

Importing Libraries

In [4]:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 %matplotlib inline
4 import seaborn as sns
5 import warnings
6 warnings.filterwarnings('ignore')
```

```
In [5]: 1 df=pd.read_csv("data files/hotel_bookings.csv")
        2 df
```

```
Out[5]:
```

	hotel	is_canceled	lead_time	arrival_date_year	arrival_date_month	arrival_date_week_number	arrival_date_day_of_month	st
0	Resort Hotel	0	342	2015	July	27	1	
1	Resort Hotel	0	737	2015	July	27	1	
2	Resort Hotel	0	7	2015	July	27	1	
3	Resort Hotel	0	13	2015	July	27	1	
4	Resort Hotel	0	14	2015	July	27	1	
...
119385	City Hotel	0	23	2017	August	35	30	
119386	City Hotel	0	102	2017	August	35	31	
119387	City Hotel	0	34	2017	August	35	31	
119388	City Hotel	0	109	2017	August	35	31	
119389	City Hotel	0	205	2017	August	35	29	

119390 rows × 32 columns



Exploratory Data Analysis and Data Cleaning

In [6]:

```
1 df.head()
```

Out[6]:

	hotel	is_canceled	lead_time	arrival_date_year	arrival_date_month	arrival_date_week_number	arrival_date_day_of_month	stays_in_week
0	Resort Hotel	0	342	2015	July	27	1	
1	Resort Hotel	0	737	2015	July	27	1	
2	Resort Hotel	0	7	2015	July	27	1	
3	Resort Hotel	0	13	2015	July	27	1	
4	Resort Hotel	0	14	2015	July	27	1	

5 rows × 32 columns



In [7]:

```
1 df.tail()
```

Out[7]:

	hotel	is_canceled	lead_time	arrival_date_year	arrival_date_month	arrival_date_week_number	arrival_date_day_of_month	stays_in_week
119385	City Hotel	0	23	2017	August	35	30	
119386	City Hotel	0	102	2017	August	35	31	
119387	City Hotel	0	34	2017	August	35	31	
119388	City Hotel	0	109	2017	August	35	31	
119389	City Hotel	0	205	2017	August	35	29	

5 rows × 32 columns



```
In [9]: 1 df.shape
```

```
Out[9]: (119390, 32)
```

```
In [10]: 1 df.columns
```

```
Out[10]: 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 [11]:

```
1 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 [18]: 1 df['reservation_status_date']=pd.to_datetime(df['reservation_status_date'])
```

In [19]: 1 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  datetime64[ns]
dtypes: datetime64[ns](1), float64(4), int64(16), object(11)
memory usage: 29.1+ MB
```

```
In [21]: 1 df.describe(include="object")
```

Out[21]:

	hotel	arrival_date_month	meal	country	market_segment	distribution_channel	reserved_room_type	assigned_room_type
count	119390	119390	119390	118902	119390	119390	119390	119390
unique	2	12	5	177	8	5	10	12
top	City Hotel	August	BB	PRT	Online TA	TATO	A	A
freq	79330	13877	92310	48590	56477	97870	85994	74053


```
In [35]: 1 for col in df.describe(include="object").columns:
          2     print(col)
          3     print(df[col].unique())
          4     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' '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 [36]: 1 df.isnull().sum()
```

```
Out[36]: 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 [38]: 1 df.drop(['company', 'agent'], axis=1, inplace=True)
         2 df.dropna(inplace=True)
```

