

```
In [15]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import psycopg2

db = psycopg2.connect(
    host='localhost',
    user='postgres',
    password='postgres',
    database='E-commerce'
)

cur = db.cursor()
```

1. List all unique cities where customers are located.

```
In [41]: query = """ select distinct(customer_city) from customers """
cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data, columns = ["Customer_cities"])
df

Out[41]:
   Customer_cities
0      bon jardin de roses
1      alto rio doce
2      alvorada do gurguaia
3      batatais
4      capao da porteira
...
4114      carbonita
4115      concordia do para
4116      independencia
4117      governador valadares
4118      baliao nova
4119 rows x 1 columns
```

2. Count the number of orders placed in 2017.

```
In [42]: query = """ select count(extract(year from cast(order_booked as TIMESTAMP))) AS years from orders
where extract(year from cast(order_booked as TIMESTAMP)) = 2017 """
cur.execute(query)
data = cur.fetchall()
["Total order placed in 2017 is", data[0][0]]

Out[42]:
('Total order placed in 2017 is', 96982)
```

3. Find the Top 5 Category by Sales.

```
In [43]: query = """ select p.product_category, cast(sum(ps.payment_value) as integer) as Total_Sales from products AS p
join order_items as o
ON o.order_id = o.p_order_id
join payments as ps
ON o.order_id = ps.order_id
group by 1
order by 2 desc
limit 5"""
cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data, columns = ["Category", "Total_Sales"])
df

Out[43]:
   Category  Total_Sales
0      bed table bath      6850216
1      HEALTH BEAUTY      6628492
2      computer accessories      6341322
3      Furniture Decoration      5720706
4      Watches present      5716887

In [44]: mtn.figure(figsize = (17,8))
ax = sns.barplot(x = "Category", y = "Total_Sales", data = df, width = 0.6, color = "brown")
ax.bar_label(ax.containers[0])
mtn.show()
```



4. Calculate the percentage of orders that were placed in installments.

```
In [46]: query = """ select round((sum(case when payment_installments >= 3 then 1 else 0 end)/1.8)) * 100, 3) as installments
from payments """
cur.execute(query)
data = cur.fetchall()
["The Percentage of Orders That were paid in installments is ", data[0][0]]

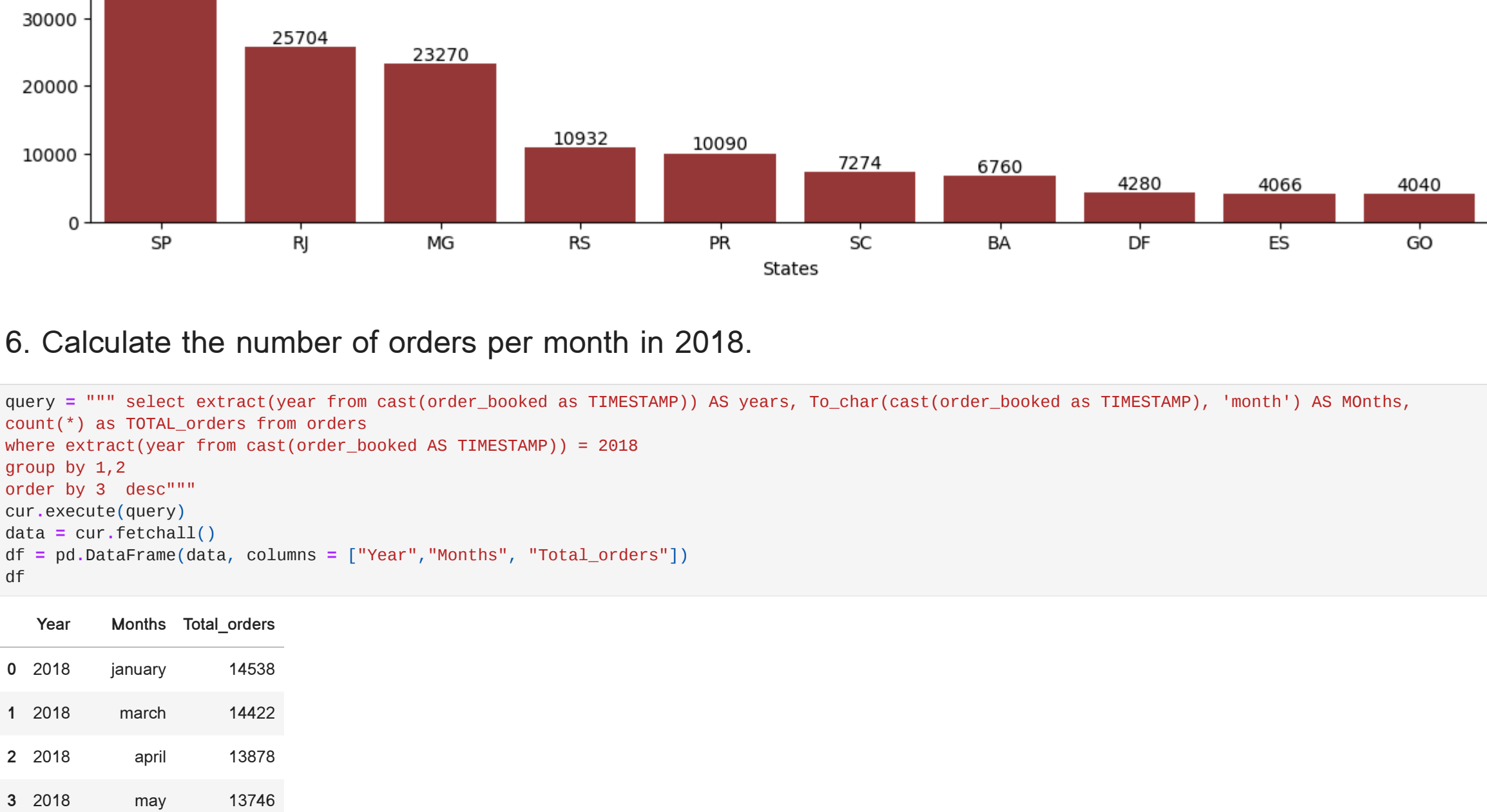
Out[46]:
('The Percentage Of Orders That were paid in installments is ',
Decimal('99.998'))
```

