

Import the libraries

```
In [186]: #Importing the libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

Import the data

```
In [191]: df=pd.read_csv('hotel_bookings.csv')
```

0	Resort Hotel	0	342	2015	July	27	1	0	0	2	...	No Deposit	NaN	NaN
1	Resort Hotel	0	737	2015	July	27	1	0	0	2	...	No Deposit	NaN	NaN
2	Resort Hotel	0	7	2015	July	27	1	0	1	1	...	No Deposit	NaN	NaN
3	Resort Hotel	0	13	2015	July	27	1	0	1	1	...	No Deposit	304.0	NaN
4	Resort Hotel	0	14	2015	July	27	1	0	2	2	...	No Deposit	240.0	NaN
...
119385	City Hotel	0	23	2017	August	35	30	2	5	2	...	No Deposit	394.0	NaN
119386	City Hotel	0	102	2017	August	35	31	2	5	3	...	No Deposit	9.0	NaN
119387	City Hotel	0	34	2017	August	35	31	2	5	2	...	No Deposit	9.0	NaN
119388	City Hotel	0	109	2017	August	35	31	2	5	2	...	No Deposit	89.0	NaN
119389	City Hotel	0	205	2017	August	35	29	2	7	2	...	No Deposit	9.0	NaN

119390 rows × 32 columns

Exploratory Data Analysis and Data Cleaning

list of first five rows

In [36]: df.head()

Exploratory Data Analysis and Data Cleaning

list of first five rows

3

Resort Hotel

0

13

2015

July

27

1

0

0

2

1

...

No Deposit

394.0

NaN

4

Resort Hotel

0

14

2015

July

27

1

0

0

2

2

...

No Deposit

240.0

NaN

5 rows x 32 columns

list of last five rows

In [191]: df.tail()

Out[191]:

	hotel	is_canceled	lead_time	arrival_date_year	arrival_date_month	arrival_date_week_number	arrival_date_day_of_month	stays_in_weekend_nights	stays_in_week_nights	adults	...	deposit_type	agent	company
119385	City Hotel	0	23	2017	August	35	30	2	5	2	...	No Deposit	394.0	NaN
119386	City Hotel	0	102	2017	August	35	31	2	5	3	...	No Deposit	9.0	NaN
119387	City Hotel	0	34	2017	August	35	31	2	5	2	...	No Deposit	9.0	NaN
119388	City Hotel	0	109	2017	August	35	31	2	5	2	...	No Deposit	89.0	NaN

list of last five rows

```
In [98]: df.shape
Out[98]: (119389, 32)
```

```
In [99]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 119389 entries, 0 to 119389
Data columns (total 32 columns):
 #   Column                Non-Null Count  Dtype  ---
 #   Column                Non-Null Count  Dtype  ---
 0   hotel                  119389 non-null  object
 1   is_cancelled           119389 non-null  bool
 2   lead_time              119389 non-null  float64
 3   arrival_date_year      119389 non-null  int64
 4   arrival_date_month     119389 non-null  object
 5   arrival_date_week_number 119389 non-null  int64
 6   arrival_date_day_of_month 119389 non-null  int64
 7   stays_in_weekend_nights 119389 non-null  float64
 8   stays_in_week_nights    119389 non-null  float64
 9   adults                 119389 non-null  int64
10   children                119389 non-null  int64
11   infants                 119389 non-null  int64
12   reservation_status      119389 non-null  object
13   reservation_status_date 119389 non-null  object
14   total_spent              119389 non-null  float64
15   total_spent_usd         119389 non-null  float64
16   total_spent_eur         119389 non-null  float64
17   total_spentGBP          119389 non-null  float64
18   total_spentJPY          119389 non-null  float64
19   total_spentAUD          119389 non-null  float64
20   total_spentCAD          119389 non-null  float64
21   total_spentCHF          119389 non-null  float64
22   total_spentCNY          119389 non-null  float64
23   total_spentHKD          119389 non-null  float64
24   total_spentINR          119389 non-null  float64
25   total_spentKRW          119389 non-null  float64
26   total_spentMXN          119389 non-null  float64
27   total_spentNZD          119389 non-null  float64
28   total_spentSGD          119389 non-null  float64
29   total_spentUSD          119389 non-null  float64
30   total_spentZAR          119389 non-null  float64
31   total_spentBRL          119389 non-null  float64
32   total_spentRUB          119389 non-null  float64
```

find the shape

