

# SQL Business Case

## Context

Target is one of the world's most recognized brands and one of America's leading retailers. Target makes itself a preferred shopping destination by offering outstanding value, inspiration, innovation and an exceptional guest experience that no other retailer can deliver.

This business case has information of 100k orders from 2016 to 2018 made at Target in Brazil. Its features allows viewing an order from multiple dimensions: from order status, price, payment and freight performance to customer location, product attributes and finally reviews written by customers.

## Data Dictionary :

Data is available in 8 csv files:

1. customers.csv
2. geolocation.csv
3. order\_items.csv
4. payments.csv
5. reviews.csv
6. orders.csv
7. products.csv
8. sellers.csv

Each feature or columns of different CSV files are described below:

The **customers.csv** contain following features:

Features	Description
customer_id	Id of the consumer who made the purchase.
customer_unique_id	Unique Id of the consumer.
customer_zip_code_prefix	Zip Code of the location of the consumer.

customer_city	Name of the City from where order is made.
customer_state	State Code from where order is made(Ex- sao paulo-SP).

The **sellers.csv** contains following features:

Features	Description
seller_id	Unique Id of the seller registered
seller_zip_code_prefix	Zip Code of the location of the seller.
seller_city	Name of the City of the seller.
seller_state	State Code (Ex- sao paulo-SP)

The **order\_items.csv** contain following features:

Features	Description
order_id	A unique id of order made by the consumers.
order_item_id	A Unique id given to each item ordered in the order.
product_id	A unique id given to each product available on the site.
seller_id	Unique Id of the seller registered in Target.
shipping_limit_date	The date before which shipping of the ordered product must be completed.
price	Actual price of the products ordered .
freight_value	Price rate at which a product is delivered from one point to another.

The **payments.csv** contain following features:

Features	Description
order_id	A unique id of order made by the consumers.
payment_sequential	sequences of the payments made in case of EMI.

payment_type	mode of payment used.(Ex-Credit Card)
payment_installments	number of installments in case of EMI purchase.
payment_value	Total amount paid for the purchase order.

The **orders.csv** contain following features:

Features	Description
order_id	A unique id of order made by the consumers.
customer_id	Id of the consumer who made the purchase.
order_status	status of the order made i.e delivered, shipped etc.
order_purchase_timestamp	Timestamp of the purchase.
order_delivered_carrier_date	delivery date at which carrier made the delivery.
order_delivered_customer_date	date at which customer got the product.
order_estimated_delivery_date	estimated delivery date of the products.

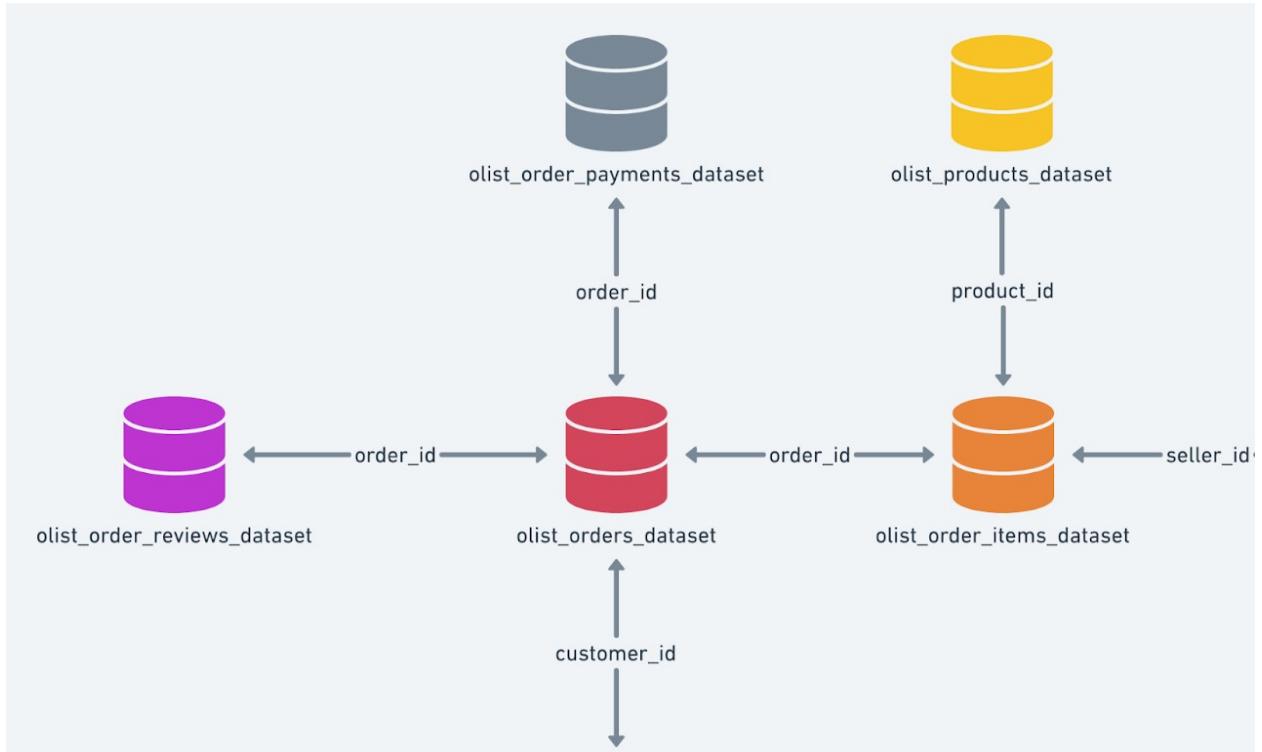
The **reviews.csv** contain following features:

Features	Description
review_id	Id of the review given on the product ordered by the order id.
order_id	A unique id of order made by the consumers.
review_score	review score given by the customer for each order on the scale of 1–5.
review_comment_title	Title of the review
review_comment_message	Review comments posted by the consumer for each order.
review_creation_date	Timestamp of the review when it is created.
review_answer_timestamp	Timestamp of the review answered.

