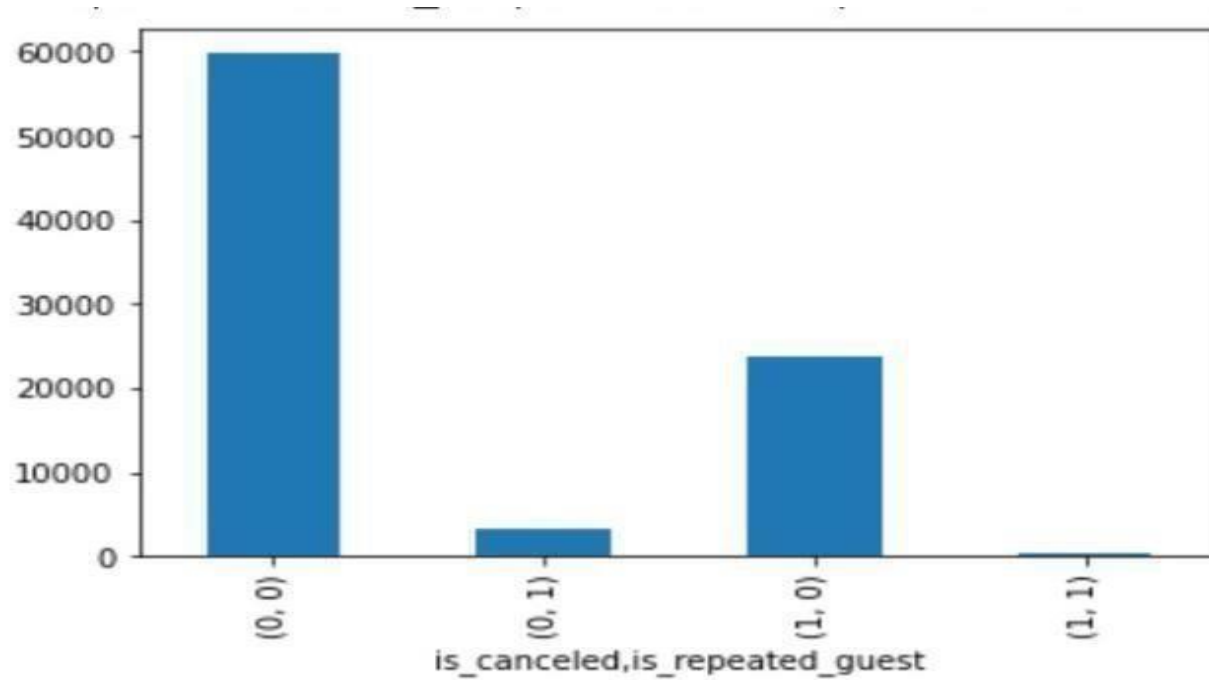# Number of kids in both hotels by countplot
