## **TARGET SQL**

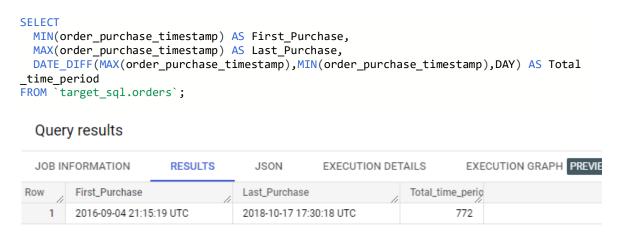
- 1. Import the dataset and do usual exploratory analysis steps like checking the structure & characteristics of the dataset.
  - 1. Data type of columns in a table

```
SELECT ddl,*
FROM target_sql.INFORMATION_SCHEMA.TABLES
WHERE table_name = 'orders';
```



In ddl column, we can get the data types of the column for the given table name.

2. Time period for which the data is given.



3. Cities and States of customers ordered during the given period



## 2. In-depth Exploration:

1. Is there a growing trend on e-commerce in Brazil? How can we describe a complete scenario? Can we see some seasonality with peaks at specific months?

```
SELECT * FROM
(SELECT
    EXTRACT(YEAR FROM order_purchase_timestamp) AS Year,
    EXTRACT(MONTH FROM order_purchase_timestamp) AS Month,
    COUNT(*) AS Month_wise_sales
FROM `target_sql.orders`
GROUP BY EXTRACT(YEAR FROM order_purchase_timestamp), EXTRACT(MONTH FROM order_purchase_timestamp)) x
ORDER BY x.Year,x.Month;
```

Quer	y results				
JOB IN	NFORMATION	RESULTS	JSON EXECUT	ION DETAILS	EXECUTION GRAPH PREVIEW
Row	Year	Month	Month_wise_sales		
1	2016	9	4		
2	2016	10	324		
3	2016	12	1		
4	2017	1	800		
5	2017	2	1780		
6	2017	3	2682		
7	2017	4	2404		
8	2017	5	3700		
9	2017	6	3245		
10	2017	7	4026		
	2247	_	1001		

#### 2. What time do Brazilian customers tend to buy (Dawn, Morning, Afternoon or Night)?

```
SELECT * FROM
(SELECT
CASE
 WHEN EXTRACT(HOUR FROM order_purchase_timestamp) BETWEEN 0 AND 6 THEN 'Dawn (0-6)'
 WHEN EXTRACT(HOUR FROM order_purchase_timestamp) BETWEEN 7 AND 12 THEN 'Morning (7-
12)'
 WHEN EXTRACT(HOUR FROM order_purchase_timestamp) BETWEEN 13 AND 16 THEN 'Afternoon (1
3-16)'
 ELSE 'Night (17-24)'
END AS Time_of_day,
COUNT(*) AS Total_sales
FROM `target_sql.orders`
GROUP BY
CASE
 WHEN EXTRACT(HOUR FROM order_purchase_timestamp) BETWEEN 0 AND 6 THEN 'Dawn (0-6)'
 WHEN EXTRACT(HOUR FROM order_purchase_timestamp) BETWEEN 7 AND 12 THEN 'Morning (7-
12)'
 WHEN EXTRACT(HOUR FROM order_purchase_timestamp) BETWEEN 13 AND 16 THEN 'Afternoon (1
3-16)'
 ELSE 'Night (17-24)'
END) x
ORDER BY x.Total_sales DESC;
```

## Query results

JOB IN	IFORMATION	RESULTS	JSON	EXECUTION DE
Row	Time_of_day	//	Total_sales	
1	Night (17-24)		40250	
2	Morning (7-12)		27733	
3	Afternoon (13-16)	)	26216	
4 Dawn (0-6)			5242	
	Row 1 2 3	1 Night (17-24) 2 Morning (7-12) 3 Afternoon (13-16)	Row Time_of_day  1 Night (17-24)  2 Morning (7-12)  3 Afternoon (13-16)	Row         Time_of_day         Total_sales           1         Night (17-24)         40250           2         Morning (7-12)         27733           3         Afternoon (13-16)         26216