The **products.csv** contain following features:

Features	Description
product_id	A unique identifier for the proposed project.
product_category_name	Name of the product category
product_name_length	length of the string which specifies the name given to the products ordered.
product_description_length	length of the description written for each product ordered on the site.
product_photos_qty	Number of photos of each product ordered available on the shopping portal.
product_weight_g	Weight of the products ordered in grams.
product_length_cm	Length of the products ordered in centimeters.
product_height_cm	Height of the products ordered in centimeters.
product_width_cm	width of the product ordered in centimeters.

High level overview of relationship between datasets:



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

**a. Get number of rows in the data**

**1. For Customers Table :**

```
1 select count(*) from `scaler-dsml-target-sql.target11.custc`
```

## Query results

JOB INFORMATION	RESULTS	JSON	EXECUTION DETA
Row 1 of 1			

2. For GeoLocation Table:

```
1 select count(*) from `scaler-dsml-target-sql.target11`
```

## Query results

JOB INFORMATION	RESULTS	JSON	EXECUTION
Row	f0		

3. For order\_items table:

```
1 select count(*) from `scaler-dsml-target-sql.target1`
```

## Query results

JOB INFORMATION	RESULTS	JSON	EXECUTION
Row	f0		

4. For order\_reviews table:

```
1 select count(*) from `scaler-dsml-target-sql.target11.order_reviews`
```

## Query results

JOB INFORMATION	RESULTS	JSON	EXECUTION
Row	f0		

5. For orders table :

```
1 select count(*) from `scaler-dsml-target-sql.target1`
```

## Query results

JOB INFORMATION

RESULTS

JSON

EXECUTI

Row	f0

6. For payments table:

```
1 select count(*) from `scaler-dsml-target-sql.target1`
```

## Query results

JOB INFORMATION

RESULTS

JSON

EXECUTIO

Row	f0_

7. For products table :

```
1 select count(*) from `scaler-dsml-target-sql.target`
```

## Query results

JOB INFORMATION	RESULTS	JSON	EXECUTI
Row	f0		

8. For sellers table :

```
1 select count(*) from `scaler-dsml-target-sql.targe
```

## Query results

JOB INFORMATION	RESULTS	JSON	EXECUT
Row	f0_		

## b. Number of null or missing values in a column

### 1. For Customers Table :

```
1 select count(*)-count(customer_id) as missing_customer_ids,
2      count(*)-count(customer_unique_id) as missing_customer_unique_ids,
3      count(*)-count(customer_zip_code_prefix) as missing_customer_zip_code_prefix,
4      count(*)-count(customer_city) as missing_customer_city,
5      count(*)-count(customer_state) as missing_customer_state,
6
7   from `scaler-dsml-target-sql.target11.customers`
```

Query results

 [SAVE RESULTS ▾](#)

JOB INFORMATION

**RESULTS**

JSON

EXECUTION DETAILS

### 2. For GeoLocation Table:

```
1 select count(*)-count(geolocation_zip_code_prefix) as missing_location_zip_code_prefix,
2      count(*)-count(geolocation_lat) as missing_geolocation_lat,
3      count(*)-count(geolocation_lng) as missing_geolocation_lng,
4      count(*)-count(geolocation_city) as missing_geolocation_city,
5      count(*)-count(geolocation_state) as missing_geolocation_state
6
7   from `scaler-dsml-target-sql.target11.geolocation`
```

Query results

 [SAVE RESULTS ▾](#)

JOB INFORMATION

**RESULTS**

JSON

EXECUTION DETAILS

### 3. For Order\_items Table:

      This

```
1 select count(*)-count(order_id) as missing_order_id,
2      count(*)-count(order_item_id) as missing_order_item_id,
3      count(*)-count(product_id) as missing_product_id,
4      count(*)-count(seller_id) as missing_seller_id,
5      count(*)-count(shipping_limit_date) as missing_shipping_limit_date,
6      count(*)-count(price) as missing_price,
7      count(*)-count(freight_value) as missing_freight_value
8
9   from `scaler-dsml-target-sql.target11.order_items`
```

Query results

 [SAVE RESULTS ▾](#)

JOB INFORMATION

**RESULTS**

JSON

EXECUTION DETAILS

#### 4. For Order\_reviews Table :

The screenshot shows a SQL query editor interface. At the top, there are buttons for RUN, SAVE, SHARE, SCHEDULE, MORE, and a checkbox labeled "This query". Below the toolbar is the SQL code:

```
1 select count(*)-count(review_id) as missing_review_id,
2      count(*)-count(order_id) as missing_order_id,
3      count(*)-count(review_score) as missing_review_score,
4      count(*)-count(review_comment_title) as missing_review_comment_title,
5      count(*)-count(review_creation_date) as missing_review_creation_date,
6      count(*)-count(review_answer_timestamp) as review_answer_timestamp
7   from `scaler-dsml-target-sql.target11.order_reviews`
```

Below the code, the text "Query results" is displayed, followed by a "SAVE RESULTS" button.

Query results

SAVE RESULTS ▾

JOB INFORMATION      RESULTS      JSON      EXECUTION DETAILS

There are 87675 missing values of Review\_comment\_title

#### 5. For Orders table :

The screenshot shows a SQL query editor interface. At the top, there are buttons for RUN, SAVE, SHARE, SCHEDULE, MORE, and a checkbox labeled "This query". Below the toolbar is the SQL code:

```
1 select count(*)-count(order_id) as missing_order_id,
2      count(*)-count(customer_id) as missing_customer_id,
3      count(*)-count(order_status) as missing_order_status,
4      count(*)-count(order_purchase_timestamp) as missing_order_purchase_timestamp,
5      count(*)-count(order_approved_at) as missing_order_approved_at,
6      count(*)-count(order_delivered_carrier_date) as missing_delivered_carrier_date,
7      count(*)-count(order_delivered_customer_date) as missing_order_delivered_customer_date,
8      count(*)-count(order_estimated_delivery_date) as missing_order_estimated_delivery_date
9   from `scaler-dsml-target-sql.target11.orders`
```

Below the code, the text "Query results" is displayed, followed by a "SAVE RESULTS" button.

