

Business Problem

In recent years, City Hotel and Resort Hotel have seen high cancellation rates. Each hotel is now dealing with a number of issues as a result, including fewer revenues and less than ideal hotel room use. Consequently, lowering cancellation rates is both hotels' primary goal in order to increase their efficiency in generating revenue and for us to offer thorough business advice to address this problem.

Research Question

1. What are the variables that affect hotel reservation cancellations?
2. How can we make hotel reservations cancellations better?
3. How will hotels be assisted in making pricing and promotional decisions?

Importing Libraries

```
In [40]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
```

Loading the Dataset

```
In [29]: df = pd.read_csv('C:/Users/DAVINA/Downloads/hotel.csv')
```

Exploratory Data Analysis and Data Cleaning

```
In [13]: df.head()
Out[13]:
```

	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	type
0	Resort Hotel	0	342	2015	July	27	1	0	0	0	2	...	No Deposit
1	Resort Hotel	0	737	2015	July	27	1	0	0	0	2	...	No Deposit
2	Resort Hotel	0	7	2015	July	27	1	0	0	1	1	...	No Deposit
3	Resort Hotel	0	13	2015	July	27	1	0	0	1	1	...	No Deposit
4	Resort Hotel	0	14	2015	July	27	1	0	0	2	2	...	No Deposit

5 rows × 32 columns

```
In [14]: df.tail()
Out[14]:
```

	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	type
119385	City Hotel	0	23	2017	August	35	30	2	5	2	...	No Deposit	
119386	City Hotel	0	102	2017	August	35	31	2	5	3	...	No Deposit	
119387	City Hotel	0	34	2017	August	35	31	2	5	3	...	No Deposit	
119388	City Hotel	0	109	2017	August	35	31	2	5	2	...	No Deposit	
119389	City Hotel	0	205	2017	August	35	29	2	7	2	...	No Deposit	

5 rows × 32 columns

```
In [15]: df.shape
Out[15]: (119390, 32)
```

```
In [16]: df.columns
Out[16]:
```

Index(['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', 'children', 'babies', 'meal', 'country', 'market_segment', 'distribution_channel', 'is_repeated_guest', 'previous_cancellations', 'previous_bookings_not_canceled', 'reserved_room_type', 'assigned_room_type', 'booking_changes', 'deposit_type', 'agent', 'company', 'days_in_waiting_list', 'customer_type', 'adr', 'required_car_parking_spaces', 'total_of_special_requests', 'reservation_status', 'reservation_status_date', 'dtype=object']

```
In [18]: df.info()
Out[18]:
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 119390 entries, 0 to 119389
Data columns (total 32 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   adults                119390 non-null int64
10  children              119390 non-null float64
11  babies                119390 non-null int64
12  meal                  119390 non-null object
13  country               118992 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 object
21  booking_changes        119390 non-null int64
22  deposit_type           119390 non-null object
23  agent                  103958 non-null float64
24  company                7797 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
Memory usage: 29.1+ MB
```

```
In [105]: df['reservation_status_date'] = pd.to_datetime(df['reservation_status_date'])
```

```
In [22]: df.describe(include = 'object')
Out[22]:
```

