Sri Lanka Institute of Information Technology

IT3021- Data Warehousing and Business Intelligence

**Assignment 1**

2022

Submitted By:
Nirmal M.D.S
IT20074340
Batch :- Y3.S1.WD.DS.05.01

# Contents

# Step 01:- Data Set Selection

**Data Set Name :** Brazilian E-Commerce Public Dataset by Olist

**Source URL :-** https://www.kaggle.com/datasets/olistbr/brazilian-ecommerce

**About Dataset:**

This Public Dataset is obtained by a Brazilian ecommerce company. This includes order information from 2016 – 2018. There are nearly 100,000 records about orders that happened at multiple marketplaces at Brazil. This dataset gives data about the customers, sellers, their geolocations, products, product category names, orders ,order items, payments, and order reviews. Original dataset only includes .csv files
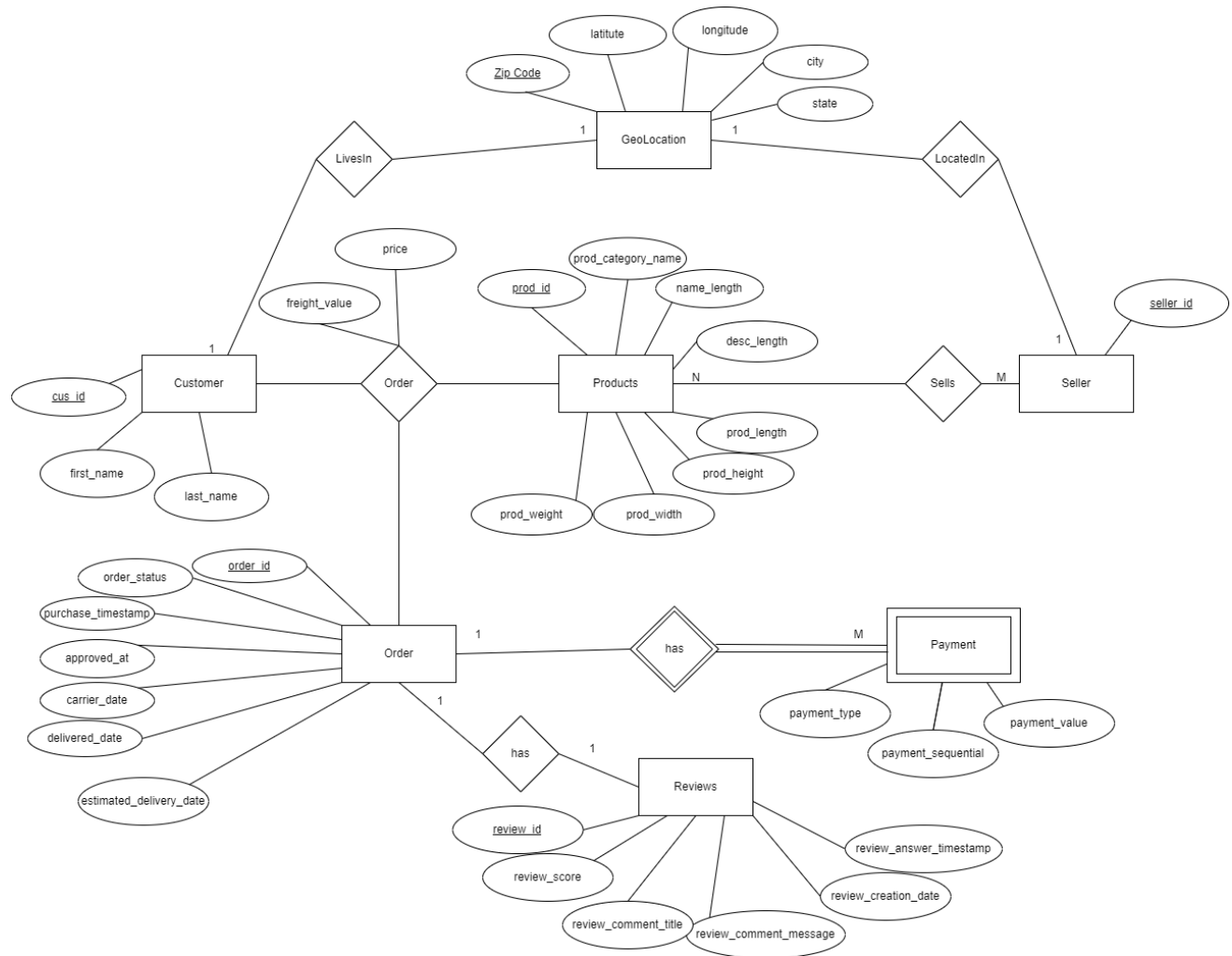
Altogether it had 9 csv files as follows.