Query results

SAVE RESULTS ▾

JOB INFORMATION      RESULTS      JSON      EXECUTION DETAILS

There are 160 missing values in order\_approved\_at column, 1783 missing values in order\_delivered\_carrier\_date and 2965 missing values in order\_estimated\_delivery\_date

## 6. For Payments table :

The screenshot shows a database query interface with the following details:

- Toolbar buttons: RUN, SAVE, SHARE, SCHEDULE, MORE.
- Query code:

```
1 select count(*)-count(order_id) as missing_order_id,
2      count(*)-count(payment_sequential) as missing_payment_sequential,
3      count(*)-count(payment_type) as missing_payment_type,
4      count(*)-count(payment_installments) as missing_payment_installments,
5      count(*)-count(payment_value) as missing_payment_value,
6
7   from `scaler-dsml-target-sql.target11.payments`
```
- Result pane: "Query results" and "SAVE RESULTS" button.

Query results

SAVE RESULTS

JOB INFORMATION

RESULTS

JSON

EXECUTION DETAILS

## 7. For Products table :

The screenshot shows a database query interface with the following details:

- Toolbar buttons: RUN, SAVE, SHARE, SCHEDULE, MORE.
- Query code:

```
1 select count(*)-count(product_id) as missing_id,
2      count(*)-count(product_category) as missing_category,
3      count(*)-count(product_name_length) as missing_name_length,
4      count(*)-count(product_description_length) as missing_desc_length,
5      count(*)-count(product_photos_qty) as missing_photos_qty,
6      count(*)-count(product_weight_g) as missing_weight_g,
7      count(*)-count(product_length_cm) as missing_length,
8      count(*)-count(product_height_cm) as missing_height,
9      count(*)-count(product_width_cm) as missing_width,
10
11   from `scaler-dsml-target-sql.target11.products`
```
- Result pane: "Query results" and "SAVE RESULTS" button.

Query results

SAVE RESULTS

There are 610 missing values in each columns of product\_category, product\_name\_length, product\_description\_length, product\_photos\_qty and 2 missing values in each of product\_length\_cm, product\_height\_cm and product\_width\_cm.

8. For Sellers table :

The screenshot shows a SQL query editor interface. At the top, there are several buttons: RUN (highlighted in black), SAVE, SHARE, SCHEDULE, and MODE. The query itself is numbered from 1 to 6:

```
1 select count(*)-count(seller_id) as missing_id,
2      |      | count(*)-count(seller_zip_code_prefix) as missing_zip_code_prefix,
3      |      | count(*)-count(seller_city) as missing_city,
4      |      | count(*)-count(seller_state) as missing_state,
5      |
6      | from `scaler-dsml-target-sql.target11.sellers`
```

Below the query editor, the text "Query results" is displayed.

### c. Data type of columns in a table

RUNSAVESHARESCHEDULEMORE

```
1 SELECT table_name,column_name,ordinal_position,data_type FROM target11.INFORMATION
```

#### Query results

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS	
Row	table_name	column_name		ordinal_posi...	data_type
1	order_items	order_id		1	STRING
2	order_items	order_item_id		2	INT64
3	order_items	product_id		3	STRING
4	order_items	seller_id		4	STRING
5	order_items	shipping_limit_date		5	TIMESTAMP
6	order_items	price		6	FLOAT64
11	sellers	seller_state		4	STRIN
12	geolocation	geolocation_zip_code_prefix		1	INT64
13	geolocation	geolocation_lat		2	FLOAT
14	geolocation	geolocation_lng		3	FLOAT
15	geolocation	geolocation_city		4	STRIN
16	geolocation	geolocation_state		5	STRIN
17	products	product_id		1	STRIN
18	products	product_category		2	STRIN

21	products	product_photos_qty	5	INT64
22	products	product_weight_g	6	INT64
23	products	product_length_cm	7	INT64
24	products	product_height_cm	8	INT64
25	products	product_width_cm	9	INT64
26	orders	order_id	1	STRING
27	orders	customer_id	2	STRING
28	orders	order_status	3	STRING

31	orders	order_delivered_carrier_date	6	TIMESTAMP
32	orders	order_delivered_customer_date	7	TIMESTAMP
33	orders	order_estimated_delivery_date	8	TIMESTAMP
34	payments	order_id	1	STRING
35	payments	payment_sequential	2	INT64
36	payments	payment_type	3	STRING
37	payments	payment_installments	4	INT64
38	payments	payment_value	5	FLOAT64

41	customers	customer_zip_code_prefix	3	INT64
42	customers	customer_city	4	STRING
43	customers	customer_state	5	STRING
44	order_reviews	review_id	1	STRING
45	order_reviews	order_id	2	STRING
46	order_reviews	review_score	3	INT64
47	order_reviews	review_comment_title	4	STRING

d. Get the time period for which the data is given

```
select
concat(
| EXTRACT(MONTH FROM
(SELECT MIN(o.order_purchase_timestamp)
FROM `scaler-dsml-target-sql.target11.orders` o)), '/',
EXTRACT(YEAR FROM
(SELECT MIN(o.order_purchase_timestamp)
FROM `scaler-dsml-target-sql.target11.orders` o)), '-'
| EXTRACT(MONTH FROM
(SELECT MAX(o.order_purchase_timestamp)
FROM `scaler-dsml-target-sql.target11.orders` o)), '/',
EXTRACT(YEAR FROM
(SELECT MAX(o.order_purchase_timestamp)
FROM `scaler-dsml-target-sql.target11.orders` o)))
)
```

ssing location: US 

---

ery results

### e. Number of cities in our dataset

```
1 select count(distinct geolocation_city) as no_of_cities from `scaler-dsml-target-sq
```

rocessing location: US 

#### Query results

---

JOB INFORMATION

RESULTS

JSON

EXECUTION DETAILS

### f. Number of states in our dataset

```
1 select count(distinct geolocation_state) as no_of_states| from `scaler-dsml-target-sql`.
```

