# Hotel bookings cancellation

## 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.

The analysis of hotel booking cancellations as well as other factors that have no bearing on their business and yearly revenue generation are the main topics of this report.

## Assumptions

1. No unusual occurrences between 2015 and 2017 will have a substantial impact on the data used.
2. The information is still current and can be used to analyze a hotel's possible plans in an efficient manner.
3. There are no unanticipated negatives to the hotel employing any advised technique.
4. The hotels are not currently using any of the suggested solutions.
5. The biggest factor affecting the effectiveness of earning income is booking cancellations.
6. Cancellations result in vacant rooms for the booked length of time.
7. Clients make hotel reservations the same year they make cancellations.

## 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?

## Hypothesis

1. More cancellations occur when prices are higher.
2. When there is a longer waiting list, customers tend to cancel more frequently.
3. The majority of clients are coming from offline travel agents to make their reservations.

# import libraries

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
```

# import data set

In [2]:

```python
data= pd.read_csv(r'Downloads/archive/hotel_booking.csv')
```
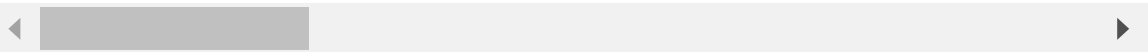
In [3]:

```python
data.head()
```

Out[3]:

| | hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | arrival_date_week_nu |
|---|---|---|---|---|---|---|
| 0 | Resort Hotel | 0 | 342 | 2015 | July | |
| 1 | Resort Hotel | 0 | 737 | 2015 | July | |
| 2 | Resort Hotel | 0 | 7 | 2015 | July | |
| 3 | Resort Hotel | 0 | 13 | 2015 | July | |
| 4 | Resort Hotel | 0 | 14 | 2015 | July | |

5 rows × 36 columns

```
1  data.tail()
```

| | hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | arrival_date_wee |
|---|---|---|---|---|---|---|
| **119385** | City Hotel | 0 | 23 | 2017 | August | |
| **119386** | City Hotel | 0 | 102 | 2017 | August | |
| **119387** | City Hotel | 0 | 34 | 2017 | August | |
| **119388** | City Hotel | 0 | 109 | 2017 | August | |
| **119389** | City Hotel | 0 | 205 | 2017 | August | |

5 rows × 36 columns

# analysis and cleaning

```
1  data.head(10)
```

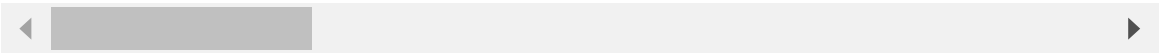| | hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | arrival_date_week_n... |
|---|---|---|---|---|---|---|
| 0 | Resort Hotel | 0 | 342 | 2015 | July | |
| 1 | Resort Hotel | 0 | 737 | 2015 | July | |
| 2 | Resort Hotel | 0 | 7 | 2015 | July | |
| 3 | Resort Hotel | 0 | 13 | 2015 | July | |
| 4 | Resort Hotel | 0 | 14 | 2015 | July | |
| 5 | Resort Hotel | 0 | 14 | 2015 | July | |
| 6 | Resort Hotel | 0 | 0 | 2015 | July | |
| 7 | Resort Hotel | 0 | 9 | 2015 | July | |
| 8 | Resort Hotel | 1 | 85 | 2015 | July | |
| 9 | Resort Hotel | 1 | 75 | 2015 | July | |

10 rows × 36 columns

In [6]:

```
1  data.tail(10)
```

Out[6]:

| | hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | arrival_date_wee |
|---|---|---|---|---|---|---|
| 119380 | City Hotel | 0 | 44 | 2017 | August | |
| 119381 | City Hotel | 0 | 188 | 2017 | August | |
| 119382 | City Hotel | 0 | 135 | 2017 | August | |
| 119383 | City Hotel | 0 | 164 | 2017 | August | |
| 119384 | City Hotel | 0 | 21 | 2017 | August | |
| 119385 | City Hotel | 0 | 23 | 2017 | August | |
| 119386 | City Hotel | 0 | 102 | 2017 | August | |
| 119387 | City Hotel | 0 | 34 | 2017 | August | |
| 119388 | City Hotel | 0 | 109 | 2017 | August | |
| 119389 | City Hotel | 0 | 205 | 2017 | August | |

10 rows × 36 columns

◀ ▶

In [7]:

