

In [1]: `!pip install matplotlib seaborn`

```

Defaulting to user installation because normal site-packages is not writeable
Looking in links: /usr/share/pip-wheels
Requirement already satisfied: matplotlib in /opt/conda/envs/anaconda-panel-2023.05-py310/lib/python3.11/site-packages (3.7.2)
Requirement already satisfied: seaborn in /opt/conda/envs/anaconda-panel-2023.05-py310/lib/python3.11/site-packages (0.12.2)
Requirement already satisfied: contourpy>=1.0.1 in /opt/conda/envs/anaconda-panel-2023.05-py310/lib/python3.11/site-packages (from matplotlib) (1.0.5)
Requirement already satisfied: cycler>=0.10 in /opt/conda/envs/anaconda-panel-2023.05-py310/lib/python3.11/site-packages (from matplotlib) (0.11.0)
Requirement already satisfied: fonttools>=4.22.0 in /opt/conda/envs/anaconda-panel-2023.05-py310/lib/python3.11/site-packages (from matplotlib) (4.25.0)
Requirement already satisfied: kiwisolver>=1.0.1 in /opt/conda/envs/anaconda-panel-2023.05-py310/lib/python3.11/site-packages (from matplotlib) (1.4.4)
Requirement already satisfied: numpy>=1.20 in /opt/conda/envs/anaconda-panel-2023.05-py310/lib/python3.11/site-packages (from matplotlib) (1.24.3)
Requirement already satisfied: packaging>=20.0 in /opt/conda/envs/anaconda-panel-2023.05-py310/lib/python3.11/site-packages (from matplotlib) (23.1)
Requirement already satisfied: pillow>=6.2.0 in /opt/conda/envs/anaconda-panel-2023.05-py310/lib/python3.11/site-packages (from matplotlib) (9.4.0)
Requirement already satisfied: pyparsing<3.1,>=2.3.1 in /opt/conda/envs/anaconda-panel-2023.05-py310/lib/python3.11/site-packages (from matplotlib) (3.0.9)
Requirement already satisfied: python-dateutil>=2.7 in /opt/conda/envs/anaconda-panel-2023.05-py310/lib/python3.11/site-packages (from matplotlib) (2.8.2)
Requirement already satisfied: pandas>=0.25 in /opt/conda/envs/anaconda-panel-2023.05-py310/lib/python3.11/site-packages (from seaborn) (2.0.3)
Requirement already satisfied: pytz>=2020.1 in /opt/conda/envs/anaconda-panel-2023.05-py310/lib/python3.11/site-packages (from pandas>=0.25->seaborn) (2023.3.post1)
Requirement already satisfied: tzdata>=2022.1 in /opt/conda/envs/anaconda-panel-2023.05-py310/lib/python3.11/site-packages (from pandas>=0.25->seaborn) (2023.3)
Requirement already satisfied: six>=1.5 in /opt/conda/envs/anaconda-panel-2023.05-py310/lib/python3.11/site-packages (from python-dateutil>=2.7->matplotlib) (1.16.0)

```

In [2]: `import pandas as pd`
`import matplotlib.pyplot as plt`
`import seaborn as sns`
`import warnings`

`warnings.filterwarnings('ignore')`

In [3]: `import os`
`print(os.getcwd())`

`/home/690e04ab-fd8c-4171-bb88-2369ac77d154`

In [4]: `os.listdir()`

```
Out[4]: ['.ipynb_checkpoints',
         '.virtualenvs',
         '.vimrc',
         '.bashrc',
         'hotel_bookings.csv',
         '.jupyter',
         '.local',
         'Untitled1.ipynb',
         'anaconda_projects',
         '.profile',
         '.pythonstartup.py',
         'README.ipynb',
         '.gitconfig',
         'Untitled2.ipynb',
         '.npm',
         '.ipython',
         '.config',
         '.anaconda',
         'Hotel_sales.ipynb',
         '.cache']
```

```
In [5]: df = pd.read_csv("hotel_bookings.csv")
```

```
In [6]: df.head()
```

```
Out[6]:
```

