# DATA WAREHOUSE IMPLEMENTATION PROJECT

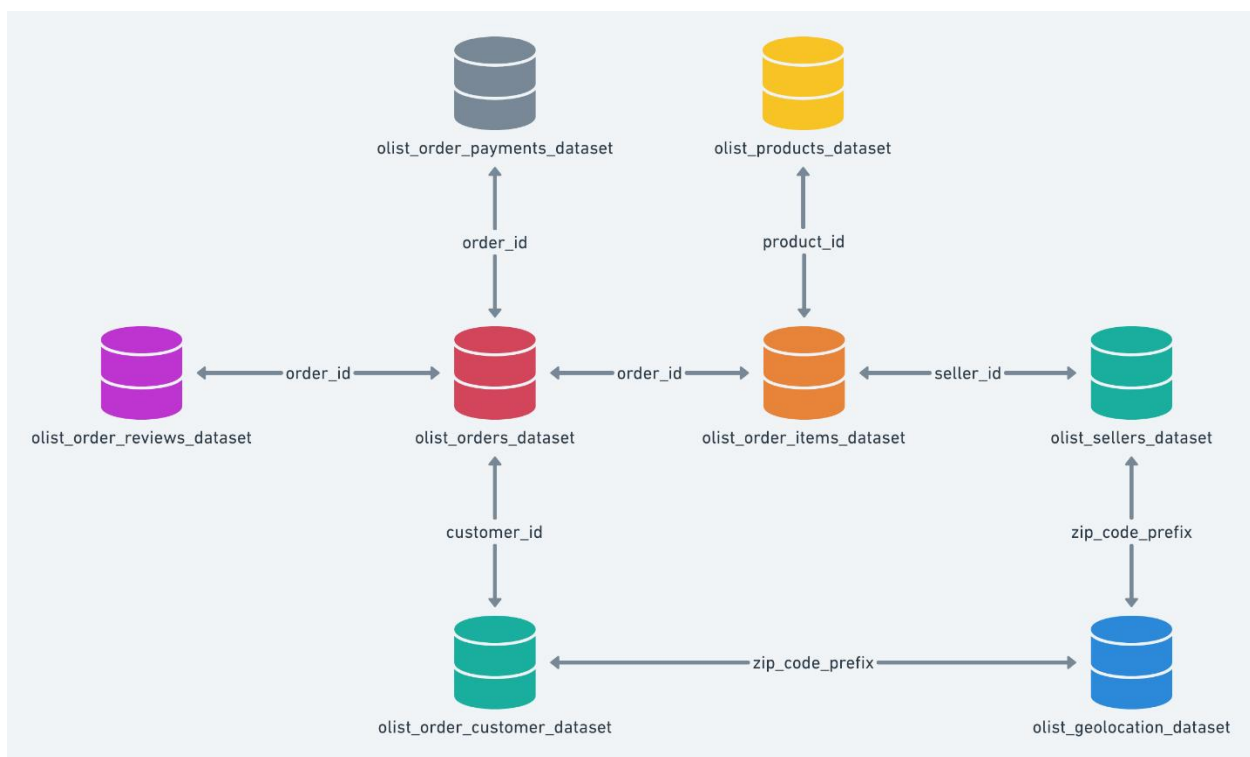Github Repo: https://github.com/divyagarigipati/datawarehouse/tree/main

Notebook: https://databricks-prod-cloudfront.cloud.databricks.com/public/4027ec902e239c93eaaa8714f173bcfc/396916982808900 9/2859625473491059/2340490761257204/latest.html

Data Source:

Raw data files

https://www.kaggle.com/datasets/olistbr/brazilian-ecommerce?select=olist_customers_dataset.csv

Initial Data Schema:



The raw e-commerce database has a relational structure, where the primary table is olist_orders_dataset which is an order lifecycle table combining different aspects. There is a olist_order_customer_dataset table where customer information is stored, and it is linked by customer_id. The olist_order_items_dataset includes order details and references olist_products_dataset` for the product details, as well as olist_sellers_dataset for seller information. In olist_geolocation_dataset, zip_code_prefix connects sellers to geographic data. Orders are linked to their payments in olist_order_payments_dataset and their reviews in olist_order_reviews_dataset. This schema provides a structured way to analyze the customer transactions, sellers' performance, delivery efficiency to fetch complete business insights.

# DATA WAREHOUSE IMPLEMENTATION PROJECT

Extracting Data:

Loading tables into data bricks hive metastore. The first row will be header, and schema will be inferred by databricks

The following images also contain the data dictionary (inferredschema) for each table.