```
1  data.shape
```

Out[7]:

(119390, 36)

# clean data / removing

In [8]:

```python
# removing personal information in data
data.drop(['name','email','phone-number','credit_card'],axis = 1 ,inplace = True )
```
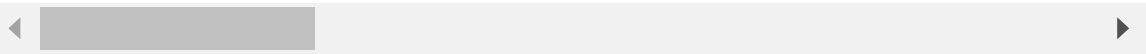
In [9]:

```python
data.head()
```

Out[9]:

| | hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | arrival_date_week_n |
|---|---|---|---|---|---|---|
| 0 | Resort Hotel | 0 | 342 | 2015 | July | |
| 1 | Resort Hotel | 0 | 737 | 2015 | July | |
| 2 | Resort Hotel | 0 | 7 | 2015 | July | |
| 3 | Resort Hotel | 0 | 13 | 2015 | July | |
| 4 | Resort Hotel | 0 | 14 | 2015 | July | |

5 rows × 32 columns

In [10]:

```python
data.shape
```

Out[10]:

```
(119390, 32)
```

In [11]:

```python
data.columns
```

Out[11]:

```
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 [12]:

```
1  data.info()
```

```
<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                        119386 non-null  float64
 11  babies                          119390 non-null  int64
 12  meal                            119390 non-null  object
 13  country                         118902 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                           103050 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(4), int64(16), object(12)
memory usage: 29.1+ MB
```

In [13]:

```
1  data['reservation_status_date'] = pd.to_datetime(data['reservation_status_date'])
```

```
1  data.info()
```

```
<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                        119386 non-null  float64
 11  babies                          119390 non-null  int64
 12  meal                            119390 non-null  object
 13  country                         118902 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                           103050 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(4), int64(16), object(11)
memory usage: 29.1+ MB
```

In [15]:

```
1  data.describe(include= 'object')
```

Out[15]:

| | hotel | arrival_date_month | meal | country | market_segment | distribution_channel |
|---|---|---|---|---|---|---|
| count | 119390 | 119390 | 119390 | 118902 | 119390 | 119390 |
| unique | 2 | 12 | 5 | 177 | 8 | 5 |
| top | City Hotel | August | BB | PRT | Online TA | TA/TO |
| freq | 79330 | 13877 | 92310 | 48590 | 56477 | 97870 |

```python
for col in data.describe(include= 'object').columns:
    print(col)
    print(data[col].unique())
    print('-'*50)
```

```
hotel
['Resort Hotel' 'City Hotel']
--------------------------------------------------
arrival_date_month
['July' 'August' 'September' 'October' 'November' 'December' 'January'
 'February' 'March' 'April' 'May' 'June']
--------------------------------------------------
meal
['BB' 'FB' 'HB' 'SC' 'Undefined']
--------------------------------------------------
country
['PRT' 'GBR' 'USA' 'ESP' 'IRL' 'FRA' nan 'ROU' 'NOR' 'OMN' 'ARG' 'POL'
 'DEU' 'BEL' 'CHE' 'CN' 'GRC' 'ITA' 'NLD' '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'
 'ISR' 'CYM' 'ZMB' 'CPV' 'ZWE' 'DZA' 'KOR' 'CRI' 'HUN' 'ARE' 'TUN' 'JAM'
 'HRV' 'HKG' 'IRN' 'GEO' 'AND' 'GIB' 'URY' 'JEY' 'CAF' 'CYP' 'COL' 'GGY'
 'KWT' 'NGA' 'MDV' 'VEN' 'SVK' 'FJI' 'KAZ' 'PAK' 'IDN' 'LBN' 'PHL' 'SEN'
 'SYC' 'AZE' 'BHR' 'NZL' 'THA' 'DOM' 'MKD' 'MYS' 'ARM' 'JPN' 'LKA' 'CUB'
 'CMR' 'BIH' 'MUS' 'COM' 'SUR' 'UGA' 'BGR' 'CIV' 'JOR' 'SYR' 'SGP' 'BDI'
 'SAU' 'VNM' 'PLW' 'QAT' 'EGY' 'PER' 'MLT' 'MWI' 'ECU' 'MDG' 'ISL' 'UZB'
 'NPL' 'BHS' 'MAC' 'TGO' 'TWN' 'DJI' 'STP' 'KNA' 'ETH' 'IRQ' 'HND' 'RWA'
 'KHM' 'MCO' 'BGD' 'IMN' 'TJK' 'NIC' 'BEN' 'VGB' 'TZA' 'GAB' 'GHA' 'TMP'
 'GLP' 'KEN' 'LIE' 'GNB' 'MNE' 'UMI' 'MYT' 'FRO' 'MMR' 'PAN' 'BFA' 'LBY'
 'MLI' 'NAM' 'BOL' 'PRY' 'BRB' 'ABW' 'AIA' 'SLV' 'DMA' 'PYF' 'GUY' 'LCA'
 'ATA' 'GTM' 'ASM' 'MRT' 'NCL' 'KIR' 'SDN' 'ATF' 'SLE' 'LAO']
--------------------------------------------------
market_segment
['Direct' 'Corporate' 'Online TA' 'Offline TA/TO' 'Complementary' 'Group
s'
 'Undefined' 'Aviation']
--------------------------------------------------
distribution_channel
['Direct' 'Corporate' 'TA/TO' 'Undefined' 'GDS']
--------------------------------------------------
reserved_room_type
['C' 'A' 'D' 'E' 'G' 'F' 'H' 'L' 'P' 'B']
--------------------------------------------------
assigned_room_type
['C' 'A' 'D' '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']
--------------------------------------------------
```

```
1  data.isnull().sum()
```

```
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
meal                              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                         16340
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
```

```
1  data.drop(['company','agent'],axis = 1 ,inplace = True )
2  data.dropna(inplace = True)
```

```
1  data.isnull().sum()
```

```
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
meal                             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
```

```
1  data.describe()
```