```
In [191]: df.shape
```

(119390, 14)

check the data info

```
In [191]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 119390 entries, 0 to 119389
Data columns (total 14 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   hotel               119390 non-null object
 1   is_canceled         119390 non-null int64
 2   lead_time          119390 non-null int64
 3   arrival_date_year   119390 non-null int64
 4   arrival_date_month  119390 non-null object
 5   arrival_date_week_number  119390 non-null int64
 6   arrival_date_day_of_month  119390 non-null int64
 7   stays_in_weekend_nights  119390 non-null int64
 8   stays_in_week_nights  119390 non-null int64
 9   children            119390 non-null float64
10   babies              119390 non-null int64
11   adults              119390 non-null int64
12   meall               118892 non-null object
13   country             119390 non-null object
14   market_segment      119390 non-null object
15   distribution_channel 119390 non-null object
16   is_repeated_guest    119390 non-null int64
17   previous_cancellations  119390 non-null int64
18   previous_bookings_not_canceled  119390 non-null int64
19   reserved_room_type   119390 non-null object
20   assigned_room_type   119390 non-null int64
21   booking_changes      119390 non-null int64
22   deposit_type         119390 non-null object
23   agent               119390 non-null float64
24   company             6797 non-null float64
25   days_in_waiting_list  119390 non-null int64
26   customer_type        119390 non-null object
27   adr                 119390 non-null float64
28   required_car_parking_spaces  119390 non-null int64
29   total_of_special_requests  119390 non-null int64
30   reservation_status   119390 non-null object
31   reservation_status_date  119390 non-null object
dtypes: float64(1), int64(10), object(12)
memory usage: 29.1+ MB
```

its convert the reservation date in to date

```
In [198]: df['reservation_status_date']=pd.to_datetime(df['reservation_status_date'])
```

```
In [191]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 119390 entries, 0 to 119389
Data columns (total 14 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   hotel               119390 non-null object
 1   is_canceled         119390 non-null int64
 2   lead_time          119390 non-null int64
 3   arrival_date_year   119390 non-null int64
 4   arrival_date_month  119390 non-null object
 5   arrival_date_week_number  119390 non-null int64
 6   arrival_date_day_of_month  119390 non-null int64
 7   stays_in_weekend_nights  119390 non-null int64
 8   stays_in_week_nights  119390 non-null int64
 9   children            119390 non-null float64
10   babies              119390 non-null int64
11   adults              119390 non-null int64
12   meall               118892 non-null object
13   country             119390 non-null object
14   market_segment      119390 non-null object
15   distribution_channel 119390 non-null object
16   is_repeated_guest    119390 non-null int64
17   previous_cancellations  119390 non-null int64
18   previous_bookings_not_canceled  119390 non-null int64
19   reserved_room_type   119390 non-null object
20   assigned_room_type   119390 non-null int64
21   booking_changes      119390 non-null int64
22   deposit_type         119390 non-null object
23   agent               119390 non-null float64
24   company             6797 non-null float64
25   days_in_waiting_list  119390 non-null int64
26   customer_type        119390 non-null object
27   adr                 119390 non-null float64
28   required_car_parking_spaces  119390 non-null int64
29   total_of_special_requests  119390 non-null int64
30   reservation_status   119390 non-null object
31   reservation_status_date  119390 non-null datetime64[ns]
dtypes: datetime64[ns](1), float64(1), int64(11), object(11)
memory usage: 29.1+ MB
```

```
In [191]: df.describe(include = "object")

Out[191]:
```

	hotel	arrival_date_month	meall	country	market_segment	distribution_channel	reserved_room_type	assigned_room_type	deposit_type	customer_type	reservation_status
count	119390	119390	119390	118902	119390	119390	119390	119390	119390	119390	119390
unique	1	12	5	177	8	5	10	12	9	4	9
top	City Hotel	August	08	Port	Online TA	TA/TO	A	A	No Deposit	Transient	Check-Out
freq	79330	13877	92310	48590	56477	97870	85994	74053	104641	89613	75166

checking for null values

```
In [192]: df.isna().sum()
```

```
Out[192]:
```

hotel	0
is_canceled	0
lead_time	0
arrival_date_year	0
arrival_date_month	0
arrival_date_week_number	0
arrival_date_day_of_month	0
stays_in_weekend_nights	0
stays_in_week_nights	0
adults	0
children	4
babies	0
meall	0
country	488
market_segment	0
distribution_channel	0
is_repeated_guest	0
previous_cancellations	0
previous_bookings_not_canceled	0
reserved_room_type	0
assigned_room_type	0
booking_changes	0
deposit_type	0
agent	26140
company	112593
days_in_waiting_list	0
customer_type	0
adr	0
required_car_parking_spaces	0
total_of_special_requests	0
reservation_status	0
reservation_status_date	0
dtype: int64	

remove the null missing values and null columns

