

# IMPORT LIBRARIES

```
In [123]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
```

```
In [124]: #LOADING THE DATASETS
```

```
In [125]: df=pd.read_csv('hotel_bookings.csv' )
```

```
In [126]: #EXPLORATORY DATA ANALYSIS AND DATA CLEANING
```

```
In [127]: df.head()
```

	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	...
0	Resort Hotel	0	342	2015	July	27	1	0	0	2	...
1	Resort Hotel	0	737	2015	July	27	1	0	0	2	...
2	Resort Hotel	0	7	2015	July	27	1	0	1	1	...
3	Resort Hotel	0	13	2015	July	27	1	0	1	1	...
4	Resort Hotel	0	14	2015	July	27	1	0	2	2	...

5 rows × 32 columns

```
In [128]: df.tail()
```

	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	...
119385	City Hotel	0	23	2017	August	35	30	2	5	2	...
119386	City Hotel	0	102	2017	August	35	31	2	5	3	...
119387	City Hotel	0	34	2017	August	35	31	2	5	2	...
119388	City Hotel	0	109	2017	August	35	31	2	5	2	...
119389	City Hotel	0	205	2017	August	35	29	2	7	2	...

5 rows × 32 columns

```
In [129]: df.shape
```

```
Out[129]: (119390, 32)
```

```
In [130]: df.columns
```

```
Out[130]: 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 [131]: 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              119390 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                 103650 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 [132]: df[['reservation_status_date']]=pd.to_datetime(df[['reservation_status_date']])
```

```
In [133]: df.describe(include='object')
```

	hotel	arrival_date_year	meal	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	2	12	5	177	8	5	10	12	3	4	3
top	City Hotel	August	BB	PRT	Online TA	TA/TO	A	A	No Deposit	Transient	Check-Out
freq	79330	13877	92310	48590	56477	97870	85994	74053	104641	89613	75166

```
In [134]: for col in df.describe(include='object').columns:
print(col)
print(df[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' 'CZ' 'GRC' 'ITA' 'IND' 'DNK' 'RUS' 'SWE' 'AUS' 'EST'
'CZE' 'BRA' 'FIN' 'NOZ' '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'
'CMR' 'HKG' 'IRN' 'GEO' 'AND' 'GIB' 'URY' 'JEP' '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'
'CHN' 'BTH' 'MUS' 'COM' 'SUR' 'UGA' 'BGR' 'CIV' 'JOR' 'SVR' 'SGP' 'BDI'
'SAU' 'VNM' 'PLW' 'GAT' 'EGY' 'PER' 'MLT' 'HMI' 'ECU' 'MDG' 'ISL' 'UZB'
'NPL' 'BHS' 'MAC' 'TGO' 'TWN' 'DJT' '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' 'MNR' 'PAN' 'BFA' 'LBY'
'MLI' 'NAM' 'BOL' 'PRY' 'GRB' 'ABW' 'ATA' 'SLV' 'DMA' 'PYF' 'GUY' 'LCA'
'ATA' 'GTM' 'ASM' 'NRT' '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 [135]: df.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	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 [136]: df.drop('agent',axis=1,inplace=True)
df.drop('company',axis=1,inplace=True)
```

```
In [141]: df.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	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
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 [86]: df.describe()
```

	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	is_canceled
count	217.000000	217.000000	217.000000	217.000000	217.000000	217.000000	217.000000	217.000000	217.000000	217.000000	217.000000
mean	0.078341	40.520737	2015.465438	38.198157	10.824885	1.56682	4.631336	1.410138	0.036866	0.0	0.0
std	0.269329	61.748375	0.720053	12.890292	7.582065	1.49270	3.552846	0.520406	0.232788	0.0	0.0
min	0.000000	0.000000	2015.000000	1.000000	1.000000	0.000000	0.000000	1.000000	0.000000	0.0	0.0
25%	0.000000	12.000000	2015.000000	33.000000	6.000000	0.000000	2.000000	1.000000	0.000000	0.0	0.0
50%	0.000000	27.000000	2015.000000	45.000000	9.000000	2.000000	4.000000	1.000000	0.000000	0.0	0.0
75%	0.000000	36.000000	2016.000000	46.000000	13.000000	2.000000	6.000000	2.000000	0.000000	0.0	0.0
max	1.000000	364.000000	2017.000000	53.000000	31.000000	9.000000	21.000000	3.000000	2.000000	0.0	0.0

```
In [142]: df.dropna()
```