|       | is_canceled   | lead_time     | arrival_date_year | arrival_date_week_number | arrival_da |
|-------|---------------|---------------|-------------------|--------------------------|------------|
| count | 118898.000000 | 118898.000000 | 118898.000000     | 118898.000000            |            |
| mean  | 0.371352      | 104.311435    | 2016.157656       | 27.166555                |            |
| std   | 0.483168      | 106.903309    | 0.707459          | 13.589971                |            |
| min   | 0.000000      | 0.000000      | 2015.000000       | 1.000000                 |            |
| 25%   | 0.000000      | 18.000000     | 2016.000000       | 16.000000                |            |
| 50%   | 0.000000      | 69.000000     | 2016.000000       | 28.000000                |            |
| 75%   | 1.000000      | 161.000000    | 2017.000000       | 38.000000                |            |
| max   | 1.000000      | 737.000000    | 2017.000000       | 53.000000                |            |

```
1  data= data[data['adr']<5000]
```

# data analysis and visualization

In [22]:

```
1  cancelled_perc = data['is_canceled'].value_counts(normalize = True)
```

In [23]:

```
1  cancelled_perc
```
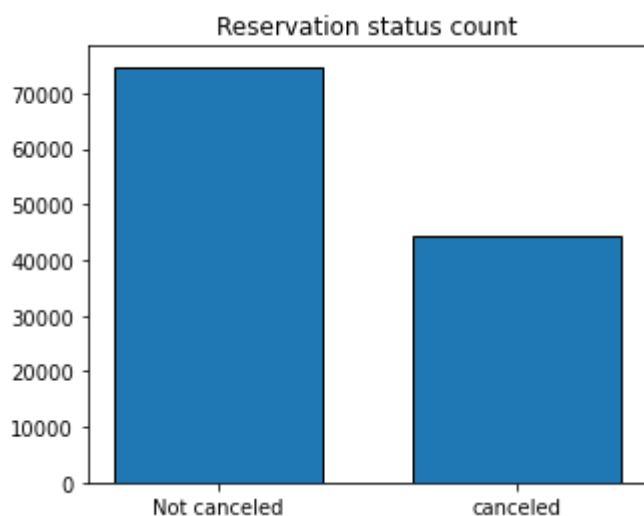
Out[23]:

