# Olist: An E-Commerce Analysis of the Brazilian Marketplace

Olist is a **Brazilian e-commerce** platform that connects small and medium-sized businesses (SMEs) with customers, allowing merchants to list and sell their products online. The platform operates as a **marketplace**, where merchants can list their products and customers can browse & purchase them online. Olist started its operation in 2015.

This particular business case focuses on analyzing the operations of Olist from **Sept-2016 to Oct-2018** in Brazil and for that we are provided with information of **approx. 100K orders** placed during this period.

The available datasets facilitate the analysis of orders across dimensions such as order status, price, payment and freight charges, customer location, product attributes, and customer reviews.

Via this analysis, we will try to shed light on key aspects of the business, including *sales distribution, revenue generation, order basket diversification, customer demographics and preferences, shipping efficiency, and customer satisfaction levels.* Through a comprehensive examination of these business aspects, we will try to derive **valuable insights** about the company's operations and performance.

```
17
18  # approx. 100K orders
19 ∨ select
20    count(order_id) as num_of_orders
21  from `olist_business_case.orders`;
22
```
✅ This script will process 827.62 MB when run.

## Query results

| Job information | **Results** | Chart |
| --- | --- | --- |

| Row | num_of_orders ▾ | |
| --- | --- | --- |
| 1 | 99441 | |

```
25  # Orders placed between Sept-2016 to Oct-2018
26 ∨ select
27    min(order_purchase_timestamp) as start_point,
28    max(order_purchase_timestamp) as end_point
29  from `olist_business_case.orders`;
30
```
✅ This script will process 827.62 MB when run.

## Query results

| Job information | **Results** | Chart | JSON | Executio |
| --- | --- | --- | --- | --- |

| Row | start_point ▾ | end_point ▾ |
| --- | --- | --- |
| 1 | 2016-09-04 21:15:19 UTC | 2018-10-17 17:30:18 UTC |

This report is organized into **two primary sections** –

- **Section 1:** Initial Data Assessment and Familiarization *(Refer to Pages 2-14)*

- **Section 2:** Deep Dive into Metrics for Business Insights *(Refer to Pages 15-37)*

# Section 1 – Initial Data Assessment and Familiarization

In this we will examine the **schemas** of all the provided tables and also figure out how these tables are related to each other i.e., the common keys. Simultaneously, we will also interpret the **business meaning** of each column. Also, we will validate the data i.e., assess the quality of data provided to us; for example, looking at **missing values**, check for **data integrity**.

```
33   # A total of 8 tables exist in the database.
34   select
35     table_name
36   from olist_business_case.INFORMATION_SCHEMA.TABLES;
37
```

This script will process 827.62 MB when run.

Query results

| Job information | Results | Chart | JSON | Ex |
|---|---|---|---|---|

| Row | table_name ▼ |
|---|---|
| 1 | order_items |
| 2 | sellers |
| 3 | geolocation |
| 4 | reviews |
| 5 | products |
| 6 | orders |
| 7 | payments |
| 8 | customers |

As shown above, we have total **8 tables/datasets** provided to us –

***orders, order_items, products, customers, sellers, geolocation, payments, reviews***

Let's skim through each table one by one –

## customers –

This table has information about customers and their locations.

```
54  select
55    column_name,
56    data_type
57  from olist_business_case.INFORMATION_SCHEMA.COLUMNS
58  where table_name = 'customers';
59
```
✅ This script will process 827.62 MB when run.

### Query results

| Job information | Results | Chart | JSON | Executic |

| Row | column_name ▾ | data_type ▾ |
|-----|---------------|-------------|
| 1 | customer_id | STRING |
| 2 | customer_unique_id | STRING |
| 3 | customer_zip_code_prefix | INT64 |
| 4 | customer_city | STRING |
| 5 | customer_state | STRING |

- Each *'customer_id'* represents a unique order of the *'orders'* table as the company's database is designed in such a way that whenever an order gets placed, a new *'customer_id'* gets created for that customer. It means that the same customer will get different *'customer_id'* for different orders they place.
- In that case what is the unique identifier for each customer - There is a column named *'customer_unique_id'* in the *'customers'* table that assigns one single unique id to each customer. So, if a customer has placed N orders, then they have N different *'customer_id'* but one single *'customer_unique_id'*.
- But what can be the intuition behind so i.e., assigning a new *'customer_id'* to the same customer for new every order they place – Maybe so that we can use the location information of the *'customers'* table as the delivery location for that order.

```
80  # As explained, the 'customer_id' column has all unique values whereas,
81  # the 'customer_unique_id' column has duplicate values in it.
82  select
83    count(customer_id),
84    count(distinct customer_id),
85    count(customer_unique_id),
86    count(distinct customer_unique_id)
87  from `olist_business_case.customers`;
oo
```
✅ This script will process 827.62 MB when run.

### Query results

| Job information | Results | Chart | JSON | Execution details | E |

| Row | f0_ ▾ | f1_ ▾ | f2_ ▾ | f3_ ▾ |
|-----|-------|-------|-------|-------|
| 1 | 99441 | 99441 | 99441 | 96096 |

- The *'customer_zip_code_prefix'* column represents the first five digits of the location zip code.
- The *'customer_city'* column represents the name of the city where the customer is located or from where the order is placed.
- The *'customer_state'* column represents the state code where the customer is located or from where the order is placed. For example, RJ represents Rio de Janeiro.

```
128   -- The customers are located across 27 different states of Brazil
129   select count(distinct customer_state) as num_of_states
130   from `olist_business_case.customers`;
131
```
✔ This script will process 827.62 MB when run.

### Query results

| Job information | Results | Chart | JSON | Execution details |
|---|---|---|---|---|

| Row | num_of_states ▾ |
|---|---|
| 1 | 27 |

- None of the columns in the *'customers'* table has missing values (NULL values) in it.

---

## sellers –

This table has information about the sellers who are listed on Olist along with their locations.

```
169   select
170     column_name,
171     data_type
172   from olist_business_case.INFORMATION_SCHEMA.COLUMNS
173   where table_name = 'sellers';
174
```
✔ This script will process 827.62 MB when run.

### Query results

| Job information | Results | Chart | JSON | Executio |
|---|---|---|---|---|

| Row | column_name ▾ | data_type ▾ |
|---|---|---|
| 1 | seller_id | STRING |
| 2 | seller_zip_code_prefix | INT64 |
| 3 | seller_city | STRING |
| 4 | seller_state | STRING |

- Each *'seller_id'* represents a unique seller listed on the platform.
- *'seller_zip_code_prefix'* column represents the first five digits of the location zip code.

- *'seller_city'* column represents the name of the city where the seller is located.
- *'seller_state'* column represents the state code where the seller is located. For example, RJ represents Rio de Janeiro.

```
206  -- The sellers are located across 23 different states of Brazil.
207  select count(distinct seller_state) as num_of_states
208  from `olist_business_case.sellers`;
209
```
This script will process 827.62 MB when run.

Query results

| Job information | Results | Chart | JSON | Execution details |

| Row | num_of_states ▼ |
|-----|-----------------|
| 1 | 23 |

- None of the columns in the *'sellers'* table has missing values (NULL values) in it.

---

# geolocation –

This table has information of Brazilian zip codes and their latitude/longitude coordinates along with the city and state names.

```
244  # This table has information of Brazilian zip codes and their
245  # latitude/longitude coordinates along with the city and state names.
246
247  select
248    column_name,
249    data_type
250  from olist_business_case.INFORMATION_SCHEMA.COLUMNS
251  where table_name = 'geolocation';
```
This script will process 827.62 MB when run.

Query results

| Job information | Results | Chart | JSON | Execution details |

| Row | column_name ▼ | data_type ▼ |
|-----|---------------|-------------|
| 1 | geolocation_zip_code_prefix | INT64 |
| 2 | geolocation_lat | FLOAT64 |
| 3 | geolocation_lng | FLOAT64 |
| 4 | geolocation_city | STRING |
| 5 | geolocation_state | STRING |

- None of the columns in the *'geolocation'* table has missing values (NULL values) in it.

## orders –

This table has information about all the orders placed at Olist from Sept-2016 to Oct-2018.

```
274  select
275    column_name,
276    data_type
277  from olist_business_case.INFORMATION_SCHEMA.COLUMNS
278  where table_name = 'orders';
270
```
✅ This script will process 827.62 MB when run.