```
In [39]: 1 df.isnull().sum()
```

```
Out[39]: 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
```

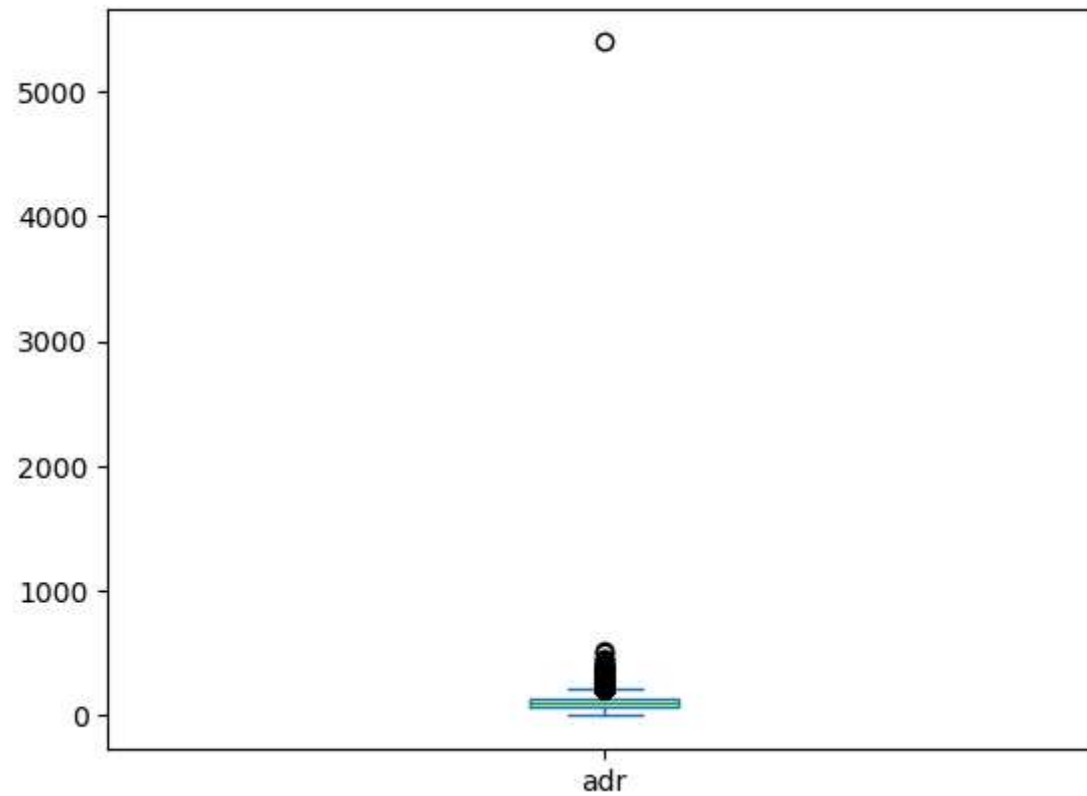
```
In [40]: 1 df.describe()
```

Out[40]:

	is_canceled	lead_time	arrival_date_year	arrival_date_week_number	arrival_date_day_of_month	stays_in_weekend_nights
count	118898.000000	118898.000000	118898.000000	118898.000000	118898.000000	118898.000000
mean	0.371352	104.311435	2016.157656	27.166555	15.800880	0.928897
min	0.000000	0.000000	2015.000000	1.000000	1.000000	0.000000
25%	0.000000	18.000000	2016.000000	16.000000	8.000000	0.000000
50%	0.000000	69.000000	2016.000000	28.000000	16.000000	1.000000
75%	1.000000	161.000000	2017.000000	38.000000	23.000000	2.000000
max	1.000000	737.000000	2017.000000	53.000000	31.000000	16.000000
std	0.483168	106.903309	0.707459	13.589971	8.780324	0.996216