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

## Analysis and Findings

In [24]:

```
1  cancelled_perc = data['is_canceled'].value_counts(normalize = True)
2  print(cancelled_perc)
3
4  plt.figure(figsize=(5,4))
5  plt.title('Reservation status count')
6  plt.bar(['Not canceled', 'canceled'],data['is_canceled'].value_counts(), edgecolor
7  plt.show()
```

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

**The accompanying bar graph shows the percentage of reservations that are canceled and those that are not. It is obvious that there are stil a significant number of reservations that have not been canceled. There are still 37% of clients who canceled their reservation, which has a significant impact on the hotels' earnings.**

In [25]:

```python
1  plt.figure(figsize= (8,4))
2  ax1 = sns.countplot(x = 'hotel', hue = 'is_canceled',data = data, palette = 'Blues'
3  legend_labels,_= ax1. get_legend_handles_labels()
4  ax1.legend(bbox_to_anchor=(1,1))
5  plt.title('Reservation status in different hotel',size = 20)
6  plt.xlabel('hotel')
7  plt.ylabel('number of reservation')
8  plt.legend(['not canceled','canceled'])
9  plt.show()
```



In [26]:

```python
1  resort_hotel = data[data['hotel'] == 'Resort Hotel']
2  resort_hotel['is_canceled'].value_counts(normalize = True)
```

Out[26]:

```
0    0.72025
1    0.27975
Name: is_canceled, dtype: float64
```

In [27]:

```python
1  city_hotel = data[data['hotel'] == 'City Hotel']
2  city_hotel['is_canceled'].value_counts(normalize = True)
```

Out[27]:

```
0    0.582918
1    0.417082
Name: is_canceled, dtype: float64
```
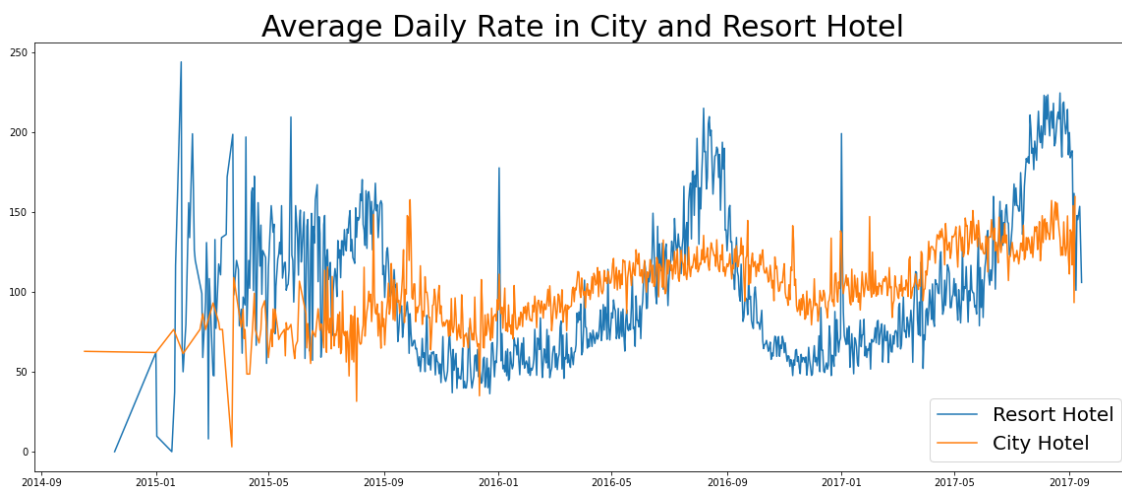
```
1  resort_hotel = resort_hotel.groupby('reservation_status_date')[['adr']].mean()
2  city_hotel = city_hotel.groupby('reservation_status_date')[['adr']].mean()
```

