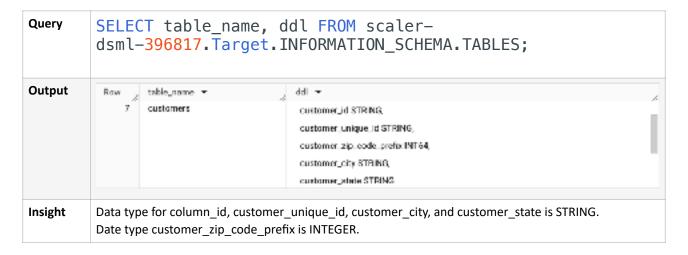
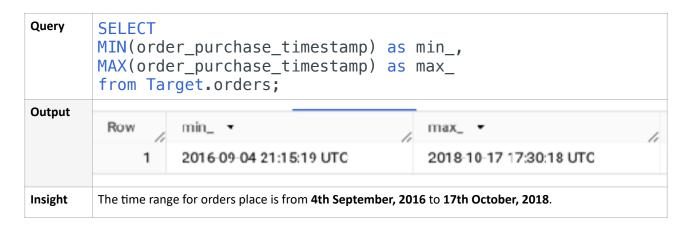
- 1. Import the dataset and do usual exploratory analysis steps like checking the structure & characteristics of the dataset:
 - 1. Data type of all columns in the "customers" table.



2. Get the time range between which the orders were placed.



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



2. In-depth Exploration:

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

Query	SELECT Year, Count FROM (SELECT EXTRA as Year, order_pur FROM Target.orders GROUP BY Year ORDER BY Year;	<pre>CT(YEAR FROM o chase_timestam</pre>	rder_purchase_ti	mestamp)
Output	Row1	Year ▼ 2016 2017 2018	Total_orders ▼ 329 45101 54011	
Insight	Considering that Target commer total order purchased in 2016 (4 In 2017 (12 months) it grew to 4 In 2018 (10 months) it further gr In conclusion, there is growing t	months) was 329. 5,101. rew to 54,011.	•	

2. Can we see some kind of monthly seasonality in terms of the no. of orders being placed?

```
Query
       SELECT * FROM (SELECT t.Month, t.YEAR, Count(*) as
       Total_orders,
       DENSE_RANK() OVER(partition by t.Year Order by Count(*)
       desc) as peak order month
       FROM (SELECT EXTRACT(MONTH FROM order purchase timestamp)
       as Month,
       EXTRACT(YEAR FROM order_purchase_timestamp) as Year,
       order purchase timestamp
       FROM scaler-dsml-396817. Target.orders) t
       GROUP BY t.Month, t.YEAR) m
       where peak order month = 1;
Output
        Row
               Month -
                              YEAR •
                                             Total_orders ▼
                                                           peak_order_month
            1
                         10
                                      2016
                                                      324
                                                                      1
            2
                          1
                                      2018
                                                     7269
                                                                      1
            3
                         11
                                      2017
                                                     7544
                                                                      1
```

After counting the total order for each month it can be observed that peak order month for each year were as follows; 2016 - October 2017 - November 2018 - January Due the limitation of data monthly seasonality cannot be determined at this point.

3. 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

6PM.

Query	SELECT CAS	SE	
		ACT(HOUR FROM order_purchase_ti	•
		OUR FROM order_purchase_timesta	amp) < <mark>7</mark> THEN
	'Dawn'	ACT/USUB FROM	
		ACT(HOUR FROM order_purchase_ti	
	'Morning'	OUR FROM order_purchase_timesta	amp) < 13 THEN
		ACT(HOUR FROM order_purchase_ti	imestamp) >= 13 AND
		OUR FROM order_purchase_timesta	•
	'Afternoo	n'	
	ELSE 'Nigl		
		me_of_day, $ extstyle{COUNT}(*)$ as $order_{COUNT}(*)$	ount
		er-dsml- <mark>396817.Target.orders</mark>	
		time_of_day	
	ORDER BY	order_count DESC;	
Output	Row	time_of_day ▼	order_count ▼
	1	Afternoon	38135
			22224
	2	Night	28331
	3	Morning	27733
	4	Dawn	5242
		Dawii	J2-72