```
In [193]: df.dropna(inplace = True)
```

```
In [194]: df.drop(columns=['agent', 'company'], inplace=True)
```

```
In [191]: df.isna().sum()

Out[191]:
```

hotel	0
is_canceled	0
lead_time	0
arrival_date_year	0
arrival_date_month	0
arrival_date_week_number	0
arrival_date_day_of_month	0
stays_in_weekend_nights	0
stays_in_week_nights	0
adults	0
children	0
babies	0
meall	0
country	0
market_segment	0
distribution_channel	0
is_repeated_guest	0
previous_cancellations	0
previous_bookings_not_canceled	0
reserved_room_type	0
assigned_room_type	0
booking_changes	0
deposit_type	0
days_in_waiting_list	0
customer_type	0
adr	0
required_car_parking_spaces	0
total_of_special_requests	0
reservation_status	0
reservation_status_date	0
dtype: int64	

how we are creating new columns for easy visualization

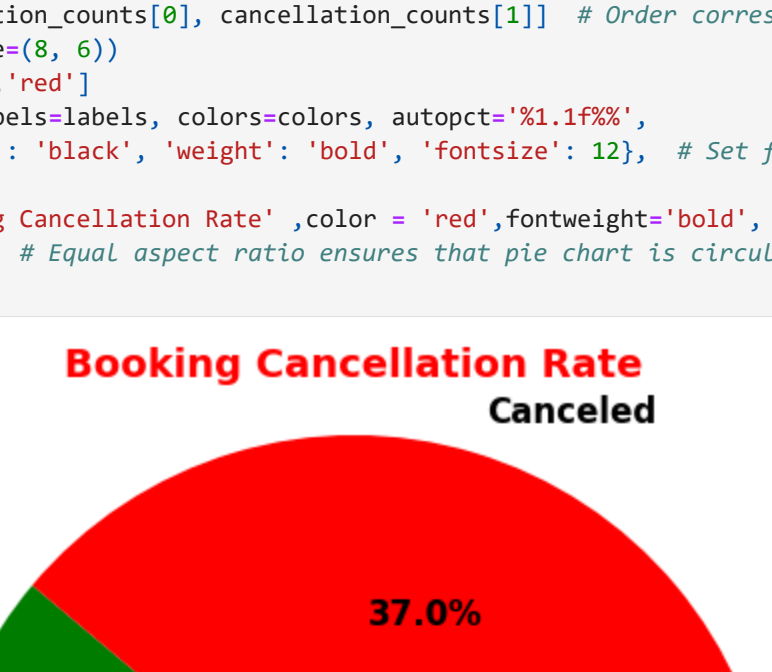
```
In [191]: df['cancellation_status'] = df['is_canceled'].astype(int)
```

Q1. What is the overall cancellation rate?

```
In [124]: cancellation_counts = df['is_canceled'].value_counts(normalize=True) * 100
```

Now we plot A pie chart