### Query results

| Job information | Results | Chart | JSON | Executio |
|---|---|---|---|---|

| Row | column_name ▼ | data_type ▼ |
|---|---|---|
| 1 | order_id | STRING |
| 2 | customer_id | STRING |
| 3 | order_status | STRING |
| 4 | order_purchase_timestamp | TIMESTAMP |
| 5 | order_approved_at | TIMESTAMP |
| 6 | order_delivered_carrier_date | TIMESTAMP |
| 7 | order_delivered_customer_date | TIMESTAMP |
| 8 | order_estimated_delivery_date | TIMESTAMP |

- *'order_id'* represents the unique identifier assigned to each order.
- *'customer_id'* represents the customer who placed that order.
- *'order_status'* represents the current status of the orders – delivered, shipped, cancelled, etc.
- *'order_purchase_timestamp'* represents the timestamp when the order was placed.
- *'order_approved_at'* represents the payment approval timestamp.
- *'order_delivered_carrier_date'* represents the timestamp when the order was handled to the logistic partner.
- *'order_delivered_customer_date'* represents the timestamp when the order was actually delivered to the customer.
- *'order_estimated_delivery_date'* represents the estimated delivery date that was informed to the customers when they placed the order.

Now, let's try to assess the missing values (NULL values) present in the *'orders'* table –

- The following columns – *'order_id'*, *'customer_id'*, *'order_status'*, *'order_purchase_timestamp'*, *'order_estimated_delivery_date'* do not have any missing values.

- The *'order_approved_at'* field is missing for 160 orders out of a total of approx. 100K orders; however, corresponding payment records exist for all these orders in the 'payments' table signifying that payment has been done for all these orders. Therefore, we can safely impute these missing values (NULL values). In order to impute these missing values, we can use the average/median approval time of all the orders.

```
319  # 160 orders out of ~100K orders have
320  # a missing "payment approval timestamp"
321  select count(*) as missing_value_count
322  from `olist_business_case.orders`
323  where order_approved_at IS NULL;
```

✅ This script will process 817.13 MB when run.

Query results

| Job information | Results | Chart | JSO |

| Row | missing_value_count |
| --- | --- |
| 1 | 160 |

```
329  # Get me the 'order_id' for which payment approval timestamp
330  # is missing and payment indeed has not been done.
331  select order_id
332  from `olist_business_case.orders`
333  where order_approved_at IS NULL
334  EXCEPT DISTINCT
335  select distinct order_id
336  from `olist_business_case.payments`;
337
```

✅ This script will process 817.13 MB when run.

Query results

| Job information | Results | Chart | JSON | Execution de |

ℹ️ There is no data to display.

- Out of approx. 100K orders placed at Olist, 1783 orders (i.e., 2%) have a missing information for when it was handled to the logistics partner for delivery and 2965 orders (i.e., 3%) have a missing information for when it was delivered to the customer.

```
347  # 1783 orders out of ~100K orders i.e., ~2% of total orders
348  # have a missing information for when it was handled to the
349  # logistics partner for delivery.
350  select count(order_id) as num_of_orders
351  from `olist_business_case.orders`
352  where order_delivered_carrier_date IS NULL;
353
```

✅ This script will process 825.8 MB when run.

## Query results

| Job information | Results | Chart | JSON | Execution de |
|---|---|---|---|---|

| Row | num_of_orders ▾ | | |
|---|---|---|---|
| 1 | 1783 | | |

```
355  # 2965 orders out of ~100K orders i.e., ~3% of total orders
356  # have a missing information for when it was delivered to
357  # the customer.
358  select count(order_id) as num_of_orders
359  from `olist_business_case.orders`
360  where order_delivered_customer_date IS NULL;
```

✅ This script will process 825.8 MB when run.

## Query results

| Job information | Results | Chart | JSON | Execution d |
|---|---|---|---|---|

| Row | num_of_orders ▾ | | |
|---|---|---|---|
| 1 | 2965 | | |

Now, let's further inspect these missing values for the *'order_delivered_carrier_date'* and *'order_delivered_customer_date'* fields so that we can think of an imputation logic for these fields.

- For the following *'order_status'* – 'created', 'approved', 'unavailable', 'processing', and 'invoiced', both the *'order_delivered_carrier_date'* and *'order_delivered_customer_date'* fields should be NULL as these orders have not yet reached to the shipment stage. Upon checking the data, it is indeed the case as well.

```
351   select distinct
352     order_delivered_carrier_date,
353     order_delivered_customer_date
354   from `olist_business_case.orders`
355   where order_status IN ('created', 'approved', 'unavailable', 'processing', 'invoiced');
356
```
✅ This script will process 817.87 MB when run.

## Query results

| Job information | **Results** | Chart | JSON | Execution details | Execution graph |
|---|---|---|---|---|---|

| Row | order_delivered_carrier_date ▾ | order_delivered_customer_date ▾ | |
|---|---|---|---|
| 1 | *null* | *null* | |

- For the orders with *'order_status'* – 'shipped', the field *'order_delivered_carrier_date'* should not be NULL and the field *'order_delivered_customer_date'* should be NULL. Upon checking the data, it is indeed the case as well.

- For the orders with *'order_status'* – 'cancelled', the fields *'order_delivered_carrier_date'* and *'order_delivered_customer_date'* may or may not be NULL as cancellation can be done at any stage.

- For the orders with *'order_status'* – 'delivered', both the fields *'order_delivered_carrier_date'* and *'order_delivered_customer_date'* should not be NULL. Upon checking the data, it is indeed the case as well (ignoring few exceptions – 8 orders as that can be a human error).

So, from the above inspection, we can comment that NULL values appearing in both the fields are completely aligned with the business logic and so, these NULL values should not be interpreted as missing values. However, if we still need to impute these NULL values, then we can do it using the average/median shipping and delivery time respectively.

- Data Integrity validation – All the customers present in the *'orders'* table are available in the parent table *'customers'*, ensuring referential integrity check.

- Anomalies detected in Delivery timestamps – Among the approximately 100K orders, 23 records present a customer delivery date earlier than the carrier handover date, while 166 orders show a carrier handover date preceding the purchase timestamp. These nonsensical timestamps require either manual correction or omission from delivery time analysis.

# Order_items –

This table contains the item-level information of each order (basically, multiple items can be purchased within the same order).



- *'order_id'* represents the unique identifier of each order.
- *'order_item_id'* is the serial number assigned to items within an order.

It is important to note the absence of a dedicated quantity field in the table. Consequently, each unit of a product within an order is stored as a separate row in the table resulting in duplicate *'order_id'* entries.

- *'product_id'* represents the unique identifier of the products included in the order.
- *'seller_id'* represents the unique identifier of the seller who will be supplying that product.

Note that different products within an order might be fulfilled by different sellers.

- *'shipping_limit_date'* represents the seller shipping limit date i.e., the deadline for the seller to hand over the order to the logistics partner.
- *'price'* represents the actual price of the associated product
- *'freight_value'* represents the delivery associated costs

The total price for an item would be represented by *'price'* + *'freight_value'*.

- Data Integrity validation – All the orders, sellers and products present in the *'order_items'* table are available in their parent tables *'orders'*, *'sellers'* and *'products'* respectively, ensuring referential integrity.

- None of the columns in the *'order_items'* table has missing values (NULL values) in it.

## products –

This table contains details of all the products listed on Olist.

```
555  select
556    column_name,
557    data_type
558  from olist_business_case.INFORMATION_SCHEMA.COLUMNS
559  where table_name = 'products';
560
```
✓ This script will process 825.8 MB when run.

### Query results

| Job information | Results | Chart | JSON | Executio |
|---|---|---|---|---|

| Row | column_name ▾ | data_type ▾ |
|---|---|---|
| 1 | product_id | STRING |
| 2 | product_category_name | STRING |
| 3 | product_name_length | INT64 |
| 4 | product_description_length | INT64 |
| 5 | product_photos_qty | INT64 |
| 6 | product_weight_g | INT64 |
| 7 | product_length_cm | INT64 |
| 8 | product_height_cm | INT64 |
| 9 | product_width_cm | INT64 |

- *'product_id'* represents the unique identifier for each product.
- *'product_category_name'* represents the product category that the product belongs to.
- Out of total 32951 products available in the table, 610 products (i.e., approx. 2%) have *'product_category_name'* missing (i.e., NULL). These missing values are replaced with the string value 'unknown'.
- There are a total of 74 distinct product categories present in the data.

```
1698  # Total 74 product categories -
1699  select
1700    count(distinct product_category_name) as num_of_product_categories
1701  from `olist_business_case.products`;
1702
```
✓ This script will process 839.17 MB when run.