- **Customers.csv -** This dataset has information about the customer and its location. Used to identify unique customers.
- **Sellers.csv -** This dataset includes data about the sellers.
- **Geolocation.csv -** This dataset has information Brazilian zip codes and its latitude and longitude coordinates
- **Products.csv -** This dataset includes data about the products sold by sellers
- **Product_category_name_translation.csv -** Translates the product category name to english.
- **Orders.csv –** This dataset includes data about orders.
- **Order_items.csv -** This dataset includes data about the items purchased within each order.
- **Order_payments.csv -** This dataset includes data about the orders payment options.
- **Order_reviews.csv -** This dataset includes data about the reviews made by the customers.

Here we have converted some of the csv files data into a Database source.

The dataset obtained was customized. It includes 5 tables in **Database** format within the '**E_CommerceSourceDB**' database. And the remaining 4 files in '**.csv**' format.

# ER Diagram



The above diagram shows the connection between the entities in the dataset.

**Assumptions:**

1. An order might have multiple items.
2. Each item might be fulfilled by a distinct seller.
3. One Seller Lives in only one location
4. One Customer lives in one location.
5. One order has one review
6. One customer can have many products in many orders
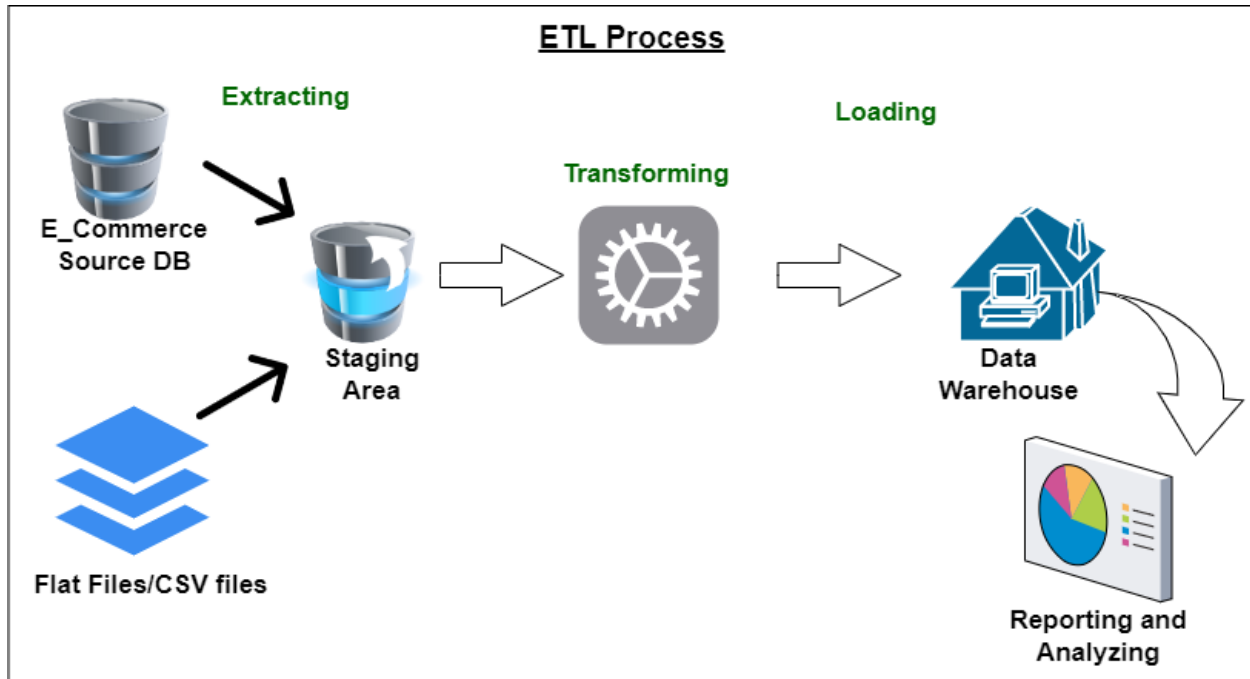
# Step 02:- Preparation of Data Sources