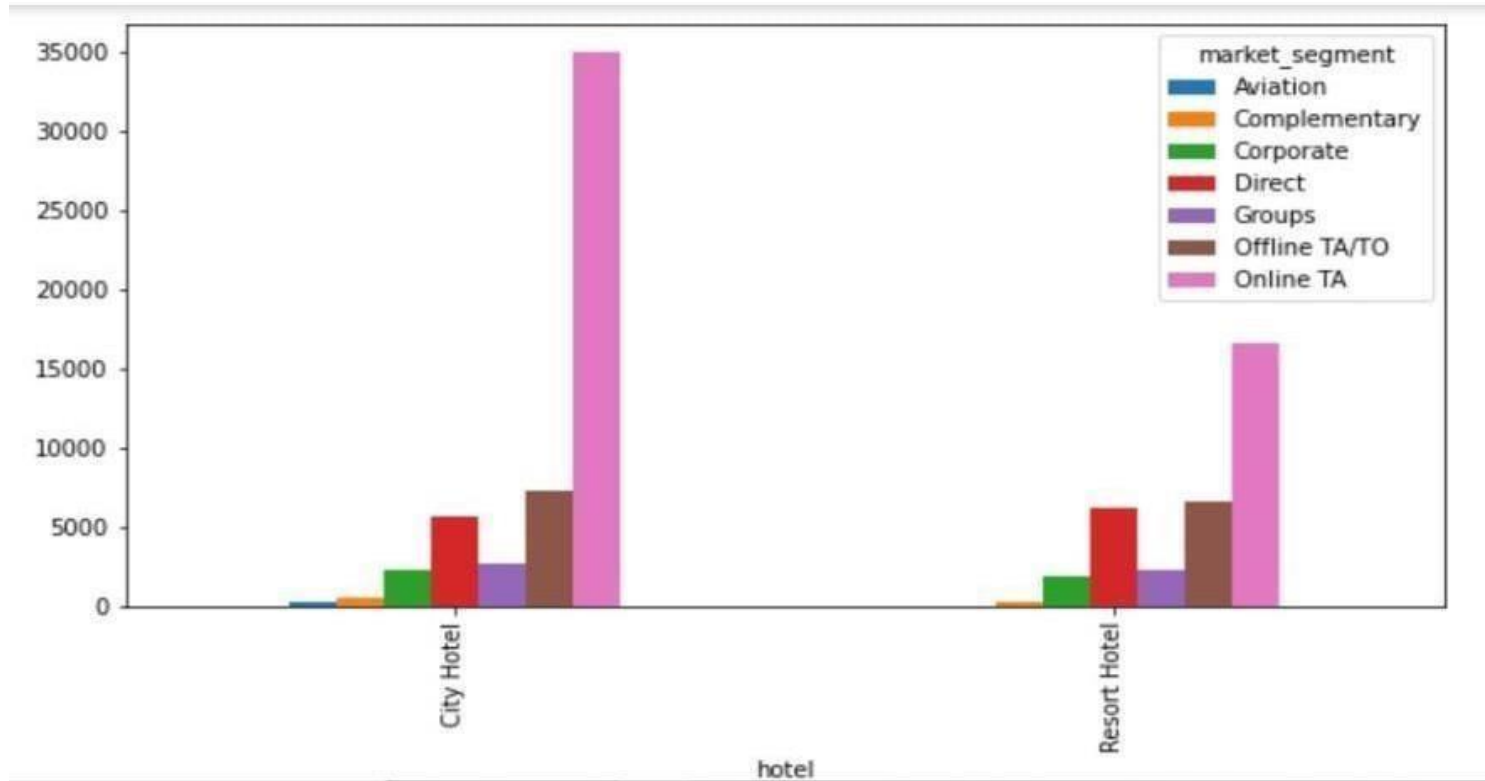
# Heatmap between type of hotel and arrival