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

1. Get the month on month no. of orders placed in each state.

Query	AS year, EXTRACT(MON c.customer	NTH FROM of state, r-dsml-3967. Target.of mer_id=0.00 ear, month,	T(YEAR FRO order_purc 5817.Targe customers customer_i ,customer_	OM order_puro chase_timesta et.orders o j c id) t _state	chase_timestamp) amp) AS month, join scaler-
		_		,	
Output	Rom	year -	month •	customer_state +	order_count =
Output		_	month •	customer_state =	order_count =
Output	Bow 1	yeer ~	1	customer_state +	2
Output	Row	yeor - 2017 2017	1 2	customer_state = AC AC	2 3
Output	Rom 1 2 3	yeer ~ 2017 2017 2017	1 2 9	AC AC AC	2 3 2
Output	1 2 3 4	yeer ~ 2017 2017 2017 2017 2017	1 2 3 4	AC AC AC	2 3 2 5
Output	Bon 2 1 2 3 4 5	yeer ~ 2017 2017 2017 2017 2017 2017	1 2 3 4	AC AC AC	2 3 2 5 8
Output	Row 2 3 4 5 6	yeer ~ 2017 2017 2017 2017 2017 2017 2017	1 2 3 4 5	AC AC AC AC	2 3 2 5 8
Output	Bow 1 2 3 4 5 6 7	yeer ~ 2017 2017 2017 2017 2017 2017 2017 2017	1 2 3 4 5 6	AC AC AC AC AC	2 3 2 5 8 4

2. How are the customers distributed across all the states?

Query	Total_custofrom scale group by co	omers r-dsml- ustome:	state, COUNT(disting - <mark>396817.Target.custo</mark> r_state ustomers desc;	_
Output		Row	customer_state ▼	Total_customers 💌
		1	93	41746
		2	RJ	12852
		3	MG	11635
		4	RS	5466
		5	PR	5045
		- 6	8G	3637
		7	BA	3380
		8	DF	2140
		9	ES	2033
		10	GD	2020

Insight Heres the no. of unique customers distributed across all the states.

SP has highest no. of unique customers.

- 4. Impact on Economy: Analyze the money movement by e-commerce by looking at order prices, freight and others.
 - 1. 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" solumn in the nayments table to get the sect of

You can use the "payment_value" column in the payments table to get the cost of orders.

```
Query
        WITH order cost AS (
        SELECT
        EXTRACT(YEAR FROM o.order purchase timestamp) AS
        order year,
        EXTRACT(MONTH FROM o.order purchase timestamp) AS
        order month,
        SUM(p.payment value) as Total cost
        FROM scaler-dsml-396817. Target. orders o JOIN scaler-
        dsml-396817. Target. payments p ON p.order id=o.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 EXTRACT(YEAR FROM o.order purchase timestamp),
        EXTRACT(MONTH FROM o.order purchase timestamp)
        ),
        Cost per year as (
        SELECT
        SUM(CASE WHEN order year=2018 THEN Total cost END) as
        CY 2018.
        SUM(CASE WHEN order year=2017 THEN Total cost END) as
        PY 2017
        FROM order cost
        SELECT *, ((CY 2018 - PY 2017) / PY 2017) * 100 AS
        percentage increase
        FROM cost per year
Output
                 CY_2018 -
                                     PY_2017 •
         Row
                                                         percentage_increase
                                     3669022.1200000113
                  8694733.8399999849
                                                            136.97687164665984
Insight
        • I have calculated the total cost of orders for each year from January to August, using the payments
         and orders table.
        • To calculate the percentage increase I have used the formula (Current year - Previous Year/ Previous
         Year)* 100
        • We can conclude that there was 136.97% increase in the cost of orders for the given period.
```

