# [A24] Data Warehousing & ETL

## Student information

- Student: Carlos Jesus Caro
- Email: [carlos.jesus-caro@edu.dsti.institute](mailto:carlos.jesus-caro@edu.dsti.institute)
- Student type: SPOC

## Project

- GitHub repository: https://github.com/carlosjesuscaro/masters_de_dwh_olist
- File: OLIST.zip
- Structure:

- o Folder "SQLQueries" - it contains all the queries used throughout the project
    - STA_Queries.sql
    - ODS_Quewries.sql
    - DWH_Queries.sql
    - DimDate_Query.sql
    - Business_Queries.sql
- o Folder "source_data" - it contains the 8 CSV files containing all the RAW data from Kaggle
- o Visual Studio data:
    - OLIST.sln
    - OLIST folder

# Dataset

- Source: Kaggle
- Link to dataset: https://www.kaggle.com/datasets/olistbr/brazilian-ecommerce
- Dataset description: The dataset is from the e-commerce business called OLIST from Brazil and it was chosen because it is very comprehensive as it includes the following tables:
    - o Customers
    - o Orders
    - o Payments
    - o Sellers
    - o OrderItems
    - o OrderPayments
    - o OrderReviews
    - o ProductCategoryNameTransation (from Portuguese to English)
- Dataset format: CSV
- Data size: ~100,000 rows

All queries are stored in the folder: "SQLQueries"

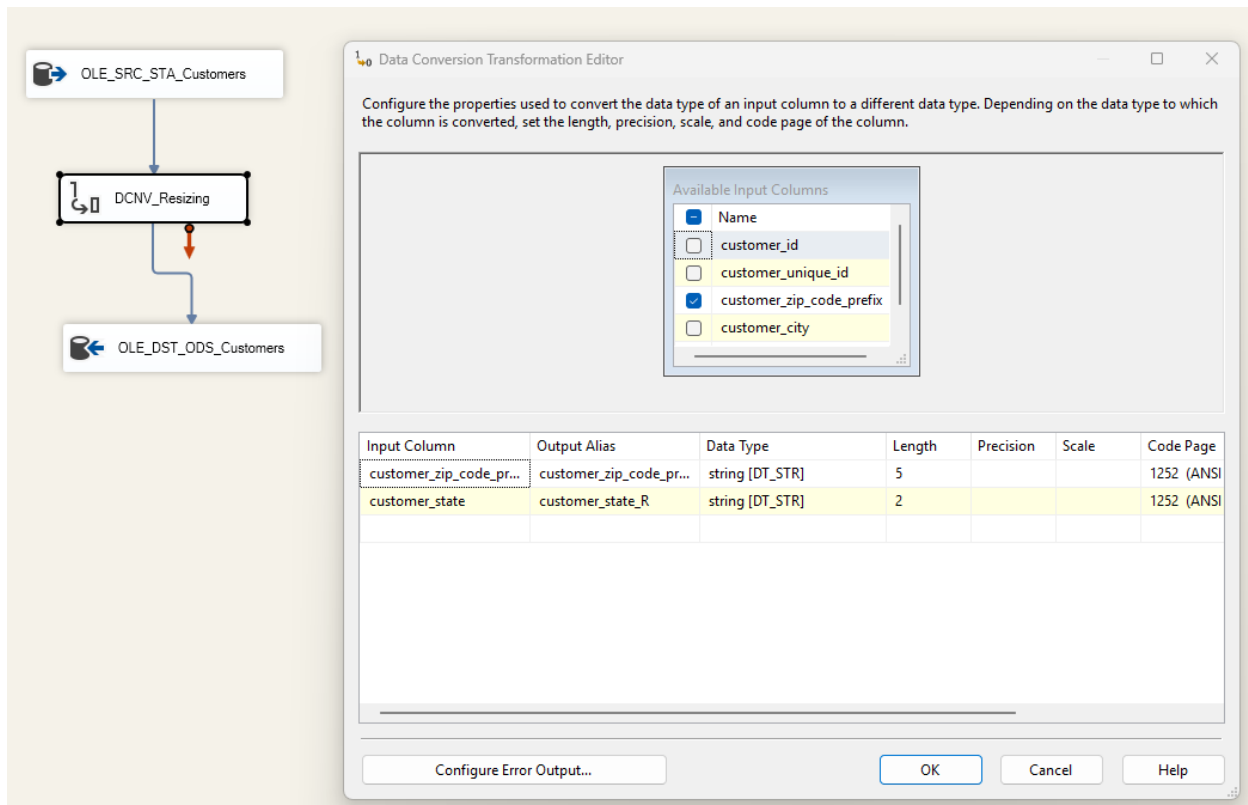# ETL Process and key transformations

## STA

All the source files were uploaded to the OLIST_STA database "as is". Here is the mapping between source CSV files and DB tables:

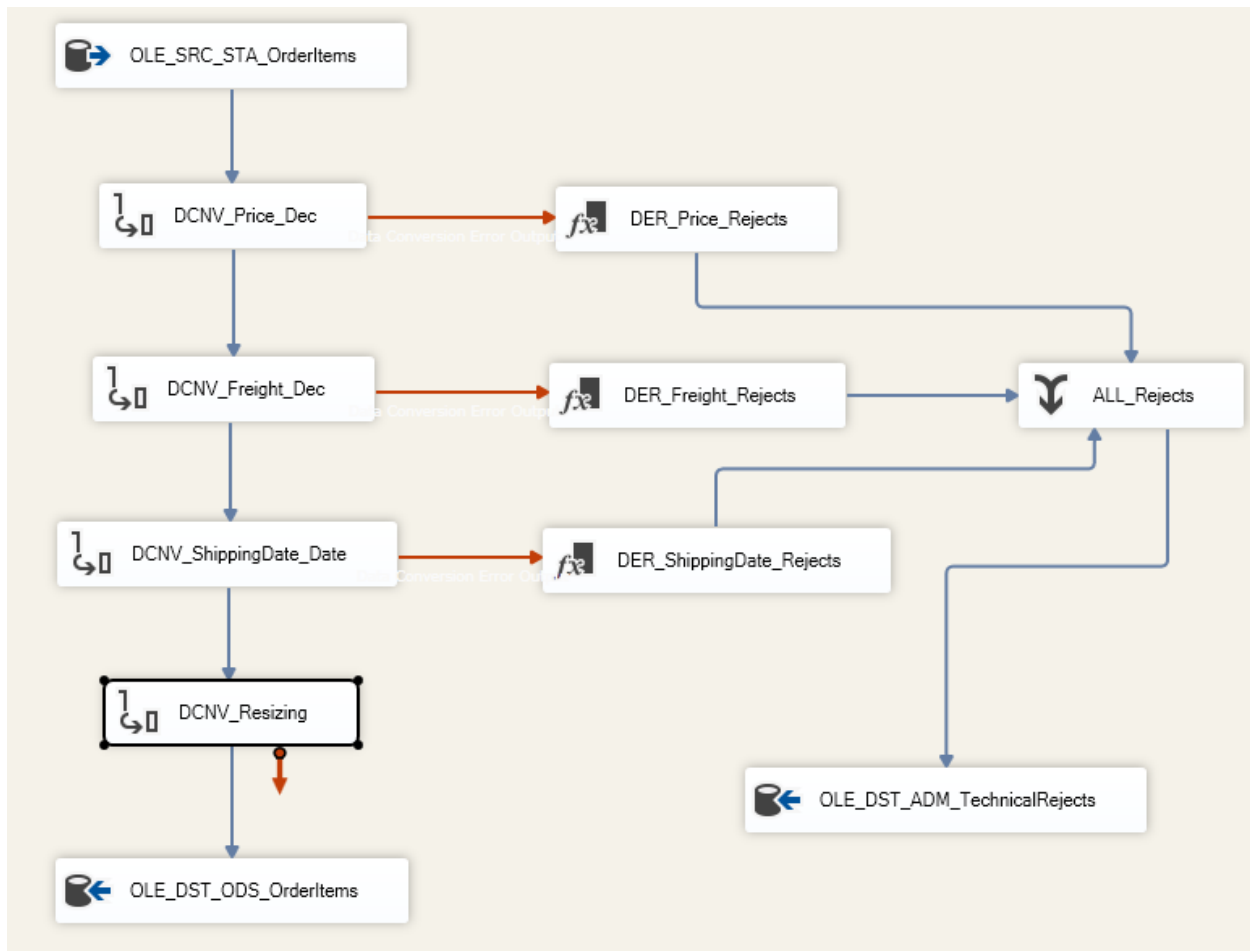| Source file | OLIST_STA table |
|---|---|
| olist_customers_dataset.csv | STA_Customers |
| olist_order_items_dataset.csv | STA_OrderItems |
| olist_order_payments_dataset.csv | STA_OrderPayments |
| olist_order_reviews_dataset.csv | STA_OrderReviews |
| olist_orders_dataset.csv | STA_Orders |
| olist_products_dataset.csv | STA_Products |
| olist_sellers_dataset.csv | STA_Sellers |
| product_category_name_transition.csv | STA_ProductTranslations |

## ODS

1. STA_Customers – ODS_Customers
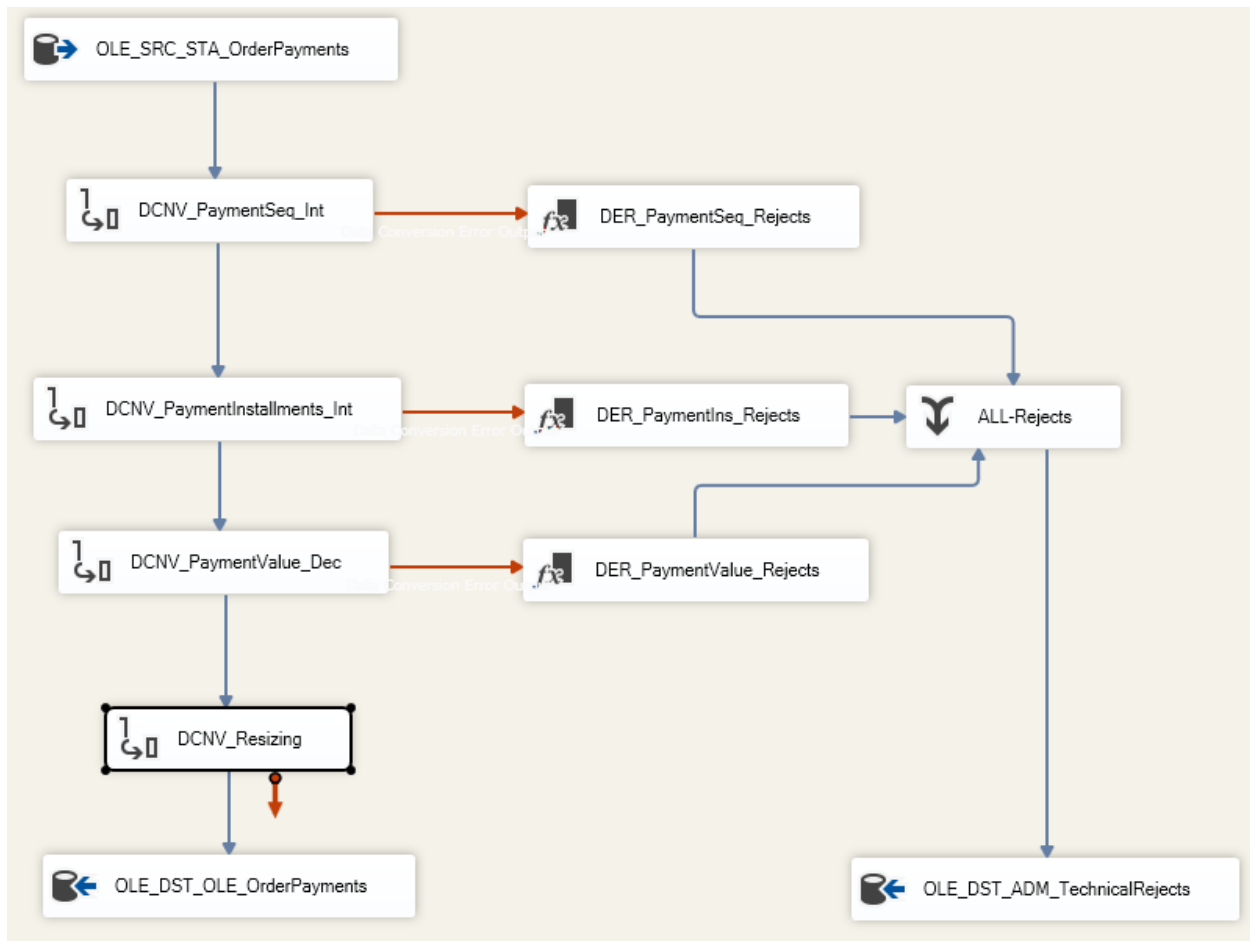   - Resizing customer_zip_code and customer_state