### Query results

| Job information | Results | Chart | JSON | Execution details |
|---|---|---|---|---|

| Row | num_of_product_categories |
|---|---|
| 1 | 74 |

- There is information available about some other attributes of the product as well like dimensions of the product, weight of the product, etc.

## payments –

This table contains details of the payments made for the orders.

```
598  select
599    column_name,
600    data_type
601  from olist_business_case.INFORMATION_SCHEMA.COLUMNS
602  where table_name = 'payments';
```

✅ Query completed

### Query results

| Job information | Results | Chart | JSON | Executic |
|---|---|---|---|---|

| Row | column_name ▾ | data_type ▾ |
|---|---|---|
| 1 | order_id | STRING |
| 2 | payment_sequential | INT64 |
| 3 | payment_type | STRING |
| 4 | payment_installments | INT64 |
| 5 | payment_value | FLOAT64 |

- *'order_id'* represents the unique identifier of the order.
- *'payment_sequential'* – A customer may pay for an order with more than one payment method (for example, credit card, voucher, debit card, etc.). If they do so, a sequence will be created to accommodate all the payments.
- *'payment_type'* represents the method of payment chosen by the customer. The different payment methods available in the data are – credit card, boleto, voucher, debit card.

  Note: "Boleto" refers to a bank slip or voucher used for payments, particularly in Brazil. It is often used to make payments for goods, services, or bills, especially for those without credit cards or online banking access.

- *'payment_value'* represents the transaction value.
- *'payment_installments'* represents the number of installments chosen by the customer to fulfill the payment.
  Note: From the data, we can observe that payment via credit card is the sole method permitting transactions to be completed in multiple installments. All other payment methods are getting processed as a single, lump-sum transaction.

```
596  # Payment via credit card is the sole method permitting transactions to be completed in multiple installments.
597  # All other payment methods are processed as a single, lump-sum transaction.
598  select
599    payment_type,
600    max(payment_installments) as max_num_of_installments
601  from `olist_business_case.payments`
602  group by payment_type;
603
```

✅ This script will process 814.56 MB when run.

## Query results

| Job information | **Results** | Chart | JSON | Execution details | Execution graph |

| Row | payment_type ▾ | max_num_of_installn |
|---|---|---|
| 1 | credit_card | 24 |
| 2 | voucher | 1 |
| 3 | not_defined | 1 |
| 4 | boleto | 1 |
| 5 | debit_card | 1 |

- Data Integrity validation – All the orders present in the *'payments'* table are available in its parent table *'orders'*, ensuring referential integrity.
- None of the columns in the *'payments'* table has missing values (NULL values) in it.

---

# reviews –

This table contains information about the reviews given by customers for their overall purchasing experience. Once the customer receives the product or the estimated delivery date is due, the customer gets a satisfaction survey by email where they can give a **review score (1 to 5)** for the purchase experience and write down some comments as well.

```
658  select
659    column_name,
660    data_type
661  from olist_business_case.INFORMATION_SCHEMA.COLUMNS
662  where table_name = 'reviews';
```

✅ Query completed

## Query results

| Job information | **Results** | Chart | JSON | Executic |

| Row | column_name ▾ | data_type ▾ |
|---|---|---|
| 1 | review_id | STRING |
| 2 | order_id | STRING |
| 3 | review_score | INT64 |
| 4 | review_comment_title | STRING |
| 5 | review_creation_date | TIMESTAMP |
| 6 | review_answer_timestamp | TIMESTAMP |

- *'review_id'* represents the identifier for the review.
- *'order_id'* represents the order associated with the review.

Multiple *'review_id'* values exist for the same *'order_id'* in the table. It may be because customer is giving more than one review for the same order.

- *'review_score'* represents a 1 to 5 integer score given by the customer for the overall purchase experience.
- *'review_comment_title'* represents the title of the review posted by customer.
- *'review_creation_date'* represents the date on which satisfaction survey was sent to the customer.
- *'review_answer_timestamp'* represents the timestamp at which the satisfaction survey was answered by the customer.

- None of the fields in the *'reviews'* table except *'review_comment_title'* has missing values (NULL values) in it. And the NULL values in the *'review_comment_title'* field should not be interpreted as missing data as customers may not always be interested in leaving a comment.

- Data Integrity validation – All the orders present in the *'reviews'* table are available in its parent table *'orders'*, ensuring referential integrity. Also, out of approx. 100K orders, 768 do not have corresponding customer reviews which is an expected behavior as providing feedback/review is not a mandatory requirement.

# Section 2 – Deep Dive into Metrics for Business Insights

Order fulfillment rate –

```
829  select
830    order_status,
831    num_of_orders,
832    CONCAT(CAST(ROUND(num_of_orders/(SUM(num_of_orders) OVER())*100, 2) AS STRING), ' %') as percent_of_orders
833  from (
834      select
835        order_status,
836        count(order_id) as num_of_orders
837      from `olist_business_case.orders`
838      group by order_status
839    ) as tbl
840  order by num_of_orders desc;
```

✅ This script will process 823.42 MB when run.

## Query results

| Job information | **Results** | Chart | JSON | Execution details | Execution graph |
|---|---|---|---|---|---|

| Row | order_status ▼ | num_of_orders ▼ | percent_of_orders ▼ | |
|---|---|---|---|---|
| 1 | delivered | 96478 | 97.02 % | |
| 2 | shipped | 1107 | 1.11 % | |
| 3 | canceled | 625 | 0.63 % | |
| 4 | unavailable | 609 | 0.61 % | |
| 5 | invoiced | 314 | 0.32 % | |
| 6 | processing | 301 | 0.3 % | |
| 7 | created | 5 | 0.01 % | |
| 8 | approved | 2 | 0 % | |

- The overwhelming majority of orders (**>97%**) reached the *'delivered'* status, indicating a **highly successful order fulfillment process** at Olist. Achieving such a high delivery rate **builds trust and satisfaction among customers**, encouraging them to return to the platform and recommend it to others (**word-of-mouth referrals**).
- Cancelled orders represent just 0.63% of the total, indicating **low order cancellation rates**.
- Overall, the data **reflects a mature and optimized e-commerce operation** with minimal friction in order processing and fulfillment.

Let's analyze the diversification of products in the order basket.

```
850  select
851    Basket_Type,
852    num_of_orders,
853    CONCAT(CAST(ROUND(num_of_orders/(SUM(num_of_orders) OVER())*100, 2) AS STRING), ' %') as percentage_of_orders
854  from (
855        select
856          'Single-item orders' as Basket_Type,
857          count(*) as num_of_orders
858        from (
859            select
860              order_id
861            from `olist_business_case.order_items`
862            group by order_id
863            having count(order_item_id) = 1
864        UNION ALL
865        select
866          'Multiple-item orders' as Basket_Type,
867          count(*) as num_of_orders
868        from (
869            select
870              order_id
871            from `olist_business_case.order_items`
872            group by order_id
873            having count(order_item_id) > 1
874      )
875  order by CASE
876            WHEN Basket_Type = 'Single-item orders' THEN 1
877            ELSE 2
878          END;
```

## Query results

| Job information | **Results** | Chart | JSON | Execution details | Exe |
|---|---|---|---|---|---|

| Row | Basket_Type ▼ | num_of_orders ▼ | percentage_of_orders ▼ |
|---|---|---|---|
| 1 | Single-item orders | 88863 | 90.06 % |
| 2 | Multiple-item orders | 9803 | 9.94 % |

- We can observe an **overwhelming majority (90%) of single-item orders** meaning that most transactions involve the purchase of just one item.
- Only 10% of orders include more than one item. This indicates **limited cross-selling or bundling activity.**
- This stark contrast between single-item orders and multiple-item orders highlights a significant opportunity to **encourage customers to purchase more items per order.**

- Olist can offer **Volume discounts** (for example, flat discounts on higher quantities or offers like Buy x Get y Free), provide **"frequently bought together" or "You may also like" product suggestions**, etc. to encourage customers to add more items to their basket.

- Also, Olist can offer **free delivery beyond a minimum order value** to incentivize customers to add more items to their basket.

- Additionally, Olist can flash "**Limited-time offers**" on related items to **create urgency** and encourage multi-unit purchases.

Let's try to investigate if price can be the reason behind these single-item orders as in are expensive products leading to single-item orders?