```
In [41]: 1 df['adr'].plot(kind='box')
```

Out[41]: <Axes: >



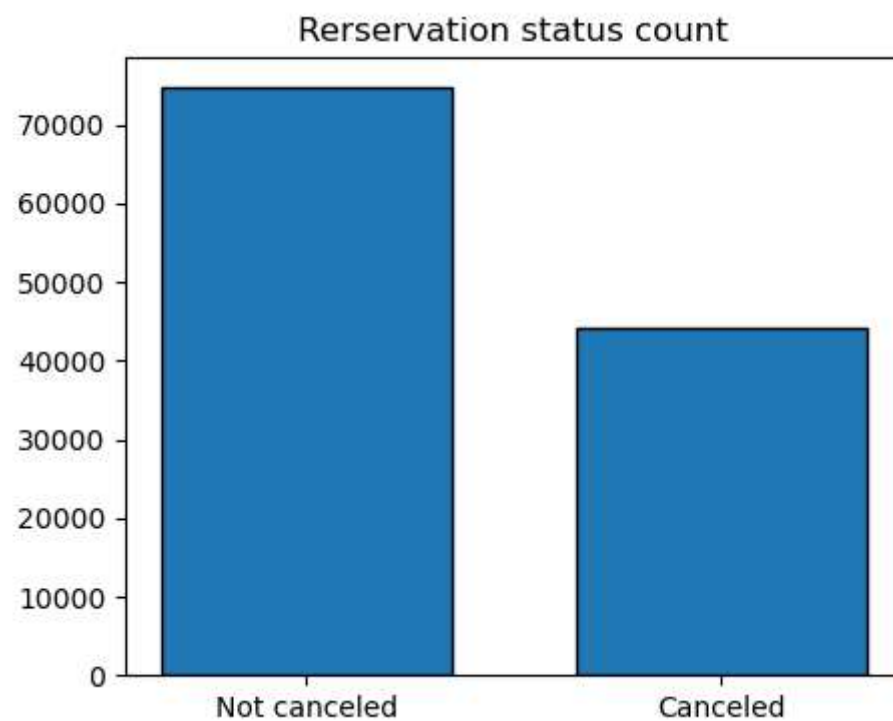
```
In [42]: 1 df=df[df['adr']<=5000]
```

Data Analysis and Visualizations

```
In [43]: 1 cancelled_perc=df["is_canceled"].value_counts(normalize=True)
```

```
In [52]: 1 print("cancelled_perc")
2 plt.figure(figsize=(5,4))
3 plt.title("Rerservation status count")
4 plt.bar(['Not canceled', 'Canceled'],df['is_canceled'].value_counts(),edgecolor='k',width=0.7)
5 plt.show()
```

cancelled_perc



```
In [56]: 1 plt.figure(figsize=(8,4))
2 ax1=sns.countplot(x='hotel',hue='is_canceled',data=df,palette='Blues')
3 legend_labels,_=ax1.get_legend_handles_labels()
4 ax1.legend(bbox_to_anchor=(1,1))
5 plt.title("Reservation status is different hotels",size=20)
6 plt.xlabel("hotel")
7 plt.ylabel("number of reservations")
8 plt.show()
```



