

Explore and Clean the Data



jupyter Hotel_Booking_Analysis Last Checkpoint: 11 days ago

File Edit View Run Kernel Settings Help Trusted

+ 🔍 📄 ▶ ⏮ ⏭ ⏪ ⏩ ⏴ ⏵ Code ▾ ▸ Open in... 🐍 Python 3 (ipykernel)

```
[1]: import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
%matplotlib inline
```

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

```
[3]: df.head(5)
```

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

```
[4]: df.shape
[4]: (119390, 32)

[5]: df.columns
[5]: 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')
```

jupyter Hotel_Booking_Analysis Last Checkpoint: 11 days ago

File Edit View Run Kernel Settings Help

Python 3 (ipykernel)

```
[6]: 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
```

jupyter

Hotel_Booking_Analysis

Last Checkpoint: 11 days ago

File

Edit

View

Run

Kernel

Settings

Help

Trusted

Code

Open in... Python 3 (ipykernel)

memory usage: 29.1+ MB

[7]:

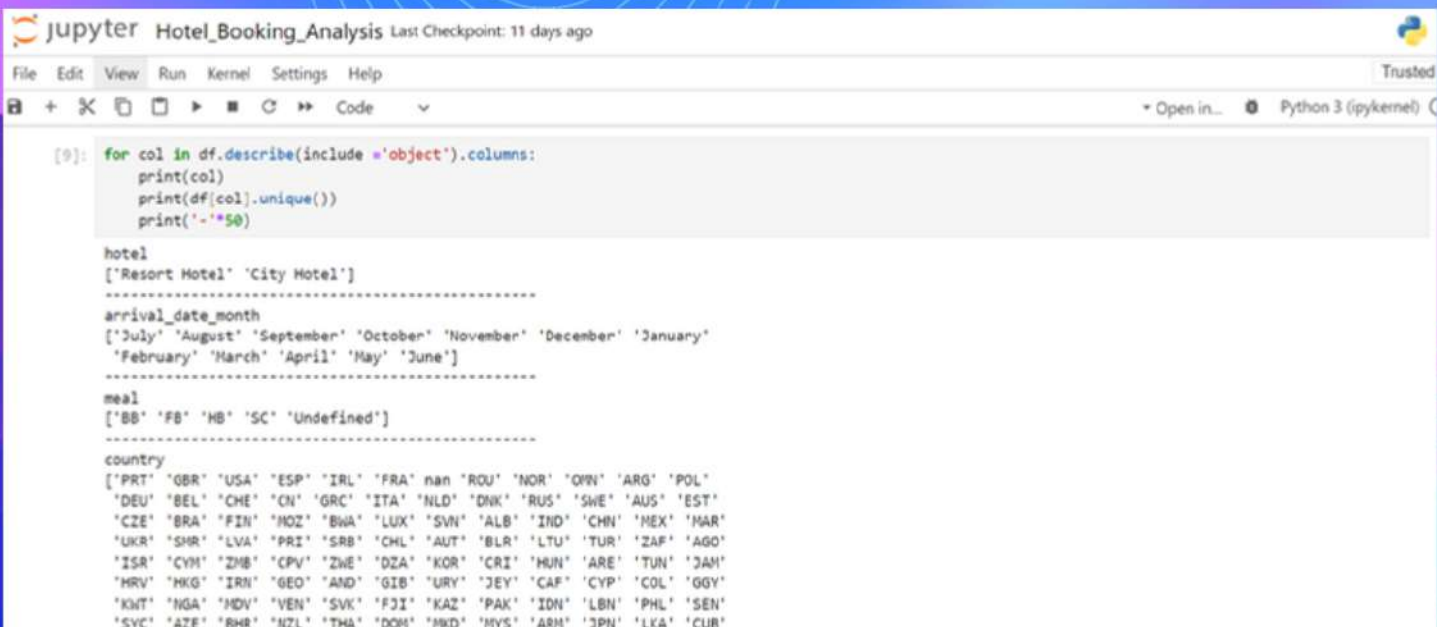
df['reservation_status_date'] = pd.to_datetime(df['reservation_status_date'], format='%d/%m/%Y')

[8]:

df.describe(include='object')