Processing location: US 

#### Query results

---

JOB INFORMATION

RESULTS

JSON

EXECUTION DETAILS

## In-depth Exploration:

- a. How many orders do we have for each order status?

```
1 select order_status, count(order_status) as no_of_orders
2 from `scaler-dsml-target-sql.target11.orders` o
3 group by order_status
```

Processing location: US 

## Query results

JOB INFORMATION		RESULTS	JSON	EXECUT
Row	order_status	no_of_orders		
1	created	5		
2	shipped	1107		
3	approved	2		
4	canceled	625		
5	invoiced	314		

b. Is there a growing trend on e-commerce in Brazil? How can we describe a complete scenario?

```
1 select *,lag(Num_of_orders_in_year,1) over (order by year) as prev_yearr_orders,
2 (Num_of_orders_in_year-lag(Num_of_orders_in_year,1) over (order by year)) as difference
3 from(
4 SELECT Extract(year from order_purchase_timestamp) as year,count(order_id) as Num_of
5 FROM `scaler-dsml-target-sql.target11.orders`
6 group by year order by year)
7 order by year
```

Query results

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS	
Row	year	Num_of_ord...	prev_yearr_...	difference_i...	
1	2016	329	null	null	

c. On what day of week brazilians customers tend to do online purchasing?

```
1 select * from
2 (select day,count(day) as no_of_orders_placed fro
3 (select FORMAT_DATE("%a",order_purchase_timestamp
4 from `scaler-dsml-target-sql.target11.orders` o)
5 group by day
6 order by count(day) desc
7 ) as X
8 limit 1
9
```

Processing location: US 

Query results

JOB INFORMATION		RESULTS	JSON	EXECU	

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

```
1 select
2 case
3 when hour>=4 and hour<=6 then 'Dawn'
4 when hour>6 and hour<12 then 'Morning'
5 when hour>=12 and hour<16 then 'Afternoon'
6 when hour>=16 and hour<=23 then 'Night'
7 else 'Midnight'
8 end as Part_of_Day, count(order_id) as no_of_orders
9
10 from (SELECT Extract(hour from order_purchase_timestamp) as hour, order_id FROM scaler-dsml-target-sc
11
12 group by Part_of_Day order by no_of_orders desc
13 limit 1
```

Processing location: US

#### Query results

TOP INFORMATION	RESULTS	ICON	EXECUTION DETAILS
	<p>e. <b>Feature Extraction:</b> Through order_purchase_timestamp in “orders” dataset extract</p> <ol style="list-style-type: none"><li>i. order_purchase_year</li><li>ii. order_purchase_month</li><li>iii. order_purchase_date</li><li>iv. order_purchase_day</li><li>v. order_purchase_dayofweek</li><li>vi. order_purchase_dayofweek_name</li><li>vii. order_purchase_hour</li><li>viii. order_purchase_time_day</li></ol> <pre>1 select 2   order_purchase_timestamp, 3     EXTRACT(year from order_purchase_timestamp) as order_purchase_year, 4     EXTRACT(month from order_purchase_timestamp) as order_purchase_month, 5     EXTRACT(date from order_purchase_timestamp) as order_purchase_date, 6     EXTRACT(day from order_purchase_timestamp) as order_purchase_day, 7     EXTRACT(dayofweek from order_purchase_timestamp) as order_purchase_dayofweek, 8     FORMAT_DATE("%a", order_purchase_timestamp) as order_purchase_dayofweek_name, 9     X.hour as order_purchase_hour, 10  case 11  when hour&gt;=4 and hour&lt;=6 then 'Dawn' 12  when hour&gt;6 and hour&lt;12 then 'Morning' 13  when hour&gt;=12 and hour&lt;16 then 'Afternoon' 14  when hour&gt;=16 and hour&lt;21 then 'Evening' 15  when hour&gt;=21 or hour=0 then 'Night' 16  else 'Midnight' 17  end as order_purcahse_time_day</pre>		

## Query results

[SAVE RESULT](#)

JOB INFORMATION		RESULTS		JSON		EXECUTION DETAILS				
Row	order_purchase_timestamp	order_purch...	order_purch...	order_purch...	order_purch...	order_purch...	order_purch...	order_purchas...	order_...	
1	2017-12-01 11:09:24 UTC	2017	12	2017-12-01	1	6	Fri			
2	2017-12-01 14:38:00 UTC	2017	12	2017-12-01	1	6	Fri			
3	2017-08-01 20:26:19 UTC	2017	8	2017-08-01	1	3	Tue			
4	2018-04-01 11:30:07 UTC	2018	4	2018-04-01	1	1	Sun			
5	2017-07-01 16:48:42 UTC	2017	7	2017-07-01	1	7	Sat			
6	2017-04-01 18:26:29 UTC	2017	4	2017-04-01	1	7	Sat			
7	2018-03-01 18:19:28 UTC	2018	3	2018-03-01	1	5	Thu			
8	2017-07-01 15:05:01 UTC	2017	7	2017-07-01	1	7	Sat			

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

### a. Get month on month orders by region

```

1 select
2 c.customer_state, EXTRACT(YEAR from o.order_purchase_timestamp) as year, EXTRACT(MONTH from o.order_purchase_timestamp)
3 no_of_orders
4 from `scaler-dsml-target-sql.target11.orders`o join `scaler-dsml-target-sql.target11.customers` c
5 on o.customer_id = c.customer_id
6 group by c.customer_state, year, month
7 order by c.customer_state,year,month

```

## Query results

[SAVE RESULTS](#)