```
In [34]: labels = ['Confirmed', 'Canceled']
sizes = [cancellation_counts[0], cancellation_counts[1]] # Order corresponds to the labels
plt.figure(figsize=(8, 6))
colors = ['#9467bd', 'red']
plt.pie(sizes, labels=labels, colors=colors, autopct='%1.1f%%',
        labeldistance=1.1, radius=1.1, # Set font color, weight, and size
        startangle=140)
plt.title('Booking Cancellation Rate', color='red', fontweight='bold', fontsize=14)
plt.axis('equal') # Equal aspect ratio ensures that pie chart is circular.
plt.show()
```

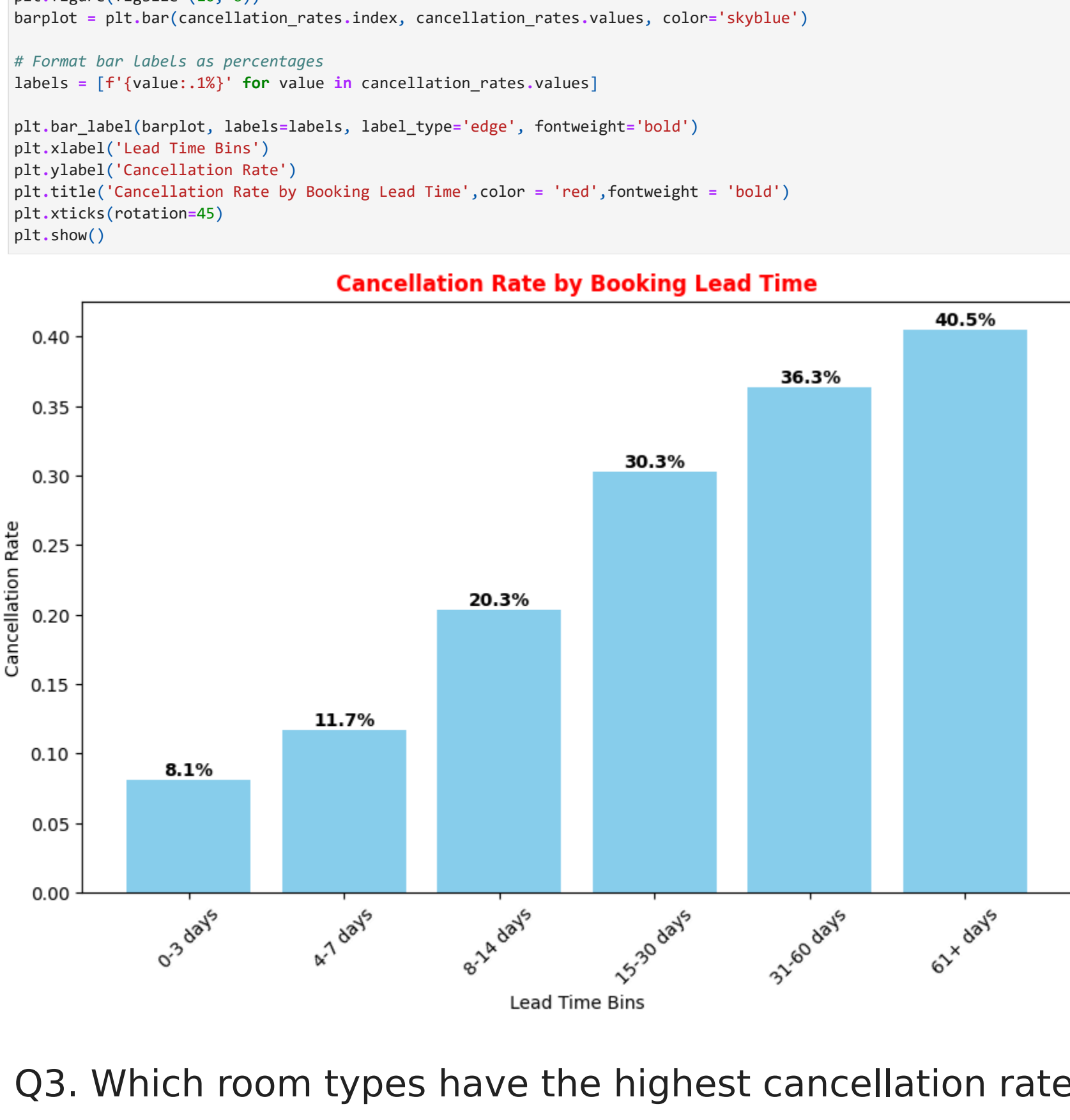


Qn2. How does the cancellation rate vary by lead time?

```
In [36]: bins = [0, 3, 7, 14, 30, 60, 180] # Define your bins
labels = ['0-3 days', '4-7 days', '8-14 days', '15-30 days', '31-60 days', '61+ days']
df['lead_time_bins'] = pd.cut(df['lead_time'], bins=bins, labels=labels, right=False)

# Calculate cancellation rates for each bin
cancellation_rates = df.groupby('lead_time_bins', observed=False)['cancellation_status'].mean()