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

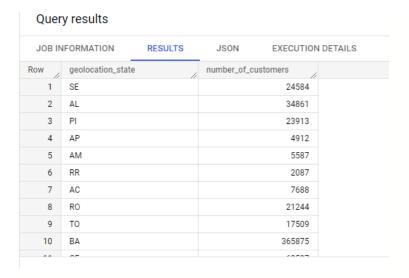
1. Get month on month orders by states

```
SELECT * FROM
(SELECT
  geolocation_state,
  EXTRACT(MONTH FROM order_purchase_timestamp) AS Month,
  COUNT(o.order_id) AS number_of_orders
FROM `target_sql.geolocation` g LEFT JOIN `target_sql.customers` c ON geolocation_zip_c
ode_prefix=customer_zip_code_prefix
JOIN `target_sql.orders` o ON o.customer_id=c.customer_id
GROUP BY geolocation_state,EXTRACT(MONTH FROM order_purchase_timestamp)) x
ORDER BY x.geolocation_state,x.Month;
```

JOB IN	IFORMATION	RESULTS	JSON	EXECUTION DETAILS	EXECUT
Row	geolocation_state	//	Month	number_of_orders	
1	AC		1	694	
2	AC		2	515	
3	AC		3	516	
4	AC		4	789	
5	AC		5	1161	
6	AC		6	563	
7	AC		7	937	
8	AC		8	1060	
9	AC		9	161	
10	AC		10	535	

#### 2. Distribution of customers across the states in Brazil

```
SELECT
   geolocation_state,
   COUNT(customer_id) AS number_of_customers
FROM `target_sql.geolocation` g JOIN `target_sql.customers` c ON geolocation_zip_code_p
refix=customer_zip_code_prefix
GROUP BY geolocation_state;
```



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

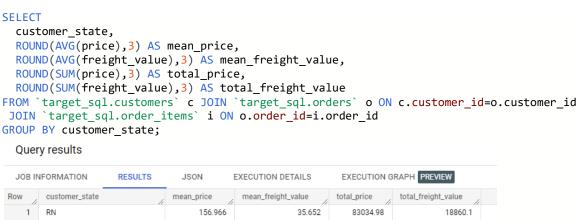
Get % increase in cost of orders from 2017 to 2018 (include months between Jan to Aug only) You can use "payment\_value" column in payments table

```
SELECT
 Year,
 Month,
 total_cost_per_month,
 IF(LAG(total_cost_per_month) OVER (ORDER BY Month ASC) = 0, 0 (total_cost_per_month)
LAG (total cost per month) OVER (ORDER BY Month ASC))/LAG (total cost per month) OVER (ORDER BY M
onth ASC)*100) AS perc_inc_in_cost
FROM
(SELECT
 EXTRACT(YEAR FROM order_purchase_timestamp) AS Year,
 EXTRACT(MONTH FROM order purchase timestamp) AS Month,
 SUM(payment_value) AS total_cost_per_month,
FROM `target_sql.payments` p JOIN `target_sql.orders` o ON p.order_id=o.order_id
GROUP BY EXTRACT(YEAR FROM order_purchase_timestamp), EXTRACT(MONTH FROM order_purchase_timestamp))
WHERE x.Year BETWEEN 2017 AND 2018 AND x.Month BETWEEN 1 AND 8
order by Year, Month;
```

#### Query results

JOB IN	FORMATION	RESULTS	JSON EXECUTION	DETAILS EXECUTION GRAPH
Row	Year	Month	total_cost_per_month	perc_inc_in_cost
1	2017	1	138488.03999999989	-87.5795945446591
2	2017	2	291908.00999999966	-70.5875271926922
3	2017	3	449863.60000000027	-61.207021291868
4	2017	4	417788.03000000032	-64.008161955988683
5	2017	5	592918.82000000111	-48.619758113243
6	2017	6	511276.38000000152	-13.769581474914128
7	2017	7	592382.92000000284	-42.143353643320104
8	2017	8	674396.32000000309	-34.039552150370795
9	2018	1	1115004.1800000065	null
10	2018	2	992463.34000000334	616.6419136266237