2. STA_OrderItems – ODS_OrderItems
   - Converting Price and Freight to numeric data type
   - Converting ShippingDate to a date data type
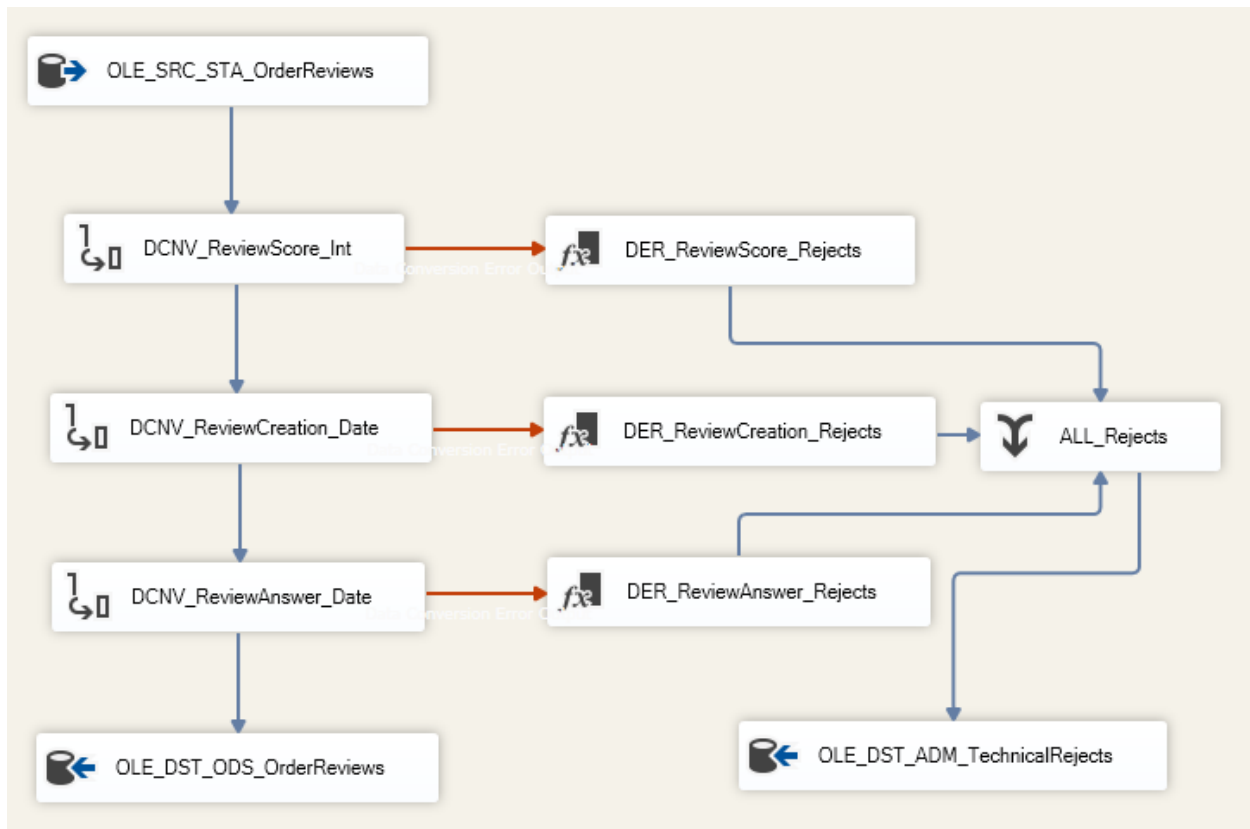   - Resizing order_item

3. STA_OrderPayments – ODS_OrderPayments
   - Converting PaymentSequential, PaymentInstallments and PaymentValue to a numeric data type
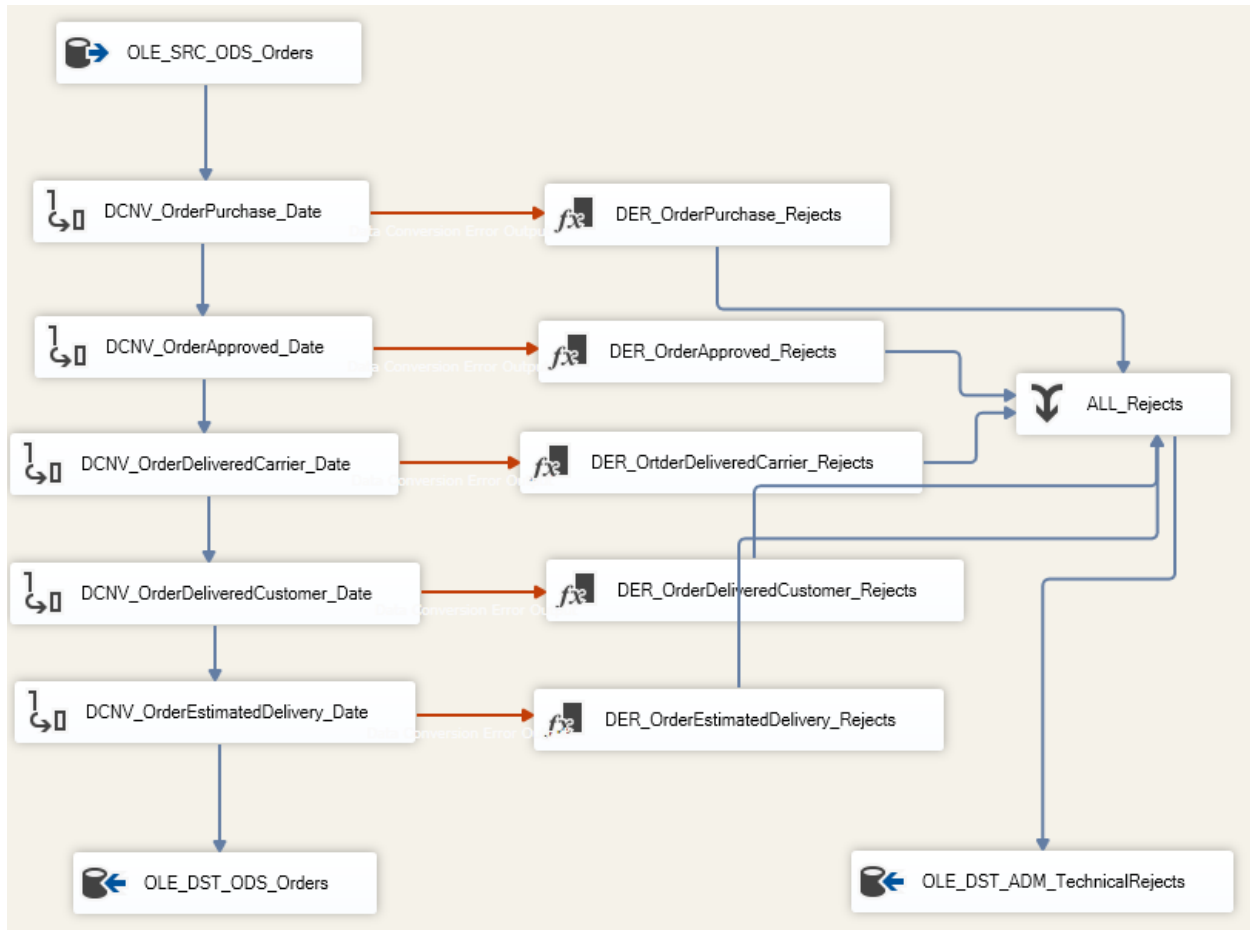   - Resizing PaymentType