JOB INFORMATION		RESULTS		JSON		EXECUTION DETAILS		
Row	customer_state	year	month	no_of_orders				
1	AC	2017	1	2				
2	AC	2017	2	3				
3	AC	2017	3	2				
4	AC	2017	4	5				
5	AC	2017	5	8				
6	AC	2017	6	4				
7	AC	2017	-	-				

## b. Total of customer orders by state

```
1 select
2 c.customer_state, count(order_id) as no_of_orders
3 from `scaler-dsml-target-sql.target11.orders`o join `scaler-dsml-target-sql.target11.customers`
4 on o.customer_id = c.customer_id
5 group by c.customer_state
6
```

### Query results

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS
Row	customer_state	no_of_orders		
1	RJ	12852		
2	RS	5466		
3	SP	41746		
4	DF	2140		
5	PR	5045		
6	MT	907		
7	MA	747		

## c. Top 10 brazilian cities most no. of orders

```
1 select
2 c.customer_city, count(order_id) as no_of_orders
3 from `scaler-dsml-target-sql.target11.orders`o join `scaler-dsml-target-sql.ta
4 on o.customer_id = c.customer_id
5 group by c.customer_city
6 order by no_of_orders desc limit 10
7
```

### Query results

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS
Row	customer_city	no_of_orders		
1	sao paulo	15540		
2	rio de janeiro	6882		
3	belo horizonte	2773		
4	brasilia	2131		
5	curitiba	1521		
6	campinas	1444		

#### d. How are customers distributed in Brazil

```
1 select
2 c.customer_state,c.customer_city, count(order_id) as no_of_orders
3 from `scaler-dsml-target-sql.target11.orders`o join `scaler-dsml-target-sql.ta
4 on o.customer_id = c.customer_id
5 group by c.customer_state,c.customer_city
6 order by c.customer_state, c.customer_city
7
```

#### Query results

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS	
Row	customer_state	customer_city		no_of_orders	
1	AC	brasileia		1	
2	AC	cruzeiro do sul		3	
3	AC	epitaciolandia		1	
4	AC	manoel urbano		1	
5	AC	porto acre		1	
6	AC	rio branco		70	

### e. City wise number of unique customers

```
1 select
2 c.customer_city, count(c.customer_unique_id) as no_of_custo
3 from `scaler-dsml-target-sql.target11.customers` c
4 group by c.customer_city
5 order by no_of_customers desc
6
```

### Query results

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAIL
Row	customer_city	no_of_custo...		
1	sao paulo	15540		
2	rio de janeiro	6882		
3	belo horizonte	2773		
4	brasilia	2131		
5	curitiba	1521		
6	campinas	1444		

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

#### Step 1: Using CTE

1. “order\_items” + “order” joined on order id where order\_purchase timestamp is already divided into month & year

2. Group data by year and month, aggregation count(order\_id), sum(price), sum(freight\_value)

3. Create new columns:

price\_per\_order = sum(price) / count(order\_id)

```
freight_per_order= sum(freight_value) / count(order_id)
```

## Step 2: Answer the following questions:

### a. Total amount sold in 2017 between Jan to August

```
1 with cte as (
2
3     select
4
5         EXTRACT(year from orders.order_purchase_timestamp) as year,
6         EXTRACT(month from orders.order_purchase_timestamp) as month,
7         count(orders.order_id) as no_of_orders, sum(items.price) as total_price,sum(items.freight_value)
8         sum(items.price)/count(orders.order_id) as price_per_order,
9         sum(items.freight_value) as freight_per_order
10    from `scaler-dsml-target-sql.target11.order_items` items join `scaler-dsml-target-sql.target11.or
11    on items.order_id = orders.order_id
12    group by year,month
13
14 )
15 select sum(total_price) from cte where year=2017 and month between 1 and 8
16
```

### Query results

### b. Total amount sold in 2018 between Jan to august

```
1 with cte as (
2
3     select
4
5         EXTRACT(year from orders.order_purchase_timestamp) as year,
6         EXTRACT(month from orders.order_purchase_timestamp) as month,
7         count(orders.order_id) as no_of_orders, sum(items.price) as total_price,sum(items.freight_value) a:
8         sum(items.price)/count(orders.order_id) as price_per_order,
9         sum(items.freight_value) as freight_per_order
10    from `scaler-dsml-target-sql.target11.order_items` items join `scaler-dsml-target-sql.target11.or
11    on items.order_id = orders.order_id
12    group by year,month
13
14 )
15 select sum(total_price) from cte where year=2018 and month between 1 and 8
16
```

### Query results

### c. % increase from 2017 to 2018

```

1 with cte as [
2   select |
3     EXTRACT(year from orders.order_purchase_timestamp) as year,
4     EXTRACT(month from orders.order_purchase_timestamp) as month,
5     count(orders.order_id) as no_of_orders, sum(items.price) as total_price,sum(items.freight_value) as total_freight_val
6     sum(items.price)/count(orders.order_id) as price_per_order,
7     sum(items.freight_value) as freight_per_order
8   from `scaler-dsml-target-sql.target11.order_items` items join `scaler-dsml-target-sql.target11.orders` orders
9   on items.order_id = orders.order_id
10  group by year,month
11 ]
12 select X.year,X.price_this_year,X.price_prev_year,X.rate_of_increase
13 from
14 (select cte.year, round(sum(cte.total_price),2) as price_this_year,
15 round((lag(sum(cte.total_price),1) over (order by cte.year )),2) as price_prev_year,
16 round(((sum(cte.total_price)-(lag(sum(cte.total_price),1) over (order by cte.year ))))/(lag(sum(cte.total_price),1) over
17 as rate_of_increase
18 from cte
19 group by cte.year
20 order by cte.year desc) as X
21 where X.year=2018

```

Query results

 [SAVE RESULTS](#)

**Step 3: Join (orders+order\_items) table from previous step with “customers” table on Customer\_id and find:**

#### a. Mean & Sum of price by customer state

```

1 with cte as (
2   select
3     customers.customer_state,
4     count(orders.order_id) as no_of_orders, sum(items.price) as total_price,sum(items.freight_value) as total_freight_'
5     sum(items.price)/count(orders.order_id) as price_per_order,
6     sum(items.freight_value) as freight_per_order
7   from `scaler-dsml-target-sql.target11.order_items` items join `scaler-dsml-target-sql.target11.orders` orders
8   on items.order_id = orders.order_id join `scaler-dsml-target-sql.target11.customers` customers on orders.customer_i
9   group by customers.customer_state
10 )
11
12 select cte.customer_state,cte.total_price,cte.price_per_order as avg_price from cte

```

Query results

 [SAVE RESULT](#)

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS
Row	customer_state			
1	MT	156453.529...	148.297184...	
2	MA	119648.219...	145.204150...	
3	AI	80314.81	180.889211...	