5. Top 10 States on the behalf of Customers

```
In [37]: query = """ select customer_state, count(*) AS Total_Customers from customers
group by 1
order by 2 desc
limit 10 """
cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data, columns = ["States", "Total_Customers"])
df

Out[37]:
   States  Total_Customers
0      SP              83492
1      RJ              25704
2      MG              23270
3      RS              10932
4      PR              10090
5      SC              7274
6      BA              6700
7      DF              4280
8      ES              4066
9      GO              4040

In [40]: mtn.figure(figsize = (14,6))
ax = sns.barplot(x = "States", y = "Total_Customers", data = df, color = "brown")
ax.bar_label(ax.containers[0])
mtn.title("Top 10 States according to customers", fontsize = 20)
mtn.show()
```



6. Calculate the number of orders per month in 2018.

```
In [42]: query = """ select extract(year from cast(order_booked as TIMESTAMP)) AS years, to_char(cast(order_booked as TIMESTAMP), 'month') AS Months,
count(*) as TOTAL_orders from orders
where extract(year from cast(order_booked as TIMESTAMP)) = 2018
group by 2,3
order by 3 desc """
cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data, columns = ["Year", "Months", "Total_orders"])
df

Out[42]:
   Year  Months  Total_orders
0  2018  january           14538
1  2018   march           14422
2  2018   april           13078
3  2018    may            13746
4  2018  february          13456
5  2018   july            13024
6  2018   june            12884
7  2018   august           12334
8  2018  september           32
9  2018  october            8

In [45]: mtn.figure(figsize = (15,7))
ax = sns.barplot(x = "Months", y = "Total_orders", data = df, color = "r")
ax.bar_label(ax.containers[0])
mtn.title("Total sales in 2018", fontsize = 20)
mtn.show()
```



January 2018 got the highest number of Orders

7. Find the average number of products per order, grouped by customer city.

```
In [72]: query = """ with one as (select o.order_id, o.customer_id, count(o.order_id) as Total_orders from orders as o
join order_items as oi
ON o.order_id = oi.order_id
group by 2,3
select o.customer_city, round(avg(one.Total_orders)) AS avg_per_order from customers as c
join one
ON c.customer_id = one.customer_id
group by 1 """
cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data, columns = ["Customer_city", "Avg_Orders"])
df

Out[72]:
   Customer_city  Avg_Orders
0      bon jardin de roses      2.75
1      alto rio doce          2.00
2      alvorada do gurguaia    2.00
3      batatais               2.21
4      capao da porteira       2.00
...
4105      concordia do para     2.00
4106      carbonita            2.00
4107      independencia        2.00
4108      governador valadares  2.19
4109      baliao nova          2.00
4110 rows x 2 columns
```

