

```
In [3]: ## Load Packages

import os
import pandas as pd
import numpy as np
import seaborn as sns
from matplotlib import pyplot as plt
```

Read and Understand the Data

```
In [3]: hotel_data = pd.read_csv('/Users/doyin/documents/Data Projects/hotel_bookings.csv')
hotel_data.head(5) # check first 5 rows of data

Out[3]:
```

	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	...	deposit_type	agent	comp
0	Resort Hotel	0	342	2015	July	27	1	0	0	2	...	No Deposit	NaN	1
1	Resort Hotel	0	737	2015	July	27	1	0	0	2	...	No Deposit	NaN	1
2	Resort Hotel	0	7	2015	July	27	1	0	1	1	...	No Deposit	NaN	1
3	Resort Hotel	0	13	2015	July	27	1	0	1	1	...	No Deposit	304.0	1
4	Resort Hotel	0	14	2015	July	27	1	0	2	2	...	No Deposit	240.0	1

5 rows × 32 columns

```
In [4]: hotel_data.dtypes # understanding the format of the data

Out[4]:
```

hotel		object
is_canceled		int64
lead_time		int64
arrival_date_year		int64
arrival_date_month		object
arrival_date_week_number		int64
arrival_date_day_of_month		int64
stays_in_weekend_nights		int64
stays_in_week_nights		int64
adults		int64
children		float64
babies		int64
meal		object
country		object
market_segment		object
distribution_channel		object
is_repeated_guest		int64
previous_cancellations		int64
previous_bookings_not_canceled		int64
reserved_room_type		object
assigned_room_type		object
booking_changes		int64
deposit_type		object
agent		float64
company		float64
days_in_waiting_list		int64
customer_type		object
adr		float64
required_car_parking_spaces		int64
total_of_special_requests		int64
reservation_status		object
reservation_status_date		object
dtype:		object

Deciding on Variables

```
In [10]: # deleting agent and company column

del hotel_data["agent"]
del hotel_data["company"]

In [5]: # looking further into unsure columns to decide if data is useful
hotel_data["market_segment"].describe()

Out[5]:
```

	count	unique	top	freq	Name:
	119390	8	Online TA	56477	market_segment, dtype: object

```
In [6]: hotel_data["market_segment"].unique()

Out[6]: array(['Direct', 'Corporate', 'Online TA', 'Offline TA/TO', 'Complimentary', 'Groups', 'Undefined', 'Aviation'], dtype=object)

In [13]: del hotel_data["market_segment"]

In [14]: del hotel_data["distribution_channel"]
```

Transforming Variables

```
In [7]: # cancelled column variable can be changed
hotel_data["is_canceled"].dtype

Out[7]: dtype('int64')
```

```
In [16]: new_is_canceled = pd.Categorical(hotel_data["is_canceled"])
new_is_canceled = new_is_canceled.rename_categories(["Cancelled", "Booked"])
new_is_canceled.describe()

Out[16]:
```

	counts	freqs
categories		
Cancelled	75166	0.629584
Booked	44224	0.370416

```
In [17]: hotel_data["Booking Status"] = new_is_canceled

In [18]: hotel_data.head(5)

Out[18]:
```

	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	...	agent	company	days_in_wait
0	Resort Hotel	0	342	2015	July	27	1	0	0	2	...	NaN	NaN	1
1	Resort Hotel	0	737	2015	July	27	1	0	0	2	...	NaN	NaN	1
2	Resort Hotel	0	7	2015	July	27	1	0	1	1	...	NaN	NaN	1
3	Resort Hotel	0	13	2015	July	27	1	0	1	1	...	304.0	NaN	1
4	Resort Hotel	0	14	2015	July	27	1	0	2	2	...	240.0	NaN	1

5 rows × 33 columns

```
In [20]: # del hotel_data[is_canceled]
```

Combining Dates Columns

hotel_data["Arrival Date"] = pd.to_datetime(hotel_data["arrival_date_year"].astype(str) + '/' + hotel_data["arrival_date_month"].astype(str) + '/' + hotel_data["arrival_date_day_of_month"].astype(str))

```
In [ ]:

In [22]: hotel_data.head(5)

Out[22]:
```