All the data sources provided in the web in .csv format. Since we need at least 2 types of data sources, here some of the .csv file data was imported to a database and some minor changes were done(some columns were removed/added in the process) to the original dataset.

Final Source Dataset was in the below formats before transforming into the Data Warehouse.

Two main sources formats listed below:

- **SQL DataBase → E_commerceSourceDB**

    - **Customer Table**
    - **Geolocation Table**
    - **Product Table**
    - **ProductCategoryNameTranslate Table**
    - **Seller Table**

- **CSV files      →**

    - **orders.csv**
    - **order_items.csv**
    - **order_payments.csv**
    - **order_reviews.csv**

# Step 03:-Solution Architecture



**Source Database:-**E_CommerceSourceDB

**Flat Files/CSV files :-** orders.csv, order_items.csv, order_payments.csv, order_reviews.csv

# Step 04:-Data Warehouse Design & Development

## Design

### Dimensional Tables

- **DimProductCategoryName -** Includes Product SubCategory names and their translation data.
- **DimProduct –** Includes Product data.
- **DimSeller –** Includes Sellers data.
- **DimCustomer-** Includes customer data.
- **DimDate –** Includes Dates and Date key.
- **DimPayment –** Includes Payment information.

## Fact Table

- **FactOrder -** Is a combination of factual data such as Surrogate keys of other dimensions for reference and the measurable data about the orders that happened in this Ecommerce company.

As mentioned above ; there are 6 Dimensional tables including the Date Dimension and the Slowly Changing Dimension which is the **DimCustomer** table. And one fact table called **FactOrder** along with them.

The Data Warehouse model is Snowflake.



➤ Hierarchies
- o DimCustomer has hierarchical attributes about Customer GeoLocation.

➤ Calculations
- o Total Price is calculated in the FactOrder table.
  TotalPrice = quantity * price

➢ Assumptions
  o Transaction table used for creating Fact Table
  o Transaction per Customer considered as the grain.

➢ Slowly Changing Dimensions
  o DimCustomer was considered as the Slowly changing Dimension

| Fixed Attributes | Historical Attributes |
|---|---|
| Customer_first_name | Customer_latitude |
| Customer_last_name | Customer_longitude |
| | Customer_State |
| | Customer_city |
| | Customer_zip_code |

# Step 05:-ETL Development

**I ) Data Extraction and Staging :** Initially in the SSIS project , I have created the Staging Database connecting to Flat File sources and the OLE DB Sources.

**<u>Control Flow</u>**

## Staging Geolocation Data



## Staging Customer Data



## Staging Seller Data



## Staging Product Category Name Data

## Staging Product Data



## Staging Order Data



## Staging Order Items Data



## Staging Payments Data

## Staging Order Review Data



## Staging Tables created and values inserted

## II ) Data Profiling

Data Profiling gives some statistics about data in our tables. Now since we have data in our Staging tables, we can run a Data Profiling task to get a better sense of data that are available in source systems. And determine what type of transformations need to be performed on the data.