	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 × 32 columns



```
In [7]: df.tail(10)
```

Out[7]:

	hotel	is_canceled	lead_time	arrival_date_year	arrival_date_month	arrival_date_w
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 × 32 columns

In [8]: `df.shape`

Out[8]: (119390, 32)

In [9]: `df.columns`

```
Out[9]: 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 [10]: `df.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 [11]: df['reservation_status_date'] = pd.to_datetime(df['reservation_status_date'], dayfi
```

```
In [12]: df.describe(include = 'object')
```

Out[12]:

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



```
In [13]: for col in df.describe(include = 'object').columns:
          print(col)
          print(df[col].unique())
```

```
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' 'Groups'
 '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']
```

```
In [14]: df.isnull().sum()
```

```

Out[14]: 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

```

```

In [15]: df.drop(['company', 'agent'], axis = 1, inplace = True)
         df.dropna(inplace = True)

```

```

In [16]: df.isnull().sum()

```

```
Out[16]: 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
```

```
In [17]: df.describe()
```

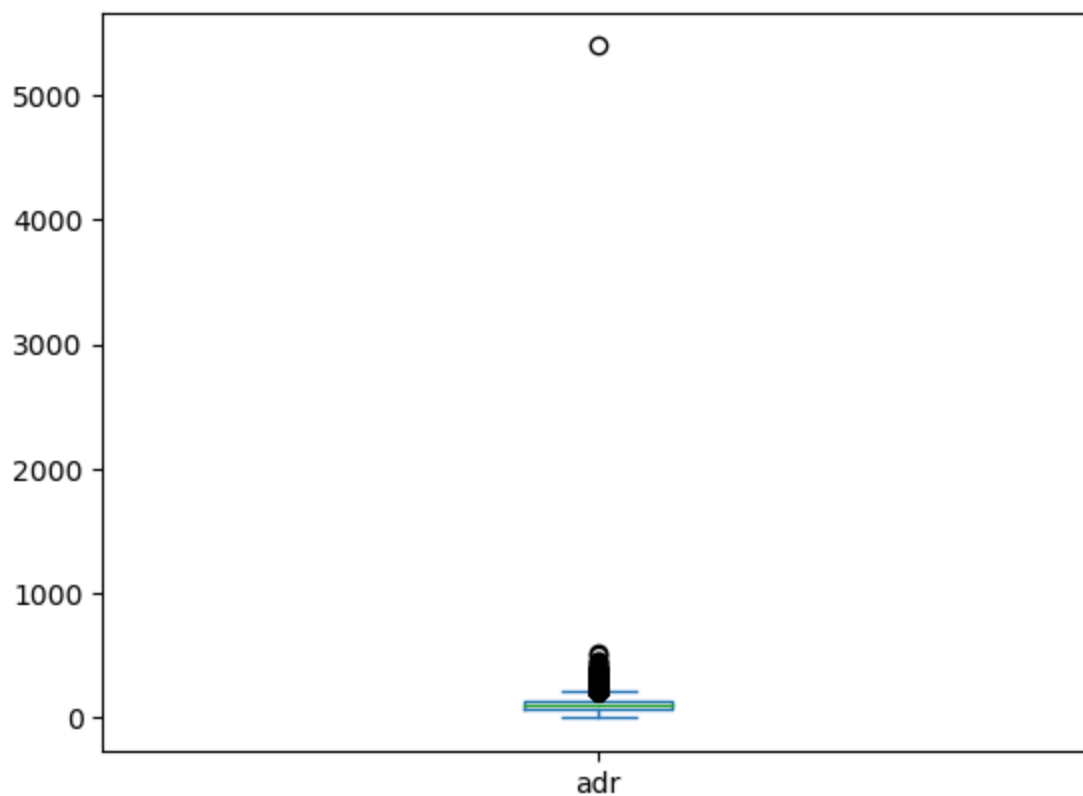
```
Out[17]:
```

	is_canceled	lead_time	arrival_date_year	arrival_date_week_number	arrival_date_day_of_month
count	118898.000000	118898.000000	118898.000000	118898.000000	118898.000000
mean	0.371352	104.311435	2016.157656	27.166555	15.859242
min	0.000000	0.000000	2015.000000	1.000000	1.000000
25%	0.000000	18.000000	2016.000000	16.000000	1.000000
50%	0.000000	69.000000	2016.000000	28.000000	15.000000
75%	1.000000	161.000000	2017.000000	38.000000	15.000000
max	1.000000	737.000000	2017.000000	53.000000	31.000000
std	0.483168	106.903309	0.707459	13.589971	10.256978



```
In [18]: df['adr'].plot(kind = 'box')
```