**In comparison to resort hotels, City hotels have more bookings. It's possible that resort hotels are more expensive than those in cities.**
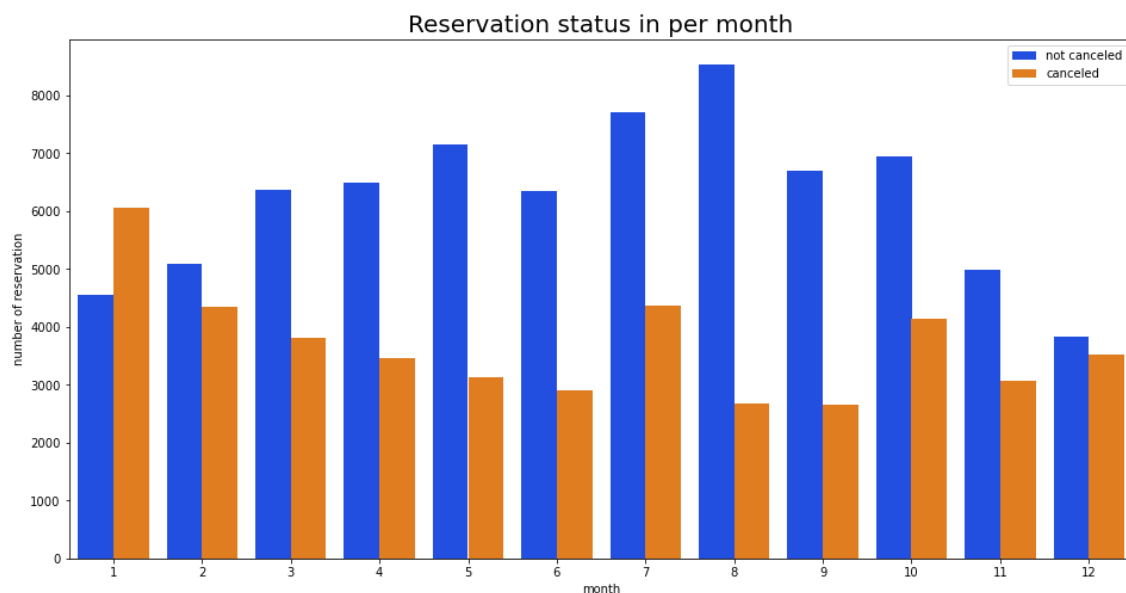
In [29]:

```
1  plt.figure(figsize= (20,8))
2  plt.title('Average Daily Rate in City and Resort Hotel',fontsize = 30)
3  plt.plot(resort_hotel.index, resort_hotel['adr'], label = 'Resort Hotel')
4  plt.plot(city_hotel.index, city_hotel['adr'], label ='City Hotel')
5  plt.legend(fontsize = 20)
6  plt.show()
```

**The line graph above shows that, on certain days, the average daily rate for a city hotel is less than that of a resort hotel, and on other days, it is even less. It goes without saying that weekends and holidays may see a rise in resort hotel rates.**
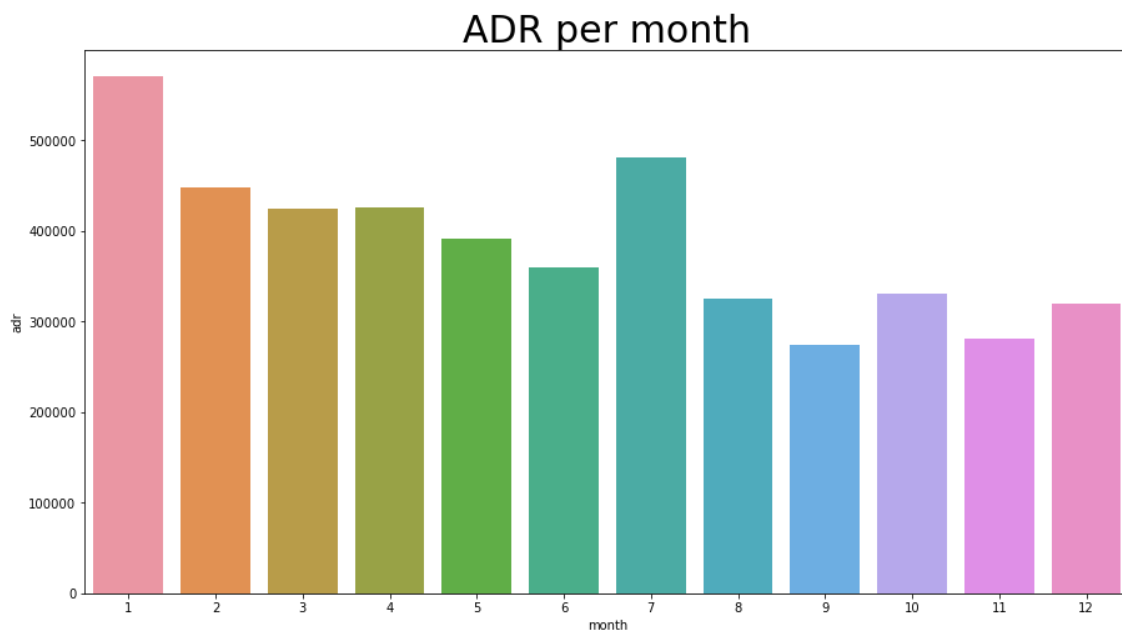
```
 1  data['month'] = data['reservation_status_date'].dt.month
 2  plt.figure(figsize= (16,8))
 3  ax1 = sns.countplot(x = 'month', hue = 'is_canceled', data = data, palette = 'brigh
 4  legend_labels,_= ax1. get_legend_handles_labels()
 5  ax1.legend(bbox_to_anchor=(1,1))
 6  plt.title('Reservation status in per month',size = 20)
 7  plt.xlabel('month')
 8  plt.ylabel('number of reservation')
 9  plt.legend(['not canceled','canceled'])
10  plt.show()
```