#### 2. Mean & Sum of price and freight value by customer state



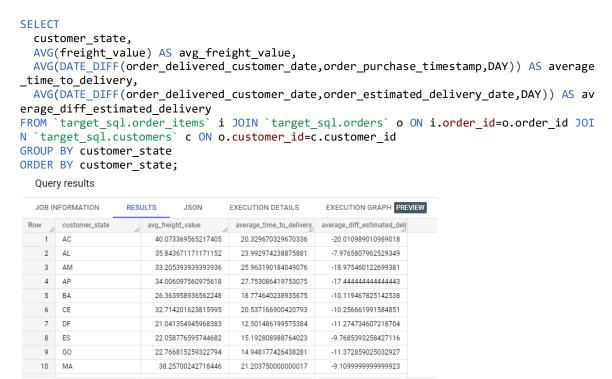
JOB IN	FORMATION	RESULTS	JSON	EXECUTION DETAILS	EXECUTION G	RAPH PREVIEW
Row	customer_state	//	mean_price	mean_freight_value	total_price	total_freight_value
1	RN		156.966	35.652	83034.98	18860.1
2	CE		153.758	32.714	227254.71	48351.59
3	RS		120.337	21.736	750304.02	135522.74
4	SC		124.654	21.47	520553.34	89660.26
5	SP		109.654	15.147	5202955.05	718723.07
6	MG		120.749	20.63	1585308.03	270853.46
7	BA		134.601	26.364	511349.99	100156.68
8	RJ		125.118	20.961	1824092.67	305589.31
9	GO		126.272	22.767	294591.95	53114.98
10	MA		145.204	38.257	119648.22	31523.77

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

Calculate days between purchasing, delivering and estimated delivery
 Find time\_to\_delivery&diff\_estimated\_delivery. Formula for the same given below:
 time\_to\_delivery = order\_purchase\_timestamp-order\_delivered\_customer\_date
 diff\_estimated\_delivery = order\_estimated\_delivery\_date-order\_delivered\_customer\_date

```
SELECT
  order_id,
  DATE_DIFF(order_delivered_customer_date,order_purchase_timestamp,DAY) AS time_to_deli
  DATE_DIFF(order_delivered_customer_date,order_estimated_delivery_date,DAY) AS diff_es
timated_delivery
FROM `target_sql.orders`;
   Query results
   JOB INFORMATION
                         RESULTS
                                       JSON
                                                   EXECUTION DETAILS
                                                                          EXECUTION GRAPH F
                                     time_to_delivery
                                                       diff_estimated_delivery
 Row
         order id
     1
         1950d777989f6a877539f5379...
                                                  30
                                                                         12
     2
         2c45c33d2f9cb8ff8b1c86cc28...
                                                  30
                                                                        -28
     3
         65d1e226dfaeb8cdc42f66542...
                                                  35
                                                                        -16
     4
         635c894d068ac37e6e03dc54e...
                                                  30
                                                                         -1
     5
         3b97562c3aee8bdedcb5c2e45...
                                                  32
                                                                         0
     6
         68f47f50f04c4cb6774570cfde...
                                                  29
                                                                         -1
     7
         276e9ec344d3bf029ff83a161c...
                                                  43
                                                                          4
     8
         54e1a3c2b97fb0809da548a59...
                                                  40
                                                                         4
         fd04fa4105ee8045f6a0139ca5...
     Q
                                                  37
                                                                          1
    10
         302bb8109d097a9fc6e9cefc5...
                                                  33
                                                                          5
```