```
In [65]: 1 resort_hotel=df[df['hotel']=="Resort Hotels"]
2 resort_hotel['is_canceled'].value_counts(normalize=True)
```

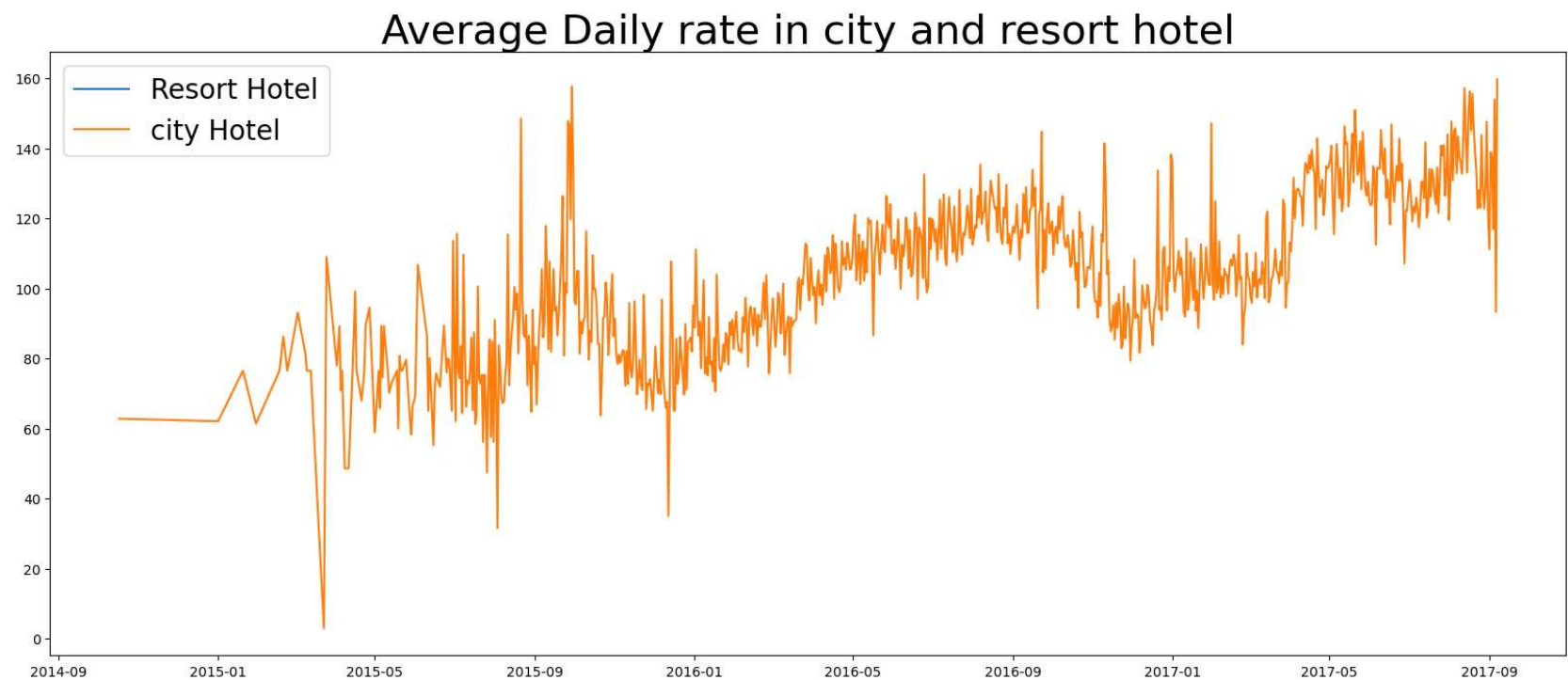
Out[65]: Series([], Name: proportion, dtype: float64)