month
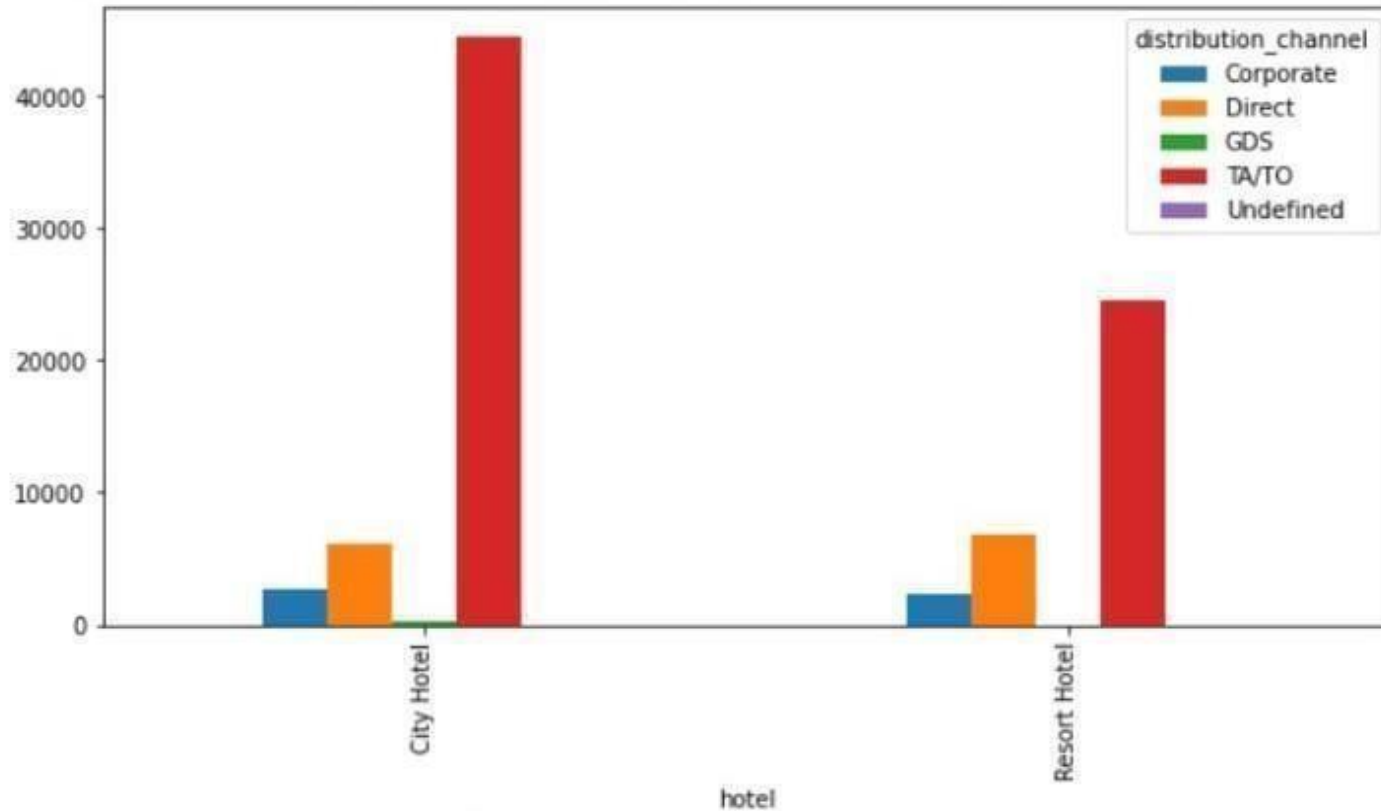
# Cancellation percentage by year