For this, we will examine if the prices involved in these single-item orders are significantly higher than the prices involved in the multiple-item orders.

```sql
893  select
894    Basket_Type,
895    num_of_orders,
896    CONCAT(CAST(ROUND(num_of_orders/(SUM(num_of_orders) OVER())*100, 2) AS STRING), ' %') as percentage_of_orders,
897    median_price,
898    mean_price
899  from (
900        select distinct
901          'Single-item orders' as Basket_Type,
902          count(order_id) over() as num_of_orders,
903          ROUND(percentile_cont(price+freight_value, 0.5) OVER(), 2) as median_price,
904          ROUND(avg(price+freight_value) OVER(), 2) as mean_price
905        from `olist_business_case.order_items`
906        where order_id IN (select
907                            order_id
908                            from `olist_business_case.order_items`
909                            group by order_id
910                            having count(order_item_id) = 1)
911      UNION ALL
912
913        select distinct
914          'Multiple-item orders' as Basket_Type,
915          count(distinct order_id) over() as num_of_orders,
916          ROUND(percentile_cont(price+freight_value, 0.5) OVER(), 2) as median_price,
917          ROUND(avg(price+freight_value) OVER(), 2) as mean_price
918        from `olist_business_case.order_items`
919        where order_id IN (select
920                            order_id
921                            from `olist_business_case.order_items`
922                            group by order_id
923                            having count(order_item_id) > 1)
924      )
925  order by CASE
926              WHEN Basket_Type = 'Single-item orders' THEN 1
927              ELSE 2
928          END;
```

## Query results

| Job information | **Results** | Chart | JSON | Execution details | Execution graph |
|---|---|---|---|---|---|

| Row | Basket_Type ▼ | num_of_orders ▼ | percentage_of_orders ▼ | median_price ▼ | mean_price ▼ |
|---|---|---|---|---|---|
| 1 | Single-item orders | 88863 | 90.06 % | 99.03 | 150.75 |
| 2 | Multiple-item orders | 9803 | 9.94 % | 73.34 | 102.91 |

- Although, median prices involved in the single-item orders is higher than the multiple-item orders, the difference is not large. So, we **can't attribute prices to be the reason behind single-item orders**.

Since price does not appear to be the main factor behind the high proportion of single-item orders, this trend may be attributed to a large number of first-time or one-time customers who have not yet developed enough trust in the platform to place large orders. To examine this, let's calculate the percentage of repeat customers on the platform –

```sql
1025  select
1026    CONCAT(CAST(ROUND(countif(num_of_orders = 1)/count(*)*100, 2) AS STRING), ' %') as percent_of_one_time_customers
1027  from (
1028      select
1029        c.customer_unique_id,
1030        count(o.order_id) as num_of_orders
1031      from `olist_business_case.orders` as o
1032      INNER JOIN `olist_business_case.customers` as c
1033      ON o.customer_id = c.customer_id
1034      group by c.customer_unique_id
1035    ) as tbl;
```

✅ Query completed

## Query results

| Job information | **Results** | Chart | JSON | Execution details | Execution graph |

| Row | percent_of_one_time_customers |
| --- | --- |
| 1 | 96.88 % |

- As anticipated, approximately **97% of the customers on the platform are first-time or one-time buyers**, and **only about 3% of the customers making repeat purchases**. This underscores a significant opportunity to **enhance customer retention** through targeted strategies such as **loyalty programs** to reward repeat customers and **post-purchase engagement**.

- In the highly competitive e-commerce market, **acquiring new customers** often involves substantial marketing and promotional expenses, resulting in a **high Customer Acquisition Cost (CAC)**. Therefore, **repeat customers** are crucial for **sustaining long-term business growth** and **profitability** as they are more likely to make additional purchases and require less investment to re-engage.

- Given that Olist was founded in 2015 and the available data spans from September 2016 to October 2018, it is reasonable to observe a high proportion of first-time customers during this growth phase. Therefore, to accurately assess the effectiveness of retention strategies, it is essential to analyze customer data over a longer period.

- Moreover, later in this report, we will also assess the customer review scores to evaluate the overall purchasing experience. Particularly, we will examine whether a high proportion of low review scores is present, as this could be a significant factor behind the low rate of repeat purchases.

Note: For analyzing the sales and revenue figures, let's not include the orders which have been *'cancelled'*.

Let's firstly calculate the total sales (number of orders) made by Olist and then analyze how it has changed over time (i.e., periodical trend) –

```
1377  # Olist fulfilled nearly 99K orders over a time period of 2 years (i.e., Sept-2016 to Oct-2018).
1378  select count(order_id) as num_of_orders
1379  from `olist_business_case.orders`
1380  where order_status != 'canceled';
```

✅ This script will process 823.42 MB when run.

## Query results

| Job information | Results | Chart | JSON | Execution details | Execution graph |
|---|---|---|---|---|---|

| Row | num_of_orders ▼ | |
|---|---|---|
| 1 | 98816 | |

Month-on-month change in sales volume (number of orders) –

```
1386  select
1387    time_period,
1388    current_month_sales,
1389    CONCAT(CAST(ROUND((current_month_sales - prev_month_sales)/prev_month_sales*100, 2) AS STRING), ' %') as percentage_change
1390  from (
1391      select
1392        time_period,
1393        num_of_orders as current_month_sales,
1394        LAG(num_of_orders, 1) OVER(ORDER BY SUBSTR(time_period, -1, 4), SUBSTR(time_period, 1, 2)) as prev_month_sales
1395      from (
1396          select
1397            FORMAT_TIMESTAMP('%m-%Y', order_purchase_timestamp) as time_period,
1398            count(order_id) as num_of_orders
1399          from `olist_business_case.orders`
1400          where order_status != 'canceled'
1401          group by FORMAT_TIMESTAMP('%m-%Y', order_purchase_timestamp)
1402      ) as tbl
1403  ) as tbl2
1404  order by SUBSTR(time_period, -1, 4), SUBSTR(time_period, 1, 2);
```

## Query results

| | Job information | Results | Chart | JSON | Execution details | Exe |
|---|---|---|---|---|---|---|

| Row | time_period ▾ | current_month_sales | percentage_change ▾ |
|---|---|---|---|
| 1 | 09-2016 | 2 | *null* |
| 2 | 10-2016 | 300 | 14900 % |
| 3 | 12-2016 | 1 | -99.67 % |
| 4 | 01-2017 | 797 | 79600 % |
| 5 | 02-2017 | 1763 | 121.2 % |
| 6 | 03-2017 | 2649 | 50.26 % |
| 7 | 04-2017 | 2386 | -9.93 % |
| 8 | 05-2017 | 3671 | 53.86 % |
| 9 | 06-2017 | 3229 | -12.04 % |
| 10 | 07-2017 | 3998 | 23.82 % |
| 11 | 08-2017 | 4304 | 7.65 % |
| 12 | 09-2017 | 4265 | -0.91 % |
| 13 | 10-2017 | 4605 | 7.97 % |
| 14 | 11-2017 | 7507 | 63.02 % |
| 15 | 12-2017 | 5662 | -24.58 % |
| 16 | 01-2018 | 7235 | 27.78 % |
| 17 | 02-2018 | 6655 | -8.02 % |
| 18 | 03-2018 | 7185 | 7.96 % |
| 19 | 04-2018 | 6924 | -3.63 % |
| 20 | 05-2018 | 6849 | -1.08 % |
| 21 | 06-2018 | 6149 | -10.22 % |
| 22 | 07-2018 | 6251 | 1.66 % |
| 23 | 08-2018 | 6428 | 2.83 % |



sales (number of orders)

- The **sales trend analysis** reveals a generally **consistent upward trajectory**, with only minor dips that are outweighed by periods of growth. Given the low percentage of repeat customers, the increasing sales figures suggest that Olist is successfully attracting a **growing number of new customers**. These trends are encouraging indicators of strong market traction and business development during the company's early years.