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

```
Query
          SELECT c.customer_state AS State,
          SUM(i.price) as total order price,
          AVG(i.price) as avg_order_price
          FROM scaler-dsml-396817. Target.order_items i JOIN scaler-
          dsml-396817. Target. orders o USING (order id)
          join scaler-dsml-396817. Target. customers c USING
          (customer id)
          GROUP BY c.customer state
          ORDER BY total order price DESC, avg order price DESC;
Output
               Row
                        State v
                                                     total_order_price <.
                                                                       avg_order_price *
                   1
                        SP
                                                     5202955.050002...
                                                                        109.6536291597...
                   2
                        RJ
                                                     1824092.669999...
                                                                       125.1178180945...
                   3
                        MG
                                                     1585308.029999...
                                                                       120.7485741488...
                   4
                        RS
                                                     750304.0200000...
                                                                       120.3374530874...
                        PR
                                                     683083.7600000...
                                                                       119.0041393728...
                   5
                   6
                        SC
                                                     520553.3400000...
                                                                       124.6535775862...
                   7
                        BA
                                                     511349.9900000...
                                                                       134.6012082126...
                   8
                        DF
                                                     302603.9399999...
                                                                       125.7705486284...
                   9
                        GO
                                                     294591.9499999...
                                                                       126.2717316759...
                   10
                        ES
                                                     275037.3099999...
                                                                       121.9137012411...
Insight
          • In order to extract the state and order price information we have joined the order item table to
           customers table.
          • I have used Group by function to club the records as per states
          • SP has the highest total order price i.e. 5202955.05 with average order price of 109.65
```

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

```
Query
SELECT c.customer_state AS State,
SUM(i.freight_value) as total_order_freight,
AVG(i.freight_value) as avg_order_freight
FROM scaler-dsml-396817.Target.order_items i JOIN scaler-
dsml-396817.Target.orders o USING (order_id)
join scaler-dsml-396817.Target.customers c USING
(customer_id)
GROUP BY c.customer_state
ORDER BY total_order_freight DESC, avg_order_freight DESC;
```

Output	Row	State ▼	total_order_freight _	avg_order_freight 🛪
	1	SP	718723.06999999	15.14727539041
	2	RJ	305589.3100000	20.96092393168
	3	MG	270853.4600000	20.63016680630
	4	RS	135522.7400000	21.73580433039
	5	PR	117851.6800000	20.53165156794
	6	BA	100156.6799999	25.36395893656
	7	SC	89660.26000000	21.47036877394
	В	PE	59449.65999999	32.91786267995
	9	GO	53114.97999999	22.76681525932
	10	DF	50625.499999999	21.04135494596

5. Analysis based on sales, freight and delivery time.

1. 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:

- time_to_deliver = order_delivered_customer_date order_purchase_timestamp
- diff_estimated_delivery = order_estimated_delivery_date order delivered customer date

```
Query

SELECT
order_id,
order_purchase_timestamp,
order_delivered_customer_date,
order_estimated_delivery_date,
DATE_DIFF(order_delivered_customer_date,
order_purchase_timestamp, DAY) AS time_to_deliver,
DATE_DIFF(order_estimated_delivery_date,
order_delivered_customer_date, DAY) AS
diff_estimated_delivery
FROM scaler-dsml-396817.Target.orders;
```