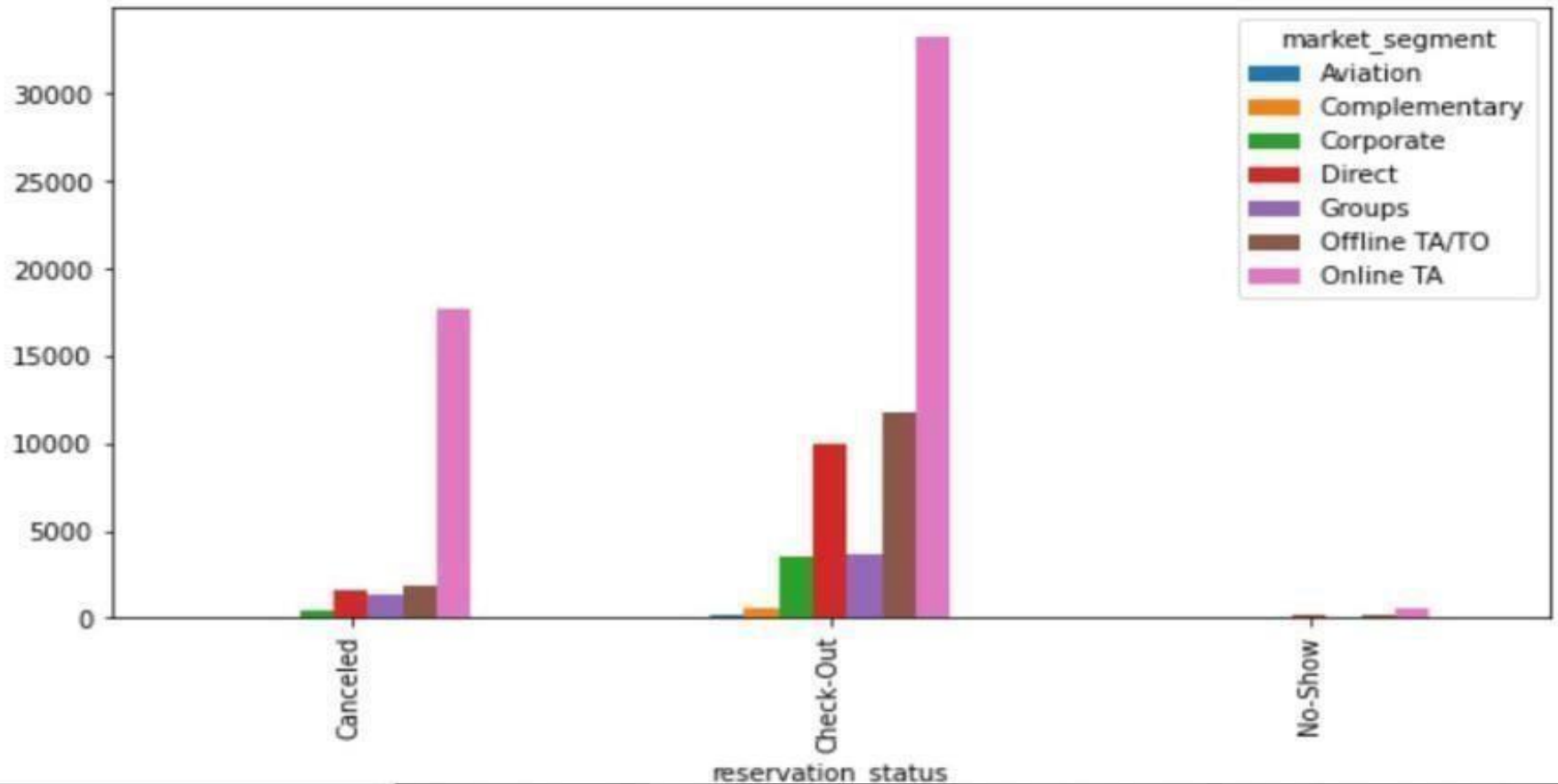# Plot between cancellation type & repeated guest