Now, let's confirm if the revenue growth is also following a similar trajectory as the sales (number of orders). Basically, we will calculate the total revenue generated and then analyze how it has changed over time (i.e., periodical trend) –

```
1420    # Olist generated a total revenue of approx. $16M (or brazilian currency)
1421    # over a time period of 2 years (i.e., Sept 2016 to Oct 2018).
1422    select ROUND(sum(payment_value), 2) as total_revenue
1423    from `olist_business_case.payments`
1424    where order_id IN (select
1425                            order_id
1426                       from `olist_business_case.orders`
1427                       where order_status != 'canceled'
1428                   );
1429
```
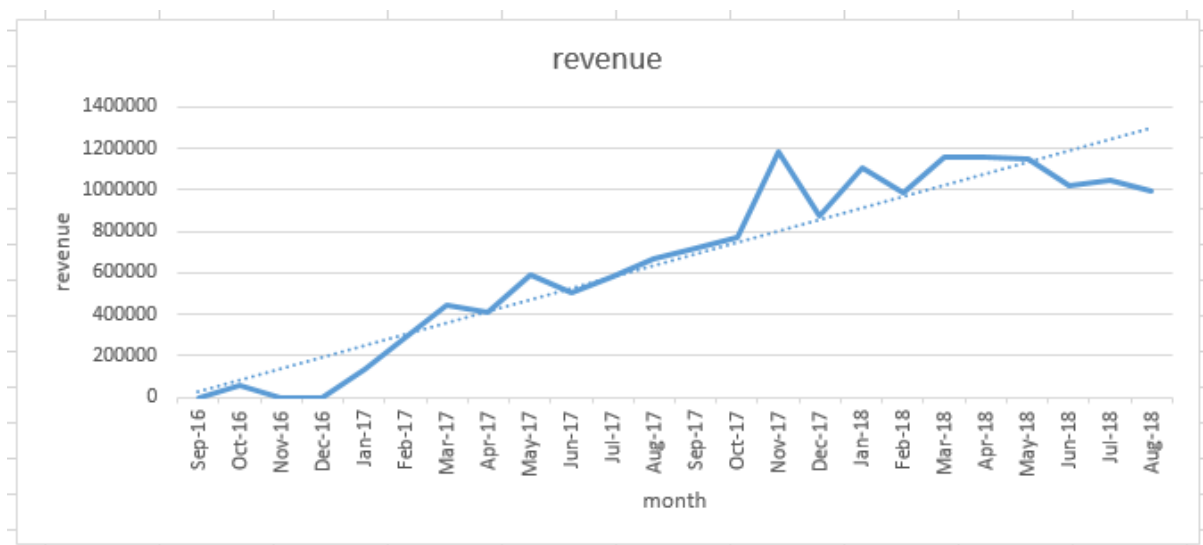⊘ This script will process 823.42 MB when run.

## Query results

| Job information | Results | Chart | JSON | Execution details | Exe |
|---|---|---|---|---|---|

| Row | total_revenue ▼ | |
|---|---|---|
| 1 | 15865616.52 | |

Month-on-month change in revenue –



- The **revenue trend analysis** closely mirrors the positive trajectory observed in sales.
- This **upward movement in revenue**, despite occasional minor declines, **reflects sustained business growth.**
- The **consistent rise in revenue, alongside growing sales volumes and an expanding customer base** during the early years, indicates that Olist is **successfully scaling its operations and strengthening its presence in the e-commerce sector**.

Note: It would be insightful to analyze **monthly seasonality** in sales or revenue numbers to determine if certain months experience significant peaks or dips. However, the current dataset of just **24 months is insufficient** to make any reliable conclusions about any such seasonal trends.

Next, let's examine the **distribution of sales across different states and cities** to identify the top and bottom performing regions –

```
1328  select
1329    *,
1330    CONCAT(CAST(ROUND(num_of_orders/sum(num_of_orders) OVER()*100, 2) AS STRING), ' %') as percent_of_orders
1331  from
1332      (select
1333        c.customer_state,
1334        count(o.order_id) as num_of_orders
1335      from `olist_business_case.orders` as o
1336      INNER JOIN `olist_business_case.customers` as c
1337      ON o.customer_id = c.customer_id
1338      where o.order_status != 'canceled'
1339      group by c.customer_state
1340      )
1341  order by num_of_orders desc;
```

Top 10 States by Sales Volume (Number of Orders) –

| Row | customer_state | num_of_orders | percent_of_orders |
|-----|----------------|---------------|-------------------|
| 1 | SP | 41419 | 41.92 % |
| 2 | RJ | 12766 | 12.92 % |
| 3 | MG | 11571 | 11.71 % |
| 4 | RS | 5441 | 5.51 % |
| 5 | PR | 5023 | 5.08 % |
| 6 | SC | 3618 | 3.66 % |
| 7 | BA | 3364 | 3.4 % |
| 8 | DF | 2133 | 2.16 % |
| 9 | ES | 2024 | 2.05 % |
| 10 | GO | 2007 | 2.03 % |

- Nearly **42%** of all orders originate from the state of **Sao Paulo (SP)** alone, highlighting its significant contribution to Olist's overall sales. The next largest shares come from **Rio de Janeiro (RJ)** and **Minas Gerais (MG)**, accounting for **13%** and **12%** of orders, respectively. The distribution reflects the demographic and economic realities of Brazil as these three states are not only the **most populous** but also the **most economically active**, together representing a substantial portion of the country's GDP.

- The **top 3 states account for 67% of total orders**, while the top 10 states contribute 90%, indicating a **strong concentration of sales within a limited number of regions**. This high concentration suggests significant market penetration in these states, but also highlights opportunities for growth in less represented regions.

- Also, note that high sales figures in these leading states mainly **represent the concentration of a substantial customer base in these regions**, as 97% of the customers on the platform are one-time purchasers.

Bottom 10 States by Sales Volume (Number of Orders) –

| 18 | PI | 491 | 0.5 % |
|----|----|-----|-------|
| 19 | RN | 485 | 0.49 % |
| 20 | AL | 412 | 0.42 % |
| 21 | SE | 349 | 0.35 % |
| 22 | TO | 279 | 0.28 % |
| 23 | RO | 250 | 0.25 % |
| 24 | AM | 148 | 0.15 % |
| 25 | AC | 81 | 0.08 % |
| 26 | AP | 68 | 0.07 % |
| 27 | RR | 45 | 0.05 % |

- The **bottom 10 states collectively contribute less than 3% to the overall sales volume**, indicating a significantly lower level of market penetration and sales activity in these regions as compared to the leading states.
- The possible reasons behind such **low sales activity in these regions** could be –
  - lower adoption of online shopping habits, particularly among rural populations
  - low digital literacy
  - lower average income levels
  - limited internet access

- Further, such low sales activity makes it cost-ineffective for companies like Olist to serve these regions efficiently.
- These combined **socio-economic and infrastructural barriers** hinder consumer's participation in e-commerce and thereby the company's ability to expand its reach in these states.

Next, we identify the top cities in terms of sales volume. And as expected, the **top-performing cities are the leading states' capitals** – sao paulo (capital of Sao Paulo, SP) leads with 16% of total sales volume, followed by rio de janeiro (capital of Rio de Janeiro, RJ) with 7%, and Belo Horizonte (capital of Minas Gerais, MG) with 3%.

Now, let's examine the **revenue distribution across states** to figure out if they follow the same pattern observed in sales distribution or not –

```sql
1371  select
1372    *,
1373    CONCAT(CAST(ROUND(revenue/sum(revenue) OVER()*100, 2) AS STRING), ' %') as percent_of_revenue
1374  from (
1375        select
1376          c.customer_state,
1377          ROUND(sum(p.payment_value), 2) as revenue
1378        from `olist_business_case.orders` as o
1379        INNER JOIN `olist_business_case.payments` as p
1380        ON o.order_id = p.order_id
1381        INNER JOIN `olist_business_case.customers` as c
1382        ON o.customer_id = c.customer_id
1383        where o.order_status != 'canceled'
1384        group by c.customer_state
1385    ) as tbl
1386  order by revenue desc;
```

Top 10 States by Revenue generated –

| Row | customer_state | revenue | percent_of_revenue |
|---|---|---|---|
| 1 | SP | 5942397.11 | 37.45 % |
| 2 | RJ | 2126444.23 | 13.4 % |
| 3 | MG | 1856375.81 | 11.7 % |
| 4 | RS | 881680.6 | 5.56 % |
| 5 | PR | 802319.18 | 5.06 % |
| 6 | SC | 613707.46 | 3.87 % |
| 7 | BA | 611796.01 | 3.86 % |
| 8 | DF | 352718.04 | 2.22 % |
| 9 | GO | 342124.8 | 2.16 % |
| 10 | ES | 324038.9 | 2.04 % |

- The analysis of revenue distribution across states reveal a pattern closely aligned with sales volume, with **Sao Paulo (SP) contributing 37.45% of total revenue**, followed by **Rio de Janeiro (RJ) at 13.4% and Minas Gerais (MG) at 11.7%.**
- Collectively, the **top 5 states account for 73%** of total revenue, while the top 10 states contribute 87%. And the **bottom 10 states collectively contribute less than 4%** of the total revenue.
- The close alignment between sales and revenue distribution across states suggests that the **average order value (AOV) is largely uniform across states.**