[8]:

	hotel	arrival_date_month	meal	country	market_segment	distribution_channel	reserved_room_type	assigned_room_type	deposit_type	customer_type	res
count	119390	119390	119390	118902	119390	119390	119390	119390	119390	119390	
unique	2	12	5	177	8	5	10	12	3	4	
top	City Hotel	August	BB	PRT	Online TA	TA/TO	A	A	No Deposit	Transient	
freq	79330	13877	92310	48590	56477	97870	85994	74053	104641	89613	



jupyter Hotel_Booking_Analysis Last Checkpoint: 11 days ago

File Edit View Run Kernel Settings Help Trusted

+ - X [] ▶ ⌂ ⌕ Code ▾

Open in... Python 3 (ipykernel)

```
[10]: df.isnull().sum()
```

```
[10]: 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
```


Jupyter Hotel_Booking_Analysis Last Checkpoint: 11 days ago

File Edit View Run Kernel Settings Help Trusted

+ - X Copy Paste Run Cell All Cells Code Open in... Python 3 (ipykernel)

```
dtype: int64
```

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

```
[12]: df.describe()
```

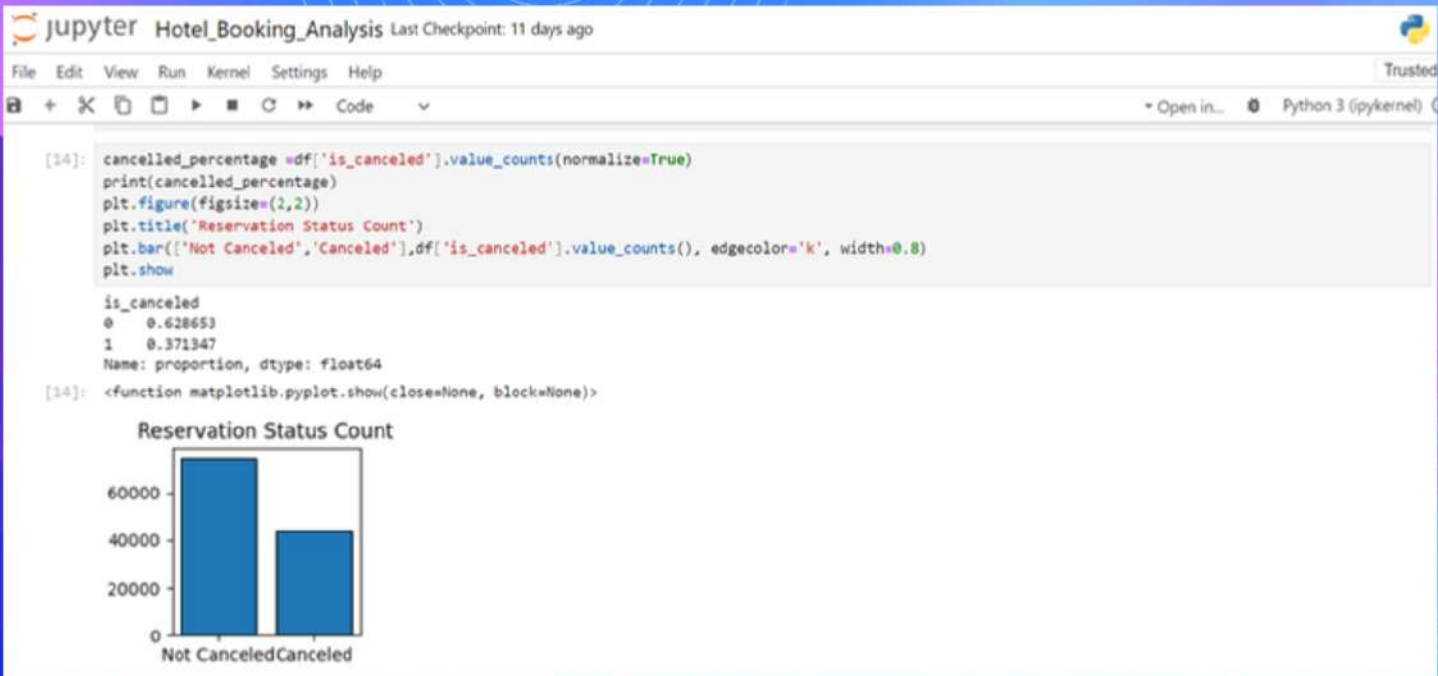