Gelocations:

# DATA WAREHOUSE IMPLEMENTATION PROJECT

Order items:

# DATA WAREHOUSE IMPLEMENTATION PROJECT

Order Payments:

# DATA WAREHOUSE IMPLEMENTATION PROJECT

Order reviews

# DATA WAREHOUSE IMPLEMENTATION PROJECT

Orders:

# DATA WAREHOUSE IMPLEMENTATION PROJECT

Products:

# DATA WAREHOUSE IMPLEMENTATION PROJECT

Sellers:

# DATA WAREHOUSE IMPLEMENTATION PROJECT

Customers:

# DATA WAREHOUSE IMPLEMENTATION PROJECT

Translation:



Loading the data into dimensions:

# DATA WAREHOUSE IMPLEMENTATION PROJECT

Displaying Schema of dimension:



Customer_dim
root
 |-- customer_id: string (nullable = true)
 |-- customer_unique_id: string (nullable = true)
 |-- customer_zip_code_prefix: integer (nullable = true)
 |-- customer_city: string (nullable = true)
 |-- customer_state: string (nullable = true)

Geolocation_dim
root
 |-- geolocation_zip_code_prefix: integer (nullable = true)
 |-- geolocation_lat: double (nullable = true)
 |-- geolocation_lng: double (nullable = true)
 |-- geolocation_city: string (nullable = true)
 |-- geolocation_state: string (nullable = true)

Order_itmes_dim
root
 |-- order_id: string (nullable = true)
 |-- order_item_id: integer (nullable = true)
 |-- product_id: string (nullable = true)
 |-- seller_id: string (nullable = true)
 |-- shipping_limit_date: timestamp (nullable = true)
 |-- price: double (nullable = true)
 |-- freight_value: double (nullable = true)

# DATA WAREHOUSE IMPLEMENTATION PROJECT

Order_payments_dim
root
 |-- order_id: string (nullable = true)
 |-- payment_sequential: integer (nullable = true)
 |-- payment_type: string (nullable = true)
 |-- payment_installments: integer (nullable = true)
 |-- payment_value: double (nullable = true)

Order_reviews_dim
root
 |-- review_id: string (nullable = true)
 |-- order_id: string (nullable = true)
 |-- review_score: string (nullable = true)
 |-- review_comment_title: string (nullable = true)
 |-- review_comment_message: string (nullable = true)
 |-- review_creation_date: string (nullable = true)
 |-- review_answer_timestamp: timestamp (nullable = true)

Orders_dim
root
 |-- order_id: string (nullable = true)
 |-- customer_id: string (nullable = true)
 |-- order_status: string (nullable = true)
 |-- order_purchase_timestamp: timestamp (nullable = true)
 |-- order_approved_at: timestamp (nullable = true)
 |-- order_delivered_carrier_date: timestamp (nullable = true)
 |-- order_delivered_customer_date: timestamp (nullable = true)
 |-- order_estimated_delivery_date: timestamp (nullable = true)


Products_dim
root
 |-- product_id: string (nullable = true)
 |-- product_category_name: string (nullable = true)
 |-- product_name_lenght: integer (nullable = true)
 |-- product_description_lenght: integer (nullable = true)
 |-- product_photos_qty: integer (nullable = true)
 |-- product_weight_g: integer (nullable = true)
 |-- product_length_cm: integer (nullable = true)
 |-- product_height_cm: integer (nullable = true)
 |-- product_width_cm: integer (nullable = true)

Sellers_dim
root
 |-- seller_id: string (nullable = true)
 |-- seller_zip_code_prefix: integer (nullable = true)
 |-- seller_city: string (nullable = true)
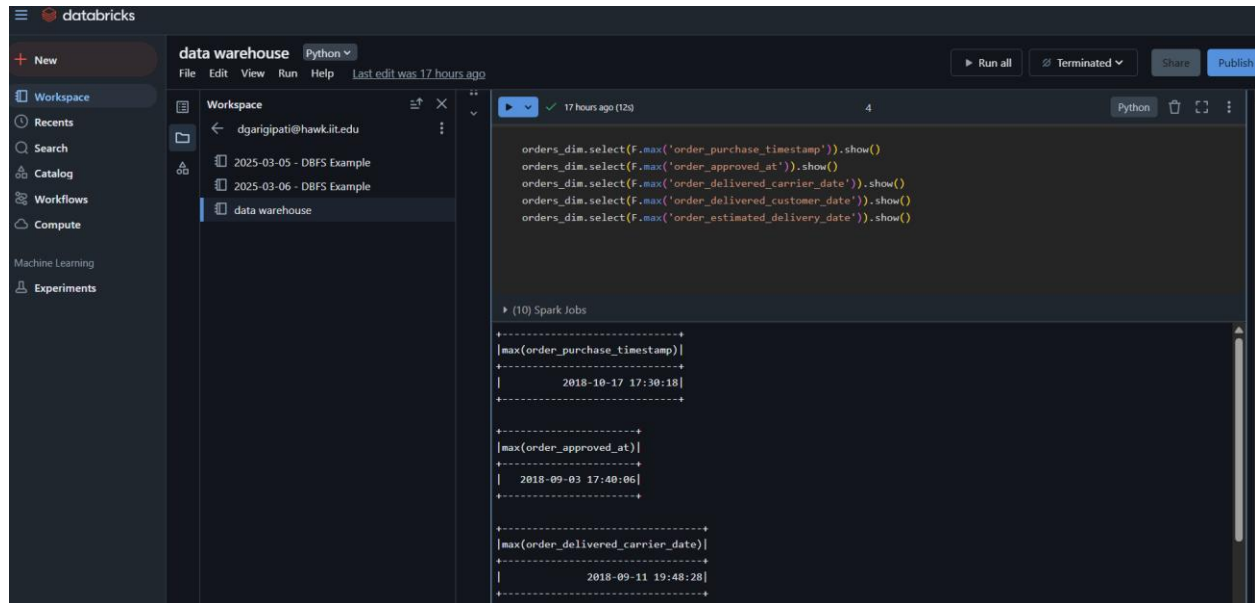 |-- seller_state: string (nullable = true)

