

```

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings("ignore")

df=pd.read_csv(r"C:\Users\MARK\1.0\OneDrive\Desktop\ML Practice\hotel_booking_data (1).csv")#load dataset
df

In [72]:
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 ... customer_type adr required_room_count

0 0 0 342 2015 7 27 1 0 0 2 ... Transient 0.00

1 1 0 737 2015 7 27 1 0 0 2 ... Transient 0.00

2 0 7 2015 7 27 1 0 1 1 ... Transient 75.00

3 0 13 2015 7 27 1 0 1 1 ... Transient 75.00

4 0 14 2015 7 27 1 0 2 2 ... Transient 98.00

... ..

119385 City 0 23 2017 8 35 30 2 5 2 ... Transient 96.14

119386 City 0 102 2017 8 35 31 2 5 3 ... Transient 225.43

119387 City 0 34 2017 8 35 31 2 5 2 ... Transient 157.71

119388 City 0 109 2017 8 35 31 2 5 2 ... Transient 104.40

119389 City 0 205 2017 8 35 29 2 7 2 ... Transient 151.20

119390 rows x 36 columns

In [73]:
df.shape

Out[73]:
(119390, 36)

In [74]:
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 119390 entries, 0 to 119389
Data columns (total 36 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 int64
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 int64
21 booking_changes 119390 non-null int64
22 deposit_type 119390 non-null object
23 agent 103058 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 object
29 total_of_special_requests 119390 non-null int64
30 reservation_status 119390 non-null object
31 reservation_status_date 119390 non-null object
32 name 119390 non-null object
33 email 119390 non-null object
34 phone-number 119390 non-null object
35 credit_card 119390 non-null object
dtypes: float64(4), int64(16), object(16)
memory usage: 32.8+ MB

In [75]:
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_repeated_guest
mean 0.370416 104.014116 2016.166554 119390.000000 119390.000000 119390.000000 119390.000000 119386.000000 119390.000000 119390.000000
std 0.482528 106.583097 0.707476 13.605138 8.780829 0.988613 1.902396 0.572861 0.98561 0.097436 0.1757
min 0.000000 0.000000 2015.000000 15.000000 1.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.0000
25% 0.000000 18.000000 2016.000000 15.000000 8.000000 0.000000 1.000000 2.000000 0.000000 0.000000 0.0000
50% 0.000000 69.000000 2016.000000 26.000000 16.000000 1.000000 2.000000 2.000000 0.000000 0.000000 0.0000
75% 1.000000 160.000000 2017.000000 38.000000 23.000000 2.000000 3.000000 2.000000 0.000000 0.000000 0.0000
max 1.000000 737.000000 2017.000000 53.000000 31.000000 19.000000 50.000000 55.000000 10.000000 10.000000 1.0000

In [76]:
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 16348
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
name 0
email 0
phone-number 0
credit_card 0
dtype: int64

In [77]:
df.drop("company",axis=1,inplace=True)

In [78]:
df.drop("agent",axis=1,inplace=True)

In [79]:
df.dropna(subset=["children","country"],inplace=True)

In [80]:
df.duplicated().sum()

0

In [81]:
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 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
name 0
email 0
phone-number 0
credit_card 0
dtype: int64

In [82]:
df.drop("is_repeated_guest",axis=1,inplace=True)

In [83]:
df.drop("babies",axis=1,inplace=True)

In [84]:
df.drop("required_car_parking_spaces",axis=1,inplace=True)

In [85]:
cat=df.select_dtypes(object).columns
cat

Index(['hotel', 'arrival_date_month', 'meal', 'country', 'market_segment',
       'distribution_channel', 'reserved_room_type', 'assigned_room_type',
       'deposit_type', 'customer_type', 'reservation_status',
       'reservation_status_date', 'name', 'email', 'phone-number',
       'credit_card'],
      dtype='object')

In [86]:
from sklearn.preprocessing import OrdinalEncoder
oe=OrdinalEncoder()
df[cat]=oe.fit_transform(df[cat])

In [87]:
df

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 ... days_in_waiting_list customer_type

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

118950 rows x 31 columns

In [88]:
d = df["is_canceled"].value_counts()
plt.figure(figsize=(8,8))
p = plt.pie(d, labels=d.index, autopct="%0.0%")
plt.title("Booking Details")

Text(0.5, 1.0, 'Booking Details')

0
67%
37%
1

In [89]:
from scipy.stats import zscore
z=np.abs(zscore(df))
z

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 ... days_in_waiting_list customer_type

0 1.415195 0.78858 2.223407 1.636365 0.138739 0.012256 1.685694 0.932429 1.316808 0.244755 ... 0.132201
1 1.415195 0.78858 5.918350 1.636365 0.138739 0.012256 1.685694 0.932429 1.316808 0.244755 ... 0.132201
2 1.415195 0.78858 0.910279 1.636365 0.138739 0.012256 1.685694 0.932429 1.790536 1.483635 ... 0.132201
3 1.415195 0.78858 0.854153 1.636365 0.138739 0.012256 1.685694 0.932429 1.790536 1.483635 ... 0.132201
4 1.415195 0.78858 0.844799 1.636365 0.138739 0.012256 1.685694 0.932429 0.264264 0.244755 ... 0.132201
... ..
119385 0.706616 0.78858 0.706010 1.190666 1.270292 0.576416 1.617159 1.075176 1.314550 0.244755 ... 0.132201
119386 0.706616 0.78858 0.021622 1.190666 1.270292 0.576416 1.731050 1.075176 1.314550 1.973145 ... 0.132201
119387 0.706616 0.78858 0.607713 1.
```