Output							
	Row	order_id -	order_purchase_timestamp	order_delivered_custory	order_estimated_delig	time_to_deliver + _	diff_estimated_delig
	1	1950d77798	2018-02-19 19:48-52 U	2018-03-21 22:03:5.	2018-03-09 00:00.	30	-12
	2	2c45c33d2f9	2014-10-09 15:29:56 U	2016-11-09 14:53:5	2016-12-08 00:00	35	26
	3	65d1e226dfa.	2016/10/03 21:01:41 U.	2016/11/08 10:58.3.	2016-11-25 00:00.	35	16
	4	625c894d06	2017-04-15 15:27:38 U	2017-05-16 14:49:5	2017-05-18 00:00	30	1
	5	3597562c2a	2017-04-14 22:21:54 U	2017-05-17 10:52-1	2017-05-18 00:00	32	D
	6	68f47f90f04e	2017-04-16 14:56:18 U.	2017/05/16 09:07:4.	2017-05-18 00:00	29	1
	7	276e9ec344d	2017-04-00 21:20:24 U	2017-05-22 14:11:3	2017-05-10 00:00.	43	-4
	0	54e1a3d7b9	2017-04-11 19:49:45 U	2017-05-22 16:10:4	2017-05-10 00:00	40	-4
	9	fd04fa4105c	2017-04-12 12:17:08 U.	2017-05-19 18:44 5.	2017-05-18 00:00	97	-1
	10	302558109d	2017-04-19 22:52:59 U	2017-05-23 14:19:4	2017-05-18 00:00	83	-5
	10	302568109d_	2017-04-19 22:52 99 U.	2017-05-23 14:19.4.	2017-05-18-00:00.	93	
Insight	timediff_daysPosit	_to_deliver c estimated_de ive integer de	diff function to calcu olumn represents to eliver represents the enotes that the order denotes that delivery	tal no. of days it to difference between was fulfilled befo	ook to fulfil an or en the estimated	der	ctual delivery

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

```
Query 1
         SELECT c.customer_state AS State,
         AVG(i.freight_value) as avg_order_freight
         FROM scaler-dsml-396817. Target.order_items i JOIN scaler-
         dsml-396817.Target.orders o USING (order_id)
         join scaler-dsml-396817.Target.customers c USING
         (customer_id)
         GROUP BY c.customer_state
         ORDER BY avg_order_freight DESC
         LIMIT 5
Output 1
                            State -
                                                           avg_order_freight
                   Row
                       1
                                                           42.98442307692...
                            RR
                       2
                            PB
                                                           42.72380398671...
                       3
                            RO
                                                           41.06971223021...
                       4
                                                           40.07336956521...
                            AC
                            ы
                                                           39.14797047970...
Insight
         • These are the top 5 states with highest average freight value
         • I have joined the customers table to order_item to get freight value and group them as per states.
         • Limit function gets the top 5 states with highest freight values.
```

```
Query 2
        SELECT c.customer state AS State,
        AVG(i.freight value) as avg order freight
        FROM scaler-dsml-396817. Target order items i JOIN scaler-
        dsml-396817. Target. orders o USING (order id)
        join scaler-dsml-396817. Target. customers c USING
        (customer id)
        GROUP BY c.customer state
        ORDER BY avg_order_freight
        LIMIT 5
Output 2
                   Row
                            State -
                                                        avg_order_freight 🤟
                       1
                            SP
                                                        15.14727539041...
                       2
                           PR
                                                        20.53165156794...
                       3
                           MG
                                                        20.63016680630.
                       4
                           RJ
                                                        20.96092393168...
                            DF
                                                        21.04135494596...
                       5
Insight
        I have ordered the average freight value in ascending order to get top 5 states with lowest freight
```

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

```
Query 1
       WITH deliverytime AS (
       SELECT
       c.customer_state,
       DATE_DIFF(order_delivered_customer_date,
       order purchase timestamp, DAY) AS time to deliver,
       AVG(DATE_DIFF(order_delivered_customer_date,
       order_purchase_timestamp, DAY)) OVER(PARTITION BY
       customer state) AS avg delivery days
       FROM scaler-dsml-396817. Target. orders o JOIN scaler-
       dsml-396817. Target. customers c
       ON o.customer id=c.customer id
       SELECT customer_state,
       ROUND(avg_delivery_days,0) AS avg_delivery_days
       FROM deliverytime
       GROUP BY customer_state,avg_delivery_days
       ORDER BY avg delivery days;
```