	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	...	deposit_type	days_in_wait
0	Resort Hotel	0	342	2015	July	27	1	0	0	2	...	No Deposit	
1	Resort Hotel	0	737	2015	July	27	1	0	0	2	...	No Deposit	
2	Resort Hotel	0	7	2015	July	27	1	0	1	1	...	No Deposit	
3	Resort Hotel	0	13	2015	July	27	1	0	1	1	...	No Deposit	
4	Resort Hotel	0	14	2015	July	27	1	0	2	2	...	No Deposit	

5 rows × 30 columns

```
In [20]: # combining children and babies
hotel_data["children"].describe()

Out[20]:
```

	count	mean	std	min	25%	50%	75%	max	Name:
	119390.000000	0.103890	0.398561	0.000000	0.000000	0.000000	0.000000	10.000000	children, dtype: float64

```
In [21]: np.sum(hotel_data.isnull())

Out[21]:
```

	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	Booking Status	Arrival Date	dtype:
		0	0	0	0	0	0	0	0	0	4	0	0	488	0	0	0	0	0	0	0	0	0	16340	112593	0	0	0	0	0	0	0	0	int64	

```
In [22]: new_children = np.where(hotel_data["children"].isnull(),
                                np.mean(hotel_data["children"]),
                                hotel_data["children"])

hotel_data["children"] = new_children
np.sum(hotel_data["children"].isnull())

Out[22]: 0

In [26]: hotel_data["Children"] = hotel_data["children"].astype(int) + hotel_data["babies"].astype(int)

In [27]: # hotel cancelled bookings
hotel_data["Booking Status"].value_counts()

Out[27]:
```

	Cancelled	Booked	Name:
	75166	44224	Booking Status, dtype: int64

```
In [23]: Cancelled = hotel_data[hotel_data["Booking Status"] == "Cancelled"]
Cancelled_per = (hotel_data["Booking Status"].shape[0]/len(Cancelled))*100

#here I have shown how to use shape[0] to get the total number of rows also use len() to get the total number of rows cancelled
Cancelled_per

Out[23]: 158.8351161961525

In [24]: hotel_data["Arrival_date_month"].value_counts()

Out[24]:
```

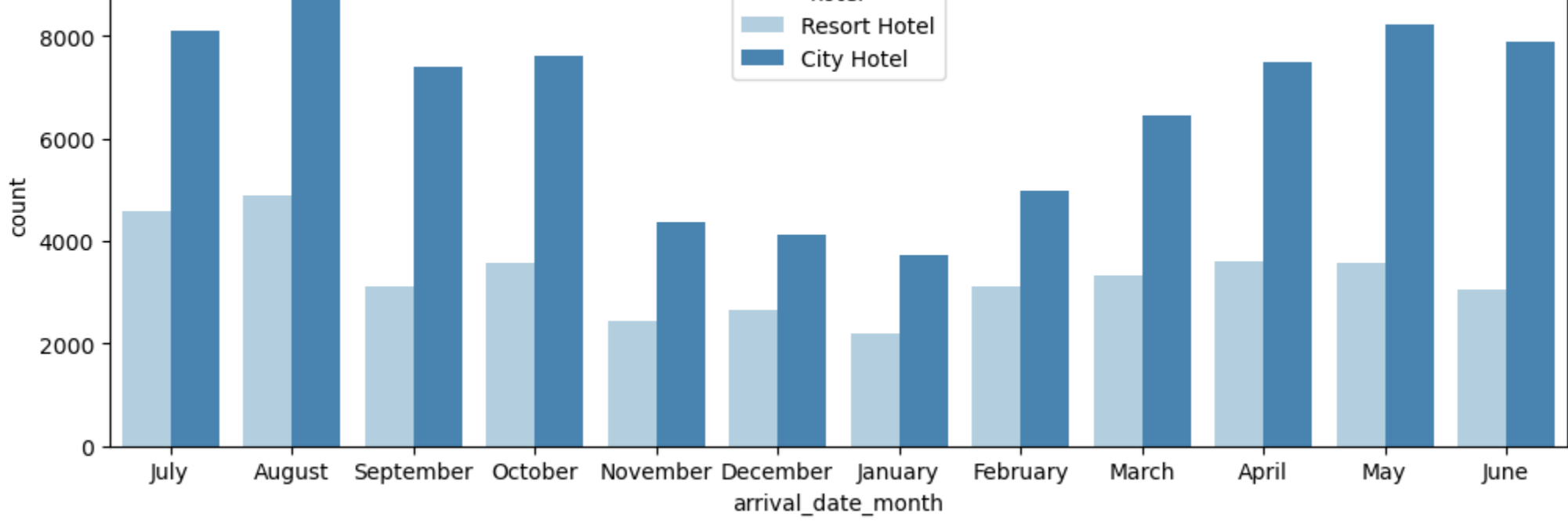
	August	July	May	October	April	June	September	March	February	November	December	January
	13877	12661	11791	11160	11089	10939	10508	9794	8668	6794	6780	5929