In [37]: plt.figure(figsize=(10, 6))
barplot = plt.bar(cancellation_rates.index, cancellation_rates.values, color='skyblue')
# Format bar labels as percentages
labels = [f'{value:.1f}' for value in cancellation_rates.values]
plt.bar_label(barplot, labels=labels, label_type='edge', fontweight='bold')
plt.xlabel('Lead Time Bins')
plt.ylabel('Cancellation Rate')
plt.title('Cancellation Rate by Booking Lead Time', color='red', fontweight='bold')
plt.xticks(rotation=45)
plt.show()
```



Q3. Which room types have the highest cancellation rates?

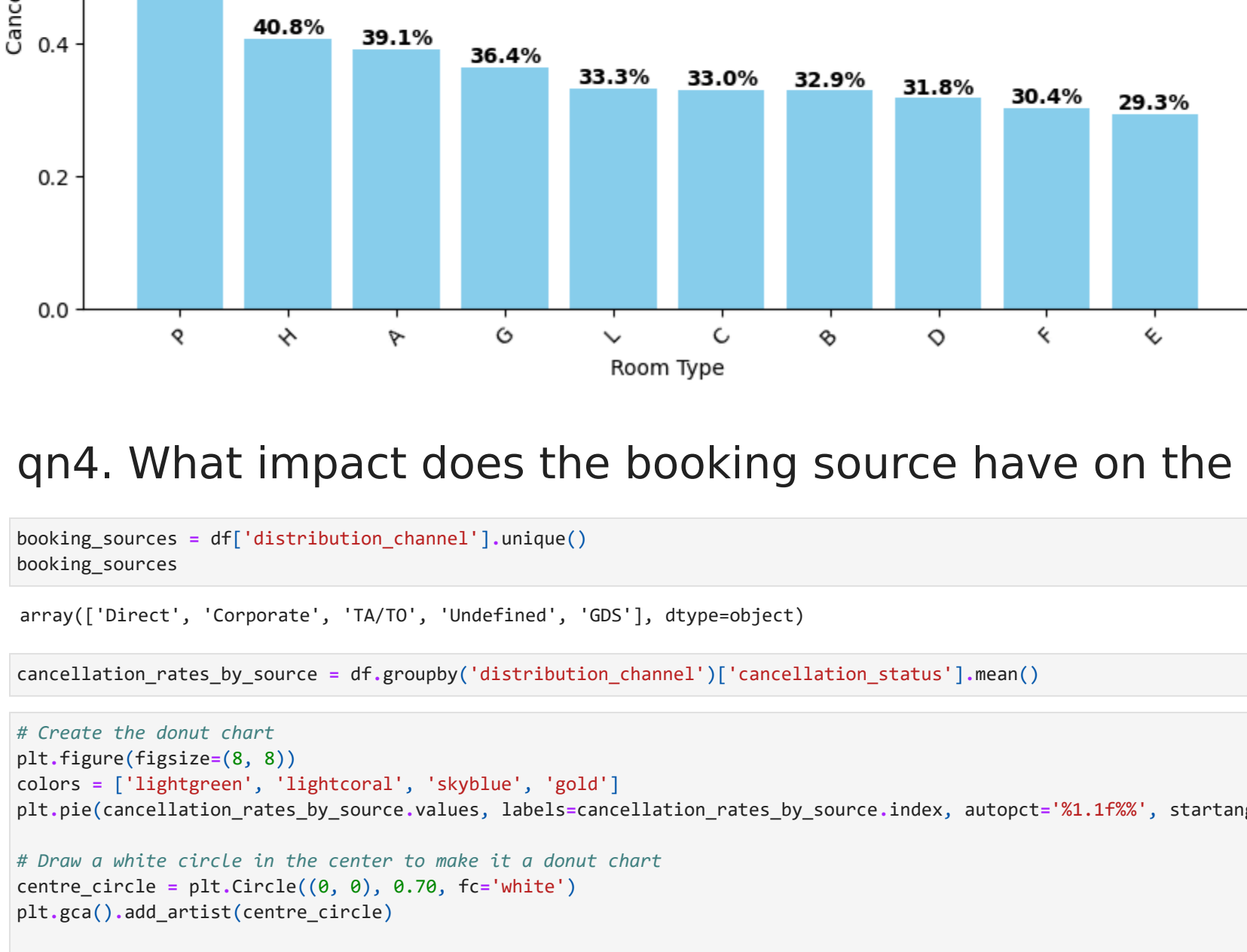
```
In [38]: cancellation_rates_by_room_type = df.groupby('reserved_room_type')['cancellation_status'].mean()

In [39]: # Sort the cancellation rates by room type in descending order
sorted_cancellation_rates = cancellation_rates_by_room_type.sort_values(ascending=False)

# Creating the bar chart with sorted data
plt.figure(figsize=(10, 6))
barplot = plt.bar(sorted_cancellation_rates.index, sorted_cancellation_rates.values, color='skyblue')

# Adding labels to the bars
plt.bar_label(barplot, labels=[f'{value:.1f}' for value in sorted_cancellation_rates.values], label_type='edge', fontweight='bold')

# Labeling the chart and axes
plt.xlabel('Room Type')
plt.ylabel('Cancellation Rate')
plt.title('Cancellation Rate by Room Type (Sorted)')
plt.xticks(rotation=45)
plt.show()
```



qn4. What impact does the booking source have on the cancellation rate?