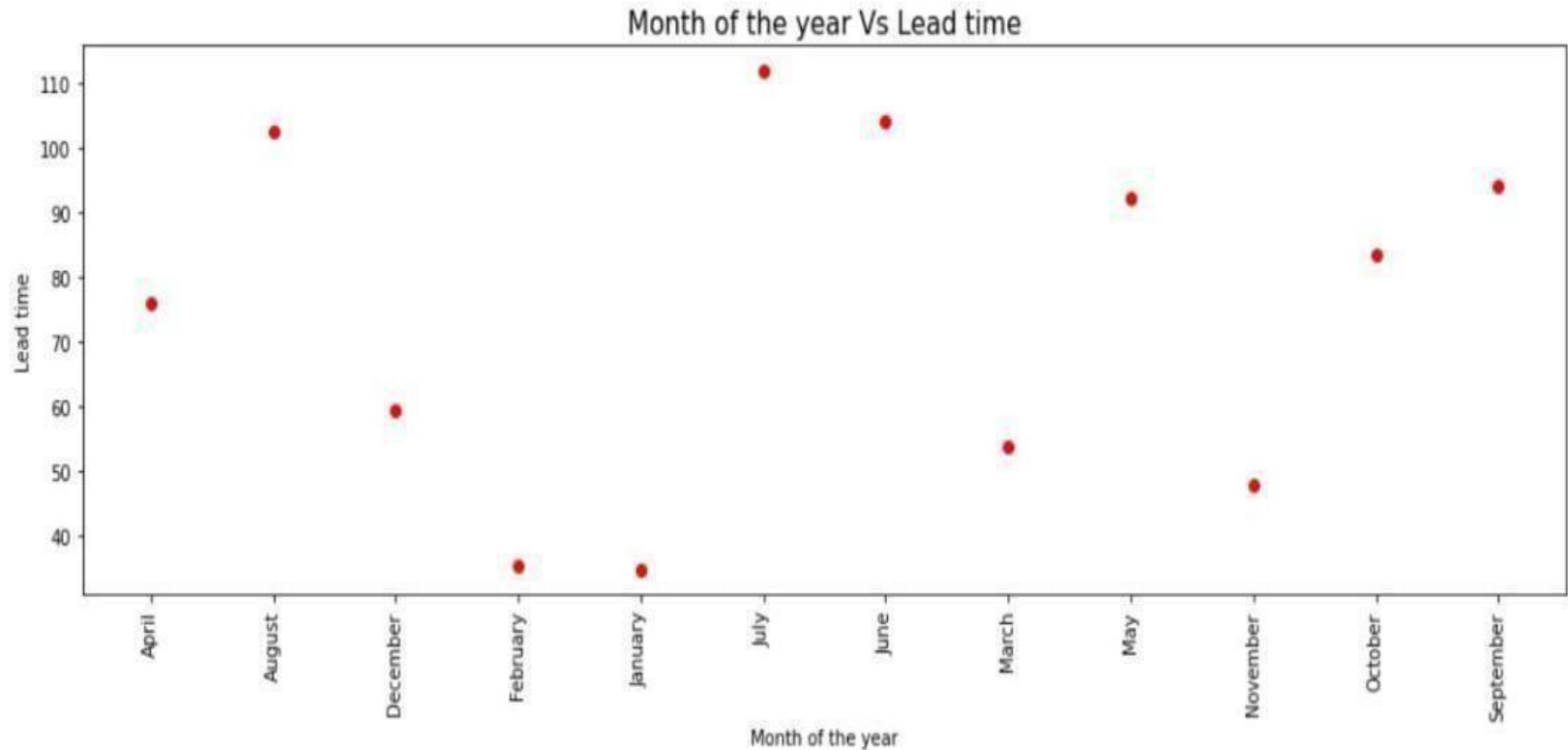
# Plot between hotel and market segment