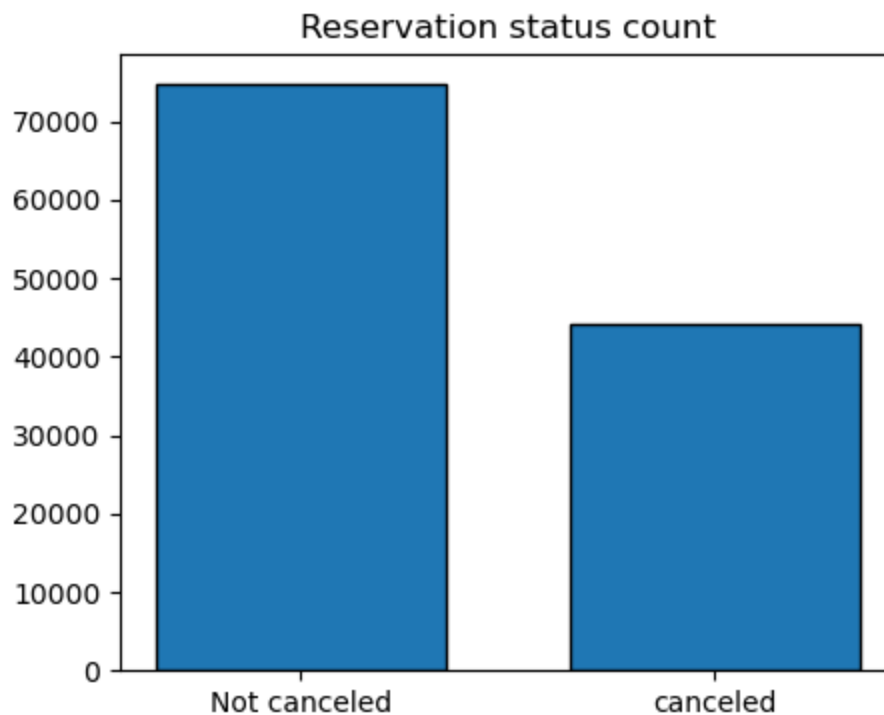
Out[18]: <Axes: >

In [19]: `df = df[df['adr'] < 5000]`

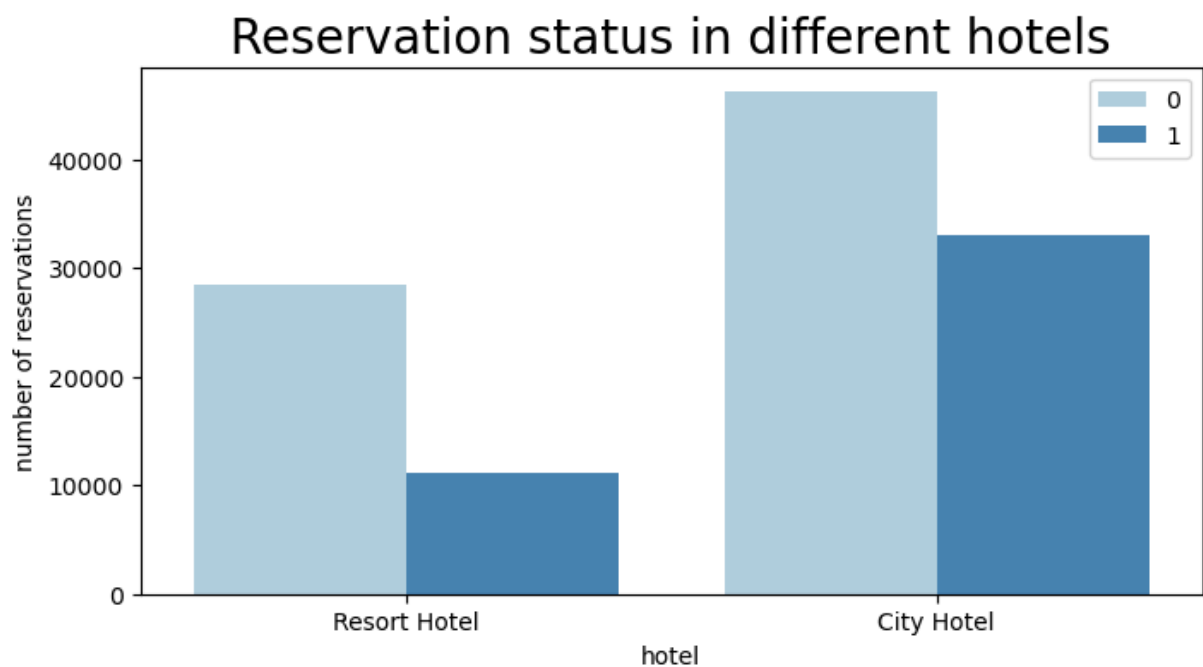
```
In [20]: cancelled_perc = df['is_canceled'].value_counts(normalize = True)
print(cancelled_perc)

plt.figure(figsize = (5,4))
plt.title('Reservation status count')
plt.bar(['Not canceled', 'canceled'], df['is_canceled'].value_counts(), edgecolor = 'r')
plt.show()
```

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

```
In [21]: plt.figure(figsize = (8,4))
ax1 = sns.countplot(x='hotel', hue='is_canceled', data=df, palette='Blues')
legend_labels, _ = ax1.get_legend_handles_labels()
ax1.legend(bbox_to_anchor=(1,1))
plt.title("Reservation status in different hotels", size=20)
plt.xlabel("hotel")
plt.ylabel("number of reservations")
plt.show()
```



```
In [22]: resort_hotel = df[df['hotel'] == 'Resort Hotel']
resort_hotel['is_canceled'].value_counts(normalize = True)
```

```
Out[22]: is_canceled
0      0.72025
1      0.27975
Name: proportion, dtype: float64
```

```
In [23]: city_hotel = df[df['hotel'] == 'City Hotel']
city_hotel['is_canceled'].value_counts(normalize = True)
```

```
Out[23]: is_canceled
0      0.582918
1      0.417082
Name: proportion, dtype: float64
```

```
In [24]: resort_hotel = resort_hotel.groupby('reservation_status_date')[['adr']].mean()
city_hotel = city_hotel.groupby('reservation_status_date')[['adr']].mean()
```