Reservation status in per month

**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 canceled reservations are largest in the month of August. whereas January is the month with the most canceled reservations.**

```
1  plt.figure(figsize=(15,8))
2  plt.title('ADR per month',fontsize = 30)
3  sns.barplot('month','adr', data =data[data['is_canceled']==1].groupby('month')[['ad
4  plt.show()
```

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

In [32]:

```
1  cancelled_data = data[data['is_canceled']== 1]
2  top_10_country = cancelled_data['country'].value_counts()[:10]
3  plt.figure(figsize = (10, 8))
4  plt.title('Top 10 country with reservation canceled ')
5  plt.pie(top_10_country, autopct = '%.2f' , labels = top_10_country.index)
6  plt.show()
7
```



Top 10 country with reservation canceled

**Now, let's see which country has the highest reservation canceled. The top country is Portugal with the highest number of cancellations.**

**Let's check the area from where guests are visiting the hotels and making reservations. Is it coming from Direct or Groups, Online or Offline Travel Agents? Around 46% of the clients come from online travel agencies, whereas 27% come from groups. Only 4% of clients book hotels directly by visiting them and making reservations.**

In [35]:

```
1  data['market_segment'].value_counts()
```

Out[35]:

```
Online TA        56402
Offline TA/TO    24159
Groups           19806
Direct           12448
Corporate         5111
Complementary      734
Aviation           237
Name: market_segment, dtype: int64
```

In [36]:

```
1  data['market_segment'].value_counts(normalize= True)
```

Out[36]:

```
Online TA        0.474377
Offline TA/TO    0.203193
Groups           0.166581
Direct           0.104696
Corporate        0.042987
Complementary    0.006173
Aviation         0.001993
Name: market_segment, dtype: float64
```
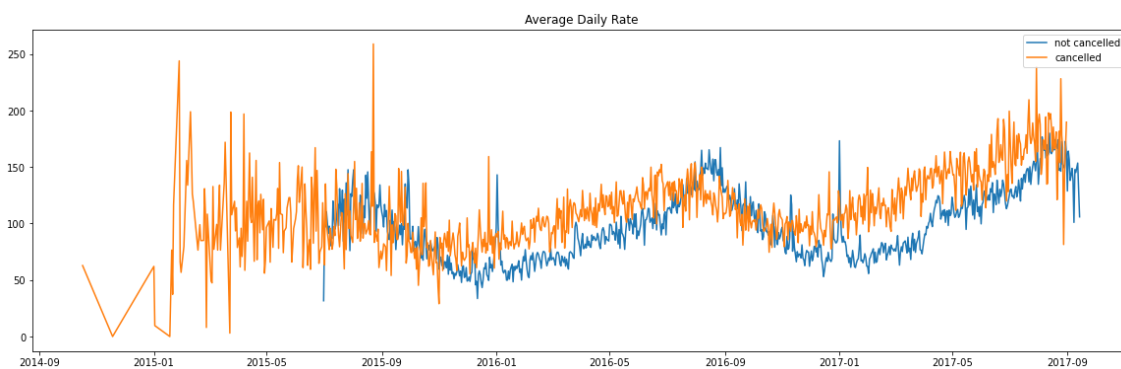
In [37]:

```
1  cancelled_data['market_segment'].value_counts(normalize= True)
```