	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	...
0	Resort Hotel	0	342	2015	July	27	1	0	0	2	...
1	Resort Hotel	0	737	2015	July	27	1	0	0	2	...
2	Resort Hotel	0	7	2015	July	27	1	0	1	1	...
3	Resort Hotel	0	13	2015	July	27	1	0	1	1	...
4	Resort Hotel	0	14	2015	July	27	1	0	2	2	...
...	...	...	...	...	...	...	...	...	...	...	...
119385	City Hotel	0	23	2017	August	35	30	2	5	2	...
119386	City Hotel	0	102	2017	August	35	31	2	5	3	...
119387	City Hotel	0	34	2017	August	35	31	2	5	2	...
119388	City Hotel	0	109	2017	August	35	31	2	5	2	...
119389	City Hotel	0	205	2017	August	35	29	2	7	2	...

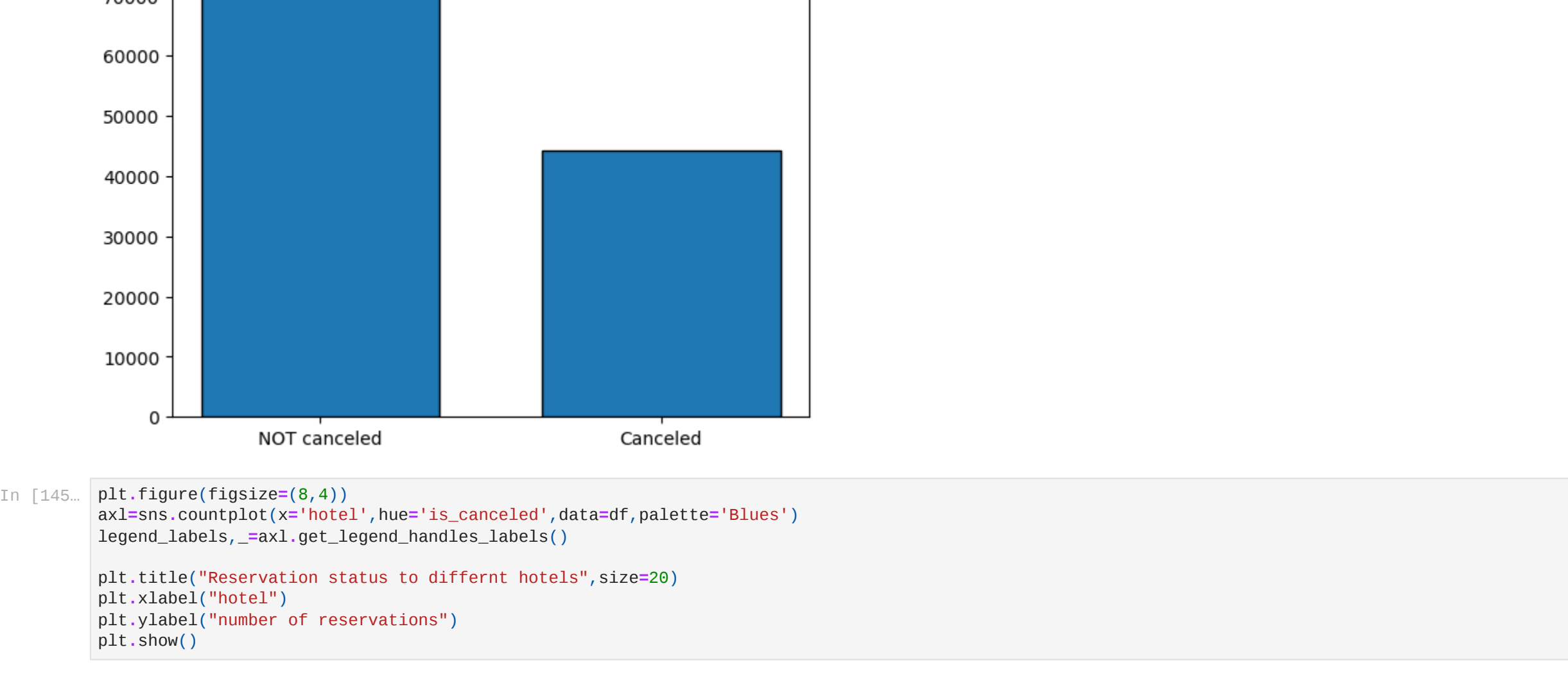
118898 rows × 30 columns

```
In [87]: df=df[df['adr']>5000]
```