	hotel	arrival_date_month	meal	country	market_segment	distribution_channel	reserved_room_type	assigned_room_type	deposit_type	customer_type	reservation_status
count	119390	119390	119390	118992	119390	119390	119390	119390	119390	119390	119390
unique	2	12	5	177	8	5	10	12	3	4	3
top	CityHotel	August	BB	PRT	Online TA	TA/TO	A	A	No Deposit	Transient	Check-Out
freq	78330	13877	92310	48590	56477	97870	85994	74053	104641	89613	75166

```
In [29]: for col in df.describe(include = 'object').columns:
print(col)
print(df[col].unique())
print("-"*56)

hotel
-----
'ResortHotel' 'CityHotel'

arrival_date_month
-----
'July' 'August' 'September' 'October' 'November' 'December' 'January'
'February' 'March' 'April' 'May' 'June'

meal
-----
['BB' 'HB' 'HB' 'SC' 'Undefined']

country
-----
'PRT' 'GBR' 'USA' 'ESP' 'IRL' 'FRA' 'nan' 'ROU' 'NOR' 'DNK' 'ARG' 'POL'
'GBR' 'GBR' 'CHE' 'CH' 'GRC' 'ITA' 'IND' 'DNK' 'RUS' 'SWE' 'AUS' 'EST'
'CZE' 'BRA' 'FIN' 'MOZ' 'BWA' 'LUX' 'SVN' 'ALB' 'IND' 'CHN' 'MEX' 'MAR'
'UKR' 'SMR' 'LVA' 'PRI' 'SRB' 'CHL' 'AUT' 'BLR' 'LTU' 'TUR' 'ZAF' 'AGO'
'CHN' 'CYP' 'ZMB' 'CPV' 'ZWE' 'DZA' 'KOR' 'CRI' 'HUN' 'ARE' 'TUR' 'JAM'
'HRV' 'HKG' 'IRN' 'GEO' 'AND' 'GIB' 'URY' 'JEY' 'CAF' 'CYP' 'COL' 'GGV'
'KWT' 'NGA' 'MDV' 'VEN' 'SVK' 'FJI' 'KAZ' 'PAK' 'IDN' 'LBN' 'PHL' 'SEN'
'AZE' 'BHR' 'MLT' 'TUN' 'DOM' 'MDG' 'MYS' 'ARM' 'JPN' 'LKA' 'CUB'
'OMR' 'BHR' 'MUS' 'COM' 'SUR' 'UGA' 'BGR' 'CIV' 'JOR' 'SYR' 'SGP' 'BOT'
'SAU' 'VNM' 'PLW' 'QAT' 'EGY' 'PER' 'MLI' 'HND' 'ECU' 'MDO' 'ISL' 'UZB'
'BGD' 'BGR' 'MWC' 'TGO' 'TWN' 'DJI' 'STP' 'KNA' 'ETM' 'IRQ' 'HND' 'RWA'
'KEN' 'MCO' 'BGD' 'JMN' 'TJK' 'NIC' 'BEN' 'VGB' 'TZA' 'GAB' 'GHA' 'TMP'
'KEN' 'LIE' 'CHN' 'HNE' 'LUX' 'HVT' 'ROU' 'HND' 'PAN' 'BSC' 'LBY'
'MLI' 'NAM' 'BOL' 'PRY' 'BRB' 'ABW' 'ALA' 'SLV' 'DMA' 'PYF' 'GUY' 'LCA'
'ATA' 'GTM' 'ASH' 'HRT' 'NCL' 'KIR' 'SDM' 'ATF' 'SLE' 'LAO'

market_segment
-----
['Direct' 'Corporate' 'Online TA' 'Offline TA/TO' 'Complementary' 'Groups'
'Undefined' 'Aviation']

distribution_channel
-----
['Direct' 'Corporate' 'TA/TO' 'Undefined' 'GDS']

reserved_room_type
-----
['C' 'A' 'P' 'E' 'G' 'F' 'H' 'L' 'P' 'B']

assigned_room_type
-----
['C' 'A' 'P' 'E' 'G' 'F' 'I' 'B' 'H' 'P' 'L' 'K']

deposit_type
-----
['No Deposit' 'Refundable' 'Non Refund']

customer_type
-----
['Transient' 'Contract' 'Transient-Party' 'Group']

reservation_status
-----
['Check-Out' 'Canceled' 'No-Show']
```

```
In [30]: df.isnull().sum()
Out[30]:
```

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	children	babies	meal	country	market_segment	distribution_channel	is_repeated_guest	previous_cancellations	previous_bookings_not_canceled	reserved_room_type	assigned_room_type	booking_changes	deposit_type	agent	company	days_in_waiting_list	customer_type	adr	required_car_parking_spaces	total_of_special_requests	reservation_status	reservation_status_date	dtype
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	

```
In [31]: df.drop(['company','agent'],axis = 1, inplace = True)
Out[31]:
```

```
In [32]: df.isnull().sum()
Out[32]:
```