Out[37]:

```
Online TA        0.469696
Groups           0.273985
Offline TA/TO    0.187466
Direct           0.043486
Corporate        0.022151
Complementary    0.002038
Aviation         0.001178
Name: market_segment, dtype: float64
```

In [52]:

```
1  cancelled_data_adr = cancelled_data.groupby('reservation_status_date')[['adr']].mea
2  cancelled_data_adr.reset_index(inplace = True )
3  cancelled_data_adr.sort_values('reservation_status_date',inplace = True )
4
5  not_cancelled_data = data[data['is_canceled']==0]
6  not_cancelled_data_adr = not_cancelled_data.groupby('reservation_status_date')[['ad
7  not_cancelled_data_adr.reset_index(inplace = True )
8  not_cancelled_data_adr.sort_values('reservation_status_date',inplace = True)
9
10 plt.figure(figsize=(20,6))
11 plt.title('Average Daily Rate')
12 plt.plot(not_cancelled_data_adr['reservation_status_date'],not_cancelled_data_adr['
13 plt.plot(cancelled_data_adr['reservation_status_date'],cancelled_data_adr['adr'], l
14 plt.legend()
15 plt.show()
```
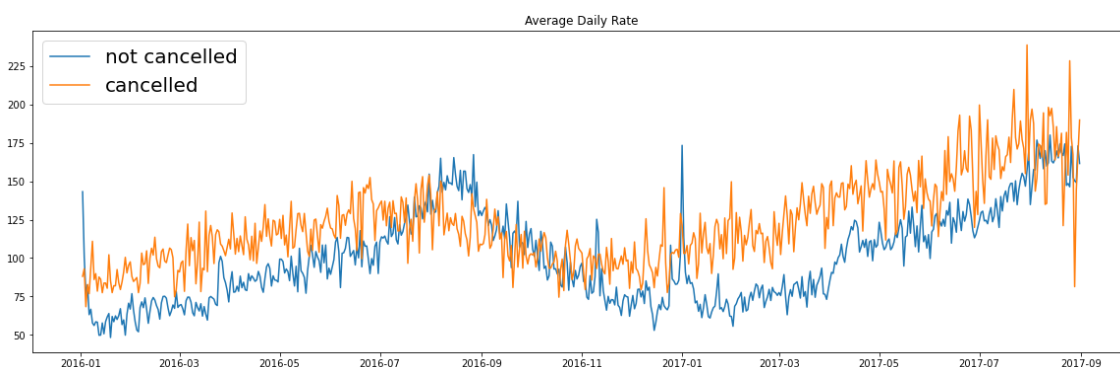


In [53]:

```
1  cancelled_data_adr = cancelled_data_adr[(cancelled_data_adr['reservation_status_dat
2  not_cancelled_data_adr = not_cancelled_data_adr[(not_cancelled_data_adr['reservatio
```

In [55]:

```
1  plt.figure(figsize=(20,6))
2  plt.title('Average Daily Rate')
3  plt.plot(not_cancelled_data_adr['reservation_status_date'],not_cancelled_data_adr['
4  plt.plot(cancelled_data_adr['reservation_status_date'],cancelled_data_adr['adr'], l
5  plt.legend(fontsize = 20)
6  plt.show()
```

**As seen in the graph, reservations are canceled when the average daily rate is higher than when it is not canceled. It clearly proves ail the above analysis, that the higher price leads to higher cancellation.**

## 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 the consumers.
2. As the ratio of the cancellation and not cancellation ofthe resort hotel i 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 reasonable amount to increase their revenue as the cancelation i 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.
5. They can charge minimum amount of booking cancellacancellation
6. Also hotels can provide a coupons for previous customer to discount on next visit

In [ ]:

```
1
```