Output 1	Row	customer_state ▼	avg_delivery_days >
	1	SP	8.0
	2	PR	12.0
	3	MG	12.0
	4	DF	13.0
	5	SC	14.0
Insight Query 2	Delivery time is calculaI have used Round fun	ction to get the no. of days as a omers table to know the orders	absolute values belong to which state and grouped then
	order_purchase AVG(DATE_DIFF(order_purchase customer_state FROM scaler-ds dsml-396817.Ta ON o.customer_) SELECT custome ROUND(avg_deli FROM deliveryt GROUP BY customer	er_delivered_custo e_timestamp, DAY) (order_delivered_c e_timestamp, DAY)) e) AS avg_delivery sml-396817.Target. arget.customers c _id=c.customer_id er_state, ivery_days,0) AS a	AS time_to_deliver, ustomer_date, OVER(PARTITION BY _days orders o JOIN scaler- vg_delivery_days ivery_days
	,		
Output 2	Row	customer_state ▼	avg_delivery_days
Output 2		customer_state ▼	avg_delivery_days 29.0
Output 2	Row	3	0
Output 2	Row 2	RR	29.0
Output 2	Row 2	RR AP	29.0 27.0

4. 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 state delivery speed as
        (
        SELECT
        c.customer state as state,
        AVG(DATE_DIFF(o.order_delivered_customer_date,order_purcha
        se timestamp,day)) OVER(PARTITION BY c.customer state) as
        avg_delivery_time,
        AVG(DATE DIFF(o.order estimated delivery date, order purcha
        se timestamp,day)) OVER(PARTITION BY c.customer state) as
        avg estimated time
        FROM scaler-dsml-396817. Target.orders o JOIN scaler-
        dsml-396817. Target.customers c
        ON o.customer id=c.customer id
        where order_delivered_customer_date is not null
        SELECT
        state,
        AVG(avg_delivery_time-avg_estimated_time) as
        delivery speed
        FROM state_delivery_speed
        GROUP BY state,(avg_delivery_time- avg_estimated_time)
        ORDER BY delivery speed
        LIMIT 5:
Output
         Row
                                                        delivery_speed ~
                    state 🔻
               1
                    AC
                                                        -20.0875000000...
               2
                    RO
                                                        -19.4732510288...
               3
                    AP
                                                        -19.1343283582...
               4
                    AM
                                                        -18.9379310344...
               5
                    RR
                                                        -16.6585365853...
        • We calculate the average delivery speed for each state by finding the difference in days
Insight
         between the actual delivery date (order delivered customer date) and the estimated
         delivery date (order_estimated_delivery_date) using the DATEDIFF function.
        • We group the result by customer state
        • To get the top 5 state we order the result in ascending order and set Limit to 5
```

6. Analysis based on the payments:

1. 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(*) AS order count
         FROM
         scaler-dsml-396817. Target.orders o
         JOIN scaler-dsml-396817. Target payments p ON o.order id =
         p.order id
         GROUP BY year, month, payment_type
         ORDER BY year, month, payment_type;
Output
                                  month v
          Row
                  Wear 3
                                                 payment_type *
                                                                           order_count •
              1
                           2016
                                                 credit_card
                                                                                       3
                                                 UP1
              2
                           2016
                                             10
                                                                                      63
                           2016
                                             10
                                                 credit_card
                                                                                     254
              3
              4
                           2016
                                             10
                                                 debit_card
                                                                                       2
                                                 voucher
                                                                                      23
              5
                           2016
                                             10
                           2016
                                             12
                                                 credit_card
              6
                                                                                       1
                                                                                     197
              7
                                                 UPI
                           2017
                                             1
                           2017
                                                 credit_card
                                                                                     583
              В
                                             1
                                                 debit_card
              9
                           2017
                                             1
                                                                                       9
             10
                           2017
                                                 voucher
                                             1
                                                                                      61
Insight
         • We extract the year and month from the order_purchase_timestamp column using the
           EXTRACT function.
         • We include the payment type column to group the results by payment type.
         • We count the number of orders for each combination of year, month, and payment type
           using COUNT(*).
         • We group the results by year, month, and payment type to get the month-on-month order
           counts for each payment type.
```

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

```
Query SELECT count(distinct order_id) as order_count
FROM scaler-dsml-396817.Target.payments
WHERE payment_installments != 0;
```

e Jakovi •		Row ord	er_count ▼	
1 1970		1	99438	
Insight	payment_installme	ents are greater than 0.	considering only orders where the distinct order IDs, which elimin	ates