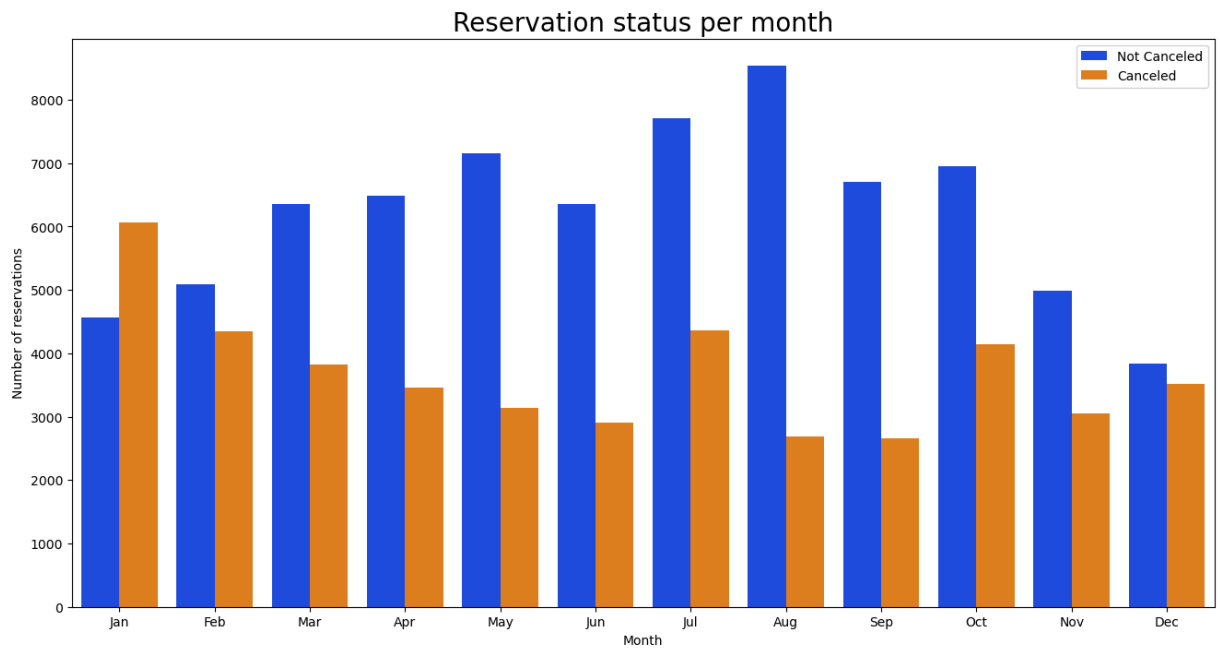
```
In [25]: import calendar
import seaborn as sns
import matplotlib.pyplot as plt

df['month'] = df['reservation_status_date'].dt.month
df['month'] = df['month'].apply(lambda x: calendar.month_abbr[x])
month_order = [calendar.month_abbr[i] for i in range(1,13)]

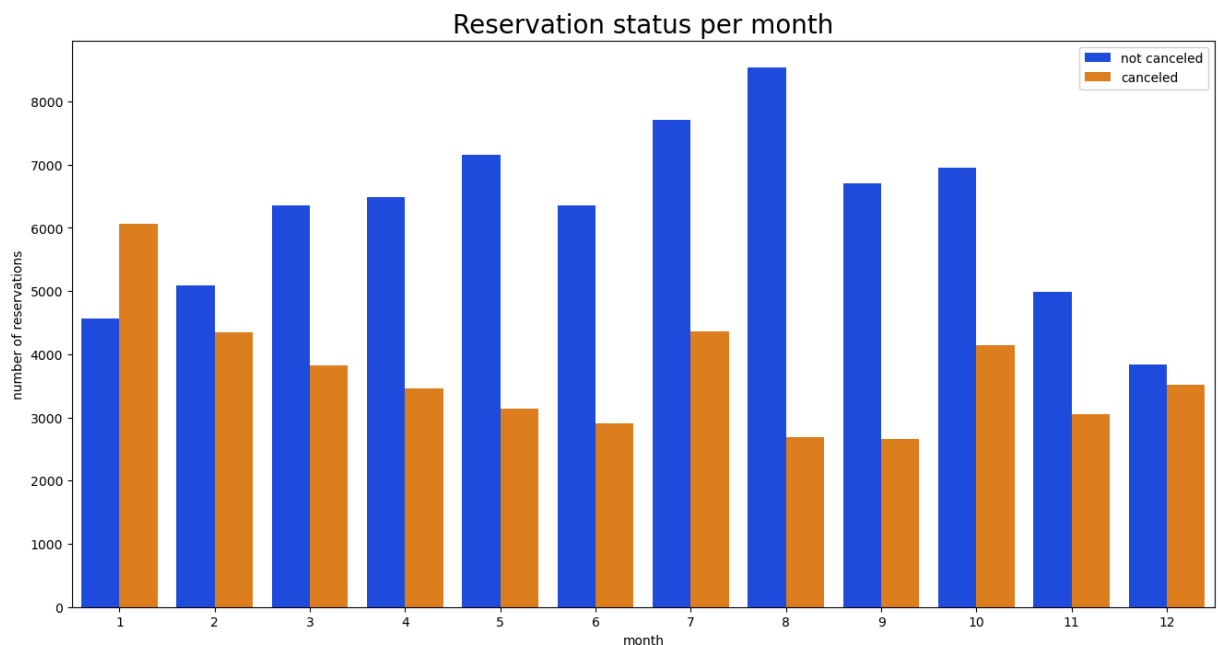
plt.figure(figsize=(16,8))
ax = sns.countplot(x='month', hue='is_canceled', data=df,
                  palette='bright', order=month_order)

handles, labels = ax.get_legend_handles_labels()
ax.legend(handles, ['Not Canceled', 'Canceled'], bbox_to_anchor=(1,1))

plt.title('Reservation status per month', size=20)
plt.xlabel('Month')
plt.ylabel('Number of reservations')
plt.show()
```