Name: arrival_date_month, dtype: int64

Visualising the data

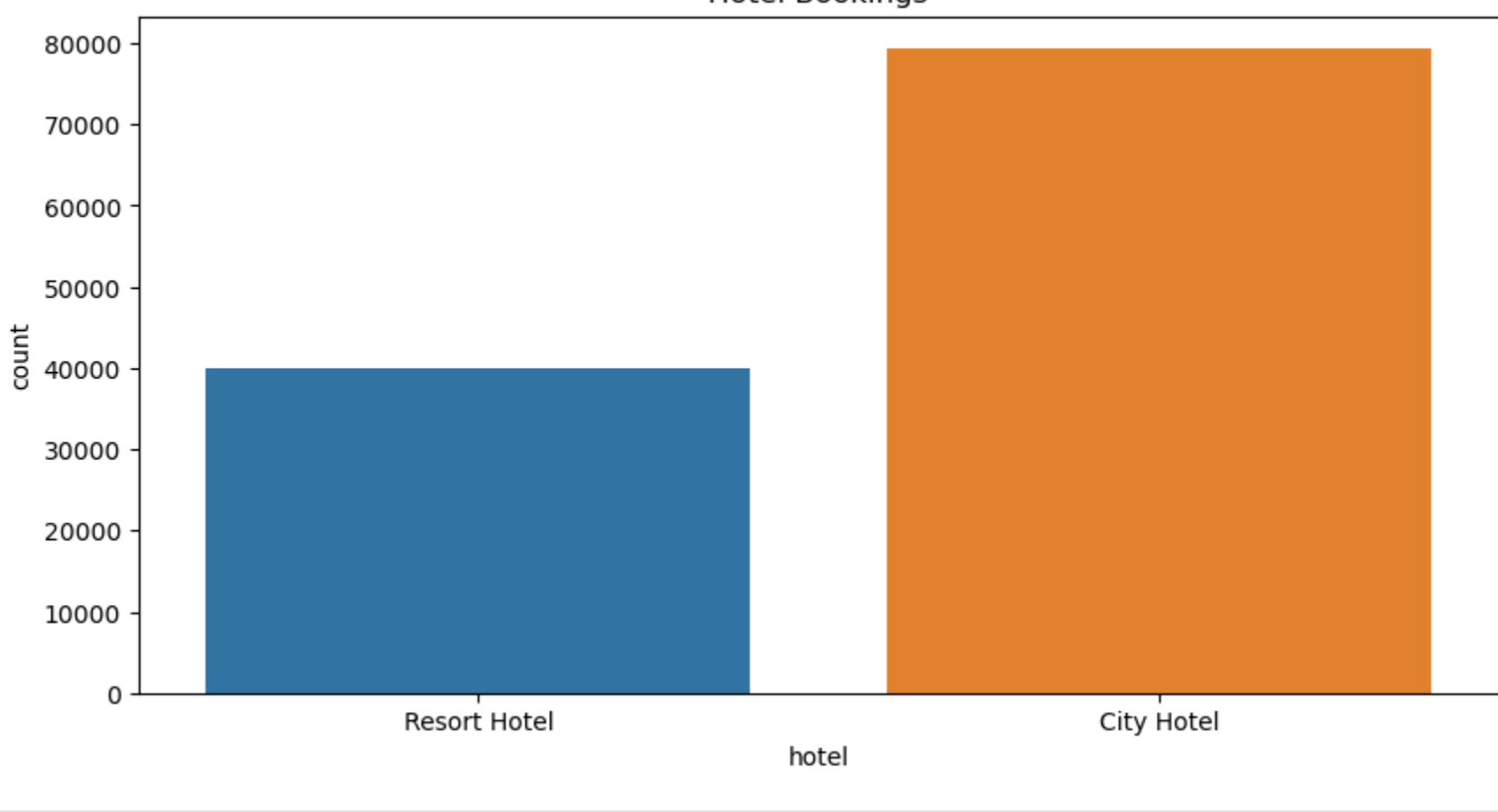
```
In [25]: plt.figure(figsize = (12,4))
sns.countplot(x = 'arrival_date_month', hue = 'hotel', data = hotel_data, palette = 'Blues')
plt.title("Hotel Bookings per Month")

Out[25]: Text(0.5, 1.0, 'Hotel Bookings per Month')
```