```
[12]:
```

	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
count	118898.000000	118898.000000	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	2.502145	1.858311
min	0.000000	0.000000	2015.000000	1.000000	1.000000	0.000000	0.000000	0.000000
25%	0.000000	18.000000	2016.000000	16.000000	8.000000	0.000000	1.000000	2.000000
50%	0.000000	69.000000	2016.000000	28.000000	16.000000	1.000000	2.000000	2.000000
75%	1.000000	161.000000	2017.000000	38.000000	23.000000	2.000000	3.000000	2.000000
max	1.000000	737.000000	2017.000000	53.000000	31.000000	16.000000	41.000000	55.000000
std	0.483168	106.903309	0.707459	13.589971	8.780324	0.996216	1.900168	0.578511

```
[13]: df=df[df['adults']<5000]
```

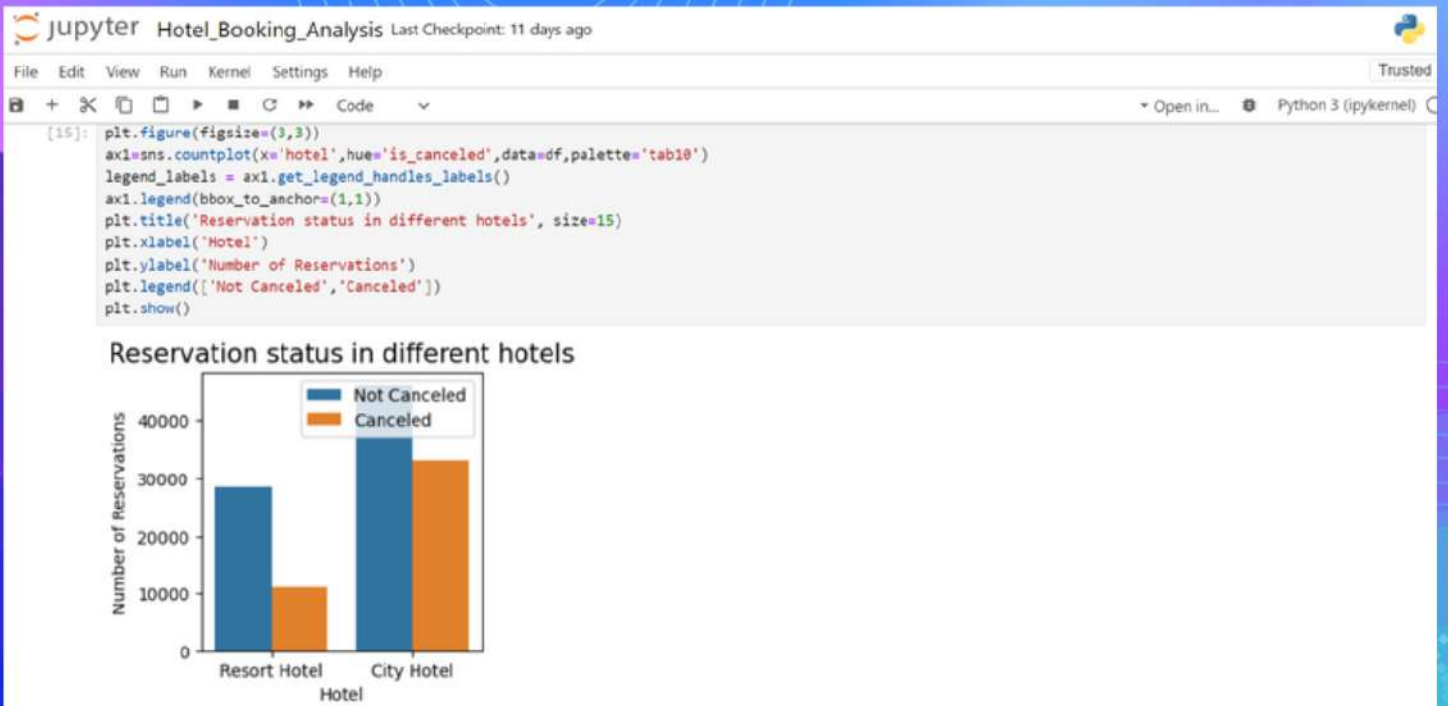
Reservation Status Count





**Reservation status in
different hotels**





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

[16]: is_canceled
      0    0.72025
      1    0.27975
      Name: proportion, dtype: float64