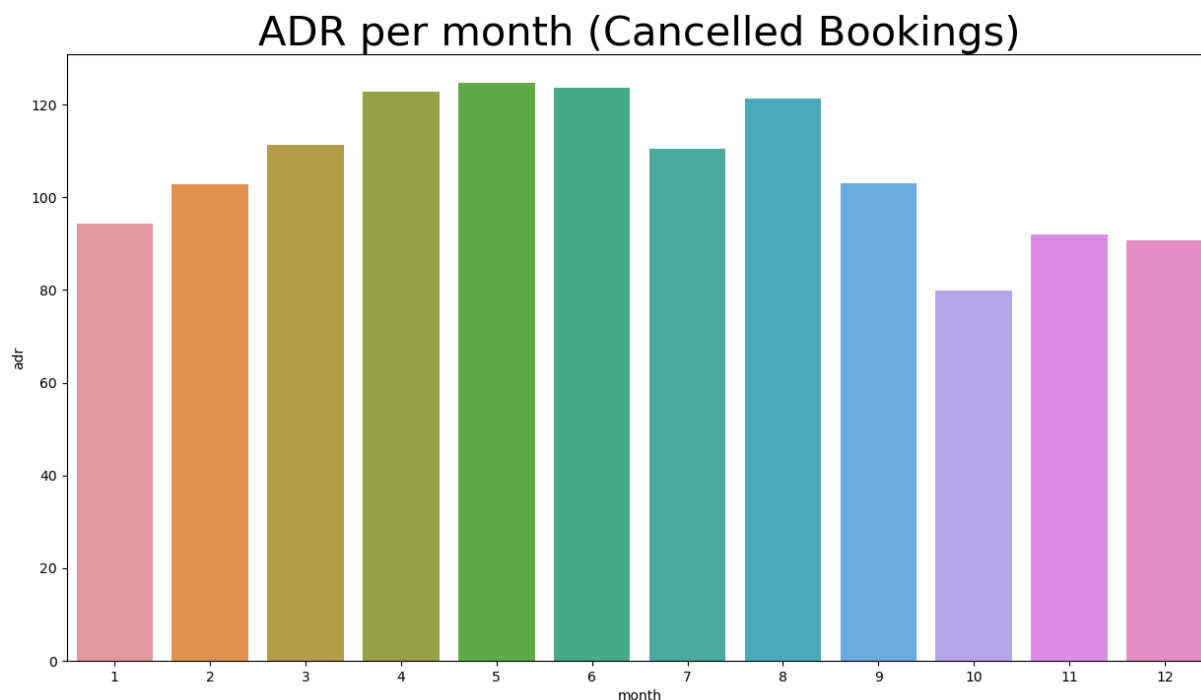
```
In [26]: df['month'] = df['reservation_status_date'].dt.month
plt.figure(figsize = (16,8))
ax1 = sns.countplot(x = 'month', hue = 'is_canceled', data = df, palette = 'bright')
legend_labels,_ = ax1.get_legend_handles_labels()
ax1.legend(bbox_to_anchor=(1,1))
plt.title('Reservation status per month', size = 20)
plt.xlabel('month')
plt.ylabel('number of reservations')
plt.legend(['not canceled', 'canceled'])
plt.show()
```



```
In [27]: cancelled_data = df[df['is_canceled'] == 1].groupby('month')[['adr']].mean().reset_index()