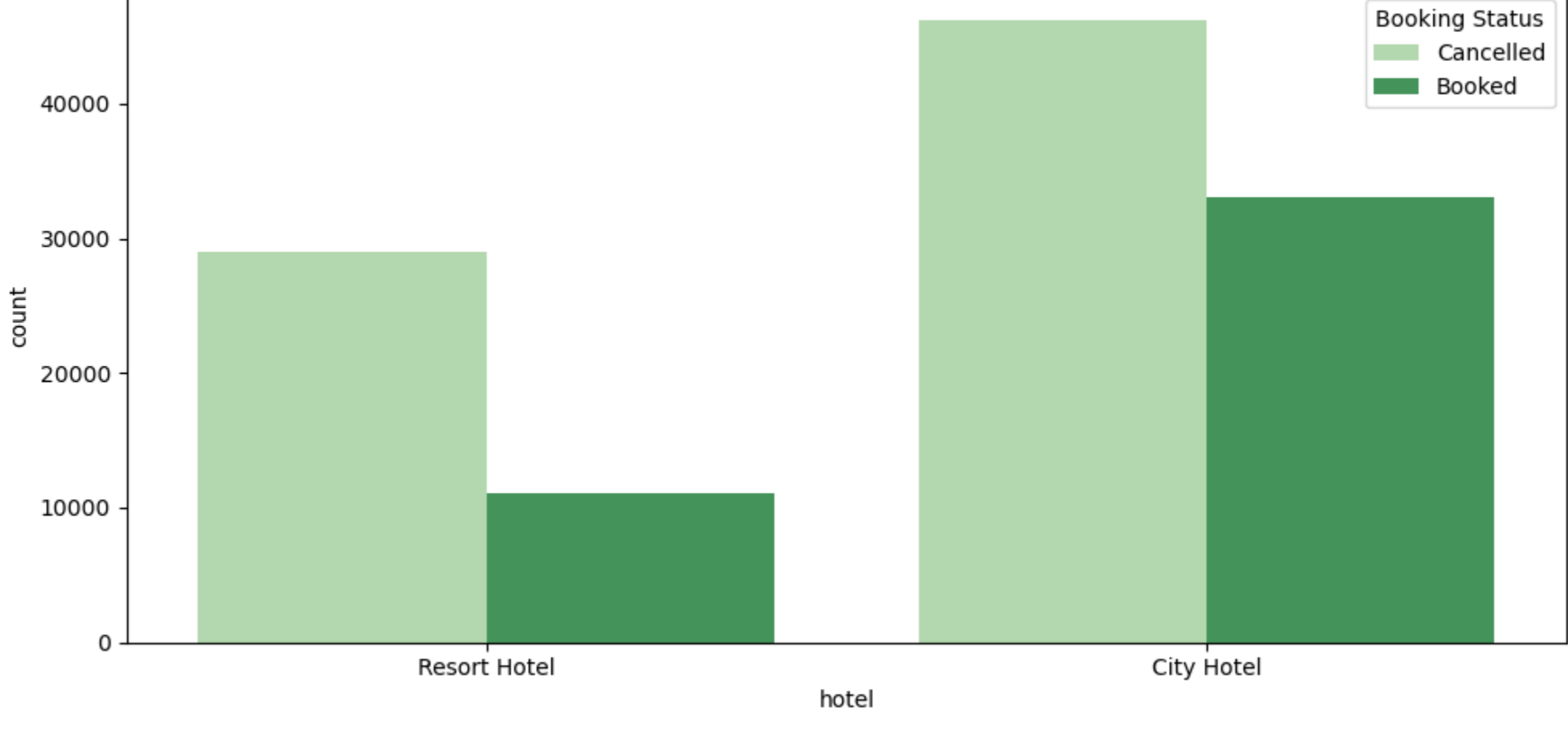


```
In [32]: plt.figure(figsize = (10,5))
sns.countplot(x = 'hotel', data = hotel_data)
plt.title("Hotel Bookings")

Out[32]: Text(0.5, 1.0, 'Hotel Bookings')
```



