SQL Target Analysis

Name: Sukanya Devi B

Problem Statement:

Assuming you are a data analyst/ scientist at Target, you have been assigned the task of analyzing the given dataset to extract valuable insights and provide actionable recommendations.

What does 'good' look like?

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

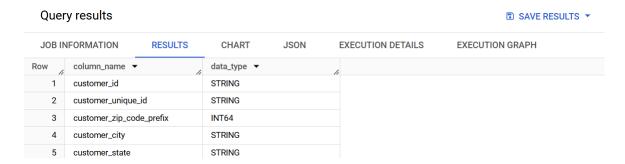
Query:

select column_name, data_type

from scaler-dsml-sql-444307. Target. INFORMATION_SCHEMA. COLUMNS

where table_name = 'customers'

Output:



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

Query:

select min(order purchase timestamp) as start time, max(order purchase timestamp) as end time

from 'scaler-dsml-sql-444307. Target.orders'

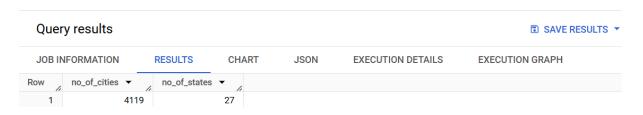


1) 3.Count the Cities & States of customers who ordered during the given period.

Query:

select count(distinct customer_city) as no_of_cities, count(distinct customer_state) as no_of_states from `scaler-dsml-sql-444307.Target.customers`

Output:



- 2) In-depth Exploration:
- 2) 1. Is there a growing trend in the no. of orders placed over the past years?

```
with final as

(select distinct extract(year from order_purchase_timestamp)as purchase_year,
count(order_id)as no_of_orders

from `scaler-dsml-sql-444307.Target.orders`

group by 1

order by 1),
previous_data as

(select * ,lag(no_of_orders,1)over(order by purchase_year)as prev

from final)

select purchase_year, no_of_orders, prev, ((no_of_orders - prev)/prev)as trend

from previous_data

order by 1
```

Query results ☐ SAVE RESULTS ▼

JOB IN	IFORMATION		RESULTS	CH	ART	JSON	EXECUTION DE	TAILS	EXECUTION GRAP
Row	purchase_year	-	no_of_orders	~ //	prev 🔻	h	trend ▼	le.	
1	201	16		329		nuli	null		
2	201	17	4	5101		329	136.0851063829		
3	201	18	5	4011		45101	0.197556595197		

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

Query:

```
with final as

(select distinct extract(month from order_purchase_timestamp)as purchase_month,

count(order_id)as no_of_orders

from `scaler-dsml-sql-444307.Target.orders`

group by 1

order by 1),

previous_data as

(select * ,lag(no_of_orders,1)over(order by purchase_month)as prev

from final)

select purchase_month, no_of_orders, prev, ((no_of_orders - prev)/prev)as trend

from previous_data

order by 1
```

Output:

Query results

■ SAVE RESULTS ▼

JOB IN	FORMATION	RESULTS	CHA	ART JS	ON	EXECUTION DETA	ILS EXECUTION	GRAPH
Row	purchase_month	no_of_orders	s v	prev 🕶	h	trend ▼		
1	1		8069		nuli	nuli		
2	2		8508		8069	0.054405750402		
3	3		9893		8508	0.162787964268		
4	4		9343		9893	-0.05559486505		
5	5		10573		9343	0.131649363159		
6	6		9412	1	0573	-0.10980800151		
7	7		10318		9412	0.096260093497		
8	8		10843	1	0318	0.050881953867		
9	9		4305	1	0843	-0.60296965784		
10	10		4959		4305	0.151916376306		
11	11		7544		4959	0.521274450494		
12	12		5674		7544	-0.24787910922		

2) 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

Query:

```
select * from

(select case when extract(hour from order_purchase_timestamp) between 0 and 6
then "Dawn"

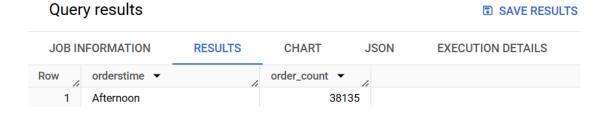
when extract(hour from order_purchase_timestamp) between 7 and 12
then "Mornings"

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 orderstime,

count(distinct order_id) as order_count
from `scaler-dsml-sql-444307.Target.orders`
group by 1
order by 2 desc
limit 1)
```



- 3) Evolution of E-commerce orders in the Brazil region
- 3) 1.Get the month on month no. of orders placed in each state.