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

[17]: is_canceled
      0    0.582918
      1    0.417082
      Name: proportion, dtype: float64

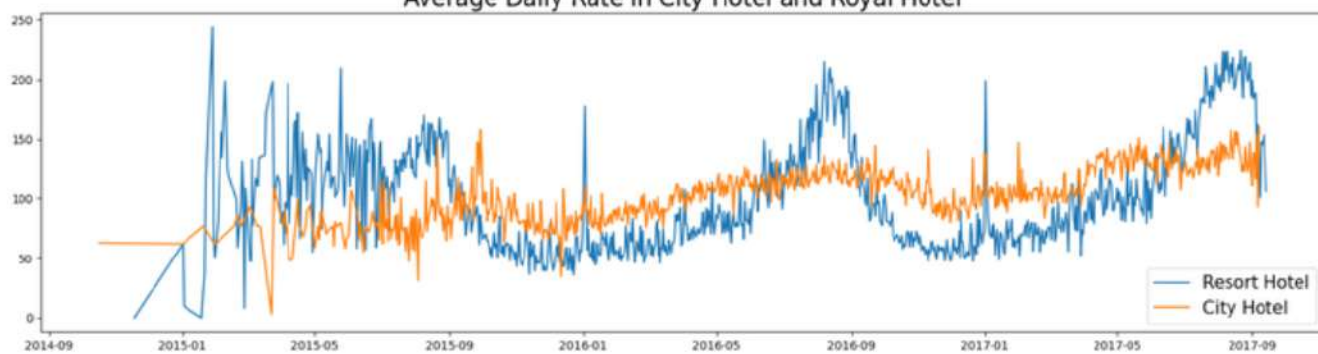
[18]: resort_hotel = resort_hotel.groupby('reservation_status_date')[['adr']].mean()
      city_hotel = city_hotel.groupby('reservation_status_date')[['adr']].mean()
```

**Average Daily Rate in City
Hotel and Royal Hotel**