## III ) Data Transformation and Loading :

Data Transformation is developed according to the dimensional modelling.

Dimension tables are loaded with data from relevant staging tables in this step.



(*The Control flow was executed successfully.I couldn't take the screenshot while the process was running.I had to take it afterwards.  )

In this step, the data in the Staging area(E_Commerce_Staging Database) is loaded to the dimension and Fact Tables in the  E_Commerce_DW.

## a)Transform and Load Product Category Name Data

```
Procedure used to identify whether to insert or update the data in DimProductCategory
                                      Table

CREATE PROCEDURE dbo.UpdateDimProductCategoryName
@ProductCategoryName nvarchar(50),
@ProductCategoryEnglish nvarchar(50)
AS BEGIN
if not exists (select ProductCategorySK
from dbo.DimProductCategoryName
where Product_category_name = @ProductCategoryName) BEGIN
insert into dbo.DimProductCategoryName
(Product_category_name, Product_category_name_english, InsertDate, ModifiedDate)
values
(@ProductCategoryName, @ProductCategoryEnglish, GETDATE(), GETDATE()) END;
if exists (select ProductCategorySK
from dbo.DimProductCategoryName
where Product_category_name = @ProductCategoryName) BEGIN
update dbo.DimProductCategoryName
set Product_category_name_english = @ProductCategoryEnglish,
ModifiedDate = GETDATE()
where Product_category_name = @ProductCategoryName END;

END;
```