Query:

```
select c.customer_state, extract (month from o.order_purchase_timestamp)as order_months, count(distinct o.order_id) as no_of_orders from `scaler-dsml-sql-444307.Target.orders` o inner join `scaler-dsml-sql-444307.Target.customers` c on o.customer_id =c.customer_id group by 1,2 order by 1,2
```

Output:

Quer	y results					SAVE RESULTS ▼
JOB IN	IFORMATION	RESULTS	CHART J	SON EXECU	TION DETAILS	EXECUTION GRAPH
Row	customer_state •		order_months ▼	no_of_orders ▼	6	
1	AC		1	8		
2	AC		2	6		
3	AC		3	4		
4	AC		4	9		
5	AC		5	10		
6	AC		6	7		
7	AC		7	9		
8	AC		8	7		
9	AC		9	5		
10	AC		10	6		
11	AC		11	5		
12	AC		12	5		
13	AL		1	39		
14	AL		2	39		

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

```
select customer_state, count(distinct customer_id)as no_of_customers from `scaler-dsml-sql-444307.Target.customers` group by 1 order by 1
```

Quer	y results					SAVE RESULTS ▼
JOB IN	FORMATION	RESULTS	CHART	JSON	EXECUTION DETAILS	EXECUTION GRAPH
Row	customer_state	· //	no_of_customers	7		
1	AC		81			
2	AL		413	3		
3	AM		148	3		
4	AP		68	3		
5	BA		3380)		
6	CE		1336	5		
7	DF		2140)		
8	ES		2033	3		
9	GO		2020)		
10	MA		747	,		
11	MG		11635	5		
12	MS		715	5		
13	MT		907	,		
14	PA		975	5		

- 4) Impact on Economy: Analyze the money movement by e-commerce by looking at order prices, freight and others.
- 4) 1.Get the % increase in the cost of orders from 2017 to 2018 (include months between Jan to Aug only).

```
WITH time_btw as
(SELECT EXTRACT(YEAR FROM o.order_purchase_timestamp) as Year,
sum(p.payment_value) as cost
FROM 'scaler-dsml-sql-444307. Target.orders' AS o
INNER JOIN 'scaler-dsml-sql-444307. Target.payments' AS p
ON o.order_id = p.order_id
WHERE EXTRACT(MONTH FROM o.order_purchase_timestamp) BETWEEN 1 AND 8 AND
EXTRACT(YEAR FROM o.order_purchase_timestamp) IN (2017,2018)
GROUP BY Year
ORDER BY Year),
lag_btw AS
(SELECT *, LAG(cost) OVER(ORDER BY Year) as lagg
FROM time_btw)
SELECT Year,
ROUND(ifnull(((cost-lagg)/lagg)*100,0),2) as Percentage_increase
FROM lag_btw
ORDER BY Year
```

Query results

SAVE RESULTS

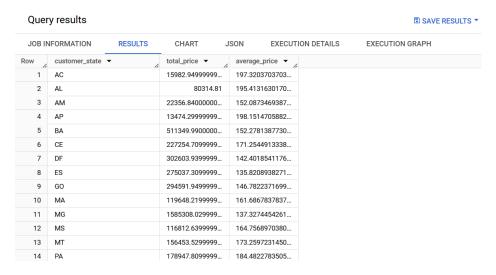
JOB IN	IFORMATION	RESULTS	CHART	JSON	EXECUTION DETAILS
Row	Year ▼	Percent	age_increase		
1	201	7	0.0		
2	201	8	136.98		

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

Query:

```
select c.customer_state ,
sum(oi.price)as total_price,
sum(oi.price)/count(distinct(o.order_id))as average_price
from `scaler-dsml-sql-444307.Target.order_items` oi
inner join `scaler-dsml-sql-444307.Target.orders` o
on oi.order_id = o.order_id
inner join `scaler-dsml-sql-444307.Target.customers` c
on o.customer_id = c.customer_id

group by c.customer_state
order by c.customer_state
```



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

Query:

```
select c.customer_state ,
sum(oi.freight_value)as total_freight_value,
sum(oi.freight_value)/count(distinct(o.order_id))as average_freight_value
from `scaler-dsml-sql-444307.Target.order_items` oi
inner join `scaler-dsml-sql-444307.Target.orders` o
on oi.order_id = o.order_id
inner join `scaler-dsml-sql-444307.Target.customers` c
on o.customer_id = c.customer_id

group by c.customer_state
order by c.customer_state
```

Output:

Query results

☐ SAVE RESULTS ▼

JOB IN	IFORMATION RE	SULTS CHART	JSON EXECUT	ION DETAILS	EXECUTION GRAPH
Row	customer_state ▼	total_freight_value	average_freight_valu		
1	AC	3686.749999999	45.51543209876		
2	AL	15914.589999999	38.72163017031		
3	AM	5478.8899999999	37.27136054421		
4	AP	2788.500000000	41.00735294117		
5	BA	100156.6799999	29.82628945801		
6	CE	48351.58999999	36.43676714393		
7	DF	50625.499999999	23.82376470588		
8	ES	49764.599999999	24.57511111111		
9	GO	53114.97999999	26.46486297957		
10	MA	31523.77000000	42.59968918918		
11	MG	270853.4600000	23.46270443520		
12	MS	19144.03000000	27.00145275035		
13	MT	29715.43000000	32.90745293466		
14	PA	38699.30000000	39.89618556701		

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

Query:

select order_id,

DATE_DIFF(order_delivered_customer_date,order_purchase_timestamp,day)as delivery_time,

date_diff(order_delivered_customer_date,order_estimated_delivery_date,day)as difference_in_days

from `scaler-dsml-sql-444307.Target.orders`

Query results				SAVE RESULTS
IOD INFORMATION	DECLUTO	CHART	NOSI	EVECUTION DETAILS

JOB IN	NFORMATION	RESULTS	CHART	JSON	EXECUTION D	ETAILS
Row	order_id ▼	6	delivery_time ▼	differ	ence_in_days	
1	1950d777989f6a	a877539f5379	30		12	
2	2c45c33d2f9cb8	8ff8b1c86cc28	30	0	-28	
3	65d1e226dfaeb8	3cdc42f66542	3:	5	-16	
4	635c894d068ac	37e6e03dc54e	30	0	-1	
5	3b97562c3aee8	bdedcb5c2e45	3:	2	0	
6	68f47f50f04c4cl	o6774570cfde	29	9	-1	
7	276e9ec344d3b	f029ff83a161c	4:	3	4	
8	54e1a3c2b97fb0)809da548a59	4	0	4	
9	fd04fa4105ee80	45f6a0139ca5	3	7	1	
10	302bb8109d097	a9fc6e9cefc5	3:	3	5	
11	66057d37308e7	87052a32828	3	В	6	
12	19135c945c554	eebfd7576c73	30	б	2	
13	4493e45e7ca10	84efcd38ddeb	34	4	0	

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

```
with cte as
(select c.customer_state,
round((sum(i.freight_value)/count(distinct i.order_id)),2) as
Average freight order price
from `Target.order_items` i
join 'Target.orders' o
on o.order_id=i.order_id
join Target.customers c
on c.customer_id=o.customer_id
group by 1
order by 1),
rank1 as
(select *,
dense_rank() over( order by Average_freight_order_price desc) as
highest_drnk,
dense_rank() over( order by Average_freight_order_price asc) as lowest_drnk
from cte)
select customer_state, Average_freight_order_price
from rank1 where lowest_drnk <=5
union all
select customer_state, Average_freight_order_price
from rank1 where highest_drnk <=5</pre>
order by Average freight order price desc
```

Query results SAVE RESULTS

JOB IN	IFORMATION	RESULTS	CHART	JSON	EXECUTION DETAILS
Row	customer_state ▼		Average_freight_	orde	
1	RR		48.5		
2	РВ		48.3	35	
3	RO		46.2	22	
4	AC		45.5	52	
5	PI		43.0	04	
6	RJ		23.9	95	
7	DF		23.8	32	
8	PR		23.5	58	
9	MG		23.4	16	
10	SP		17.3	37	

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

```
with AverageDelivertime as
(SELECT
c.customer\_state, ROUND(AVG(DATE\_DIFF(DATE(o.order\_delivered\_customer\_date),
DATE(o.order_purchase_timestamp),Day)),2) AS Average_delivery_time
from `Target.orders` o
join 'Target.customers' c
on c.customer_id=o.customer_id
group by 1
order by 2),
rank1 as
(select customer_state,Average_delivery_time,
     dense_rank() over(order by Average_delivery_time desc) as
Highest_Average_delivery_time,
     dense_rank() over(order by Average_delivery_time asc) as
Lowest_Average_delivery_time
     from AverageDelivertime)
select customer_state,Average_delivery_time
    from rank1
```

```
where Highest_Average_delivery_time <=5 or
Lowest_Average_delivery_time <=5
order by Average_delivery_time desc
```

Query results SAVE RESULTS

JOB IN	IFORMATION	RESULTS	CHART	JSON	EXECUTION DETAILS
Row	customer_state	•	Average_delivery_t	im	
1	RR		29.34		
2	AP		27.18		
3	AM		26.36		
4	AL		24.5		
5	PA		23.73		
6	SC		14.91		
7	DF		12.9		
8	MG		11.95		
9	PR		11.94		
10	SP		8.7		