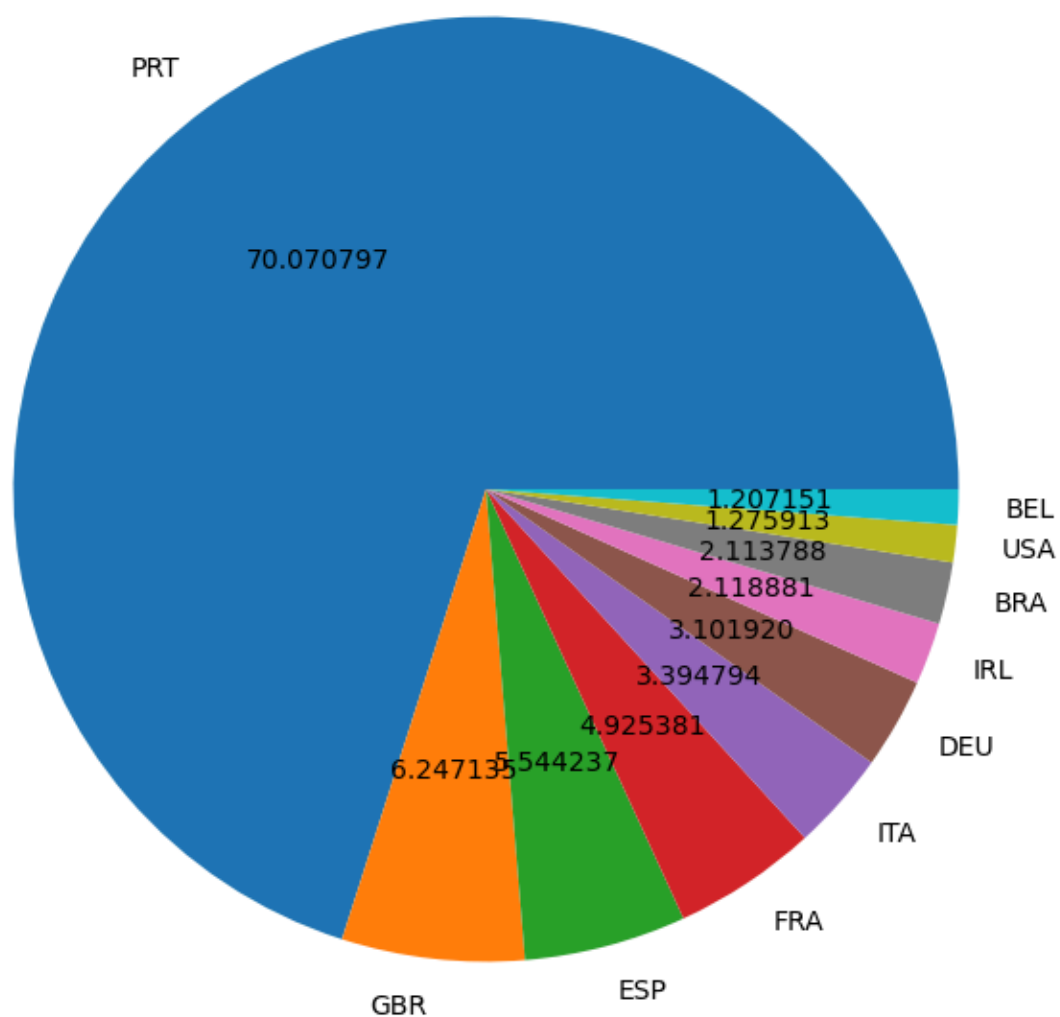
plt.figure(figsize=(15,8))
plt.title('ADR per month (Cancelled Bookings)', fontsize=30)
```

```
sns.barplot(x='month', y='adr', data=cancelled_data)
plt.show()
```



```
In [28]: cancelled_data = df[df['is_canceled'] == 1]
top_10_country = cancelled_data['country'].value_counts()[:10]
plt.figure(figsize = (8,8))
plt.title('Top 10 countries with reservation canceled')
plt.pie(top_10_country, autopct = '%2f', labels = top_10_country.index)
plt.show()
```

