

```
import numpy as np
import pandas as pd
import seaborn as sea
import matplotlib.pyplot as plt
```

```
In [1]: df = pd.read_csv('C:\Users\chris\Desktop\CSV files\Hotel Booking CSV\hotel_booking_data.csv', encoding = "unicode_escape")
```

```
In [4]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 119369 entries, 0 to 119368
Data columns (total 32 columns):
 #   Column                Non-Null Count  Dtype  
---  --
 0   hotel                 non-null      object 
 1   is_cancelled          non-null      int64  
 2   lead_time             non-null      int64  
 3   arrival_date_year     non-null      int64  
 4   arrival_date_month    non-null      int64  
 5   arrival_date_week_number non-null      int64  
 6   arrival_date_day_of_month non-null      int64  
 7   stays_in_weekend_nights non-null      int64  
 8   stays_in_week_nights  non-null      int64  
 9   adults               non-null      float64 
10   children              non-null      int64  
11   babies               non-null      int64  
12   meal                 non-null      object 
13   country              non-null      object 
14   market_segment       non-null      object 
15   distribution_channel  non-null      int64  
16   is_repeated_guest     non-null      int64  
17   previous_cancellations non-null      int64  
18   previous_bookings_not_cancelled non-null      int64  
19   reserved_room_type    non-null      object 
20   assigned_room_type     non-null      int64  
21   booking_changes       non-null      object 
22   deposit_type          non-null      object 
23   agent                non-null      float64 
24   company              non-null      int64  
25   days_in_waiting_list  non-null      int64  
26   customer_type         non-null      object 
27   adr                  non-null      float64 
28   required_car_parking_spaces non-null      int64  
29   total_of_special_requests non-null      int64  
30   reservation_status    non-null      object 
31   reservation_status_date non-null      object 
dtypes: float64(4), int64(16), object(12)
memory usage: 28.1+ MB
```

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

```
Out[5]:
```

	is_cancelled	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	...	deposit_type	agent	company	days_in_waiting_list	customer_type	adr	required_car_parking_spaces	total_of_special_requests
0	Resort Hotel	0	342	2015	July	27	1	0	0	2	No Deposit	NaN	NaN	0	Transient	0.0	0	
1	Resort Hotel	0	737	2015	July	27	1	0	0	2	No Deposit	NaN	NaN	0	Transient	75.0	0	
2	Resort Hotel	0	7	2015	July	27	1	0	1	1	No Deposit	NaN	NaN	0	Transient	75.0	0	
3	Resort Hotel	0	13	2015	July	27	1	0	1	1	No Deposit	304.0	NaN	0	Transient	75.0	0	
4	Resort Hotel	0	14	2015	July	27	1	0	2	2	No Deposit	240.0	NaN	0	Transient	98.0	0	

5 rows x 32 columns

```
In [7]: df["reservation_status_date"] = pd.to_datetime(df["reservation_status_date"], format = "%m/%d")
```

```
In [9]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 119369 entries, 0 to 119368
Data columns (total 32 columns):
 #   Column                Non-Null Count  Dtype  
---  --
 0   hotel                 non-null      object 
 1   is_cancelled          non-null      int64  
 2   lead_time             non-null      int64  
 3   arrival_date_year     non-null      int64  
 4   arrival_date_month    non-null      int64  
 5   arrival_date_week_number non-null      int64  
 6   arrival_date_day_of_month non-null      int64  
 7   stays_in_weekend_nights non-null      int64  
 8   stays_in_week_nights  non-null      int64  
 9   adults               non-null      float64 
10   children              non-null      int64  
11   babies               non-null      int64  
12   meal                 non-null      object 
13   country              non-null      object 
14   market_segment       non-null      object 
15   distribution_channel  non-null      int64  
16   is_repeated_guest     non-null      int64  
17   previous_cancellations non-null      int64  
18   previous_bookings_not_cancelled non-null      int64  
19   reserved_room_type    non-null      object 
20   assigned_room_type     non-null      int64  
21   booking_changes       non-null      object 
22   deposit_type          non-null      object 
23   agent                non-null      float64 
24   company              non-null      int64  
25   days_in_waiting_list  non-null      int64  
26   customer_type         non-null      object 
27   adr                  non-null      float64 
28   required_car_parking_spaces non-null      int64  
29   total_of_special_requests non-null      int64  
30   reservation_status    non-null      object 
31   reservation_status_date non-null      object 
dtypes: float64(4), int64(16), float64(4), int64(16), object(11)
memory usage: 29.3+ MB
```