4. STA_OrderReviews – ODS_OrderReviews
   - Converting ReviewScore to a numeric data type
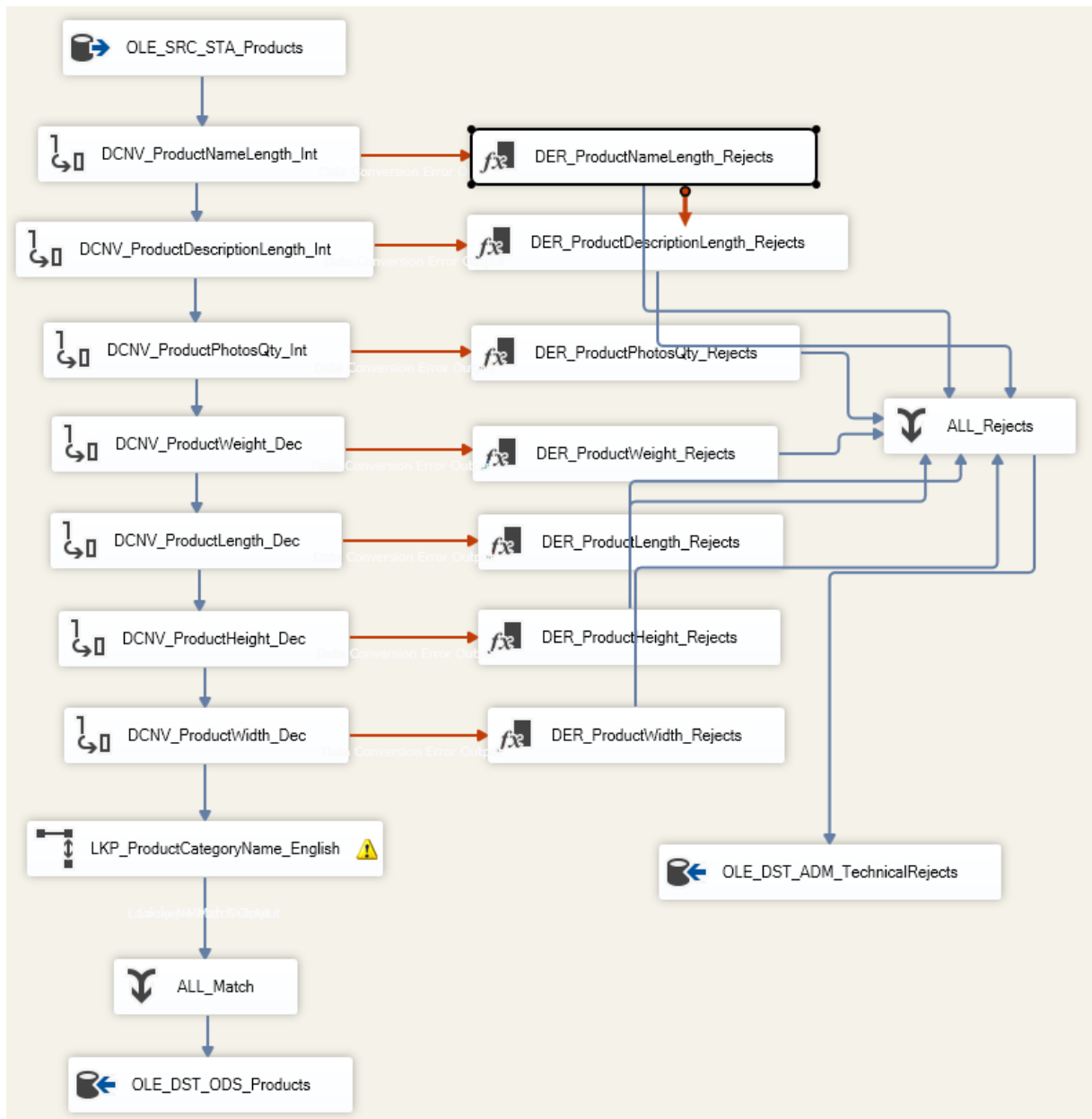   - Converting ReviewCreation and ReviewAnswer to a date data type

5. STA_Orders – ODS_Orders
   - Converting OrderPurchase, OrderApproved, OrderDeliveredCarrier, OrderDeliveredCustomer and OrderEstimatedDelivery to a date data type

6. STA_Products – ODS_Products
   - Converting ProductNameLength, ProductDescriptionLength, ProductPhotosQty, ProductWeight, ProductHeight and ProductWidth to a numeric data type
   - Lookup between the opriginal product_category_name (in Portuguese) to product_category_name_english. This is why the table STA_ProductTranslations is used on this lookup but not later used anymore

7. STA_Sellers – ODS_Sellers
   • Resizing seller_zip_code and seller_state

In all cases of data conversions, the technical rejects have been populated to the TechnicalRejects table in the OLIST_ADM database.