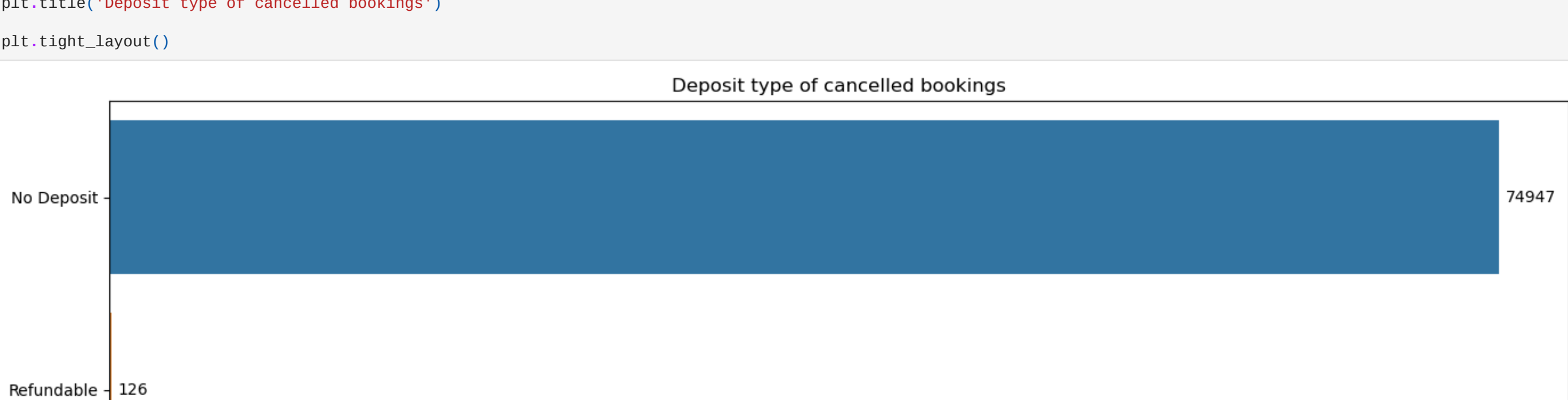
```
In [33]: plt.figure(figsize = (10,5))
sns.countplot(x = 'hotel', data = hotel_data, hue = 'Booking Status', palette = 'Greens')
plt.title("Canceleation vs Confirmed Bookings")
plt.tight_layout()
```



```
In [20]: # if null values in column is more than 70% of overall rows
for col in hotel_data.columns:
    if np.sum(hotel_data[col].isnull()) > (hotel_data.shape[0]*0.70):
        print("print above 70%")

print above 70%
```

```
In [30]: fig, ax = plt.subplots(figsize = (14, 6))
sns.countplot(ax = ax,
              y = 'deposit_type',
              data = Cancelled,
              orient = "h",
              )
ax.bar_label(ax.containers[0], padding = 4)
plt.title('Deposit type of cancelled bookings')
plt.tight_layout()
```



```
In [ ]:
```