#### 2. Group data by state, take mean of freight value, time to delivery, diff estimated delivery



- 3. Sort the data to get the following:
- 4. Top 5 states with highest/lowest average freight value sort in desc/asc limit 5

```
SELECT * FROM
(SELECT
    customer_state,
    AVG(freight_value) AS avg_freight_value
FROM `target_sql.order_items` i JOIN `target_sql.orders` o ON i.order_id=o.order_id JOI
N `target_sql.customers` c ON o.customer_id=c.customer_id
GROUP BY customer_state) x
ORDER BY x.avg_freight_value DESC LIMIT 5;
```

JOB INFORMATION		RESULTS	JSON	EXECUTION	ON DETAIL
Row	customer_state	//	avg_freight_val	lue //	
1	RR		42.984423	076923093	
2	PB		42.723803	986710941	
3	RO		41.069712	230215842	
4	AC		40.073369	565217405	
5	PI		39.1479704	479704767	

#### 5. Top 5 states with highest/lowest average time to delivery

```
SELECT * FROM
(SELECT
    customer_state,
    AVG(DATE_DIFF(order_delivered_customer_date,order_purchase_timestamp,DAY)) AS average
    _time_to_delivery
FROM `target_sql.order_items` i JOIN `target_sql.orders` o ON i.order_id=o.order_id JOI
N `target_sql.customers` c ON o.customer_id=c.customer_id
GROUP BY customer_state ) x
ORDER BY x.average_time_to_delivery DESC LIMIT 5;
```

#### Query results

JOB IN	IFORMATION	RESULTS	JSON	EXECUTION DETAILS
Row	customer_state	//	average_time	_to_delivery
1	RR		27.826	6086956521738
2	AP		27.753	3086419753075
3	AM		25.963	3190184049076
4	AL		23.992	2974238875881
5	PA		23.301	1707779886126

#### 6. Top 5 states where delivery is really fast/ not so fast compared to estimated date

```
SELECT * FROM
(SELECT
   customer_state,
   AVG(DATE_DIFF(order_delivered_customer_date,order_estimated_delivery_date,DAY)) AS av
erage_diff_estimated_delivery
FROM `target_sql.order_items` i JOIN `target_sql.orders` o ON i.order_id=o.order_id JOI
N `target_sql.customers` c ON o.customer_id=c.customer_id
GROUP BY customer_state ) x
ORDER BY x.average_diff_estimated_delivery DESC LIMIT 5;
```

Quer	y results			
JOB INFORMATION		FORMATION RESULTS		EXECUTION DETAI
Row	customer_state	Į,	average_diff_	estimated_deliy
1	AL		-7.976580	7962529349
2	MA		-9.109999	9999999923
3	SE		-9.165333	3333333276
4	ES		-9.768539	3258427116
5	BA		-10.11946	7825142538

## 6. Payment type analysis:

1. Month over Month count of orders for different payment types

```
SELECT * FROM
(SELECT
    EXTRACT(YEAR FROM order_purchase_timestamp) AS Year,
    EXTRACT(MONTH FROM order_purchase_timestamp) AS Month,
    payment_type,
    COUNT(p.order_id) AS count_of_sales_per_payment
FROM `target_sql.payments` p JOIN `target_sql.orders` o ON p.order_id=o.order_id
GROUP BY EXTRACT(YEAR FROM order_purchase_timestamp), EXTRACT(MONTH FROM order_purchase
    timestamp),payment_type)
ORDER BY Year,Month,payment_type;
```

JOB IN	IFORMATION	RESULTS	JSON	EXECUTION DET	AILS EXE	CUTION
Row	Year	Month	payment_type	//	count_of_sales_	
1	2016	9	credit_card		3	
2	2016	10	UPI		63	
3	2016	10	credit_card		254	
4	2016	10	debit_card		2	
5	2016	10	voucher		23	
6	2016	12	credit_card		1	
7	2017	1	UPI		197	
8	2017	1	credit_card		583	
9	2017	1	debit_card		9	
10	2017	1	voucher		61	

## 2. Count of orders based on the no. of payment installments

```
SELECT
  payment_installments,
  COUNT(order_id) AS count_of_orders
FROM `target_sql.payments`
GROUP BY payment_installments;
```