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	children	babies	meal	country	market_segment	distribution_channel	is_repeated_guest	previous_cancellations	previous_bookings_not_canceled	reserved_room_type	assigned_room_type	booking_changes	deposit_type	days_in_waiting_list	customer_type	adr	required_car_parking_spaces	total_of_special_requests	reservation_status	reservation_status_date	dtype
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

```
In [33]: df.describe()
Out[33]:
```

	is_canceled	lead_time	arrival_date_year	arrival_date_week_number	arrival_date_day_of_month	stays_in_weekend_nights	stays_in_week_nights	adults	children	babies
count	118898.000000	118898.000000	118898.000000	118898.000000	118898.000000	118898.000000	118898.000000	118898.000000	118898.000000	118898.0000
mean	0.371352	104.311435	2016.157656	27.166555	15.800890	0.928897	2.502145	1.858391	0.104207	0.0079
std	0.483168	106.903309	0.707459	13.589971	8.780324	0.996216	1.900168	0.578576	0.399172	0.0973
min	0.000000	0.000000	2015.000000	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0000
25%	0.000000	18.000000	2016.000000	16.000000	8.000000	0.000000	1.000000	2.000000	0.000000	0.0000
50%	0.000000	69.000000	2016.000000	28.000000	16.000000	1.000000	2.000000	2.000000	0.000000	0.0000
75%	1.000000	161.000000	2017.000000	38.000000	23.000000	2.000000	3.000000	2.000000	0.000000	0.0000
max	1.000000	737.000000	2017.000000	53.000000	31.000000	16.000000	41.000000	55.000000	10.000000	10.0000

```
In [41]: df['adr'].plot(kind = 'box')
Out[41]:
```

<Axes: >

```
In [42]: df[df['adr'] <= 5000]
```

Data Analysis and Visualizations

```
In [52]: cancelled_perc = df['is_canceled'].value_counts(normalize = True)
print(cancelled_perc)
plt.figure(figsize = (5,4))
plt.title('Reservation Status count')
plt.bar(['not cancelled', 'cancelled'],df['is_canceled'].value_counts(),edgecolor = 'k',width = 0.7)
plt.show()
```

```
0    0.628653
1    0.371347
Name: is_canceled, dtype: float64
```

Reservation Status count

```
In [83]: plt.figure(figsize = (8,4))
ax1 = sns.countplot(x = 'hotel',hue = 'is_canceled',data = df,palette = 'Blues')
plt.title('Reservation Status count')
ax1.legend(bbox_to_anchor=(1,1))
plt.xlabel('Reservation Status in different Hotels', size = 20)
plt.ylabel('Number of Reservations')
plt.legend(['not cancelled', 'cancelled'])
plt.show()
```

Reservation Status in different Hotels

From the bar graph, it shows the percentage of reservations that are canceled and those that are not. It is obvious that there are still a significant number of reservations that have not been cancelled. There are still 37% of clients who cancelled their reservation, which has a significant impact on the hotel's earnings.

In comparison to resort hotels, city hotels have more bookings. It's possible that resort hotels are more expensive than those in cities.

```
In [71]: resort_hotel = df[df['hotel']=='Resort Hotel']
resort_hotel['is_canceled'].value_counts(normalize = True)
Out[71]:
```

	is_canceled	dtype
0	0.72825	float64
1	0.27175	float64

```
In [72]: city_hotel = df[df['hotel']=='City Hotel']
city_hotel['is_canceled'].value_counts(normalize = True)
Out[72]:
```

	is_canceled	dtype
0	0.582928	float64
1	0.417072	float64

```
In [79]: resort_hotel = resort_hotel.groupby('reservation_status_date')['adr'].mean()
city_hotel = city_hotel.groupby('reservation_status_date')['adr'].mean()
In [106]: df['month'] = df['reservation_status_date'].dt.month
plt.figure(figsize = (16,8))
ax1 = sns.countplot(x = 'month',hue = 'is_canceled',data = df, palette = 'bright')
legend_labels = ax1.get_legend_handles_labels()
ax1.legend(bbox_to_anchor = (1,1))
plt.title('Reservation Status per Month', size = 20)
plt.xlabel('Month')
plt.ylabel('number of reservation')
plt.legend(['not cancelled', 'cancelled'])
plt.show()
```