# Datawarehouse

## STAR Schema



## Dimension tables
- The chosen option for "Slow Changing Dimension – SDC" was type 1through the use of built-in feature in the SSIS toolbox
- A surrogate key has been added to all the dimension tables
- Tables:
    - DWH_DimDate
    - DWH_DimCustomer

- DWH_DimProducts

- DWH_DimSellers

- DWH_DimOrderReviews

- DWH_DimOrderItems



- DWH_DimOrderPayments

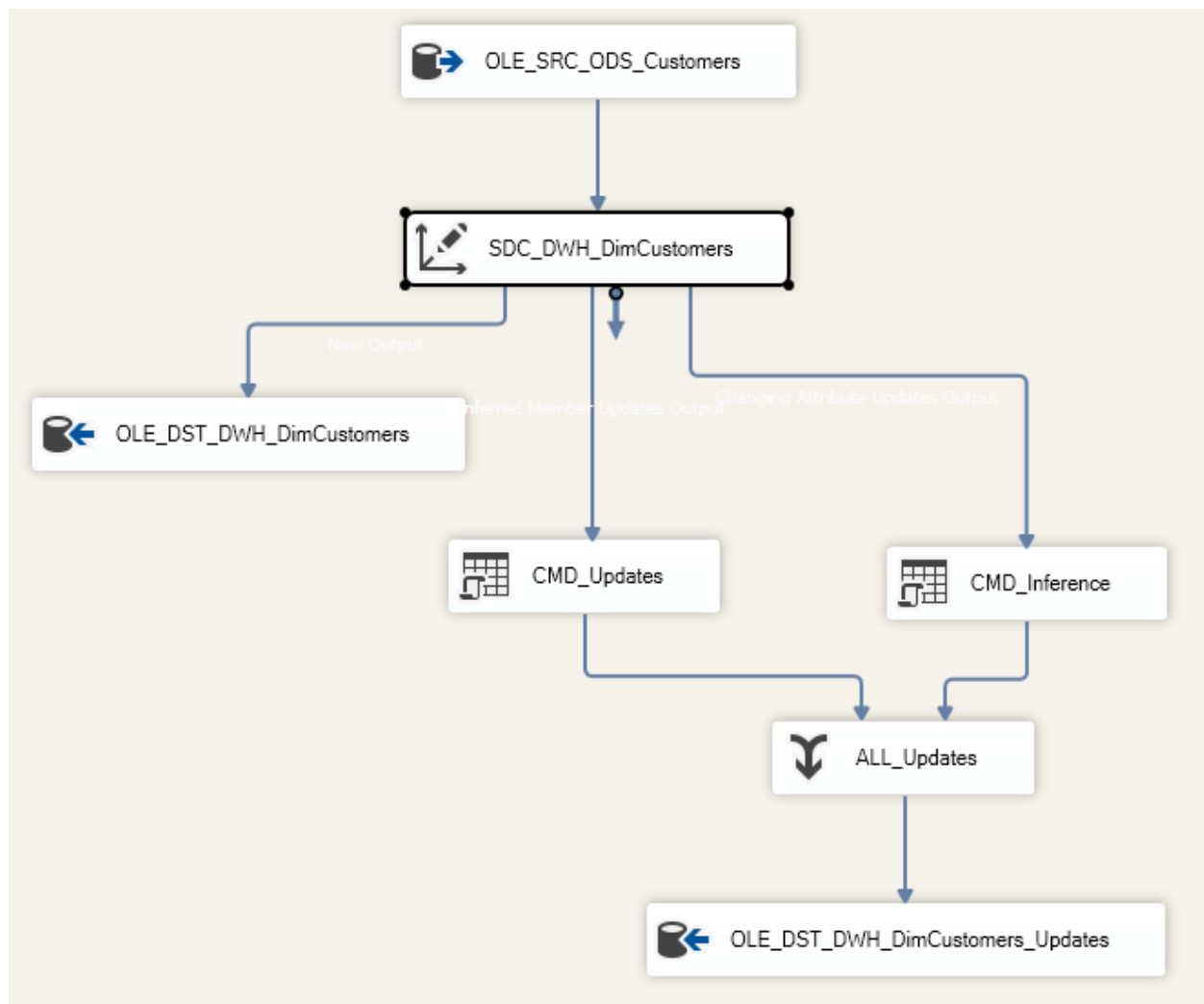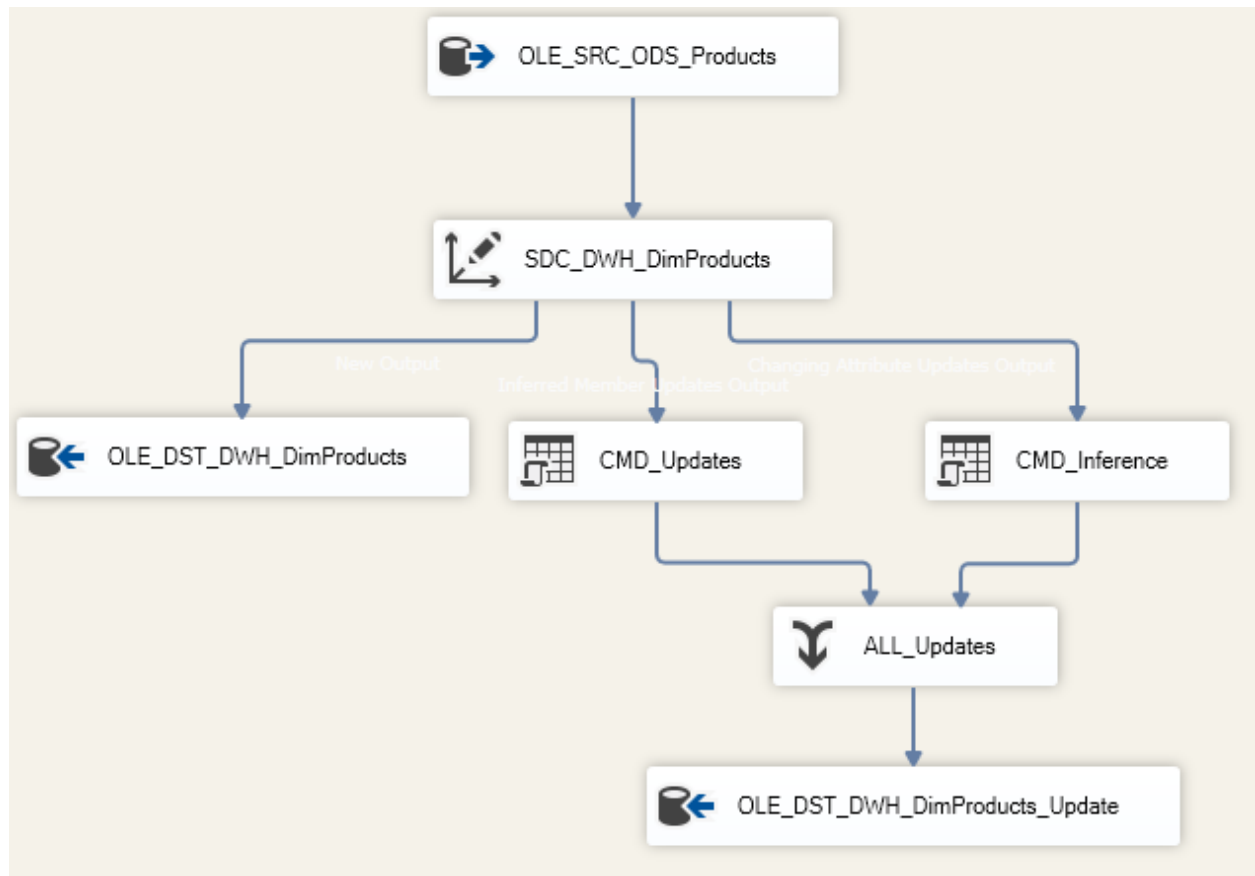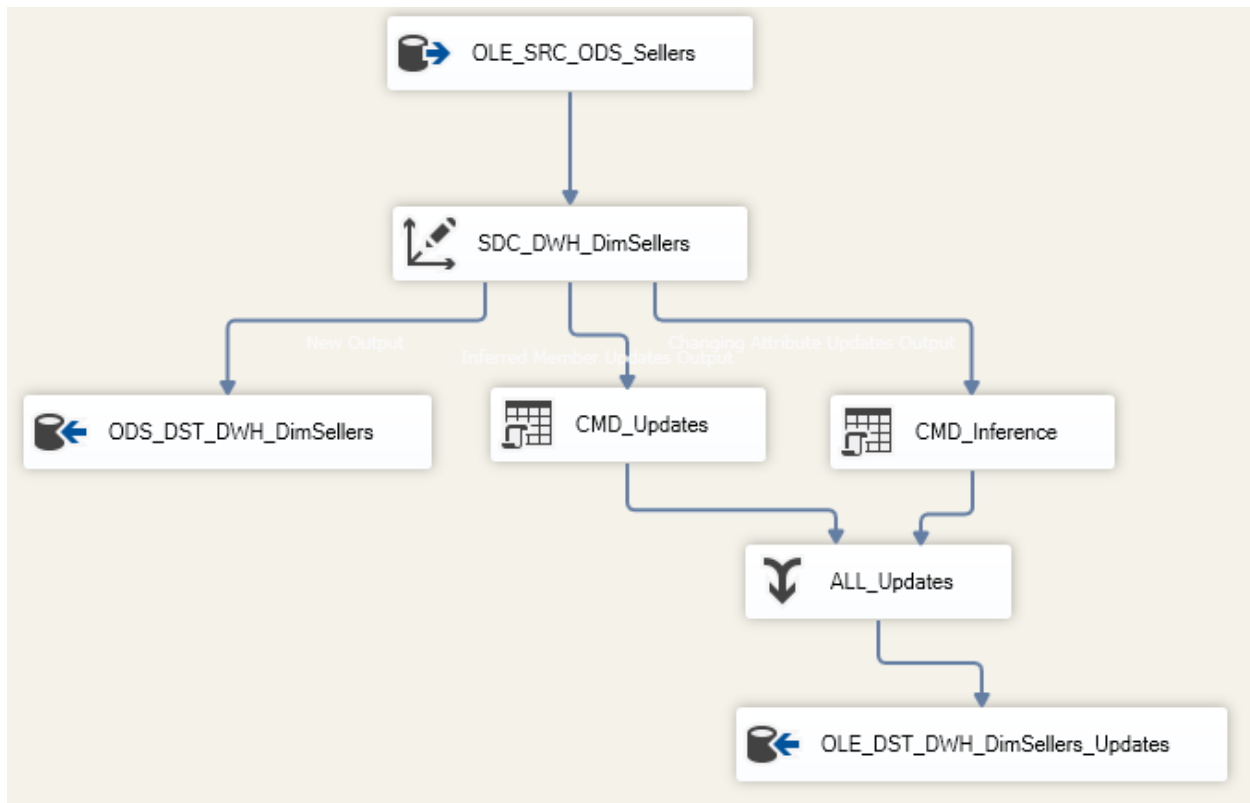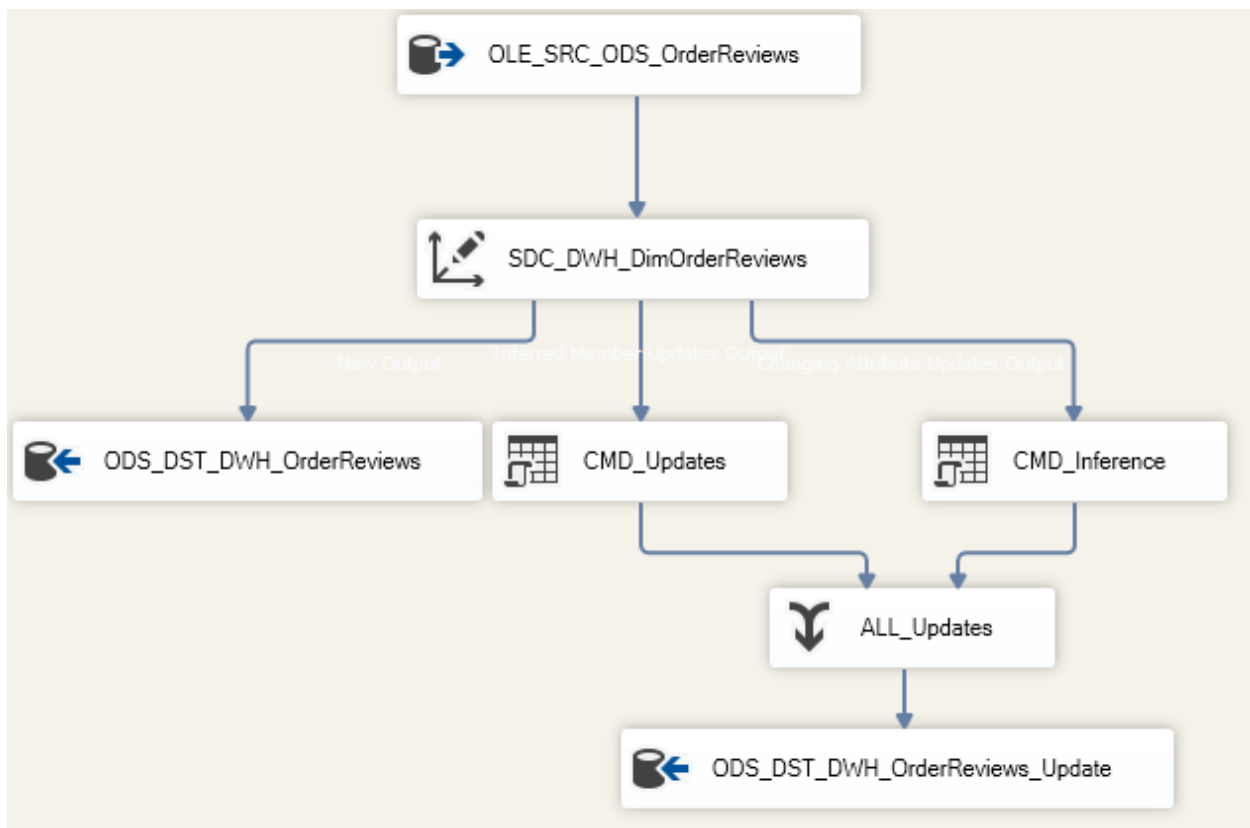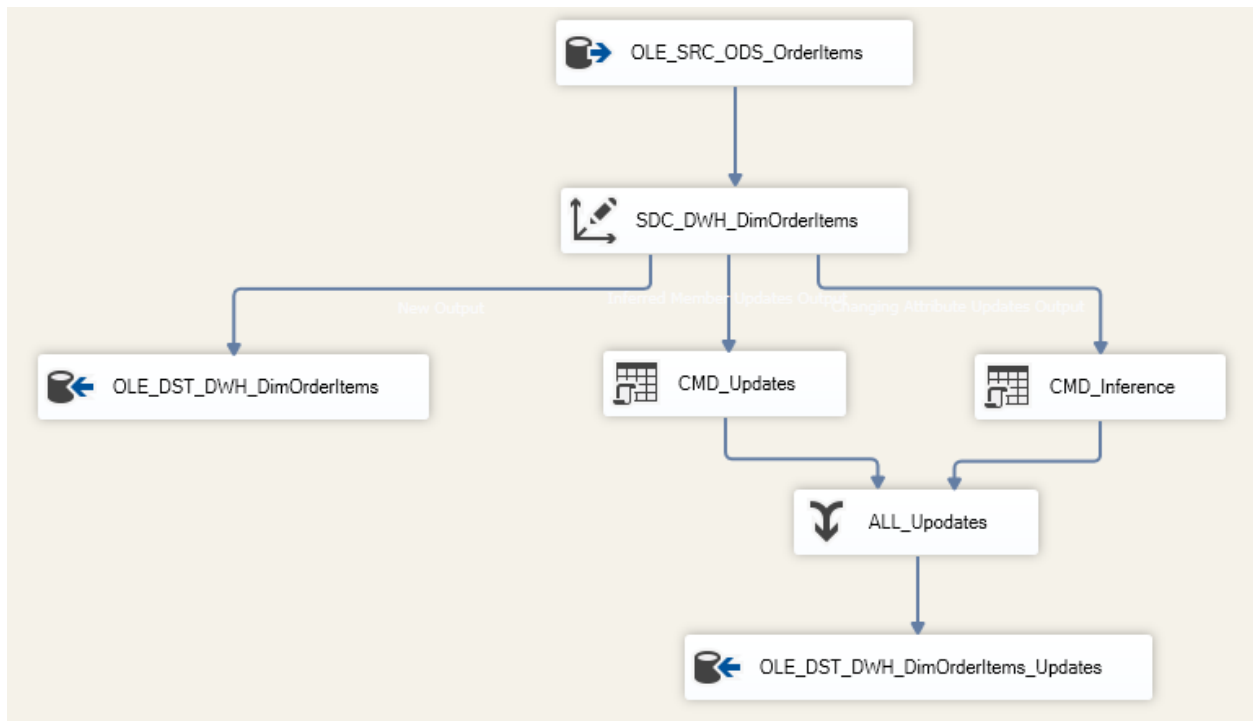The source table to the FACT table is ODS_Orders which contains the following fields:

- Order_id
- Customer_id
- Order_status
- Order_purchase_timestamp (timestamp)
- Order_approved_at (timestamp)
- Order_delivered_carrier_date (timestamp)
- Order_delivered_customer_date (timestamp)
- Order_estimated_delivery_date (timestamp)

The FACT table is DWH_FactOrders contains the following fields:

- Order_id
- Customer_key (surrogate key from the DWH_DimCustomers)
- Order_status
- Order_purchase_datekey (date) - removing the timestamp to match the date_key INT from the DimDate table
- Order_approved_datekey (date) - removing the timestamp to match the date_key INT from the DimDate table
- Order_delivered_carrier_datekey (date) - removing the timestamp to match the date_key INT from the DimDate table
- Order_delivered_customer_datekey (date) - removing the timestamp to match the date_key INT from the DimDate table
- Order_estimated_delivery_datekey (date) - removing the timestamp to match the date_key INT from the DimDate table
- Order_purchase_timestamp (timestamp)
- Order_approved_at (timestamp)
- Order_delivered_carrier_date (timestamp)
- Order_delivered_customer_date (timestamp)
- Order_estimated_delivery_date (timestamp)

In the cases where the surrogate key does not exist because the record does not exist, the surrogate key is created as a –1 and added to the table FunctionalRejects in the OLIST_ADM database

## Execution

Execute the SSIS package "OLIST_Exec.dtsx" to execute the entire pipeline in the required order



# Data insights

Note: all the queries are in the file Business_Queries.sql

1. Top 10 Brazilian states with the highest number of orders

```
 2
 3        -- 1. Rank the states based on number of orders
 4
 5     ∨  SELECT TOP 10
 6            dc.customer_state,
 7            COUNT(DISTINCT fo.order_id) AS orders_count
 8        FROM DWH_FactOrders AS fo
 9        LEFT JOIN DWH_DimCustomers AS dc
10            ON fo.customer_key = dc.customer_key
11        LEFT JOIN DWH_DimOrderPayments AS dop
12            ON fo.order_id = dop.order_id
13        GROUP BY customer_state
14        ORDER BY orders_count DESC;
```