### b. Mean & Sum of freight value by customer state

```
1 with cte as (
2   select
3     customers.customer_state,
4     count(orders.order_id) as no_of_orders, sum(items.price) as total_price,sum(items.freight_value) as total_freight
5     sum(items.price)/count(orders.order_id) as price_per_order,
6     sum(items.freight_value)/count(orders.order_id) as freight_per_order
7   from `scaler-dsml-target-sql.target11.order_items` items join `scaler-dsml-target-sql.target11.orders` orders
8   on items.order_id = orders.order_id join `scaler-dsml-target-sql.target11.customers` customers on orders.customer
9   group by customers.customer_state
10 )
11
12 select cte.customer_state,cte.total_freight_value,round(cte.freight_per_order,2) as avg_freight_value from cte
```

Query results

 [SAVE RESI](#)

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS	
Row	customer_state	total_freight_value	avg_freight_value		
1	MT	29715.4300...	28.1662843...		
2	MA	31523.7700...	38.2570024...		
3	DE	15914.5800	25.8436711		

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

1. Calculating days between purchasing, delivering and estimated delivery

2. Create columns:

time\_to\_delivery = order\_purchase\_timestamp-order\_delivered\_customer\_date

diff\_estimated\_delivery = order\_estimated\_delivery\_date-order\_delivered\_customer\_date

3. Grouping data by state, take mean of freight\_value, time\_to\_delivery, diff\_estimated\_delivery

4. Sort the data to get the following:

## a. Top 5 states with highest/lowest average freight value

### States with Highest freight value :

```
1 select
2 X.customer_state, avg(X.freight_value) as avg_freight_value
3 from (select
4 c.customer_state,
5 DATE_DIFF(o.order_delivered_customer_date, o.order_purchase_timestamp, DAY) as time_to_delivery,
6 DATE_DIFF(o.order_estimated_delivery_date, o.order_delivered_customer_date, DAY) as diff_estimated_delivery,
7 oi.freight_value
8
9 from `scaler-dsml-target-sql.target11.orders` o join `scaler-dsml-target-sql.target11.order_items` oi
10 on o.order_id = oi.order_id
11 join `scaler-dsml-target-sql.target11.customers` c
12 on o.customer_id = c.customer_id
13 where o.order_delivered_customer_date is not null and o.order_purchase_timestamp is not null and o.order_estimate
14 ) as X
15 group by X.customer_state
16 order by avg_freight_value desc limit 5
```

### Query results

[SAVE RES](#)

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS
Row	customer_state	avg_freight...		
1	PB	43.0916894...		

### States with Lowest freight Value :

```
1 select
2 X.customer_state, avg(X.freight_value) as avg_freight_value
3 from (select
4 c.customer_state,
5 DATE_DIFF(o.order_delivered_customer_date, o.order_purchase_timestamp, DAY) as time_to_delivery,
6 DATE_DIFF(o.order_estimated_delivery_date, o.order_delivered_customer_date, DAY) as diff_estimated_delivery,
7 oi.freight_value
8
9 from `scaler-dsml-target-sql.target11.orders` o join `scaler-dsml-target-sql.target11.order_items` oi
10 on o.order_id = oi.order_id
11 join `scaler-dsml-target-sql.target11.customers` c
12 on o.customer_id = c.customer_id
13 where o.order_delivered_customer_date is not null and o.order_purchase_timestamp is not null and o.order_estimate
14 ) as X
15 group by X.customer_state
16 order by avg_freight_value asc limit 5
```

### Query results

[SAVE R](#)

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS
Row	customer_state	avg_freight...		
1	SP	15.1149940...		

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

States with highest avg delivery time :

```
1 select
2 X.customer_state, avg(X.time_to_delivery) as avg_delivery_time
3 from (select
4 c.customer_state,
5 DATE_DIFF(o.order_delivered_customer_date, o.order_purchase_timestamp, DAY) as time_to_delivery,
6 DATE_DIFF(o.order_estimated_delivery_date, o.order_delivered_customer_date, DAY) as diff_estimated_delivery,
7 oi.freight_value
8 from `scaler-dsml-target-sql.target11.orders` o join `scaler-dsml-target-sql.target11.order_items` oi
9 on o.order_id = oi.order_id
10 join `scaler-dsml-target-sql.target11.customers` c
11 on o.customer_id = c.customer_id
12 where o.order_delivered_customer_date is not null and o.order_purchase_timestamp is not null and o.order_estimated_d
13 ) as X
14 group by X.customer_state
15 order by avg_delivery_time desc limit 5
```

Query results

 [SAVE RESULT](#)

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS
Row	customer_state	avg_delivery...		
1	RR	27.8260869...		

```
1 select
2 X.customer_state, avg(X.time_to_delivery) as avg_delivery_time
3 from (select
4 c.customer_state,
5 DATE_DIFF(o.order_delivered_customer_date, o.order_purchase_timestamp, DAY) as time_to_delivery,
6 DATE_DIFF(o.order_estimated_delivery_date, o.order_delivered_customer_date, DAY) as diff_estimated_delivery,
7 oi.freight_value
8 from `scaler-dsml-target-sql.target11.orders` o join `scaler-dsml-target-sql.target11.order_items` oi
9 on o.order_id = oi.order_id
10 join `scaler-dsml-target-sql.target11.customers` c
11 on o.customer_id = c.customer_id
12 where o.order_delivered_customer_date is not null and o.order_purchase_timestamp is not null and o.order_estimated.
13 ) as X
14 group by X.customer_state
15 order by avg_delivery_time asc limit 5
```

Query results

 [SAVE RESU](#)

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS
Row	customer_state	avg_delivery...		
1	SP	8.25960855...		