5) 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:

Select

```
c.customer_state,ROUND(AVG(DATE_DIFF(DATE(order_estimated_delivery_date),DATE(order_deliverd_customer_date),Day)),2) AS avg_delivery_days

FROM `scaler-dsml-sql-444307.Target.orders`as o

INNER JOIN `scaler-dsml-sql-444307.Target.customers` AS c

ON o.customer_id = c.customer_id

GROUP BY c.customer_state

ORDER BY avg_delivery_days

LIMIT 5
```

Query results

■ SAVE RESULTS ▼

JOB IN	NFORMATION	RESULTS	CHART	JSON	EXECUTION DETAILS	E
Row	customer_state	▼	avg_delivery_days	i.		
1	AL		8.71			
2	MA		9.57	,		
3	SE		10.02	!		
4	ES		10.5	i		
5	BA		10.79)		

6) Analysis based on the payments

6) 1. Find the month on month no. of orders placed using different payment types.

Query:

```
select p.payment_type, extract (month from o.order_purchase_timestamp)as order_months,
count(distinct o.order_id) as no_of_orders
from `scaler-dsml-sql-444307.Target.orders` o
inner join `scaler-dsml-sql-444307.Target.payments` p
on o.order_id = p.order_id
group by 1,2
order by 2,1
```

Query results SAVE RESULTS							
JOB IN	IFORMATION RESULTS	CHART J	SON EXECUT	ION DETAILS			
Row	payment_type ▼	order_months ▼	no_of_orders ▼				
1	UPI	1	1715				
2	credit_card	1	6093				
3	debit_card	1	118				
4	voucher	1	337				
5	UPI	2	1723				
6	credit_card	2	6582				
7	debit_card	2	82				
8	voucher	2	288				
9	UPI	3	1942				
10	credit_card	3	7682				
11	debit_card	3	109				
12	voucher	3	395				
13	UPI	4	1783				
14	credit_card	4	7276				

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

```
Query:
WITH cte_table AS
(SELECT
c.customer state AS state,
SUM(price) AS total_price,
COUNT(DISTINCT(o.order id)) AS num orders
FROM 'scaler-dsml-sql-444307. Target.orders' o
INNER JOIN 'scaler-dsml-sql-444307. Target.order_items' i
ON o.order_id= i.order_id
INNER JOIN 'scaler-dsml-sql-444307. Target.customers' c
ON o.customer id=c.customer id
GROUP BY state)
SELECT state,
total_price,
num_orders,
(total_price/num_orders) AS avg_price
FROM cte table
ORDER BY total_price DESC;
SELECT payment_installments AS installments,
COUNT(order_id) AS num_orders
FROM 'scaler-dsml-sql-444307. Target.payments'
WHERE payment installments >= 1
GROUP BY payment_installments
ORDER BY num_orders DESC;
```

Query results

SAVE RESULTS ▼

JOB INFORMATION		RESULTS		CHART		JSON	
ow /	installments 🔻	h	num_orders	~			
1		1		52546			
2		2		12413			
3		3		10461			
4		4		7098			
5	1	10		5328			
6		5		5239			
7		8		4268			
8		6		3920			
9		7		1626			
10		9		644			
11	1	12		133			
12	1	15		74			
13	1	18		27			
14	1	11		23			

Insights & Business recommendations

- 1. Orders are placed within the time range of 2 years and 1 month (September 4, 2016 October 17, 2018) and customers from 4119 districts and 27 states have placed their order during the given time period.
- 2.Growing trend can be seen in no.of.orders placed over the past years and it's fluctuating if we consider the monthly basis. Brazilian customers mostly placed their orders in the Afternoon.
- 3. The percentage of increase in cost of orders from 2017 to 2018 between Jan to Aug is 136.98%
- 4. AL, MA, SE, ES & BA are the top 5 states where the delivery is fast compared to estimated delivery date.

Avg delivery time is quite high for most of those states from where the company is receiving quite less volume of orders, detailed study is needed further for checking the other reasons behind such a low volume of orders from the majority of states. Huge delivery time can be one of the reasons and we need to work on it.

5. Also only 3 states contribute to maximum volume, and the rest of the states need to be focused on improving the business.

From the analysis, We can see how the orders trajectory is showing a very abrupt increase in orders volume within a very short time. Looking at the overall trend, it is seen that business is picking up very fast in Brazil so the company has to be ready with an extra workforce. To avoid high risk, it can consider hiring contractual employees.