## b) Transform and Load Product Data

```
    Procedure used to identify whether to insert or update the data in DimProduct Table

CREATE PROCEDURE dbo.UpdateDimProduct
@AlternateId nvarchar(50),
@ProductCategoryNameKey int,
@ProductNameLength int,
@ProductDescriptionLength int,
@ProductPhotosQty int,
@ProductWeightG int,
@ProductLengthCm int,
@ProductHeightCm int,
@ProductWidthCm int
AS BEGIN
if not exists (select ProductSK
from dbo.DimProduct
where Alternate_product_id = @AlternateId) BEGIN
insert into dbo.DimProduct
(Alternate_product_id, Product_category_name_key, Product_name_length,
Product_description_length,
Product_photos_qty, Product_weight_g, Product_length_cm, Product_height_cm,
Product_width_cm, InsertDate, ModifiedDate)
values
(@AlternateId, @ProductCategoryNameKey, @ProductNameLength , @ProductDescriptionLength,
@ProductPhotosQty,
@ProductWeightG, @ProductLengthCm, @ProductHeightCm, @ProductWidthCm,  GETDATE(),
GETDATE()) END;
if exists (select ProductSK
from dbo.DimProduct
where Alternate_product_id = @AlternateId) BEGIN
update dbo.DimProduct
set ModifiedDate = GETDATE()
where Alternate_product_id = @AlternateId END;

END;
```
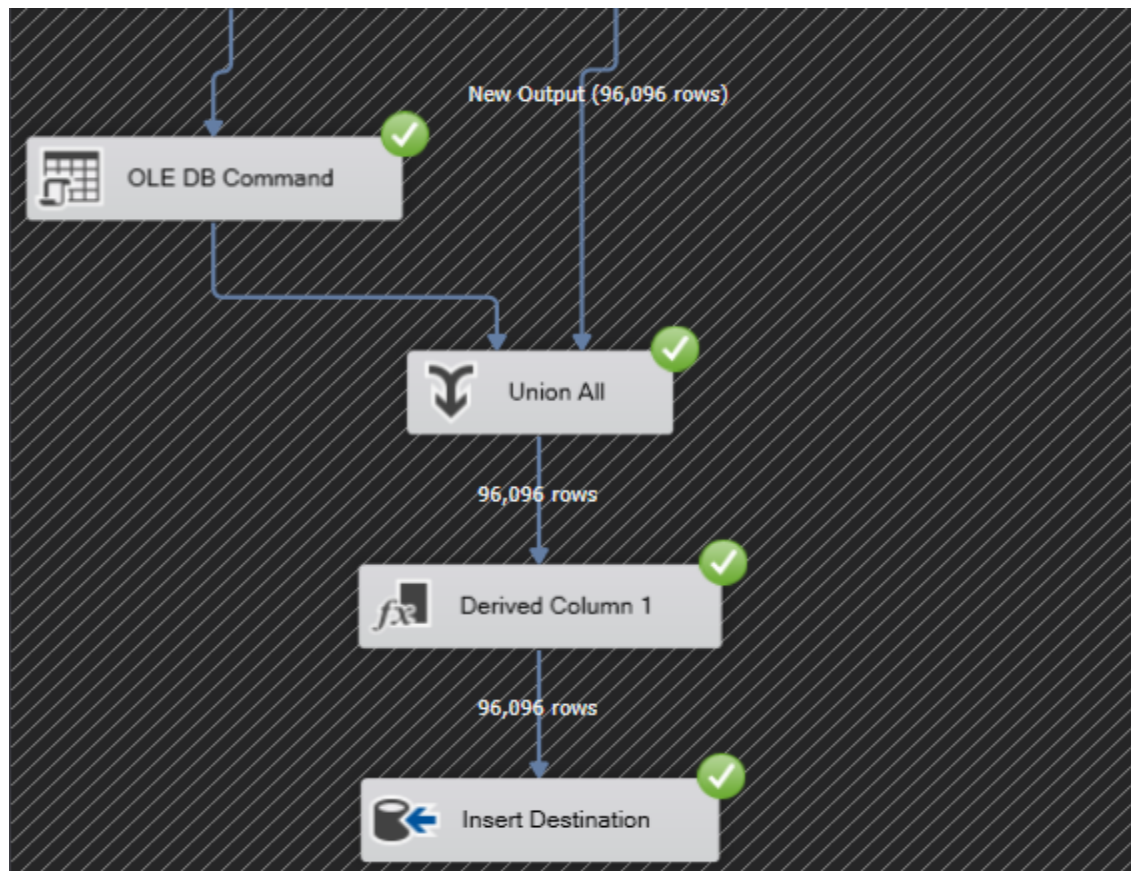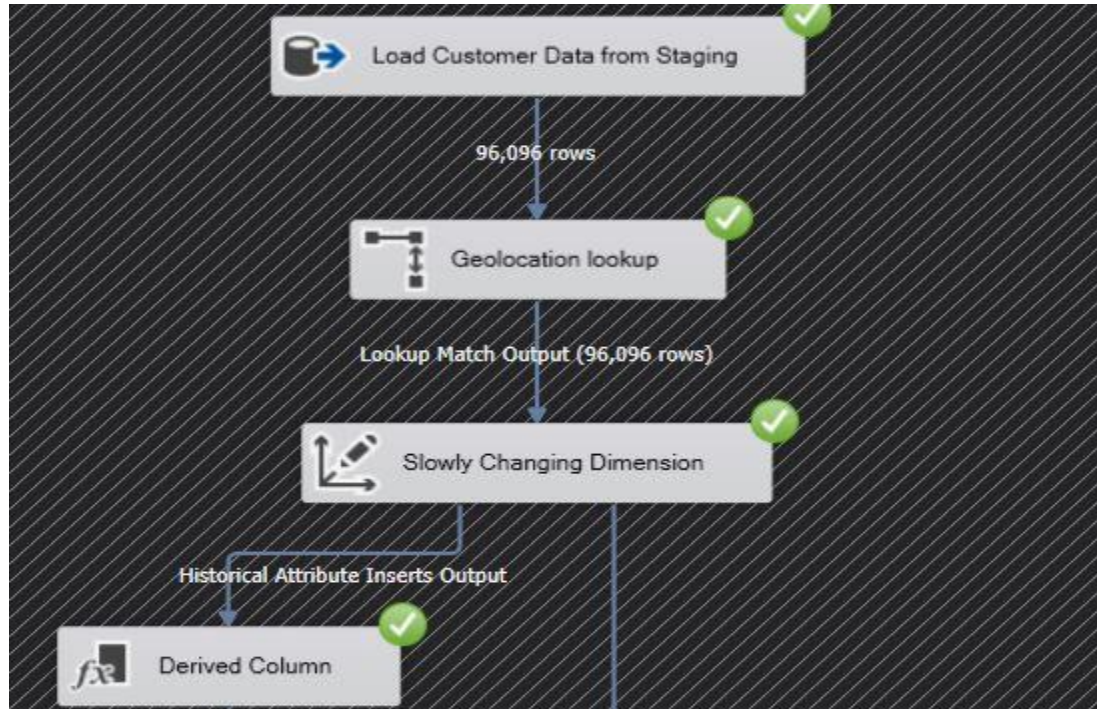
## c) Transform and Load Sellers Data

```
    Procedure used to identify whether to insert or update the data in DimSeller Table


CREATE PROCEDURE dbo.UpdateDimSeller
@SellerID nvarchar(50),
@SellerZip nvarchar(50),
@SellerCity nvarchar(50),
@SellerState nvarchar(10),
@Latitude float,
@Longitude float

AS BEGIN
if not exists (select SellerSK
from dbo.DimSeller
where Seller_id = @SellerID) BEGIN
insert into dbo.DimSeller
(Seller_id, Seller_zip_code_prefix, Seller_city, Seller_state, Seller_latitude,
Seller_longitude, InsertDate, ModifiedDate)
values
(@SellerID, @SellerZip , @SellerCity , @SellerState, @Latitude, @Longitude, GETDATE(),
GETDATE()) END;
if exists (select SellerSK
from dbo.DimSeller
where Seller_id = @SellerID) BEGIN
update dbo.DimSeller
set Seller_zip_code_prefix = @SellerZip, Seller_city = @SellerCity, Seller_state =
@SellerState, ModifiedDate = GETDATE()
where Seller_id = @SellerID END;

END;
```