# Plot between hotel and distribution channel

# Plot between reservation status and market segment

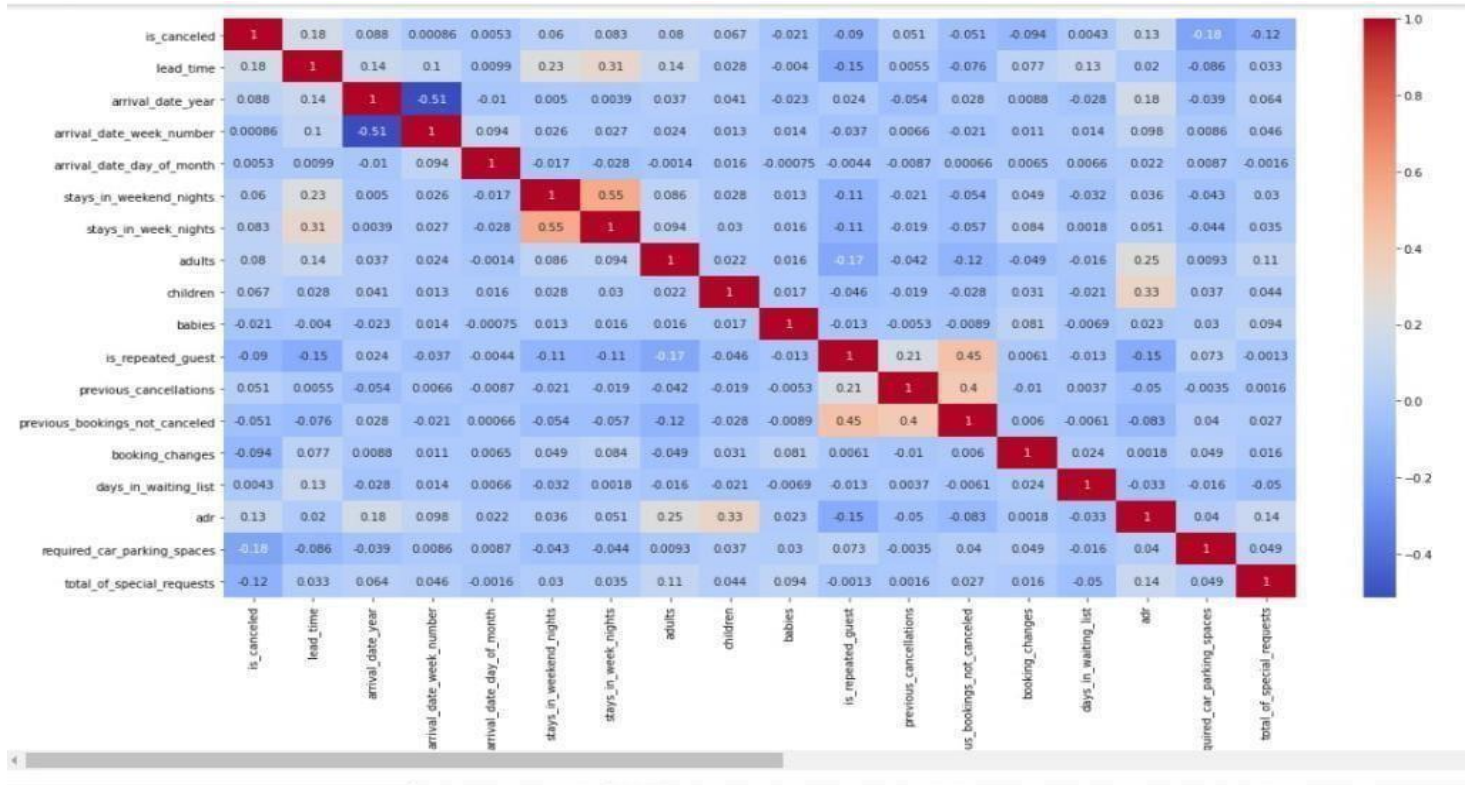# Plot between Month of year & Lead time



Month of the year Vs Lead time

# Plot between Month of year & Days in the waiting list



Month of the year Vs Days in the waiting list

# Heatmap to show the relation between variables

# Challenges

- The name of the countries was not in the proper format, because of which weare not able to plot the geomap plot

- Company and agent column has lots of duplicate value

- There were many rows with almost similar data

- Lots of null values in the dataset

# CONCLUSION

- Month of August and July receives most no. of booking.
- Booking for city hotels is twice as for resort hotels.
- Repeated costumers cancel their hotel in very rare cases.
- Customers coming from aviation industry has very less time i.e. they book urgently
- People with no kid prefer to choose city hotel over resort hotel

# Strategies to counter high cancellations at Hotel

- Since we see, our repetitive customers are most loyal customers
- To maintain them we can provide them with some bonus points, which can be redeem in the next booking
- Month of January and December receives less no. of booking, hotels can offer discounted packages for these months.
- Family with kids prefer resorts, we can provide with holiday family packages.
- Great number of the bookings are coming from travel agents, so we can provide them some commission.

# Thank You