```
[31]: plt.figure(figsize=(20,5))
plt.title('Average Daily Rate in City Hotel and Royal Hotel', size=20)
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=15)
plt.show()
```

Average Daily Rate in City Hotel and Royal Hotel



Reservation status per month

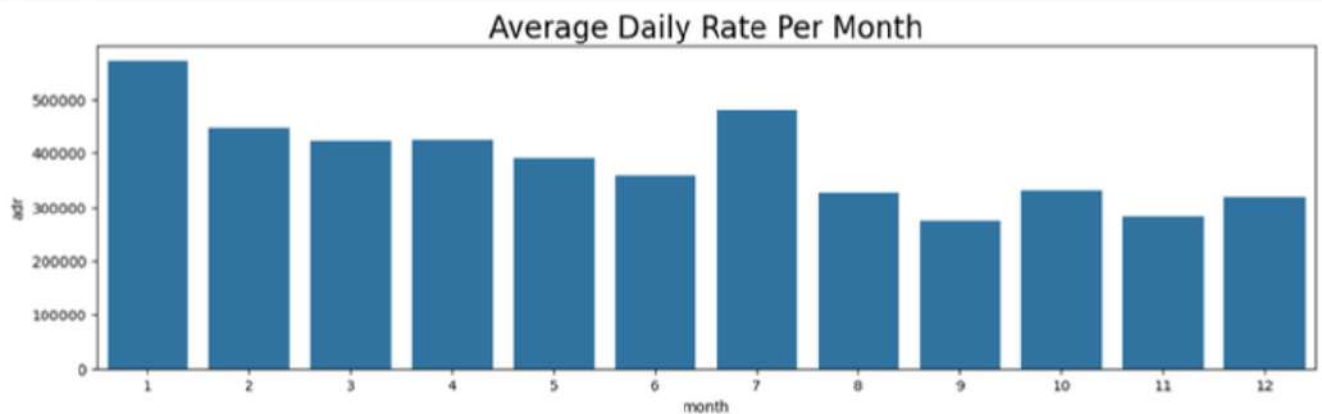




Average Daily Rate Per
Month

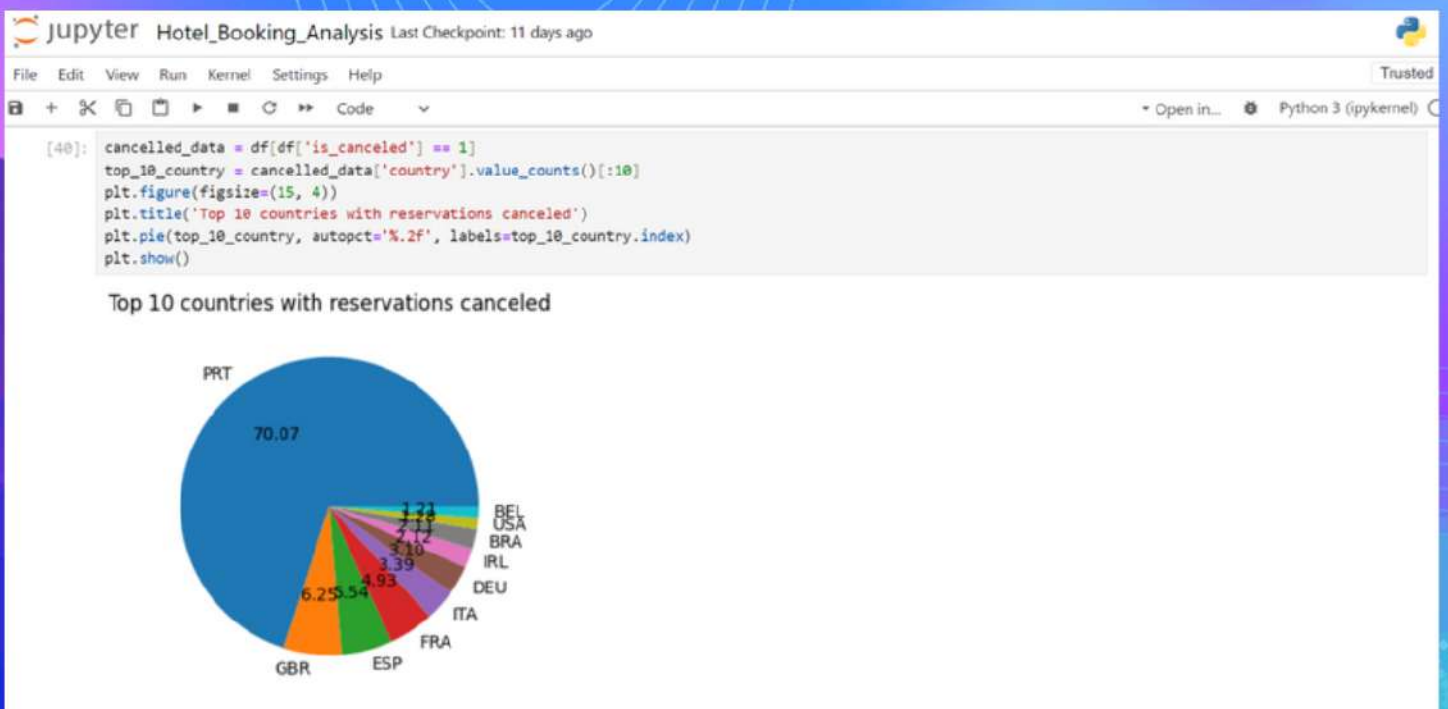


```
[34]: plt.figure(figsize=(15,4))
plt.title('Average Daily Rate Per Month', fontsize=20)
sns.barplot(x='month', y='adr', data=df[df['is_canceled']==1].groupby('month')['adr'].sum().reset_index())
plt.show()
```



**Top 10 countries with
reservations canceled**