8. Calculate the percentage of total revenue contributed by each product category.

```
In [74]: query = """ select product_category, round(cast((total_revenue/(select sum(payment_value) from payments)))100 AS Decimal),2) as total_perlon from
(select p.product_category, cast(sum(p.payment_value) as integer) as total_revenue from products AS p
join order_items as oi
ON p.product_id = o.product_id
join payments as ps
ON o.order_id = ps.order_id
group by 1
order by 2 desc """
cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data, columns = ["Category", "Sales%"])
df

Out[74]:
   Category  Sales%
0      bed table bath      21.40
1      HEALTH BEAUTY      20.71
2      computer accessories      19.81
3      Furniture Decoration      17.87
4      Watches present       17.86
...
69      PC Gamer            0.03
70      House Comfort 2      0.02
...
72      Fashion Children's Clothing      0.01
73      insurance and services      0.00
74 rows x 2 columns
```

9. Identify the correlation between product price and the number of times a product has been purchased.

```
In [77]: query = """ select p.product_category, round(cast(avg(oq.price) AS decimal),2) as avg_price, count(oq.order_id) as Total_orders from products as p
join order_items as oi
ON oi.order_id = p.product_id
join oq
ON oi.order_id = oq.order_id
group by 1
order by 3 desc """
cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data, columns = ["Category", "Price", "Total_orders"])
df

Out[77]:
   Category  Price  Total_orders
0      bed table bath      93.30      22230
1      HEALTH BEAUTY      130.16      18940
2      sport leisure       114.34      17282
3      Furniture Decoration      67.56      16668
4      computer accessories      116.51      15654
...
69      cds music discs       52.14         28
70      La Cuisine           146.79         28
71      PC Gamer            171.77         18
72      Fashion Children's Clothing      71.23         16
73      insurance and services      141.65          4
74 rows x 3 columns
```

10. Identify the Top 6 Rank holders by total sales

```
In [36]: query = """ select *, rank()over(order by total_revenue desc) AS Rank from
(select s.seller_id, cast(sum(p.payment_value) as INTEGER) AS Total_revenue from sellers as s
join payments as p
ON s.seller_id = p.seller_id
join payments as ps
ON o.order_id = ps.order_id
group by 1
order by 2 desc """
cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data, columns = ["Seller_id", "Total_revenue", "Rank"])
df

Out[36]:
   Seller_id  Total_revenue  Rank
0  7ed7a144b0205e963d30c0a6010ab      2028665      1
1  10203a244467041d0b30b05a050a6cfa      1232885      2
2  4eb30515b174a040b6374361493884a      1204981      3
3  150f030178b18a58a9849023252100      1161014      4
4  6324358a14d6a2430216180348905      1139612      5
5  de8f023146b17ba28114ed3b05a48da      1038877      6

In [36]: mtn.figure(figsize = (12,7))
ax = sns.barplot(x = "Seller_id", y = "Total_revenue", data = df, width = 0.5, color = "#808080")
ax.bar_label(ax.containers[0])
mtn.title("Top 6 Rank holders by Sales", fontsize = 20)
mtn.xticks(rotation = 45)
mtn.show()
```



11. Calculate the moving average of order values for each customer over their order history.

```
In [47]: query = """select customer_id, order_booked, payment,
avg(payment)over(partition by customer_id order by order_booked rows between 2 preceding and current row) from
(select o.customer_id, o.order_booked, p.payment_value as payment from orders as o
join payments as p
ON o.order_id = p.order_id) """
cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data, columns = ["Customer_id", "Date", "Payment", "Moving_avg"])
df

Out[47]:
   Customer_id  Date  Payment  Moving_avg
0  00012a2e0ff80a20a050a9849f703  2017-11-14  16.0826      114.74
1  00012a2e0ff80a20a050a9849f703  2017-11-14  16.0826      114.74
2  00012a2e0ff80a20a050a9849f703  2017-11-14  16.0826      114.74
3  00012a2e0ff80a20a050a9849f703  2017-11-14  16.0826      114.74
4  0001f4e05b00a020910070a4427140  2017-07-16  09.4032      67.41      67.41
...
415539  ff031725277f656a70036a7e63aae8  2017-09-02  11.5332      45.50      45.50
415540  ff06050ea30b76d5a076a7f60b0b99  2017-09-29  14.0703      18.37      18.37
415541  ff06050ea30b76d5a076a7f60b0b99  2017-09-29  14.0703      18.37      18.37
415542  ff06050ea30b76d5a076a7f60b0b99  2017-09-29  14.0703      18.37      18.37
415543  ff06050ea30b76d5a076a7f60b0b99  2017-09-29  14.0703      18.37      18.37
415544 rows x 4 columns
```