✔ No issues found

Results | Messages

| customer_state | orders_count |
|---|---|
| SP | 40489 |
| RJ | 12351 |
| MG | 11352 |
| RS | 5342 |
| PR | 4923 |
| SC | 3547 |
| BA | 3256 |
| DF | 2080 |
| ES | 1995 |
| GO | 1957 |

2. Top 10 Brazilian states with the highest average payment value per order

```
16    -- 2. Rank the states based on the average order value
17
18  ∨ SELECT TOP 10
19        dc.customer_state,
20        AVG(dop.payment_value) AS avg_payment
21    FROM DWH_FactOrders AS fo
22    LEFT JOIN DWH_DimCustomers AS dc
23        ON fo.customer_key = dc.customer_key
24    LEFT JOIN DWH_DimOrderPayments AS dop
25        ON fo.order_id = dop.order_id
26    GROUP BY customer_state
27    ORDER BY avg_payment DESC;
28
```

00 %    ▼        ✅ No issues found        ◀

⊞ Results  📄 Messages

|    | customer_state | avg_payment |
|----|----------------|-------------|
| 1  | PB             | 250.301996  |
| 2  | AP             | 234.536231  |
| 3  | AC             | 234.488795  |
| 4  | AL             | 230.022116  |
| 5  | RO             | 226.147490  |
| 6  | RR             | 220.476097  |
| 7  | PA             | 215.110835  |
| 8  | PI             | 210.711011  |
| 9  | SE             | 206.839331  |
| 10 | TO             | 204.662779  |

3. Listing the sellers based on number of order and the stated where they belong to

```sql
-- 3. Seller ranking based on number of order and including the seller's state

SELECT
    DISTINCT doi.seller_id,
    ds.seller_state,
    COUNT(doi.order_id) OVER(PARTITION BY doi.seller_id) AS count_order_per_seller
FROM DWH_DimOrderItems AS doi
LEFT JOIN DWH_DimSellers AS ds
ON doi.seller_id = ds.seller_id
ORDER BY count_order_per_seller DESC;
```

%    ✓ No issues found                                                    Ln: 38    Ch: 38

Results    Messages

| seller_id | seller_state | count_order_per_seller |
|---|---|---|
| 6560211a19b47992c3666cc44a7e94c0 | SP | 2033 |
| 4a3ca9315b744ce9f8e9374361493884 | SP | 1987 |
| 1f50f920176fa81dab994f9023523100 | SP | 1931 |
| cc419e0650a3c5ba77189a1882b7556a | SP | 1775 |
| da8622b14eb17ae2831f4ac5b9dab84a | SP | 1551 |
| 955fee9216a65b617aa5c0531780ce60 | SP | 1499 |
| 1025f0e2d44d7041d6cf58b6550e0bfa | SP | 1428 |
| 7c67e1448b00f6e969d365cea6b010ab | SP | 1364 |
| ea8482cd71df3c1969d7b9473ff13abc | SP | 1203 |
| 7a67c85e85bb2ce8582c35f2203ad736 | SP | 1171 |
| 4869f7a5dfa277a7dca6462dcf3b52b2 | SP | 1156 |
| 3d871de0142ce09b7081e2b9d1733cb1 | SP | 1147 |

4. Top 10 product categories with the highest review scores

```sql
-- 4. Top 10 product_categories with the higghest reviews

SELECT TOP 10
    dp.product_category_name_english,
    AVG(dor.review_score) AS product_category_review
FROM DWH_FactOrders AS fo
LEFT JOIN DWH_DimOrderItems AS doi
ON fo.order_id = doi.order_id
LEFT JOIN DWH_DimProducts AS dp
ON doi.product_id = dp.product_id
LEFT JOIN DWH_DimOrderReviews AS dor
ON dor.order_id = doi.order_id
WHERE
    dp.product_category_name_english IS NOT NULL
    AND dor.review_score IS NOT NULL
GROUP BY product_category_name_english
ORDER BY product_category_review DESC;
```

✔ No issues found

**Results** | Messages

| product_category_name_english | product_category_review |
| --- | --- |
| fashion_childrens_clothes | 5 |
| fashion_sport | 4 |
| consoles_games | 4 |
| cds_dvds_musicals | 4 |
| small_appliances | 4 |
| garden_tools | 4 |
| fashion_underwear_beach | 4 |
| arts_and_craftmanship | 4 |
| home_appliances | 4 |
| housewares | 4 |

5. Top 10 sellers with the highest number of late deliveries

```sql
      -- 5. Top 10 sellers with the highest number of late deliveries
      WITH delivery AS (
          SELECT
          fo.order_id,
          CASE
              WHEN fo.order_delivered_customer_date > fo.order_estimated_delivery_date THEN 'Late'
                  ELSE 'On time'
              END AS deadline
          FROM DWH_FactOrders AS fo)
      SELECT TOP 10
          ds.seller_id,
          COUNT(delivery.deadline) AS late_count
      FROM DWH_DimOrderItems AS doi
      LEFT JOIN delivery
          ON doi.order_id = delivery.order_id
      LEFT JOIN DWH_DimSellers AS ds
          ON ds.seller_id = doi.seller_id
      WHERE
          delivery.deadline = 'Late'
      GROUP BY ds.seller_id
      ORDER BY late_count DESC;
```

No issues found                                                                Ln: 79   Ch: 1   TABS

Results  Messages

| seller_id | late_count |
| --- | --- |
| 4a3ca9315b744ce9f8e9374361493884 | 214 |
| 1f50f920176fa81dab994f9023523100 | 182 |
| 4869f7a5dfa277a7dca6462dcf3b52b2 | 133 |
| 1025f0e2d44d7041d6cf58b6550e0bfa | 131 |
| 7c67e1448b00f6e969d365cea6b010ab | 130 |
| 6560211a19b47992c3666cc44a7e94c0 | 124 |
| ea8482cd71df3c1969d7b9473ff13abc | 123 |
| 955fee9216a65b617aa5c0531780ce60 | 119 |
| da8622b14eb17ae2831f4ac5b9dab84a | 113 |
| cc419e0650a3c5ba77189a1882b7556a | 103 |