```
[23]: df['market_segment'].value_counts()
```

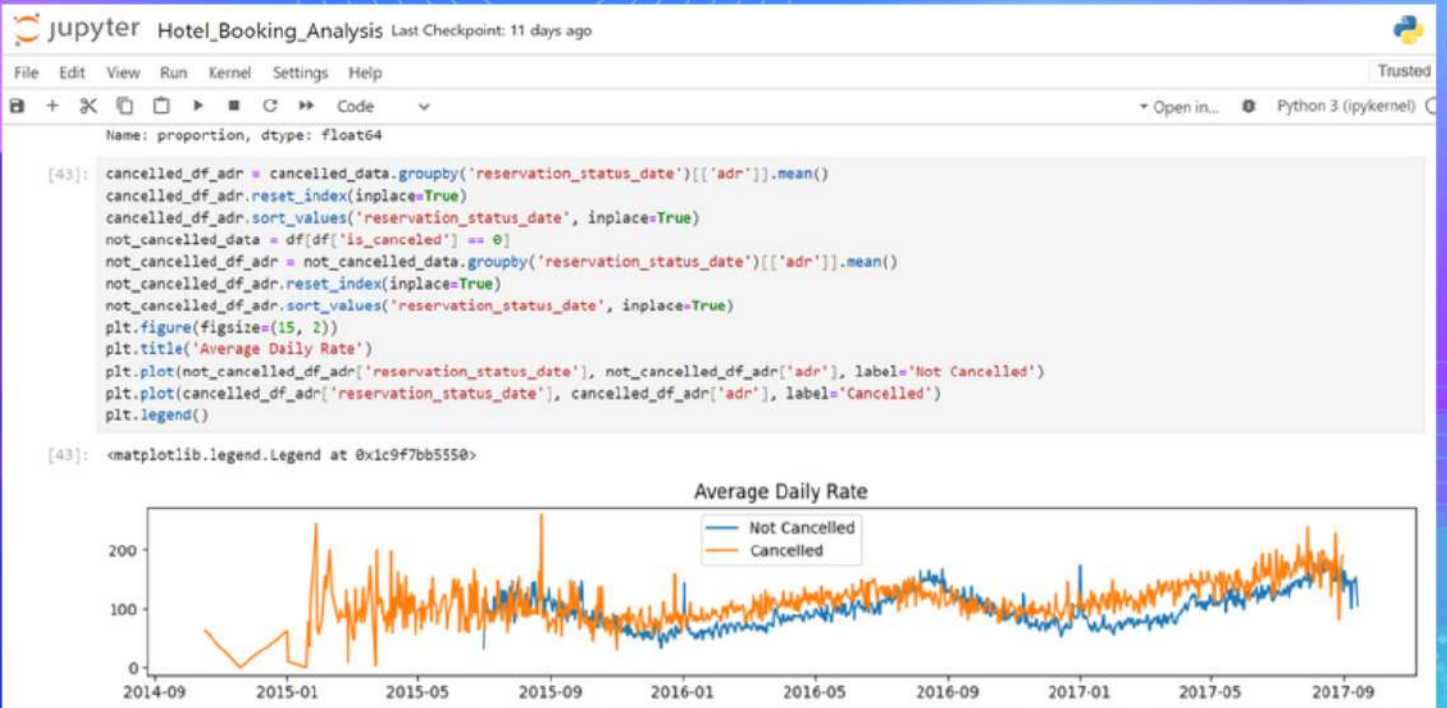
```
[23]: market_segment
Online TA      56402
Offline TA/TO  24159
Groups         19806
Direct         12448
Corporate       5111
Complementary   734
Aviation        237
Name: count, dtype: int64
```

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

```
[24]: 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
```

Average Daily Rate





```
[47]: cancelled_df_adr = cancelled_df_adr[(cancelled_df_adr['reservation_status_date'] > '2016') & (cancelled_df_adr['reservation_status_date'] < '2017')]
not_cancelled_df_adr = not_cancelled_df_adr[(not_cancelled_df_adr['reservation_status_date'] > '2016') & (not_cancelled_df_adr['reservation_status_date'] < '2017')]
plt.figure(figsize=(15, 3))
plt.title('Average Daily Rate')
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.legend(fontsize=10)
plt.show()
```