Next, let's examine the **distribution of sales across different product categories** to identify the top and bottom performing categories –

```sql
1224 select
1225   *,
1226   CONCAT(CAST(ROUND(num_of_orders/(sum(num_of_orders) over())*100, 2) AS STRING), ' %') as percent_of_orders
1227 from (
1228     select
1229       p.product_category_name,
1230       count(order_id) as num_of_orders,
1231     from
1232     (select *
1233     from `olist_business_case.order_items`
1234     where order_id IN (select order_id from `olist_business_case.orders` where order_status != 'canceled')) as oi
1235     INNER JOIN `olist_business_case.products` as p
1236     ON oi.product_id = p.product_id
1237     group by p.product_category_name
1238   )
1239 order by num_of_orders desc;
```

Top 10 Product categories by Sales Volume (Number of Orders) –

| Row | product_category_name | num_of_orders | percent_of_orders |
|---|---|---|---|
| 1 | bed table bath | 11097 | 9.9 % |
| 2 | HEALTH BEAUTY | 9634 | 8.59 % |
| 3 | sport leisure | 8590 | 7.66 % |
| 4 | Furniture Decoration | 8298 | 7.4 % |
| 5 | computer accessories | 7781 | 6.94 % |
| 6 | housewares | 6915 | 6.17 % |
| 7 | Watches present | 5970 | 5.33 % |
| 8 | telephony | 4527 | 4.04 % |
| 9 | Garden tools | 4328 | 3.86 % |
| 10 | automotive | 4205 | 3.75 % |

- The **top 10 product categories account for a substantial 64% share of total sales volume**. Many of these categories, such as household goods and health and beauty products represent **everyday essentials, explaining their high demand.**
- Also, the top 20 product categories (out of a total of 74) accounts for 88% of total sales volume.
- These figures indicate that **customer demand is heavily focused on a limited selection of categories.** This insight is valuable for **inventory management**, as it enables the company to strategically prioritize stocking and replenishment for these high-demand categories.

Now that we have identified the product categories with high and low demand, let's examine the **pricing** within these categories. This will provide us valuable insights into the **purchasing capacity/behavior of the customers** –

```
1325  select distinct
1326    p.product_category_name,
1327    ROUND(PERCENTILE_CONT(oi.price, 0.5) OVER(PARTITION BY p.product_category_name), 2) as median_price,
1328    ROUND(AVG(oi.price) OVER(PARTITION BY p.product_category_name), 2) as avg_price,
1329    COUNT(oi.order_id) OVER(PARTITION BY p.product_category_name) as num_of_orders,
1330    ROUND(COUNT(oi.order_id) OVER(PARTITION BY p.product_category_name)/(COUNT(oi.order_id) OVER())*100, 2)
1331    as percent_of_orders
1332  from
1333  (select *
1334  from `olist_business_case.order_items`
1335  where order_id IN (select order_id from `olist_business_case.orders` where order_status != 'canceled')) as oi
1336  INNER JOIN `olist_business_case.products` as p
1337  ON oi.product_id = p.product_id
1338  order by num_of_orders desc;
```

Mostly, **product categories having high Sales Volume have median prices below 100 dollars** (or Brazil currency) –

| Row | product_category_name | median_price | avg_price | num_of_orders | percentage_of_orders |
|---|---|---|---|---|---|
| 1 | bed table bath | 79.9 | 93.36 | 11097 | 9.9 % |
| 2 | HEALTH BEAUTY | 79.9 | 130.34 | 9634 | 8.59 % |
| 3 | sport leisure | 78.0 | 114.06 | 8590 | 7.66 % |
| 4 | Furniture Decoration | 65.49 | 87.67 | 8298 | 7.4 % |
| 5 | computer accessories | 81.99 | 116.22 | 7781 | 6.94 % |
| 6 | housewares | 59.8 | 90.65 | 6915 | 6.17 % |
| 7 | Watches present | 129.0 | 200.7 | 5970 | 5.33 % |
| 8 | telephony | 29.99 | 71.2 | 4527 | 4.04 % |
| 9 | Garden tools | 59.9 | 111.14 | 4328 | 3.86 % |
| 10 | automotive | 84.9 | 139.51 | 4205 | 3.75 % |

Mostly, **product categories having higher price shows negligible sales** –

| Row | product_category_name | median_price | avg_price | num_of_orders | percentage_of_orders |
|---|---|---|---|---|---|
| 1 | PCs | 1100.0 | 1098.34 | 203 | 0.18 % |
| 2 | HOUSE PASTALS OVEN AND C... | 587.0 | 624.29 | 76 | 0.07 % |
| 3 | Agro Industria e Comercio | 258.65 | 342.12 | 212 | 0.19 % |
| 4 | ELECTRICES 2 | 227.99 | 470.85 | 235 | 0.21 % |
| 5 | Furniture | 179.0 | 183.75 | 109 | 0.1 % |
| 6 | Furniture office | 144.99 | 161.88 | 1690 | 1.51 % |
| 7 | insurance and services | 141.64 | 141.65 | 2 | 0 % |
| 8 | climatization | 139.99 | 185.5 | 295 | 0.26 % |
| 9 | La Cuisine | 137.0 | 146.78 | 14 | 0.01 % |
| 10 | Cool Stuff | 129.99 | 164.24 | 3780 | 3.37 % |

To get further insights into the purchasing capacity/behavior of the customers, let's examine the **distribution of sales across different price segments** –

```sql
1731  with master_table as (
1732    select
1733      product_id,
1734      count(order_id) as sales_volume,
1735      avg(price) as price
1736    from `olist_business_case.order_items`
1737    group by product_id
1738  ),
1739
1740  price_categories as (
1741    select
1742      *,
1743      case
1744        when price < 100 then '<100'
1745        when price between 100 and 200 then '100-200'
1746        when price between 200 and 500 then '200-500'
1747        else '>500'
1748      end as price_range
1749    from master_table
1750  )
1751
1752  select distinct
1753    price_range,
1754    sum(sales_volume) over(partition by price_range) as total_sales_volume,
1755    CONCAT(CAST(ROUND(100*sum(sales_volume) over(partition by price_range)/(sum(sales_volume) over()), 2) AS STRING), ' %')
1756    as percent_of_sales_volume
1757  from price_categories
1758  order by case
1759          when price_range = '<100' then 1
1760          when price_range = '100-200' then 2
1761          when price_range = '200-500' then 3
1762          else 4
1763        end;
```

| Row | price_range | total_sales_volume | percent_of_sales_volume |
|---|---|---|---|
| 1 | <100 | 70914 | 62.95 % |
| 2 | 100-200 | 27958 | 24.82 % |
| 3 | 200-500 | 10502 | 9.32 % |
| 4 | >500 | 3276 | 2.91 % |

- The **vast majority of orders (nearly 63%) come from products priced under 100 dollars** (or Brazil currency). Products priced between 100 and 200 dollars account for about 25% of total sales volume. So, **overall, 88% of total sales volume is coming from products priced below 200 dollars**.

- Almost all **high-volume categories** that we identified before, such as household goods and health and beauty products, **fall into the lower price range of under 100 dollars**.

- These observations highlight **largely price-sensitive and utility-driven shopping behavior of customers**, suggesting that demand is largely driven by affordability and everyday needs, rather than premium or niche products. However, **mid-range products are also getting significant traction.**

- Additionally, the fact that **97% of customers on the platform are one-time buyers (first-time customers)** suggests **limited trust or confidence in the platform which could be discouraging the customers from purchasing premium (expensive) products**, further reinforcing the preference for lower/moderate priced items.

Furthermore, since **high-volume categories are priced similarly**, the revenue distribution will largely follow sales patterns, which means these product categories will be the **primary revenue drivers for the company.**

```
1543  select
1544    *,
1545    CONCAT(CAST(ROUND(num_of_orders/sum(num_of_orders) OVER()*100, 2) AS STRING), ' %') as percent_of_orders
1546  from (
1547      select
1548        payment_type,
1549        count(order_id) as num_of_orders
1550      from `olist_business_case.payments`
1551      group by payment_type
1552    )
1553  order by num_of_orders desc;
```

| Row | payment_type | num_of_orders | percent_of_orders |
|-----|-------------|---------------|-------------------|
| 1 | credit_card | 76795 | 73.92 % |
| 2 | boleto | 19784 | 19.04 % |
| 3 | voucher | 5775 | 5.56 % |
| 4 | debit_card | 1529 | 1.47 % |

- Nearly **74% of all transactions are completed via credit cards** – meaning almost three out of every four transactions are done using credit cards, making it the most preferred payment method among customers. This also reflects a **purchasing behavior that is largely credit-driven.**
- Olist could **collaborate with leading credit card providers** to facilitate seamless card payments for its customers while offering attractive **perks such as reward points, cashback offers, special discounts**, etc. This would help in enhancing the overall **user experience**.