## d)Transform and Load Payment Data

```
    Procedure used to identify whether to insert or update the data in DimPayment Table


CREATE PROCEDURE dbo.UpdateDimOrderPayment
@AlternateId int,
@OrderID nvarchar(50),
@PaymentSeq int,
@PaymentType nvarchar(50),
@PaymentIns int,
@PaymentVal float
AS BEGIN
if not exists (select PaymentSK
from dbo.DimPayment
where AlternatePaymentId = @AlternateId) BEGIN
insert into dbo.DimPayment
(AlternatePaymentId, OrderId, PaymentSequent, PaymentType, PaymentInstallment,
PaymentValue, InsertDate, ModifiedDate)
values
(@AlternateId, @OrderID, @PaymentSeq , @PaymentType , @PaymentIns, @PaymentVal,
GETDATE(), GETDATE()) END;
if exists (select PaymentSK
from dbo.DimPayment
where AlternatePaymentId = @AlternateId) BEGIN
update dbo.DimPayment
set PaymentSequent = @PaymentSeq, PaymentType = @PaymentType, PaymentInstallment =
@PaymentIns, PaymentValue = @PaymentVal, ModifiedDate = GETDATE()
where AlternatePaymentId = @AlternateId END;

END;
```

**e) Transform and Load Customer Data(Slowly Changing Dimension)**

## f) Load Order Fact Table

# Step 06:-ETL Development - Accumulating Fact Tables

In order to create the Accumulated Fact table I created a new SSIS Package called 'Brazillian_Ecom_AccumulatedFact.dtsx' and updated the accm_txn_complete_time and txn_process_time_hours as follows.

**\*\* END \*\***