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

states with delivery fast compared to estimated delivery :

```
1 select
2 X.customer_state, avg(X.time_to_delivery-X.diff_estimated_delivery) as avg_delay
3 from (select
4 c.customer_state,
5 DATE_DIFF(o.order_delivered_customer_date, o.order_purchase_timestamp, DAY) as time_to_delivery,
6 DATE_DIFF(o.order_estimated_delivery_date, o.order_delivered_customer_date, DAY) as diff_estimated_delivery,
7 oi.freight_value
8 from `scaler-dsml-target-sql.target11.orders` o join `scaler-dsml-target-sql.target11.order_items` oi
9 on o.order_id = oi.order_id
10 join `scaler-dsml-target-sql.target11.customers` c
11 on o.customer_id = c.customer_id
12 where o.order_delivered_customer_date is not null and o.order_purchase_timestamp is not null and o.order_estimated_
13 ) as X
14 group by X.customer_state
15 order by avg_delay asc limit 5
```

Query results

 [SAVE RESULT](#)

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS
Row	customer_state	avg_delay		
1	SP	-2.0059858...		

states with late delivery compared to estimated delivery :

```
1 select
2 X.customer_state, avg(X.time_to_delivery-X.diff_estimated_delivery) as avg_delay
3 from (select
4 c.customer_state,
5 DATE_DIFF(o.order_delivered_customer_date, o.order_purchase_timestamp, DAY) as time_to_delivery,
6 DATE_DIFF(o.order_estimated_delivery_date, o.order_delivered_customer_date, DAY) as diff_estimated_delivery,
7 oi.freight_value
8 from `scaler-dsml-target-sql.target11.orders` o join `scaler-dsml-target-sql.target11.order_items` oi
9 on o.order_id = oi.order_id
10 join `scaler-dsml-target-sql.target11.customers` c
11 on o.customer_id = c.customer_id
12 where o.order_delivered_customer_date is not null and o.order_purchase_timestamp is not null and o.order_estimated_
13 ) as X
14 group by X.customer_state
15 order by avg_delay desc limit 5
```

Query results

 [SAVE RES](#)

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS
Row	customer_state	avg_delay		
1	AL	16.0163934...		

## 6. Payment type analysis: Join “payments” dataset with the existing data on order\_id

### a. Count of orders for different payment types

```
1 select
2 X.payment_type, count(X.order_id) as no_of_orders
3 from (select
4 c.customer_state,
5 p.payment_type,o.order_id,p.payment_installments,
6 DATE_DIFF(o.order_delivered_customer_date, o.order_purchase_timestamp, DAY) as time_to_delivery,
7 DATE_DIFF(o.order_estimated_delivery_date,o.order_delivered_customer_date, DAY) as diff_estimated_delivery,
8 oi.freight_value
9 from `scaler-dsml-target-sql.target11.orders` o join `scaler-dsml-target-sql.target11.order_items` oi
10 on o.order_id = oi.order_id
11 join `scaler-dsml-target-sql.target11.customers` c
12 on o.customer_id = c.customer_id
13 join `scaler-dsml-target-sql.target11.payments` p
14 on p.order_id = o.order_id
15 where o.order_delivered_customer_date is not null and o.order_purchase_timestamp is not null and o.order_estimated
16 ) as X
```

Query results

 SAVE RES

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS
Row	payment_type	no_of_orders		
1	credit card	0.1006		

### b. Distribution of payment installments and count of orders

```
1 select
2 X.payment_installments, count(X.order_id) as no_of_orders
3 from (select
4 c.customer_state,p.payment_type,o.order_id,p.payment_installments,
5 DATE_DIFF(o.order_delivered_customer_date, o.order_purchase_timestamp, DAY) as time_to_delivery,
6 DATE_DIFF(o.order_estimated_delivery_date,o.order_delivered_customer_date, DAY) as diff_estimated_delivery,
7 oi.freight_value
8 from `scaler-dsml-target-sql.target11.orders` o join `scaler-dsml-target-sql.target11.order_items` oi
9 on o.order_id = oi.order_id
10 join `scaler-dsml-target-sql.target11.customers` c
11 on o.customer_id = c.customer_id
12 join `scaler-dsml-target-sql.target11.payments` p
13 on p.order_id = o.order_id
14 where o.order_delivered_customer_date is not null and o.order_purchase_timestamp is not null and o.order_estimated
15 ) as X
16 group by X.payment_installments order by no_of_orders desc
```

Query results

 SAVE RES

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS
Row	payment_in...	no_of_orders		
1	1	0.1006		

### c. Count of orders for different payment types Month over Month

```

1 select
2 X.year,X.MONTH, X.payment_type,count(X.order_id) as no_of_orders
3 from (select
4 c.customer_state,p.payment_type,o.order_id,p.payment_installments,
5 EXTRACT(MONTH from o.order_purchase_timestamp) as month,
6 EXTRACT(YEAR from o.order_purchase_timestamp) as year,
7 oi.freight_value
8 from `scaler-dsml-target-sql.target11.orders` o join `scaler-dsml-target-sql.target11.order_items` oi
9 on o.order_id = oi.order_id
10 join `scaler-dsml-target-sql.target11.customers` c
11 on o.customer_id = c.customer_id
12 join `scaler-dsml-target-sql.target11.payments` p
13 on p.order_id = o.order_id
14 where o.order_delivered_customer_date is not null and o.order_purchase_timestamp is not null and o.order_estimated_delivery_date is not null
15 ) as X
16 group by X.year,X.month,X.payment_type order by X.year,X.month asc

```

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS	
Row	year	MONTH	payment_type	no_of_orders	
1	2016	10	credit_card	253	
2	2016	10	UPI	61	
3	2016	10	voucher	20	
4	2016	10	debit_card	2	
5	2016	12	credit_card	1	
6	2017	1	credit_card	659	
7	2017	1	voucher	68	
8	2017	1	UPI	226	
9	2017	1	debit_card	10	
..	..	..	..	..	..