- The **'Boleto' payment method**, a popular alternative in Brazil, accounts for almost one in five orders (**19%**), **highlighting its continued relevance** for customers who may not have access to credit cards or prefer not to use them.

| Row | payment_type | revenue | percent_of_total_revenue |
|-----|-------------|---------|--------------------------|
| 1 | credit_card | 12542084.19 | 78.34 % |
| 2 | boleto | 2869361.27 | 17.92 % |
| 3 | voucher | 379436.87 | 2.37 % |
| 4 | debit_card | 217989.79 | 1.36 % |

- The **distribution of revenue across payment methods closely follows the pattern observed in transaction volume.** Credit cards account for 78% of total revenue, which aligns with their 74% share of total transactions. Similarly, 'Boleto' contributes 18% of revenue, similar to its 19% share of transactions.

- This close alignment between transaction volume and revenue share for each payment method suggests a largely **uniform Average Order Value (AOV) across payment methods.**

Next, we will examine the preferences for installment options to make the payment.

```
1573  select
1574     distinct payment_type
1575  from `olist_business_case.payments`
1576  where payment_installments > 1;
1577
```
✅ This script will process 851.81 MB when run.

Query results

| Job information | Results | Chart |
|---|---|---|

| Row | payment_type ▾ | |
|---|---|---|
| 1 | credit_card | |

- Among all payment methods, only credit cards offer the option of paying in more than one installment.

```
1580  select
1581     CONCAT(CAST(ROUND(100*countif(payment_installments > 1)/count(*), 2) AS STRING),' %')
1582     as percent_of_installment_transactions
1583  from `olist_business_case.payments`;
1584
```
✅ This script will process 851.81 MB when run.

Query results

| Job information | Results | Chart | JSON | Execution details | Execution graph |
|---|---|---|---|---|---|

| Row | percent_of_installment_transactions |
|---|---|
| 1 | 49.42 % |

- **Nearly half (49.4%) of all transactions are made using installment plans**, **indicating a strong customer preference** for flexible payment options.

- Given this strong preference, Olist should prioritize making **installment plans readily available and prominently displayed during the checkout process**. By doing so, the company can **enhance customer convenience**, reduce financial barriers to purchase – especially for higher-value items—and **potentially increase conversion rates** and overall sales.

Review scores are a direct feedback from the customers representing their overall purchasing experience and that's why review scores can be seen as a valuable **measure of customer satisfaction**.

Let's begin by examining the **distribution of review scores** to gauge overall customer sentiment –

```
1051  select
1052    *,
1053    CONCAT(CAST(ROUND(num_of_orders/(sum(num_of_orders) OVER())*100, 2) AS STRING), ' %') as percent_of_orders
1054  from (
1055        select
1056          review_score,
1057          count(distinct order_id) as num_of_orders
1058        from `olist_business_case.reviews`
1059        group by review_score
1060    ) as tbl
1061  order by num_of_orders desc;
```

| Row | review_score | num_of_orders | percent_of_orders |
|---|---|---|---|
| 1 | 5 | 57076 | 57.73 % |
| 2 | 4 | 19098 | 19.32 % |
| 3 | 1 | 11393 | 11.52 % |
| 4 | 3 | 8160 | 8.25 % |
| 5 | 2 | 3148 | 3.18 % |

- Nearly **58% of all orders have received the highest rating of 5** and overall, approximately **77% of all orders have ratings of 4 or 5.** These figures **indicate a strong level of customer satisfaction and positive purchasing experience.**

- Less than **12% of all the orders received the lowest rating of 1** and overall, approximately **15% of all orders have ratings of 1 or 2**. These figures suggest that while most customers are satisfied with their purchasing experience, a **significant minority had a poor experience**.

- Also, seeing such high proportion of positive review scores, we can comment that **purchasing experience is not the reason behind such low rate of repeat customers (3%)**.

Let's further deep-dive into these review scores to identify the best and worst performing product categories –

To assess the best and worst-performing product categories based on customer reviews, we will calculate the proportion of **low (review score <= 3)** and **high (review score > 3)** rated reviews for each product category. This will help us **identify categories having large proportion of "low rated reviews" (priority for improvement) or "high rated reviews" (potential bestsellers).**

Note: To ensure statistical reliability, we will consider only those categories for our analysis that have received more than 100 reviews (median 'number of reviews' for product categories was calculated to be 244).

```
1145  WITH master_table as (
1146
1147    select distinct
1148      oi.order_id,
1149      p.product_category_name,
1150      r.review_score
1151    from `olist_business_case.order_items` as oi
1152    INNER JOIN `olist_business_case.products` as p
1153    ON oi.product_id = p.product_id
1154    INNER JOIN `olist_business_case.reviews` as r
1155    ON oi.order_id = r.order_id
1156  )
1157
1158  select
1159    product_category_name,
1160    ROUND(100*countif(review_score <= 3)/count(review_score), 2) as percentage_of_low_ratings
1161  from master_table
1162  group by product_category_name
1163  having count(order_id) > 100
1164  order by percentage_of_low_ratings desc;
```

5 best performing product categories based on customer reviews –

| Row | product_category_name ▼ | percentage_of_high_ratings |
|---|---|---|
| 1 | General Interest Books | 88.39 |
| 2 | technical books | 86.38 |
| 3 | Bags Accessories | 83.98 |
| 4 | Drink foods | 83.7 |
| 5 | foods | 83.6 |

5 worst performing product categories based on customer reviews –

| Row | product_category_name ▼ | percentage_of_low_ratings |
|---|---|---|
| 1 | Furniture office | 37.28 |
| 2 | House comfort | 32.49 |
| 3 | Fashion Men's Clothing | 32.43 |
| 4 | audio | 31.61 |
| 5 | CONSTRUCTION SECURITY TO... | 28.31 |

- Among the 52 product categories with sufficient reviews (over 100 reviews), 29 categories outperformed the overall platform benchmark i.e., more than 77% of their reviews received high (review score >3) ratings. Furthermore, among these, the following **2 categories – *'General Interest Books'* and *'technical books'* performed exceptionally well, with more than 85% of their reviews receiving high ratings.**

- Performance of the following **4 categories – *'Furniture Office'*, *'House comfort'*, *'Fashion Men's Clothing'*, and *'audio'*** signals significant customer dissatisfaction as more than 30% of their reviews received low (review score <=3) ratings.

Furthermore, let's assess the **review scores of the key product categories** that we previously identified as the main contributors to Olist's overall sales and revenue –

| Row | product_category_name | percent_of_orders | percentage_of_low_ratings |
|---|---|---|---|
| 1 | bed table bath | 9.9 % | 26.88 |
| 2 | HEALTH BEAUTY | 8.59 % | 20.49 |
| 3 | sport leisure | 7.66 % | 20.31 |
| 4 | Furniture Decoration | 7.4 % | 25.62 |
| 5 | computer accessories | 6.94 % | 24.31 |
| 6 | housewares | 6.17 % | 21.73 |
| 7 | Watches present | 5.33 % | 23.83 |
| 8 | telephony | 4.04 % | 25.48 |
| 9 | Garden tools | 3.86 % | 21.45 |
| 10 | automotive | 3.75 % | 22.28 |

- There are **no significant performance concerns among these high-volume product categories**, which reflects an overall positive customer response to Olist's most in-demand product segments.

- Nevertheless, given their substantial impact on total sales, it remains crucial for the company to continually strive for service improvements and higher customer ratings. For instance, while the 'bed table bath' category leads in sales accounting for nearly 10% of total sales volume, it still receives 27% low ratings (review scores ≤3).

**Timely delivery is often considered as a major factor influencing customer experience in online shopping.** So, we will now explore whether there is a correlation between order delivery times and the review scores customers have provided.