12. Calculate the cumulative sales per month for each year.

```
In [4]: query = """select *, sum(total_revenue)over(partition by years order by months) as cumulative_revenue from
(select extract(year from cast(o.order_booked as TIMESTAMP)) AS years,
extract(month from cast(o.order_booked as TIMESTAMP)) as months,
round(cast(sum(p.payment_value) as decimal),2) as Total_revenue
from orders as o
join payments as p
ON o.order_id = p.order_id
group by 2,3 """
cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data, columns = ["years", "Months", "Total_sales", "Cumulative_sales"])
df

Out[4]:
   Years  Months  Total_sales  Cumulative_sales
0  2016         9      1008.96      1008.96
1  2016        10      232861.92      232972.88
2  2016        12       78.46      232954.36
3  2017         1      553962.16      553962.16
4  2017         2      1167032.04      1721094.20
5  2017         3      1798564.60      3521058.80
6  2017         4      1671512.12      5192570.92
7  2017         5      2374976.28      7567547.00
8  2017         6      2045055.52      9609602.52
9  2017         7      2386931.68      11996534.20
10  2017         8      2697045.38      14693579.48
11  2017         9      2911549.80      17605129.28
12  2017        10      318711.52      20792240.80
13  2017        11      477831.20      25470551.80
14  2017        12      3513805.92      28984357.72
15  2018         1      1460016.72      44080514.44
16  2018         2      3958053.36      84368047.80
17  2018         3      4638026.48      130648314.56
18  2018         4      4643141.92      177179556.48
19  2018         5      4619328.60      223372845.08
20  2018         6      4090322.00      264276067.08
21  2018         7      4306163.00      30683764.08
22  2018         8      4088701.28      34776535.36
23  2018         9      17768.16      34794303.52
24  2018        10      2356.68      34796660.20
```

13. Calculate the year-over-year growth rate of total sales.

```
In [51]: query = """ with one as (select extract(year from cast(o.order_booked as TIMESTAMP)) AS years, round(cast(sum(p.payment_value)as DECIMAL),2) as Total_revenue from orders as o
join payments as p
ON o.order_id = p.order_id
group by 1)
select years,
total_revenue,
round((total_revenue - lag(total_revenue,1)over(order by years))/lag(total_revenue,1)over(order by years)) * 100,2) as YoY_Growth
from one """
cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data, columns = ["Years", "Total_sales", "YoY_Growth"])
df

Out[51]:
   Years  Total_sales  YoY_Growth
0  2016      237449.36      None
1  2017  3899896.82      1212.70
2  2018  34796620.20      20.00
```

14. Calculate the retention rate of customers, defined as the percentage of customers who make another purchase within 6 months of their first purchase.