Quer	y results			
JOB IN	IFORMATION RE	SULTS	JSON	EXECUTION DETAILS
Row	payment_installments/	count_of_	orders	
1	0		2	
2	1		52546	
3	2		12413	
4	3		10461	
5	4		7098	
6	5		5239	
7	6		3920	
8	7		1626	
9	8		4268	
10	9		644	

## **Additional Insights:**

1. Survey based on customer reviews – Product Category with lowest rating

```
SELECT product_category,COUNT(*) AS order_count
FROM target_sql.products WHERE
product_id IN(SELECT product_id FROM target_sql.order_items WHERE
order_id IN (SELECT order_id FROM target_sql.reviews WHERE review_score =1))
GROUP BY product_category
ORDER BY 2 DESC;
```

#### Query results

JOB IN	IFORMATION RESULTS	JSON	EXECUTION
Row	product_category	order_count	
1	bed table bath	852	
2	Furniture Decoration	659	
3	sport leisure	603	
4	HEALTH BEAUTY	510	
5	housewares	470	
6	computer accessories	465	
7	automotive	352	
8	telephony	327	
9	Watches present	308	
10	toys	283	

## **Actionable Insights:**

- i. Data type of a particular column of any table can be found using the data present in the 'ddl' column fetched.
- ii. The 'Target' dataset contains order-data belonging to the below mentioned period of time: 2016-09-04 To 2018-10-17, which is a total of 772 days.
- iii. This Dataset provides information of customer orders from 27 states and 4119 cities across Brazil.
- iv. Since the dataset doesn't contain complete data for the year 2016 and 2018, this 'seasonality' analysis is based on the year 2017. We can see that the no.of.orders gradually increasefrom March to October, reaching its peak during the month of November with 7.54k orders (Christmas festival), then it

- reaches another peak in the month of January with 7.26k orders (New Year purchases).
- v. Brazilians tend to purchase the most during the night (17hrs 24hrs) and the least during dawn (0hrs 6 hrs).
- vi. From the fetched result, we can infer that the state 'SP' has the highest no.of.orders (660764) in the month of August and the state 'AM' has the lowest no.of.orders (49) in the month of October.
- vii. The state 'RR' has the lowest number of customers (2087) and the state 'SP' has the highest number of customers (5620430) among all the 27 states present in the dataset.
- viii. Compared to 2017 there is a linear percentage increase in the cost of orders per month in 2018.
  - ix. Top 5 states with highest average freight value RR>PB>RO>AC>PI And Top 5 states with lowest average freight value SP<PR<MG<RJ<DF
  - x. Top 5 states with highest average time to delivery RR>AP>AM>AL>PA And Top 5 states with lowest average time to delivery SP<PR<MG<DF<SC
- xi. Top 5 states where delivery is really fast compared to estimated date AC,RO,AM,AP,RR And Top 5 states where delivery is not so fast compared to estimated date AL,MA,SE,ES,BA
- xii. The most preferred payment type among customers are credit card and UPI.
- xiii. The most preferred EMI installment period among customers are 1, 2 and 3.

#### Recommendations:

I. We can see there is a peak in number of orders during the month November, December and January, we can continue that streak by providing vouchers/coupons to the customers which should be claimed within the next month, making them a regular.

- II. In the months of May and June, there is a considerable decrease in the number of orders. This situation can be rectified by initiating some clearance sales and providing more discounts, so that no.of.orders will increase.
- III. From the extracted data, we can see that the deliveries takes longer than the estimated time. So, if we improve the time taken to delivery and match the estimated time provided to the customer during the order process, we can improve the overall customer experience and thus it will result in increased orders.
- IV. Since 1, 2 and 3 months installment is popular among the customers, we can provide some offers and cashback during EMI process, so that the customer would opt the other installment periods also.
- V. Debit card and vouchers seems to be the opted by very less customers compared to credit card. We can offer some bank offers to encourage customers to opt for the above payment methods also.
- VI. The following product categories are given low rating by customers: bed table bath, Furniture Decoration, sport leisure. By improving the product quality we can ensure customer satisfaction and improve sales.