Note: For this analysis of order delivery times, we will only consider those 97% orders that have been *'delivered'* as for the remaining 3% orders, the delivery timestamp is not available yet.

```
1074  select
1075    r.review_score,
1076    ROUND(avg(d.delivery_time_in_days), 2) as avg_delivery_time_in_days
1077  from
1078      (select
1079        order_id,
1080        TIMESTAMP_DIFF(order_delivered_customer_date, order_purchase_timestamp, SECOND)/86400
1081        as delivery_time_in_days
1082      from `olist_business_case.orders`
1083      where order_status = 'delivered') as d
1084  INNER JOIN
1085      (select
1086        order_id,
1087        review_score
1088      from `olist_business_case.reviews`) as r
1089  ON d.order_id = r.order_id
1090  group by r.review_score
1091  order by avg_delivery_time_in_days desc;
```

| Row | review_score | avg_delivery_time_in_days |
|---|---|---|
| 1 | 1 | 21.31 |
| 2 | 2 | 16.66 |
| 3 | 3 | 14.26 |
| 4 | 4 | 12.31 |
| 5 | 5 | 10.69 |

- The data suggests a **clear and strongly negative correlation between average delivery time and customer ratings** – longer delivery times correspond to lower customer ratings.

- **Orders with the lowest ratings (1 and 2) have significantly longer average delivery times** (21.3 and 16.7 days, respectively), while orders with the highest ratings (4 and 5) are delivered much faster (12.3 and 10.7 days, respectively).

- Moreover, Orders that received a 1-star rating experienced an average delivery time of 21.3 days, which is **almost double the delivery time for orders rated 4 or 5** (12.3 and 10.7 days, respectively). This stark difference underscores the **crucial role of timely delivery in shaping customer satisfaction.**

In the previous section, we saw a clear and strongly negative correlation between average delivery time and customer satisfaction – longer delivery times correspond to lower customer ratings. Therefore, let's examine the order delivery times to identify potential areas for improvement and opportunities to shorten delivery times –

To better understand and optimize the order fulfillment process, we have identified **three key stages that each order passes through**:

- the **approval stage** (from order placement to approval),
- the **waiting stage** (from approval to handover to the logistics partner),
- and, the **shipping stage** (from carrier handover to final delivery to the customer).

Let's examine the average time taken at each of these stages –

```sql
select distinct
  'approval_time (in hours)' as different_stages,
  ROUND(AVG(TIMESTAMP_DIFF(order_approved_at, order_purchase_timestamp, SECOND)/3600) OVER(), 2) as mean,
  ROUND(PERCENTILE_CONT(TIMESTAMP_DIFF(order_approved_at, order_purchase_timestamp, SECOND)/3600, 0.5) OVER(), 2) as median
from `olist_business_case.orders`

UNION ALL

select distinct
  'waiting_time (in days)' as different_stages,
  ROUND(AVG(TIMESTAMP_DIFF(order_delivered_carrier_date, order_approved_at, SECOND)/86400) OVER(), 2) as mean,
  ROUND(PERCENTILE_CONT(TIMESTAMP_DIFF(order_delivered_carrier_date, order_approved_at, SECOND)/86400, 0.5) OVER(), 2) as median
from `olist_business_case.orders`

UNION ALL

select distinct
  'shipping_time (in days)' as different_stages,
  ROUND(AVG(TIMESTAMP_DIFF(order_delivered_customer_date, order_delivered_carrier_date, SECOND)/86400) OVER(), 2) as mean,
  ROUND(PERCENTILE_CONT(TIMESTAMP_DIFF(order_delivered_customer_date, order_delivered_carrier_date, SECOND)/86400, 0.5) OVER(), 2) as median
from `olist_business_case.orders`
where order_status = 'delivered'

order by CASE
            WHEN different_stages = 'approval_time (in hours)' THEN 1
            WHEN different_stages = 'waiting_time (in days)' THEN 2
            ELSE 3
         END;
```

| Row | different_stages | mean | median |
|-----|------------------|------|--------|
| 1 | approval_time (in hours) | 10.42 | 0.34 |
| 2 | waiting_time (in days) | 2.81 | 1.82 |
| 3 | shipping_time (in days) | 9.33 | 7.1 |

- The **median approval time is exceptionally fast at just 20 minutes**, indicating that **approvals might be automated** and thus doesn't require any manual intervention.

- However, the large gap between the median and mean approval times (20 minutes vs. 10.4 hours) highlights that large number of orders experience substantial delays – over 27% of orders take more than 10 hours for approval and **over 17% of orders take more than 24 hours for approval**. These significant delays **may be due to requirement of manual reviews** or exceptional circumstances.

- Once approved, **sellers typically hand over orders to logistics partners within two days**, which is competitive by industry standards and **suggests efficient seller operations**.

- The **shipping stage**, with a **median duration of 7.1 days** and a mean of 9.3 days, emerges as the **primary contributor to overall delivery time**.

- To further enhance customer satisfaction and reduce delivery times, **efforts should focus on investigating and optimizing** both the outlier approval cases and, most importantly, the shipping process.

Next, we will evaluate Olist's **delivery efficiency** by **examining its ability to deliver the order on time**. For this we will analyze the proportion of orders delivered early, on time, and late based on the estimated delivery date the customer was informed during order placement –

```
1673  select
1674    'Early delivery' as delivery,
1675    CONCAT(CAST(ROUND(100*countif(extract(date from order_delivered_customer_date) < extract(date from order_estimated_delivery_date))/count(*), 2) AS STRING), ' %')
1676    as percent_of_orders
1677  from `olist_business_case.orders`
1678  where order_status = 'delivered'
1679
1680  UNION ALL
1681
1682  select
1683    'On-time delivery' as delivery,
1684    CONCAT(CAST(ROUND(100*countif(extract(date from order_delivered_customer_date) = extract(date from order_estimated_delivery_date))/count(*), 2) AS STRING), ' %')
1685    as percent_of_orders
1686  from `olist_business_case.orders`
1687  where order_status = 'delivered'
1688
1689  UNION ALL
1690
1691  select
1692    'Late delivery' as delivery,
1693    CONCAT(CAST(ROUND(100*countif(extract(date from order_delivered_customer_date) > extract(date from order_estimated_delivery_date))/count(*), 2) AS STRING), ' %')
1694    as percent_of_orders
1695  from `olist_business_case.orders`
1696  where order_status = 'delivered'
1697
1698  order by case
1699          when delivery = 'Early delivery' then 1
1700          when delivery = 'On-time delivery' then 2
1701          else 3
1702        end;
```

| Row | delivery | percent_of_orders |
|---|---|---|
| 1 | Early delivery | 91.88 % |
| 2 | On-time delivery | 1.34 % |
| 3 | Late delivery | 6.77 % |

- The data reveals an **outstanding delivery efficiency**, with nearly **92% of orders getting delivered earlier than the promised delivery date**. This positions Olist as a reliable and customer-centric platform in terms of order fulfillment.

- However, the **disproportionately high "early deliveries" suggests that Olist's estimated delivery windows may be overly cautious/conservative**, resulting in most orders arriving ahead of schedule. Let's validate it by examining the difference between *'average actual order delivery times'* and *'average estimated delivery times'* –

```
1702  select distinct
1703    'actual_delivery_time' as delivery_time_in_days,
1704    ROUND(AVG(TIMESTAMP_DIFF(order_delivered_customer_date, order_purchase_timestamp, SECOND)/86400) OVER(), 2) as mean,
1705    ROUND(PERCENTILE_CONT(TIMESTAMP_DIFF(order_delivered_customer_date, order_purchase_timestamp, SECOND)/86400, 0.5) OVER(), 2)as median
1706  from `olist_business_case.orders`
1707  where order_status = 'delivered'
1708
1709  UNION ALL
1710
1711  select distinct
1712    'estimated_delivery_time' as delivery_time_in_days,
1713    ROUND(AVG(TIMESTAMP_DIFF(order_estimated_delivery_date, order_purchase_timestamp, SECOND)/86400) OVER(), 2) as mean,
1714    ROUND(PERCENTILE_CONT(TIMESTAMP_DIFF(order_estimated_delivery_date, order_purchase_timestamp, SECOND)/86400, 0.5) OVER(), 2) as median
1715  from `olist_business_case.orders`
1716  where order_status = 'delivered'
1717
1718  order by CASE
1719         WHEN delivery_time_in_days = 'actual_delivery_time' THEN 1
1720         ELSE 2
1721       END;
```

| Row | delivery_time_in_days | mean | median |
|---|---|---|---|
| 1 | actual_delivery_time | 12.56 | 10.22 |
| 2 | estimated_delivery_time | 23.74 | 23.23 |

- As anticipated, there is a **significantly large gap – nearly double – between Olist's estimated and actual delivery times**, confirming that the company's **delivery estimates are highly conservative**.
- Olist should consider refining its forecasting of estimated delivery times as longer estimated delivery times informed to customers may have unintended consequences. **Overly long delivery estimates at checkout could deter potential customers from completing their purchases.**