```
In [41]: booking_sources = df['distribution_channel'].unique()
booking_sources

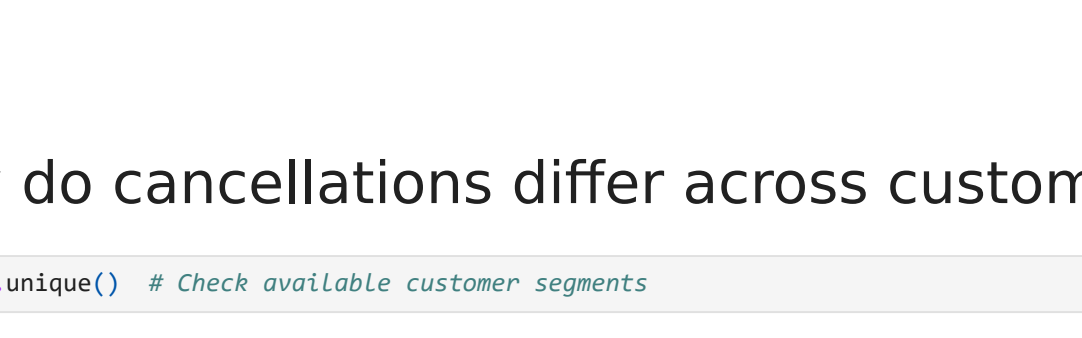
Out[41]: array(['Direct', 'Corporate', 'TA/TO', 'Undefined', 'GDS'], dtype=object)

In [42]: cancellation_rates_by_source = df.groupby('distribution_channel')['cancellation_status'].mean()

In [43]: # Create the donut chart
plt.figure(figsize=(8, 6))
colors = ['lightgreen', 'lightcoral', 'skyblue', 'gold']
plt.pie(cancellation_rates_by_source.index, labels=cancellation_rates_by_source.index, autopct='%1.1f%%', startangle=90, colors=colors)

# Show a white circle in the center to make it a donut chart
centre_circle = plt.Circle((0, 0), 0.70, fc='white')
plt.gca().add_artist(centre_circle)

plt.title('Cancellation Rate by Booking Source')
plt.show()
```



Qn5. How do cancellations differ across customer segments?

```
In [51]: df['customer_type'].unique() # Check available customer segments
Out[51]: array(['Transient', 'Contract', 'Transient-Party', 'Group'], dtype=object)

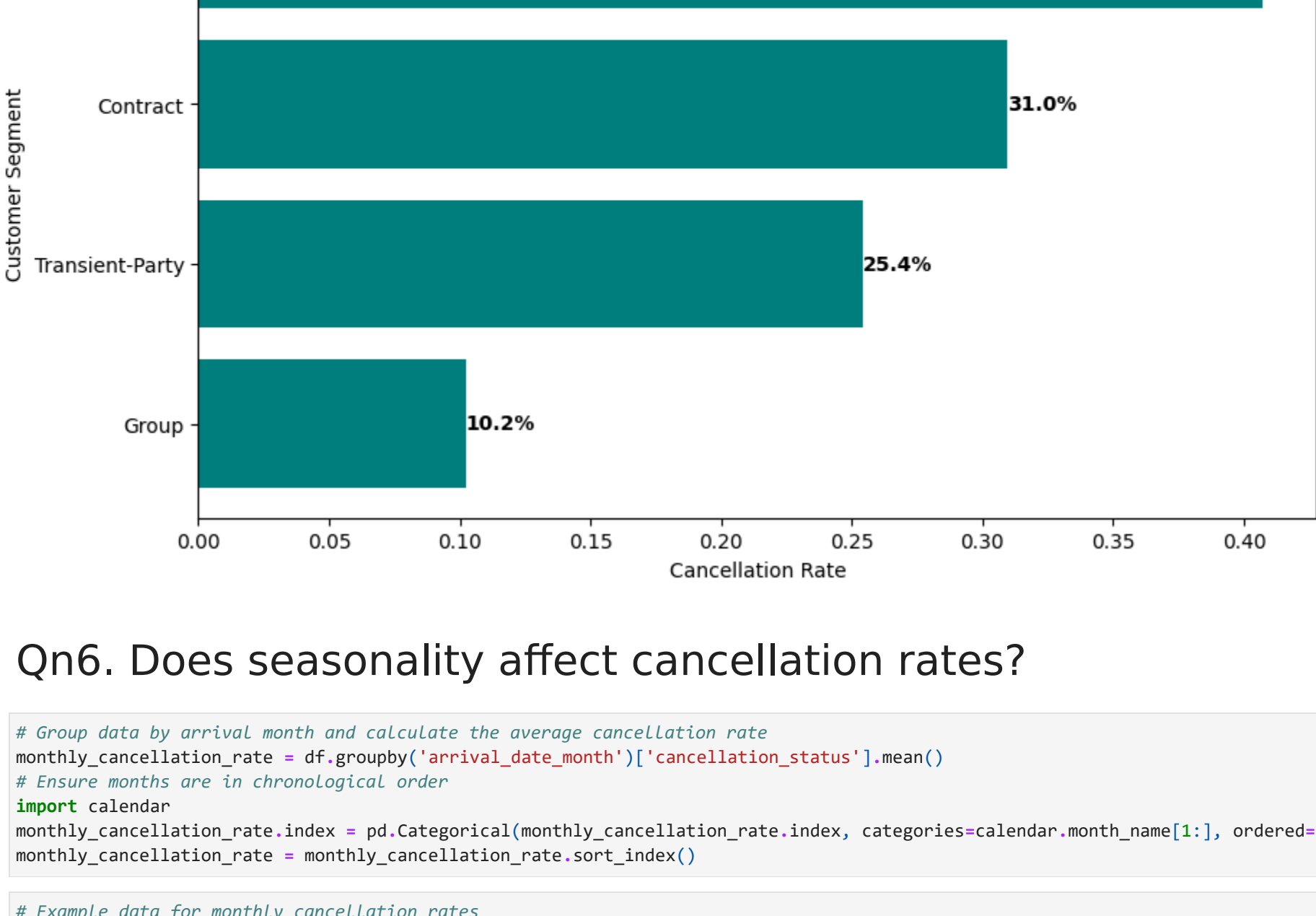
In [54]: cancellation_rates_by_customer_type = df.groupby('customer_type')['cancellation_status'].mean()