Reservation Status per Month

From above, we have developed the grouped bar graph to analyze the months with the highest and lowest reservation levels according to reservation status. As can be seen, both the number of confirmed reservations and the number of cancelled reservations are largest in the month of August. Whereas January is the month with the most cancelled reservations.

```
In [117]: plt.figure(figsize=(15, 8))
plt.title('ADR per month for Canceled Reservations', fontsize=30)
data_to_plot = df[df['is_canceled'] == 1].groupby('month')['adr'].reset_index()
sns.barplot(x='month', y='adr', data=data_to_plot)
plt.xlabel('Month', fontsize=20)
plt.ylabel('ADR', fontsize=20)
plt.show()
```

ADR per month for Canceled Reservations

This bar graph demonstrates that cancellations are most common when prices are greatest and are least when they are lowest. Therefore, the cost of the accommodation is solely responsible for the cancellations.

```
In [118]: cancelled_data = df[df['is_canceled']==1]
top_10_country = cancelled_data['country'].value_counts()[0:10]
plt.figure(figsize = (8,8))
plt.title('Top 10 Countries with Reservation Cancelled')
plt.pie(top_10_country, autopct = '%.2f', labels= top_10_country.index)
plt.show()
```

Top 10 Countries with Reservation Cancelled

Taking a look at this pie chart, we can see which country has the highest reservation cancelled. The top country is Portugal with the highest number of cancellations.

Let's also look at the area from where guests are visiting the hotels and making reservations. Is it coming from Direct or Groups, Online or Official travel agents? Around 46% of the clients book hotels directly by visiting them and making reservations.

```
In [119]: df['market_segment'].value_counts()
Out[119]:
```

	market_segment	count
Online TA	56082	
Offline TA/TO	24159	
Groups	19806	
Direct	12448	
Corporate	5111	
Complementary	734	
Aviation	237	
Name: market_segment, dtype: int64		

```
In [120]: df['market_segment'].value_counts(normalize = True)
Out[120]:
```

	market_segment	count
Online TA	0.474377	
Offline TA/TO	0.203193	
Groups	0.166581	
Direct	0.104696	
Corporate	0.042987	
Complementary	0.006172	
Aviation	0.001993	
Name: market_segment, dtype: float64		

```
In [121]: cancelled_data['market_segment'].value_counts(normalize = True)
Out[121]:
```

	market_segment	count
Online TA	0.469696	
Groups	0.273985	
Offline TA/TO	0.167446	
Direct	0.043486	
Corporate	0.022151	
Complementary	0.002038	
Aviation	0.00178	
Name: market_segment, dtype: float64		

```
In [ ]:
In [134]: cancelled_df_adr = cancelled_data.groupby('reservation_status_date')['adr'].mean()
cancelled_df_adr.reset_index(inplace = True)
cancelled_df_adr['reservation_status_date'] = pd.to_datetime(cancelled_df_adr['reservation_status_date'])
not_cancelled_df_adr = not_cancelled_data.groupby('reservation_status_date')['adr'].mean()
not_cancelled_df_adr.reset_index(inplace = True)
not_cancelled_df_adr['reservation_status_date'] = pd.to_datetime(not_cancelled_df_adr['reservation_status_date'])
plt.figure(figsize = (20, 6))
plt.title('Average Daily Rate', fontsize=30)
plt.plot(not_cancelled_df_adr['reservation_status_date'],not_cancelled_df_adr['adr'], label = 'not cancelled')
plt.plot(cancelled_df_adr['reservation_status_date'],cancelled_df_adr['adr'], label = 'cancelled')
plt.xlabel('ADR', fontsize=20)
plt.ylabel('Date', fontsize=20)
plt.legend(fontsize = 8)
plt.show()
```

Average Daily Rate

As seen in the graph, reservations are cancelled when the average daily rate is higher than when it is not cancelled. It clearly proves all the above analysis, that the higher the price the more the cancellations.

Suggestions

1. Cancellation rates rise as the price does. In order to prevent cancellations of reservations, hotels could work on their pricing strategies and try to lower the rates for specific hotels based on locations. They can also provide some discounts to the consumers.

2. As the ratio of the cancellation and not cancellation of the resort hotel is higher in the resort hotel than the city hotels. So the hotels should provide a reasonable discount on the room prices on weekends or on holidays.

3. In the month of January, hotels can start campaigns or marketing with a holiday amount to increase their revenue as the cancellation is the highest in this month

4. They can also increase the quality of their hotels and their services mainly in Portugal to reduce the cancellation rate.

```
In [ ]:
In [ ]:
```