```
In [66]: 1 city_hotel=df[df['hotel']=='City Hotel']
2 city_hotel['is_canceled'].value_counts(normalize=True)
```

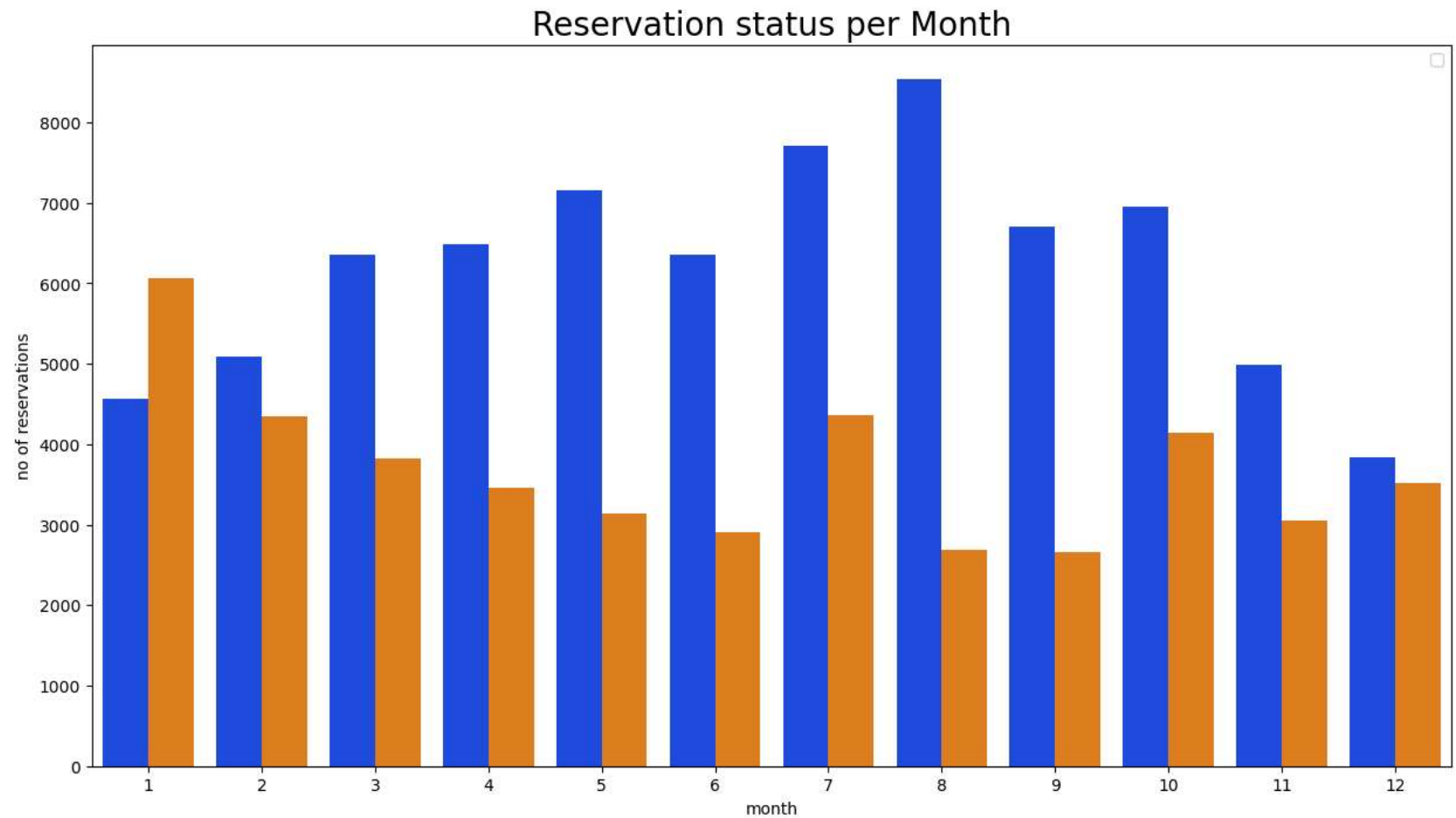
```
Out[66]: is_canceled
0      0.582918
1      0.417082
Name: proportion, dtype: float64
```

```
In [67]: 1 resort_hotel=resort_hotel.groupby('reservation_status_date')[['adr']].mean()
2 city_hotel=city_hotel.groupby('reservation_status_date')[['adr']].mean()
```

```
In [74]: 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()
```



```
In [81]: 1 df['month']=df['reservation_status_date'].dt.month
2 plt.figure(figsize=(15,8))
3 ax1=sns.countplot(x='month',hue='is_canceled',data=df,palette='bright')
4 ax1.legend(bbox_to_anchor=(1,1))
5 plt.title("Reservation status per Month",size=20)
6 plt.xlabel("month")
7 plt.ylabel("no of reservations")
8 plt.legend("not canceled",'canceled')
9 plt.show()
```



```
In [87]: 1 df['market_segment'].value_counts()
```

```
Out[87]: market_segment
Online TA      56402
Offline TA/TO  24159
Groups         19806
Direct         12448
Corporate       5111
Complementary   734
Aviation        237
Name: count, dtype: int64
```

```
In [88]: 1 df['market_segment'].value_counts(normalize=True)
```

```
Out[88]: 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 [ ]: 1
```

```
In [ ]: 1
```