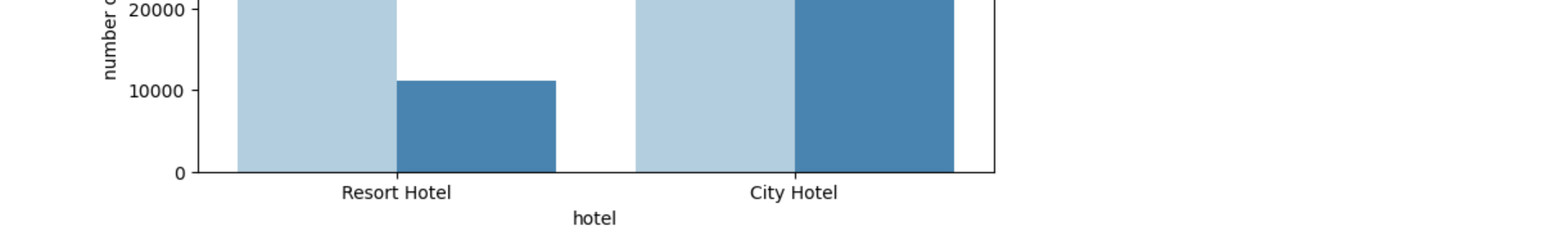
# ANALYSIS AND VISUALIZATION

```
In [143]: cancelled_perc=df['is_canceled'].value_counts(normalize=True)
print(cancelled_perc)
0    0.629584
1    0.370416
Name: is_canceled, dtype: float64
```

```
In [144]: plt.figure(figsize=(8,4))
plt.title("Reservation status count")
plt.bar(['NOT canceled','Canceled'],df['is_canceled'].value_counts(),edgecolor='k',width=0.7)
plt.show()
```



```
In [145]: plt.figure(figsize=(8,4))
ax=sns.countplot(x='hotel',hue='is_canceled',data=df,palette='Blues')
legend_labels,_=ax.get_legend_handles_labels()
plt.title("Reservation status to different hotels",size=20)
plt.xlabel("hotel")
plt.ylabel("number of reservations")
plt.show()
```



```
In [146]: resort_hotel=df[df['hotel']=='Resort Hotel']
resort_hotel['is_canceled'].value_counts(normalize=True)
```

```
Out[146]: 0    0.723366
1    0.277634
Name: is_canceled, dtype: float64
```

```
In [147]: city_hotel=df[df['hotel']=='City Hotel']
city_hotel['is_canceled'].value_counts(normalize=True)
```

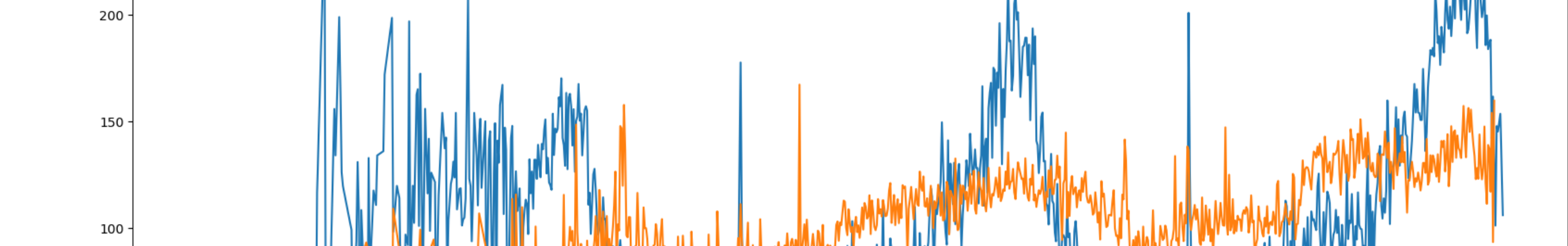
```
Out[147]: 0    0.58273
1    0.41727
Name: is_canceled, dtype: float64
```

```
In [148]: resort_hotel=resort_hotel.groupby("reservation_status_date")[['adr']].mean()
city_hotel=city_hotel.groupby("reservation_status_date")[['adr']].mean()
```

```
In [149]: plt.figure(figsize=(20,8))
plt.title("Average rate in Resort and City hotel",fontsize=30)
plt.plot(resort_hotel.index,resort_hotel['adr'],label="Resort Hotel")
plt.plot(city_hotel.index,city_hotel['adr'],label="City Hotel")
plt.legend(fontsize=20)
plt.show()
```