### Insights :

- There are 99441 orders placed in the given data
- There are 32951 different types of products which the customers ordered
- There are 87675 reviews where title is missing for review\_comments.
- There are 160 missing values in order\_approved\_at, 1783 missing values in order\_delivered\_customer\_date and 2965 missing values in order\_estimated\_delivery\_date in Orders table
- There are 610 missing values in each columns of product\_category, product name product\_description\_length, product photos qty and 2 missing values in each of product length\_cm, product\_height\_cm and product\_width\_cm.
- The data is between September 2016 and October 2018.
- There are 27 states and 8011 cities in the data

- There is growing trend in ecommerce usage over the years.

```

1 select *,lag(Num_of_orders_in_year,1) over (order by year) as prev_yearr_orders,
2 (Num_of_orders_in_year-lag(Num_of_orders_in_year,1) over (order by year)) as difference
3 from(
4 |  SELECT Extract(year from order_purchase_timestamp) as year,count(order_id) as Num_of_
5 |  FROM `scaler-dsml-target-sql.target11.orders`_
6 |  group by year order by year)
7 order by year

```

## Query results

JOB INFORMATION		RESULTS		JSON	EXECUTION DETAILS	
Row	year	Num_of_ord...	prev_yearr...	difference_i...	nuli	nuli
1	2016	329				

- More number of orders are placed on Monday

```

1 select * from
2 (select day,count(day) as no_of_orders_placed fro
3 (select FORMAT_DATE("%a",order_purchase_timestamp
4 from `scaler-dsml-target-sql.target11.orders` o)
5 group by day
6 order by count(day) desc
7 ) as X
8 limit 1
9

```

Processing location: US

## Query results

JOB INFORMATION		RESULTS		JSON	EXECU'	

- The states Sp, RJ and Mg have more number of orders while the states AC,AP and AR have least number of orders
- The “sau palo” city has more number of orders
- There is 19.99% of revenue increase from 2017 to 2018 even when the data of 2018 is till October only
- The states PB and SP have highest and lowest avg freight rates respectively.
- The states RR and SP have highest and lowest avg delivery times.
- SP state has the lowest delay where as AL state has highest order delay between estimated and delivered time

- More number of orders are placed through credit card

```

1 select
2 X.payment_type, count(X.order_id) as no_of_orders
3 from (select
4 c.customer_state,
5 p.payment_type,o.order_id,p.payment_installments,
6 DATE_DIFF(o.order_delivered_customer_date, o.order_purchase_timestamp, DAY) as time_to_delivery,
7 DATE_DIFF(o.order_estimated_delivery_date,o.order_delivered_customer_date, DAY) as diff_estimated_delivery,
8 oi.freight_value
9 from `scaler-dsml-target-sql.target11.orders` o join `scaler-dsml-target-sql.target11.order_items` oi
10 on o.order_id = oi.order_id
11 join `scaler-dsml-target-sql.target11.customers` c
12 on o.customer_id = c.customer_id
13 join `scaler-dsml-target-sql.target11.payments` p
14 on p.order_id = o.order_id
15 where o.order_delivered_customer_date is not null and o.order_purchase_timestamp is not null and o.order_estimated
16 ) as X

```

### Query results

[SAVE RES](#)

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS
Row	payment_type	no_of_orders		
1	credit card	0.0006		

- More number of orders are paid in one installment

```

1 select
2 X.payment_type, count(X.order_id) as no_of_orders
3 from (select
4 c.customer_state,
5 p.payment_type,o.order_id,p.payment_installments,
6 DATE_DIFF(o.order_delivered_customer_date, o.order_purchase_timestamp, DAY) as time_to_delivery,
7 DATE_DIFF(o.order_estimated_delivery_date,o.order_delivered_customer_date, DAY) as diff_estimated_delivery,
8 oi.freight_value
9 from `scaler-dsml-target-sql.target11.orders` o join `scaler-dsml-target-sql.target11.order_items` oi
10 on o.order_id = oi.order_id
11 join `scaler-dsml-target-sql.target11.customers` c
12 on o.customer_id = c.customer_id
13 join `scaler-dsml-target-sql.target11.payments` p
14 on p.order_id = o.order_id
15 where o.order_delivered_customer_date is not null and o.order_purchase_timestamp is not null and o.order_estimated
16 ) as X

```

### Query results

[SAVE RES](#)

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS
Row	payment_type	no_of_orders		
1	credit card	0.0006		

## **Recommendations :**

- Every time an order is placed by a customer, each instance is treated as different customer but there will be a high possibility that there are repeated customers. If each instance of order is treated as different customer, it is difficult to analyze the customer behavior.
- There are 2965 missing values in estimated\_delivery\_time which is not a good experience for a customer to trace or expect the order
- The data is available for only four months i.e.(September to December) of data in 2016 and only ten months i.e.(January to October) in 2018. It's better to have data for whole year so that we can have proper estimation of trends over years.
- The states AC, AP and AR has least number of orders. We can attract more customers by having more campaigns and offers specific to those states.
- The RR state has more avg time taken for delivery. We can decrease the time by having more physical stores and delivery partners.
- The AL state has more difference in delay between the expected delivery and actual delivery time. Need to trace out the reasons for that.
- More number of orders are placed in single installments. We can have more number of orders if we offer no cost EMI