Top 10 countries with reservation canceled



```
In [29]: df['market_segment'].value_counts()
```

```
Out[29]: market_segment
Online TA      56402
Offline TA/TO  24159
Groups         19806
Direct         12448
Corporate       5111
Complementary   734
Aviation        237
Name: count, dtype: int64
```

```
In [30]: df['market_segment'].value_counts(normalize = True)
```

```
Out[30]: market_segment
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: proportion, dtype: float64
```

```
In [31]: cancelled_data['market_segment'].value_counts(normalize = True)
```

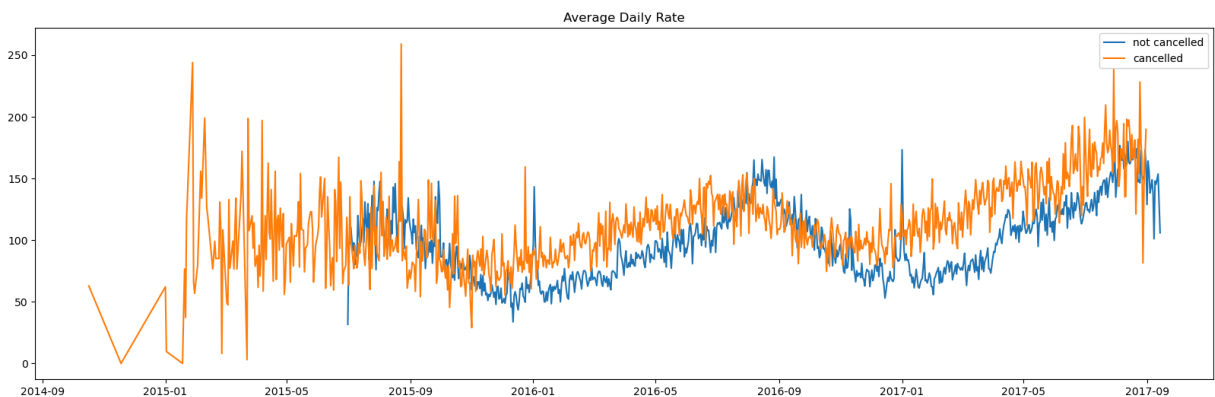
```
Out[31]: market_segment
Online TA      0.469696
Groups         0.273985
Offline TA/TO  0.187466
Direct         0.043486
Corporate      0.022151
Complementary  0.002038
Aviation       0.001178
Name: proportion, dtype: float64
```

```
In [32]: cancelled_df_adr = cancelled_data.groupby('reservation_status_date')[['adr']].mean(
cancelled_df_adr.reset_index(inplace = True)
cancelled_df_adr.sort_values('reservation_status_date', inplace = True)

not_cancelled_data = df[df['is_canceled'] == 0]
not_cancelled_df_adr = not_cancelled_data.groupby('reservation_status_date')[['adr']]
not_cancelled_df_adr.reset_index(inplace = True)
not_cancelled_df_adr.sort_values('reservation_status_date', inplace = True)

plt.figure(figsize = (20,6))
plt.title('Average Daily Rate')
plt.plot(not_cancelled_df_adr['reservation_status_date'],not_cancelled_df_adr['adr'])
plt.plot(cancelled_df_adr['reservation_status_date'],cancelled_df_adr['adr'], label='cancelled')
plt.legend()
```

```
Out[32]: <matplotlib.legend.Legend at 0x7e688334d890>
```



```
In [33]: import calendar
df['month'] = df['reservation_status_date'].dt.month
df['month'] = df['month'].apply(lambda x: calendar.month_abbr[x])
```

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []: