# Problem Statement
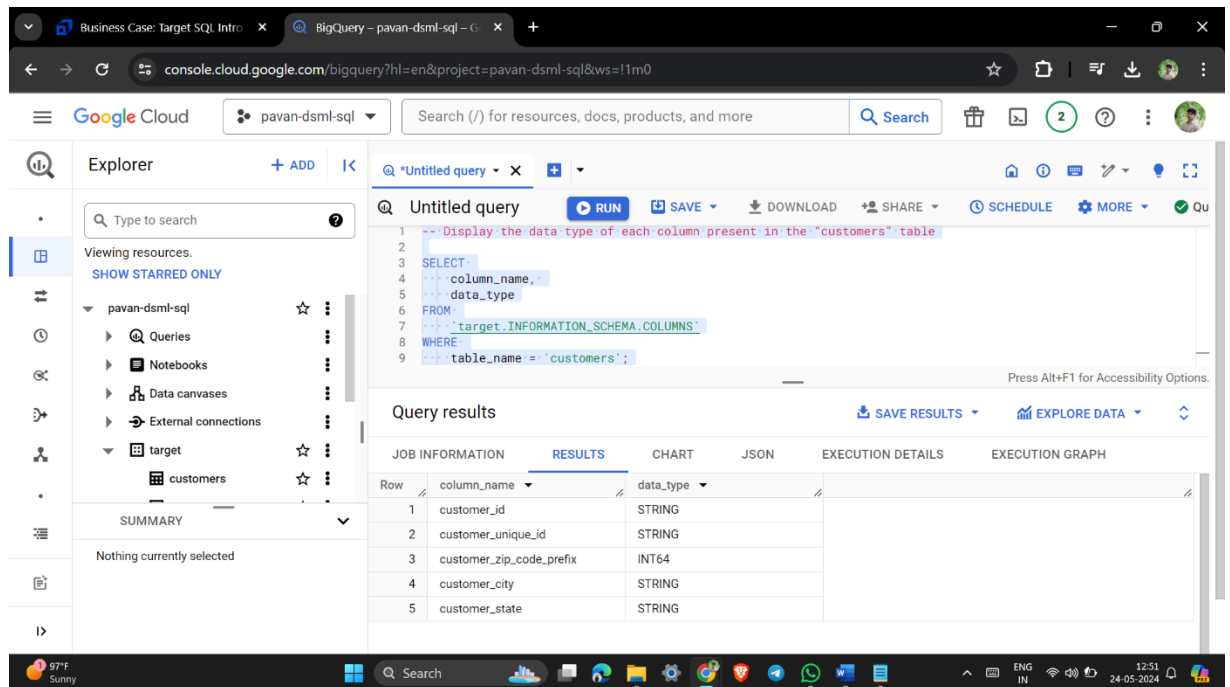
1. **Import the dataset and do usual exploratory analysis steps like checking the structure & characteristics of the dataset:**

## A. Data type of all columns in the "customers" table.

**QUERY:**
```
SELECT
    column_name,
    data_type
FROM
    `target.INFORMATION_SCHEMA.COLUMNS`
WHERE
    table_name = 'customers';
```



## B. Get the time range between which the orders were placed.

**QUERY:**
```
SELECT
    MIN(order_purchase_timestamp) AS first_order_date,
    MAX(order_purchase_timestamp) AS last_order_date
```

FROM
  `target.orders`;



## C. Count the Cities & States of customers who ordered during the given period.

**QUERY:**
SELECT
  COUNT(DISTINCT customer_city) AS unique_cities,
  COUNT(DISTINCT customer_state) AS unique_states
FROM
  `target.customers`;

## 2. In-depthExploration:

**A. Is there a growing trend in the no. of orders placed over the past years?**

**QUERY:**

```
SELECT
  EXTRACT(YEAR FROM order_purchase_timestamp) AS year,
  EXTRACT(MONTH FROM order_purchase_timestamp) AS month,
  COUNT(order_id) AS number_of_orders
FROM
  `target.orders`
GROUP BY
  year, month
ORDER BY
  year, month;
```
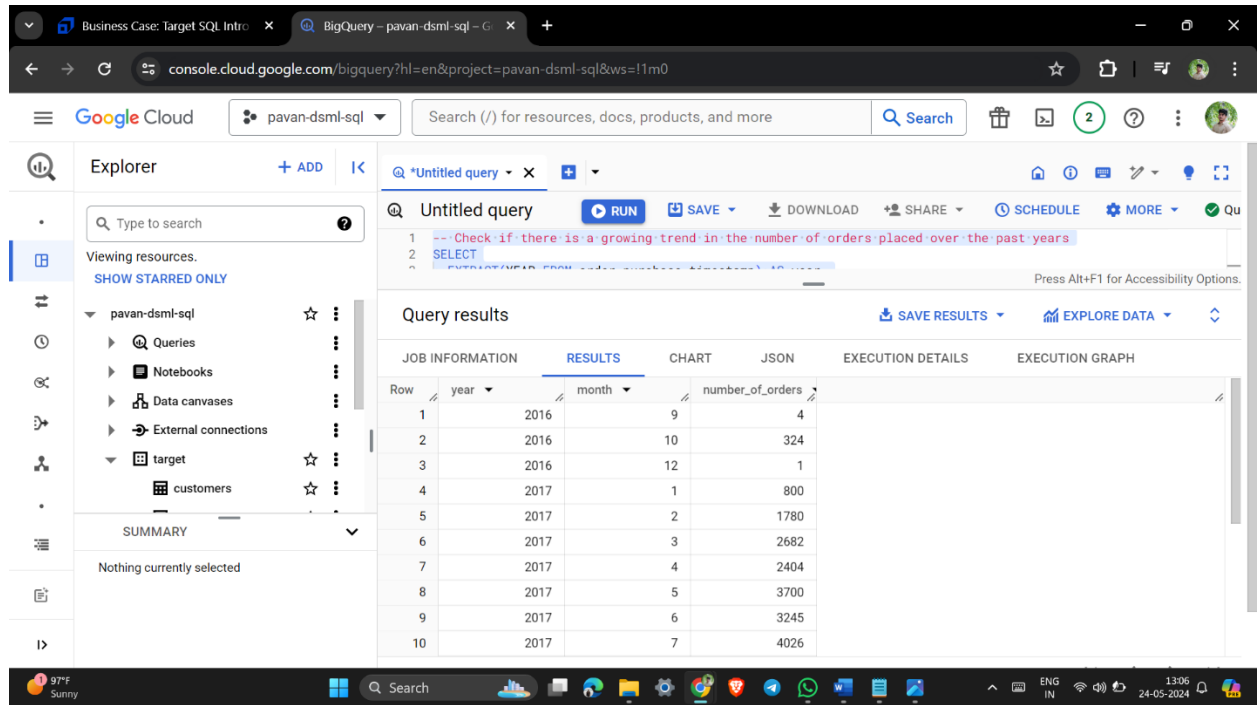
**B.** **Can we see some kind of monthly seasonality in terms of the no. of orders being placed?**

**QUERY:**

SELECT
  EXTRACT(MONTH FROM order_purchase_timestamp) AS month,
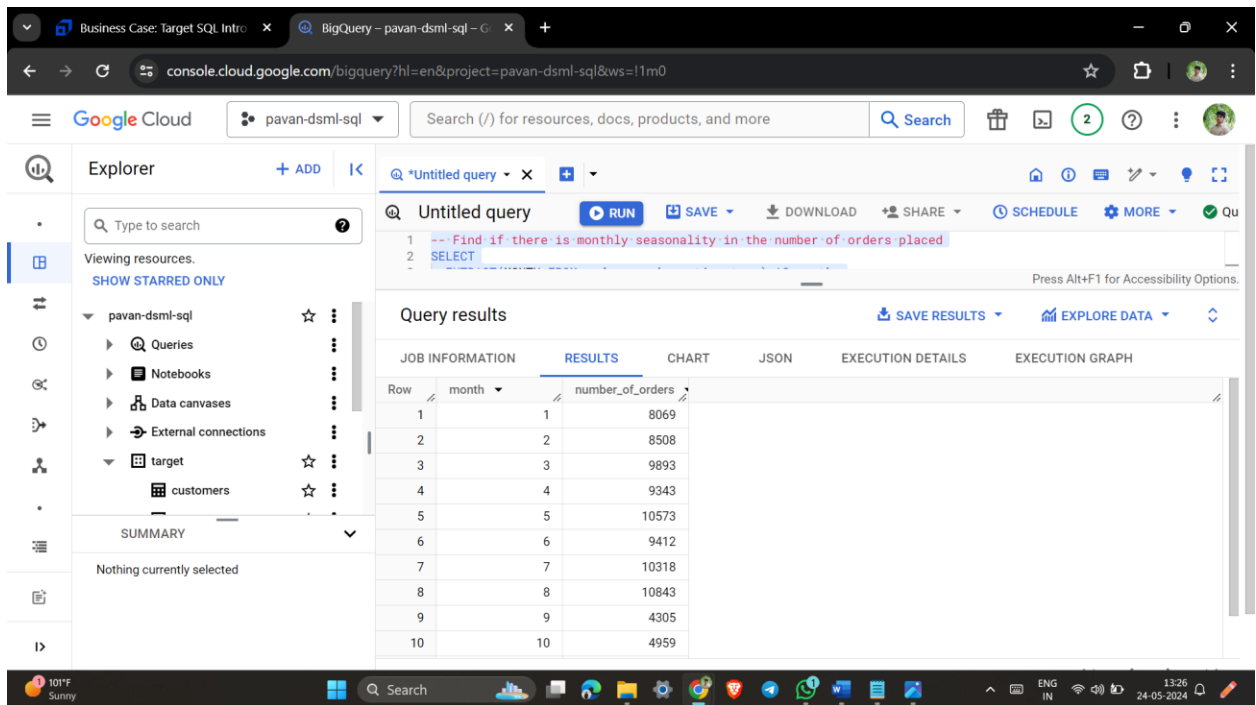  COUNT(order_id) AS number_of_orders
FROM
  `target.orders`
GROUP BY
  month
ORDER BY
  month;

## C. During what time of the day, do the Brazilian customers mostly place their orders? (Dawn, Morning, Afternoon or Night)

- **0-6 hrs : Dawn**
- **7-12 hrs : Mornings**
- **13-18 hrs : Afternoon**
- **19-23 hrs : Night**

**QUERY:**
```
SELECT
  CASE
    WHEN EXTRACT(HOUR FROM order_purchase_timestamp) BETWEEN 0 AND 6 THEN 'Dawn'
    WHEN EXTRACT(HOUR FROM order_purchase_timestamp) BETWEEN 7 AND 12 THEN 'Morning'
    WHEN EXTRACT(HOUR FROM order_purchase_timestamp) BETWEEN 13 AND 18 THEN 'Afternoon'
    WHEN EXTRACT(HOUR FROM order_purchase_timestamp) BETWEEN 19 AND 23 THEN 'Night'
  END AS time_of_day,
  COUNT(order_id) AS number_of_orders
FROM
  `target.orders`
GROUP BY
```

time_of_day
ORDER BY
number_of_orders DESC;



## 3.  Evolution of E-commerce orders in the Brazil region:

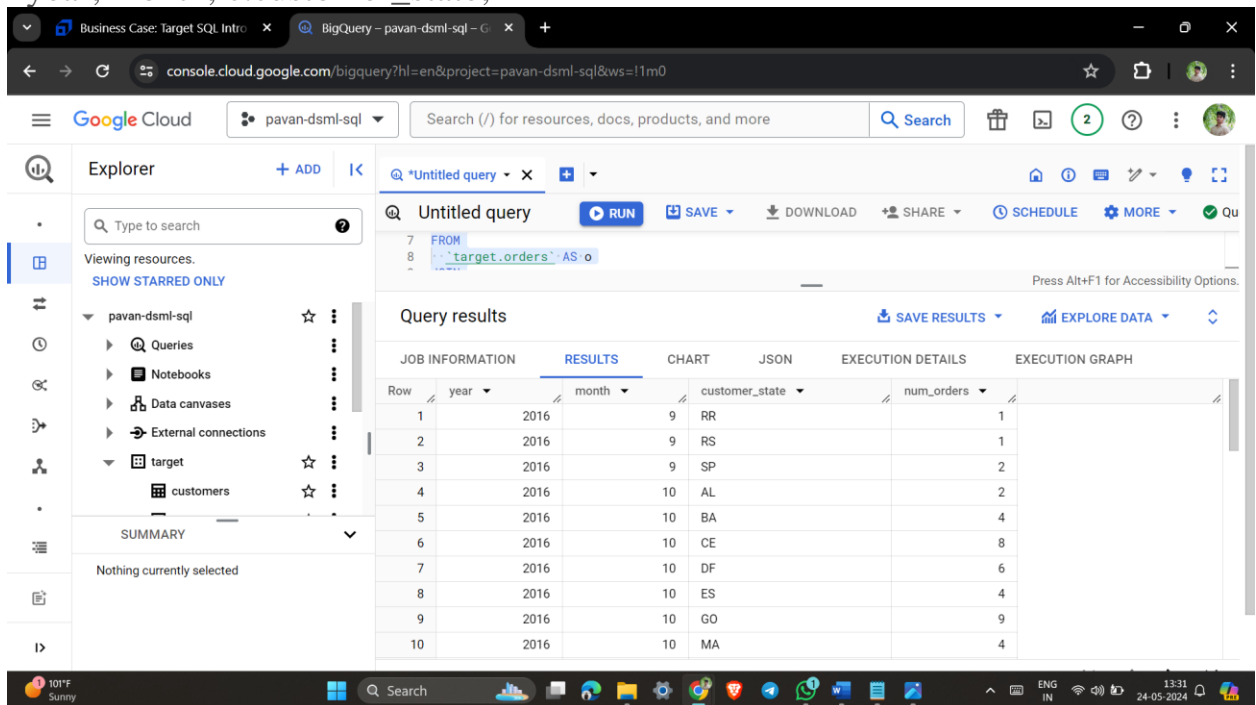**A. Get the month-on-month no. of orders placed in each state.**
**QUERY:**
SELECT
 EXTRACT(YEAR FROM o.order_purchase_timestamp) AS year,
 EXTRACT(MONTH FROM o.order_purchase_timestamp) AS month,
 c.customer_state,
 COUNT(o.order_id) AS num_orders
FROM
 `target.orders` AS o
JOIN
 `target.customers` AS c
ON
 o.customer_id = c.customer_id
GROUP BY
 year, month, c.customer_state
ORDER BY

year, month, c.customer_state;



## B. How are the customers distributed across all the states?
**QUERY:**
```
SELECT
  customer_state,
  COUNT(DISTINCT customer_id) AS num_customers
FROM
  `target.customers`
GROUP BY
  customer_state
ORDER BY
  customer_state;
```

**4. Impact on Economy: Analyze the money movement by e-commerce by looking at order prices, freight and others.**

**A. Get the % increase in the cost of orders from year 2017 to 2018 (include months between Jan to Aug only).**
**You can use the "payment_value" column in the payments table to get the cost of orders.**

**QUERY:**

```
WITH order_costs AS (
 SELECT
  EXTRACT(YEAR FROM o.order_purchase_timestamp) AS year,
  EXTRACT(MONTH FROM o.order_purchase_timestamp) AS month,
  SUM(p.payment_value) AS total_payment_value
 FROM
  `target.orders` AS o
 JOIN
  `target.payments` AS p
 ON
  o.order_id = p.order_id
 WHERE
  EXTRACT(YEAR FROM o.order_purchase_timestamp) IN (2017, 2018)
  AND EXTRACT(MONTH FROM o.order_purchase_timestamp)
BETWEEN 1 AND 8
```

```
  GROUP BY
    year, month
)
SELECT
  (oc_2018.total_payment_value - oc_2017.total_payment_value) /
oc_2017.total_payment_value * 100 AS percent_increase
FROM
  order_costs oc_2017
JOIN
  order_costs oc_2018
ON
  oc_2017.month = oc_2018.month
WHERE
  oc_2017.year = 2017 AND oc_2018.year = 2018;
```



**B. Calculate the Total & Average value of order price for each state.**

**QUERY:**
```
SELECT
  c.customer_state,
  SUM(oi.price) AS total_order_price,
  AVG(oi.price) AS avg_order_price
FROM
  `target.order_items` AS oi
```

JOIN
 `target.orders` AS o
ON
 oi.order_id = o.order_id
JOIN
 `target.customers` AS c
ON
 o.customer_id = c.customer_id
GROUP BY
 c.customer_state
ORDER BY
 c.customer_state;



## C. Calculate the Total & Average value of order freight for each state.

**QUERY:**
SELECT
 c.customer_state,
 SUM(oi.freight_value) AS total_freight_value,
 AVG(oi.freight_value) AS avg_freight_value
FROM
 `target.order_items` AS oi
JOIN
 `target.orders` AS o

```
ON
  oi.order_id = o.order_id
JOIN
  `target.customers` AS c
ON
  o.customer_id = c.customer_id
GROUP BY
  c.customer_state
ORDER BY
  c.customer_state;
```



5. **Analysis based on sales, freight and delivery time.**
A. **Find the no. of days taken to deliver each order from the order's purchase date as delivery time.**
   **Also, calculate the difference (in days) between the estimated & actual delivery date of an order.**
   **Do this in a single query.**

   **You can calculate the delivery time and the difference between the estimated & actual delivery date using the given formula:**
      i. **time_to_deliver = order_delivered_customer_date - order_purchase_timestamp**

ii. **diff_estimated_delivery = order_delivered_customer_date - order_estimated_delivery_date**

**QUERY:**

SELECT
 order_id,
 DATE_DIFF(order_delivered_customer_date, order_purchase_timestamp, DAY) AS delivery_time,
 DATE_DIFF(order_estimated_delivery_date, order_delivered_customer_date, DAY) AS diff_estimated_delivery
FROM
 `target-project.target_brazil.orders`
WHERE
 order_status = 'delivered';



**B. Find out the top 5 states with the highest & lowest average freight value.**

**QUERY:**

with highest_freight as(
select customer_state,avg(freight_value) high_av,
row_number() over(order by avg(freight_value) desc) rhf

from `target.order_items` oi join `target.orders` o on oi.order_id=o.order_id
join `target.customers` c on c.customer_id=o.customer_id
group by customer_state
order by high_av desc
limit 5),
lowest_freight as(
select customer_state,avg(freight_value) low_av,
row_number() over(order by avg(freight_value) asc) rlf
from `target.order_items` oi join `target.orders` o on oi.order_id=o.order_id
join `target.customers` c on c.customer_id=o.customer_id
group by customer_state
order by low_av asc
limit 5)
select hf.customer_state,hf.high_av,lf.customer_state,lf.low_av from
highest_freight hf join lowest_freight lf
on hf.rhf=lf.rlf;



**C. Find out the top 5 states with the highest & lowest average delivery time.**

**QUERY:**

with cte as(select
customer_state,

avg(datetime_diff(order_delivered_customer_date,order_purchase_timesta
mp,day ) ) av_del_time,
from
`target.orders`o join  `target.customers`c on
o.customer_id=c.customer_id
group by customer_state),
cte2 as(
select *,
row_number() over(order by av_del_time desc) row_h,
from cte),
cte3 as(
  select *,
  row_number() over(order by av_del_time asc) row_l
  from cte
)
select
cte2.customer_state,cte2.av_del_time,cte3.customer_state,cte3.av_del_time
from
cte2 join cte3
on cte2.row_h=cte3.row_l
limit 5;

**D. Find out the top 5 states where the order delivery is really fast as compared to the estimated date of delivery.**
**You can use the difference between the averages of actual & estimated delivery date to figure out how fast the delivery was for each state.**

**QUERY:**
WITH fast_delivery_states AS (
 SELECT
  c.customer_state,
  AVG(DATE_DIFF(o.order_delivered_customer_date,
o.order_estimated_delivery_date, DAY)) AS avg_delivery_difference
 FROM
  `target.orders` AS o
 JOIN
  `target.customers` AS c
 ON
  o.customer_id = c.customer_id
 WHERE
  o.order_status = 'delivered'
 GROUP BY
  c.customer_state
)
SELECT
 customer_state,
 avg_delivery_difference
FROM
 fast_delivery_states
ORDER BY
 avg_delivery_difference ASC
LIMIT
 5;

## 6. Analysis based on the payments:

**A. Find the month-on-month no. of orders placed using different payment types.**

**QUERY:**

SELECT
 EXTRACT(YEAR FROM o.order_purchase_timestamp) AS year,
 EXTRACT(MONTH FROM o.order_purchase_timestamp) AS month,
 p.payment_type,
 COUNT(DISTINCT o.order_id) AS num_orders
FROM
 `target.orders` AS o
INNER JOIN
 `target.payments` AS p
ON
 o.order_id = p.order_id
GROUP BY
 year, month, payment_type
ORDER BY
 year, month, payment_type;

## B. Find the no. of orders placed on the basis of the payment installments that have been paid.

**QUERY:**
SELECT
  payment_installments,
  COUNT(DISTINCT order_id) AS num_orders
FROM
  `target.payments`
WHERE
  payment_value > 0
GROUP BY
  payment_installments
ORDER BY
  payment_installments;