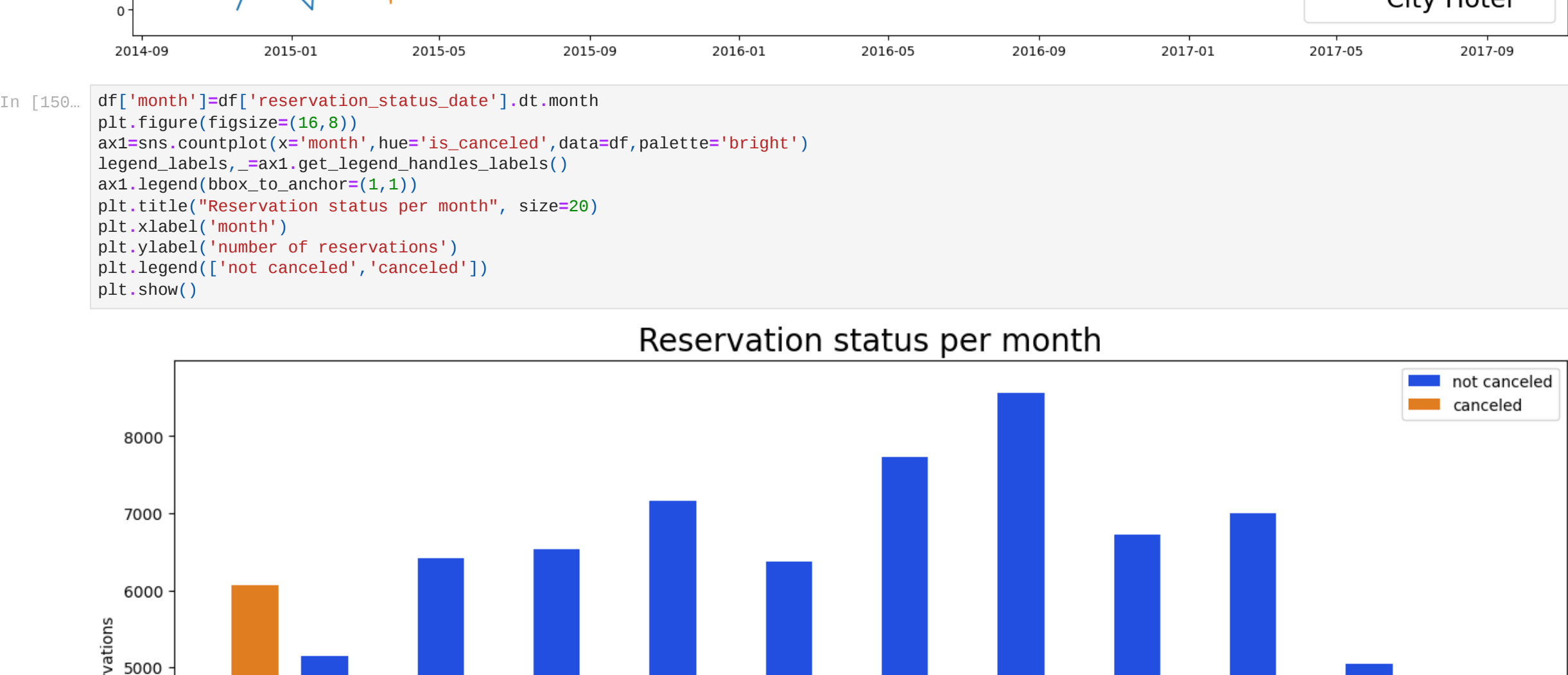
```
In [150]: df['month']=df['reservation_status_date'].dt.month
plt.figure(figsize=(16,8))
ax=sns.countplot(x='month',hue='is_canceled',data=df,palette='bright')
legend_labels,_=ax.get_legend_handles_labels()
ax.legend(bbox=(0.5,0.5))
plt.title("Reservation status per month",size=20)
plt.xlabel('month')
plt.ylabel('number of reservations')
plt.legend(['not canceled','canceled'])
plt.show()
```



```
In [162]: plt.figure(figsize=(15, 8))
plt.title('ADR PER MONTH', fontsize=30)
sns.barplot(x='month', y='adr', data=df[df['is_canceled'] == 1].groupby('month')[['adr']].sum().reset_index())
plt.legend(fontsize=20)
plt.show()
```

No artists with labels found to put in legend. Note that artists whose label start with an underscore are ignored when legend() is called with no argument.

# ADR PER MONTH



```
In [ ]: 
```

```
In [ ]: 
```