1. unique values of object datatype

```
In [10]: df.describe(include = "object")
```

```
Out[10]:
```

	hotel	arrival_date_month	meal	country	market_segment	distribution_channel	reserved_room_type	assigned_room_type	deposit_type	customer_type	reservation_status	
count	119369	11	12	5	177	8	5	10	12	3	4	3
unique	2	12	5	177	8	5	10	12	3	4	3	3
top	City Hotel	August	BB	PBE	Online TA	TATTO	A	A	No Deposit	Transient	Check-Out	
freq	79330	13877	82310	4680	56477	97870	60994	74053	104641	89613	75146	

```
In [11]: for col in df.describe(include = "object"):
print(col)
print(df[col].unique())
print("\n--")
```

```
hotel
['Resort Hotel' 'City Hotel']

arrival_date_month
['July' 'August' 'September' 'October' 'November' 'December' 'January'
 'February' 'March' 'April' 'May' 'June']

meal
['B' 'D' 'HB' 'SC' 'Undefined']

country
['PRT' 'GBR' 'USA' 'ESP' 'ITA' 'FRA' 'ROU' 'NOR' 'CPM' 'ARG' 'POL'
 'CHL' 'BEL' 'CHE' 'CAN' 'GRC' 'ITA' 'ROM' 'ISR' 'SWE' 'JAP' 'EST'
 'CZE' 'BRA' 'FIN' 'HUN' 'RUS' 'UKR' 'DNK' 'ALB' 'IND' 'CHW' 'PER' 'PAK'
 'VEN' 'DNK' 'USA' 'PRT' 'ESP' 'COL' 'AUT' 'BLU' 'LTU' 'TUR' 'ZAF' 'AGO'
 'ISL' 'CWM' 'CPV' 'CPV' 'QZA' 'KOR' 'CRI' 'HUN' 'ARG' 'TUR' 'JAM'
 'MLI' 'ARM' 'AND' 'AND' 'GER' 'URY' 'ECU' 'SRB' 'CYP' 'COL' 'SVK'
 'KWT' 'MDA' 'HRV' 'VEN' 'SRB' 'FJI' 'KAZ' 'PAK' 'IDN' 'LBN' 'PHL' 'SEN'
 'EGY' 'AZE' 'ARM' 'UZB' 'TUR' 'DOM' 'PRY' 'SRB' 'ARM' 'JPN' 'KAZ' 'CUB'
 'CMR' 'BEN' 'RUS' 'CPM' 'SUR' 'UGA' 'BDI' 'CYP' 'JOM' 'SYR' 'ZMB' 'BOT'
 'SDN' 'CHN' 'VUA' 'KWT' 'KWT' 'MLT' 'MAR' 'SRB' 'HND' 'ESA' 'COS'
 'NPL' 'BHS' 'PAR' 'YGO' 'TMR' 'DJT' 'STP' 'KNA' 'ETH' 'IRQ' 'ABO' 'PAK'
 'ARM' 'HND' 'IRN' 'TJK' 'KOL' 'BEN' 'TUR' 'KAZ' 'DNK' 'THP'
 'GUP' 'KEN' 'LIE' 'KOR' 'RME' 'JMT' 'MYT' 'FRO' 'ARM' 'PAK' 'BFA' 'LBY'
 'MLI' 'ARM' 'BOL' 'VEN' 'BGR' 'ARM' 'AZA' 'SLV' 'HND' 'PRY' 'GER' 'LEA'
 'ATA' 'CHN' 'ARM' 'PRT' 'MLI' 'KIL' 'DMR' 'ATP' 'GLE' 'LAO']

market_segment
['Direct' 'Corporate' 'Online TA' 'Offline TATTO' 'Complementary' 'Groups'
 'Undefined' 'Aviation']

distribution_channel
['Direct' 'Corporate' 'TA/TTO' 'Undefined' 'GDS']

reserved_room_type
['C' 'A' 'D' 'E' 'G' 'F' 'P' 'B' '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' 'Cancelled' 'No-show']
```

```
In [12]: pd.isnull(df).sum()
```

```
Out[12]:
```

	hotel	is_cancelled	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	is_repeated_guest	previous_cancellations	previous_bookings_not_cancelled	booking_changes	deposit_type	customer_type	reservation_status
count	119369	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

```
Out[12]:
```

	hotel	is_cancelled	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	is_repeated_guest	previous_cancellations	previous_bookings_not_cancelled	booking_changes	deposit_type	customer_type	reservation_status
mean	0.371342	104.312018	2016.157657	27.166674	15.808082	0.928965	2.621445												