In [55]: # Sort cancellation rates by customer type
sorted_cancellation_rates_by_customer_type = cancellation_rates_by_customer_type.sort_values(ascending=False)

# Creating the bar chart
plt.figure(figsize=(10, 6))
barplot = plt.bar(sorted_cancellation_rates_by_customer_type.index, sorted_cancellation_rates_by_customer_type.values, color='teal')

# Adding labels to the bars
plt.bar_label(barplot, labels=[f'{value:.1f}' for value in sorted_cancellation_rates_by_customer_type.values], label_type='edge', fontweight='bold')

# Labeling the chart and axes
plt.xlabel('Customer Segment')
plt.ylabel('Cancellation Rate by Customer Segment')
plt.title('Cancellation Rate by Customer Segment')
plt.xticks(rotation=45)
plt.show()
```



Qn6. Does seasonality affect cancellation rates?

```
In [56]: # Group data by arrival month and calculate the average cancellation rate
monthly_cancellation_rate = df.groupby('arrival_date_month')['cancellation_status'].mean()
# Ensure months are in chronological order
import calendar
monthly_cancellation_rate.index = pd.CategoricalIndex(monthly_cancellation_rate.index, categories=calendar.month_name[1:], ordered=True)
monthly_cancellation_rate = monthly_cancellation_rate.sort_index()
```

```
In [57]: # Example data for monthly cancellation rates
data = {
    'Month': ['January', 'February', 'March', 'April', 'May', 'June', 'July', 'August', 'September', 'October', 'November', 'December'],
    'Rate': [0.11, 0.10, 0.09, 0.08, 0.07, 0.06, 0.05, 0.04, 0.03, 0.02, 0.01, 0.00]
}
monthly_cancellation_rate = pd.Series(data)