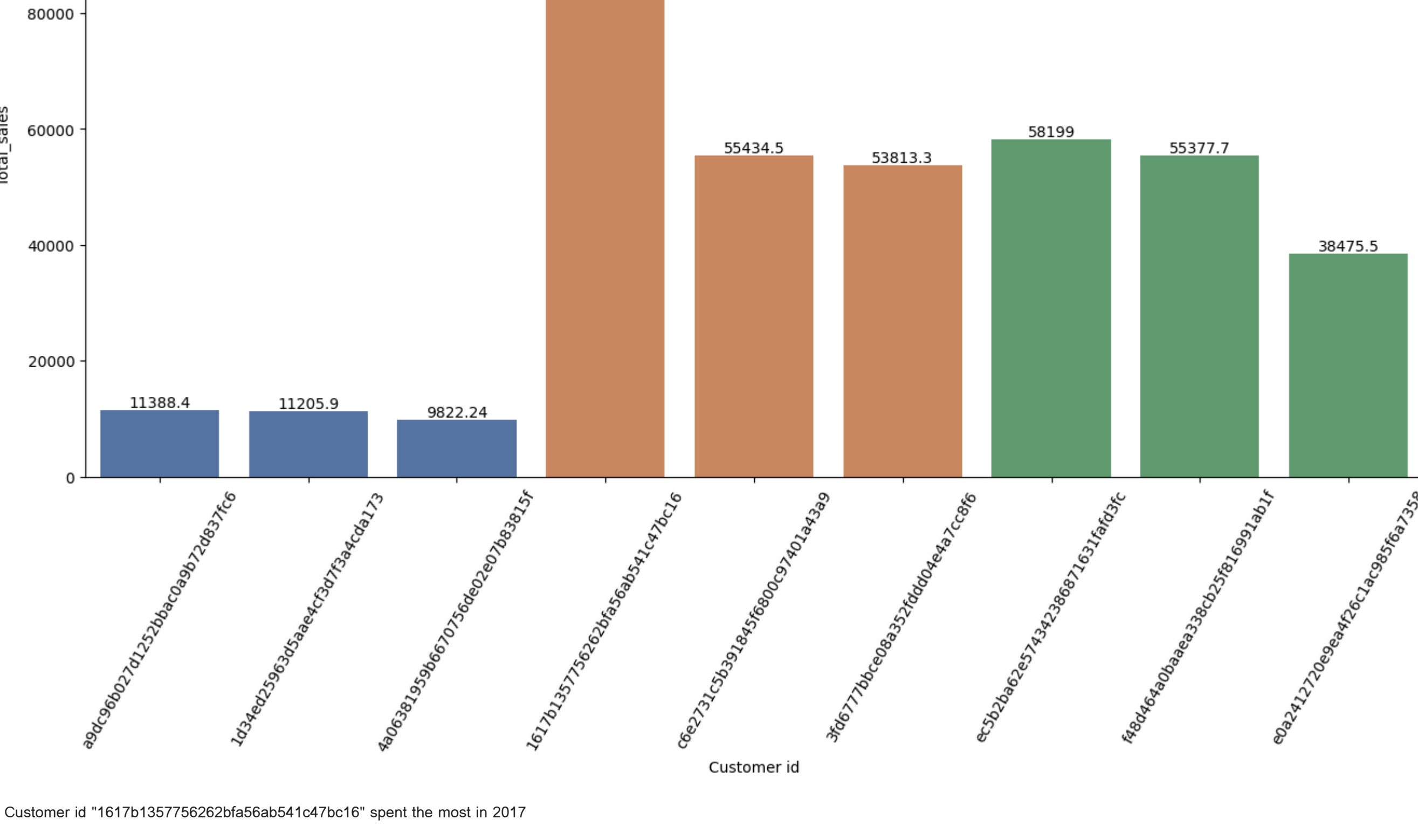
```
In [34]: query = """ with a as (select c.customer_id, min(CAST(o.order_booked as TIMESTAMP)) as first_order from customers c
join orders as o
ON c.customer_id = o.customer_id
select a.customer_id, count(distinct cast(o.order_booked as TIMESTAMP)) AS order from a
join orders as o
ON a.customer_id = o.customer_id
and cast(o.order_booked as TIMESTAMP) > first_order + interval '6 months'
and cast(o.order_booked as TIMESTAMP) < first_order + interval '6 months'
group by 1 """
cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data)
df

Out[34]:
   No customers had purchased the order within 6 month of their 1st order

In [34]: query = """select * from
(select s.seller_id, rank()over(partition by years order by total_revenue desc) as rank from
(select extract(year from cast(order_booked as TIMESTAMP)) AS years, c.customer_id, round(cast(sum(p.payment_value) as decimal),2) as Total_revenue from customers AS c
join payments as p
ON s.seller_id = p.seller_id
join payments as ps
ON o.order_id = ps.order_id
group by 2,3 """
where rank < 3"""
cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data, columns = ["years", "Customer id", "Total_sales", "Rank"])
df

Out[34]:
   years  Customer id  Total_sales  Rank
0  2016  e9a036a027f1252ba0da0b72a8376b      11388.40      1
1  2016  1e34ed59638f6a0e43f7f34da173      11205.92      2
2  2016  4a0381958607056a02a07b0381f9      9822.24      3
3  2016  1617b1357705670626a56a541cf7bc16      9822.24      1
4  2016  df6273125b39184b8830a074614e3a0b      55434.48      2
5  2017  3a857778b0da93255a304e4a7c0d8e      52813.28      3
6  2016  ec2b2a62a74342388671631a1836      58189.04      1
7  2016  f4846a40a0e0329a209f169f1a011      55377.68      2
8  2016  efba412720ba04261ac9a080a7358      38476.52      3

In [34]: mtn.figure(figsize = (16,8))
ax = sns.barplot(x = "Customer id", y = "Total_sales", hue = "years", data = df, palette = "deep")
ax.bar_label(ax.containers[0])
ax.bar_label(ax.containers[1])
mtn.title("Top 3 Customers by their Amount Spent in each Year", fontsize = 20)
mtn.xticks(rotation = 60)
mtn.show()
```



Customer id '1617b1357705670626a56a541cf7bc16' spent the most in 2017