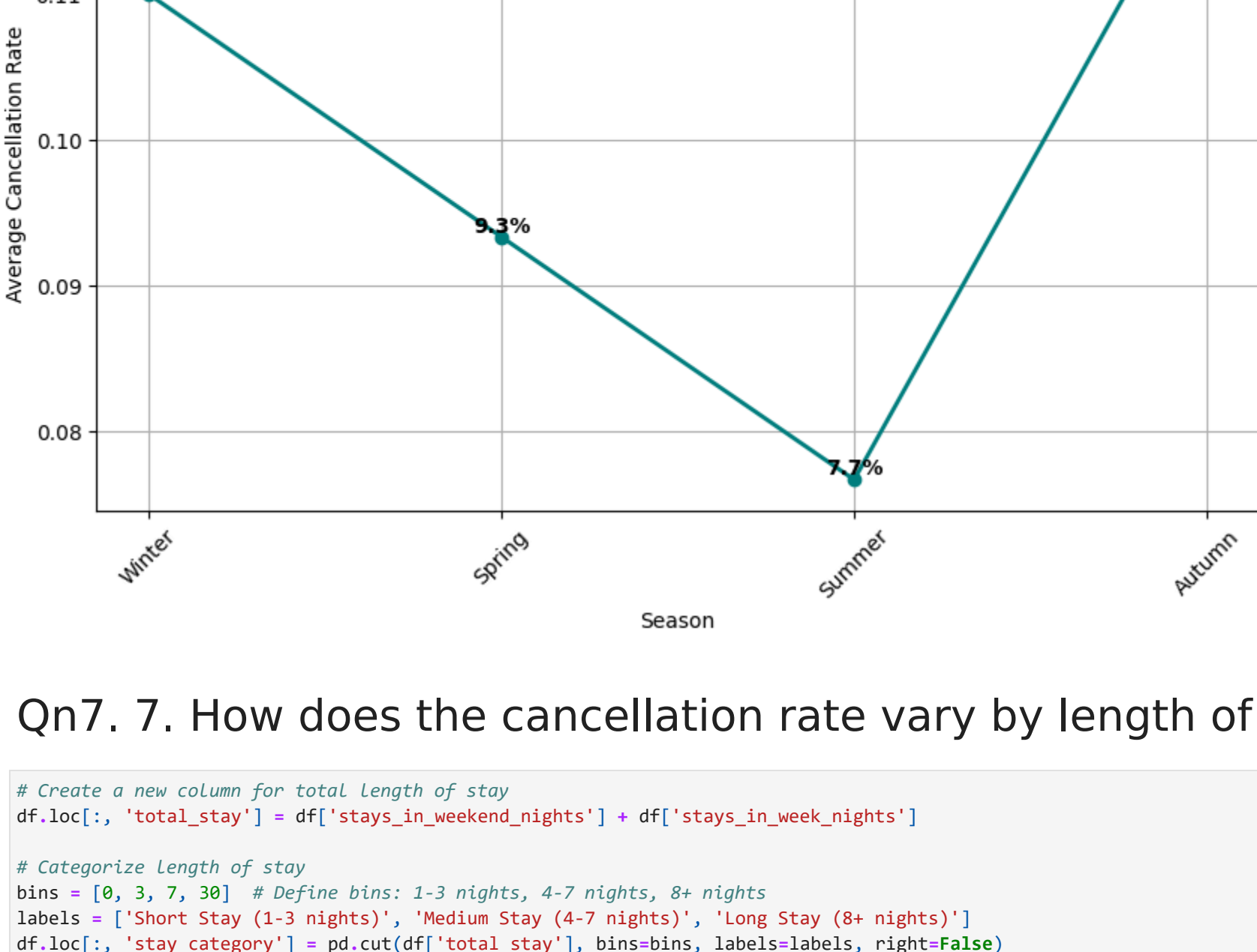
# Create a mapping of months to seasons
season_map = {
    'January': 'Winter', 'February': 'Winter', 'March': 'Spring',
    'April': 'Spring', 'May': 'Spring', 'June': 'Summer',
    'July': 'Summer', 'August': 'Summer', 'September': 'Autumn',
    'October': 'Autumn', 'November': 'Autumn', 'December': 'Winter'
}

# Map months to seasons and calculate average cancellation rate by season
monthly_cancellation_rate.index = pd.Index(monthly_cancellation_rate.index.map(season_map))
seasonal_cancellation_rate = monthly_cancellation_rate.groupby(monthly_cancellation_rate.index).mean()
```

```
# Sorting seasons for proper plotting order
season_order = ['Winter', 'Spring', 'Summer', 'Autumn']
seasonal_cancellation_rate = seasonal_cancellation_rate.reindex(season_order)

# Plotting
plt.figure(figsize=(10, 6))
plt.plot(seasonal_cancellation_rate.index, seasonal_cancellation_rate.values, marker='o', color='teal', linestyle='-', linewidth=2)

# Adding labels at each point
for i, value in enumerate(seasonal_cancellation_rate.values):
    plt.text(i, value, f'{value:.1f}', ha='center', va='bottom', fontweight='bold')
```



Qn7. 7. How does the cancellation rate vary by length of stay?

```
In [70]: # Create a new column for total length of stay
df['loc', 'total_stay'] = df['stays_in_weekend_nights'] + df['stays_in_week_nights']

# Categorize length of stay
bins = [0, 3, 7, 30] # Define bins: 1-3 nights, 4-7 nights, 8+ nights
labels = ['Short Stay (1-3 nights)', 'Medium Stay (4-7 nights)', 'Long Stay (8+ nights)']
df['loc', 'stay_category'] = pd.cut(df['total_stay'], bins=bins, labels=labels, right=False)

# Calculate cancellation rates by length of stay category with observed parameter
cancellation_rates = df.groupby('stay_category', observed=True)['is_canceled'].mean() * 100 # Multiply by 100 for percentage
```

```
# Plotting
plt.figure(figsize=(10, 6))
bars = plt.bar(cancellation_rates.index, cancellation_rates.values, color='teal')
plt.bar_label(bars, labels=[f'{val:.1f}' for val in cancellation_rates.values], label_type='edge', fontweights='bold')

# Labeling the chart
plt.xlabel('Length of Stay Category')
plt.ylabel('Cancellation Rate (%)')
plt.title('Cancellation Rate by Length of Stay')
plt.xticks(rotation=18)
plt.ylim(100) # Set y-axis limit from 0% to 100%
plt.grid(axis='y')

# Show the chart
plt.show()
```