# DATA WAREHOUSE IMPLEMENTATION PROJECT

Prod_transalation_dim
root
 |-- product_category_name: string (nullable = true)
 |-- product_category_name_english: string (nullable = true)
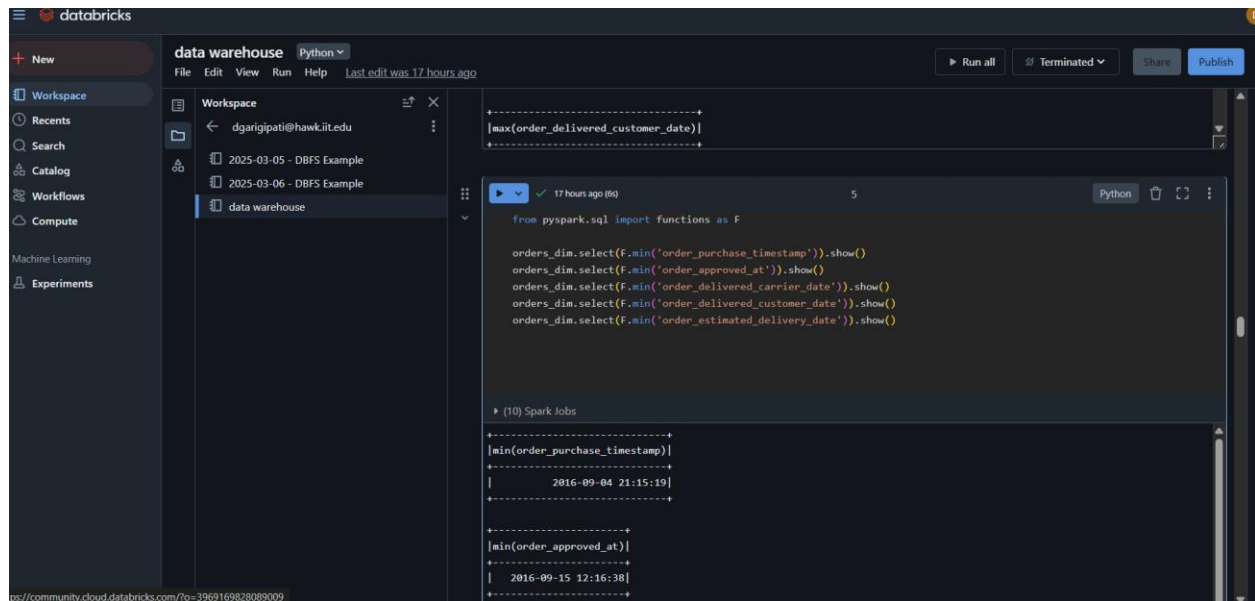
Getting Max and Min dates to generate date dimension:

# DATA WAREHOUSE IMPLEMENTATION PROJECT

Generating Date Dim:



Creating fact table by joining all the dimensions:

# DATA WAREHOUSE IMPLEMENTATION PROJECT

Fact and dimension ERD in star Schema:



The ERD represents a star schema designed for an e-commerce sales analysis system. At the center is the **fact_table** table, which records transactional data, including order details, prices, freight costs, payment values, and review scores. It is linked to multiple dimension tables that provide descriptive context.

The **customer_dim** stores customer details, while **sellers_dim** contains seller information. **products_dim** holds product attributes, and its relationship with **prod_translation_dim** allows category name translations. **orders_dim** tracks various order-related timestamps, enabling order lifecycle analysis. **date_dim** facilitates time-based reporting, and **geolocation_dim** maps geographic details using zip code prefixes. The structure optimizes analytical queries by reducing data redundancy and improving performance.

# DATA WAREHOUSE IMPLEMENTATION PROJECT

Analytical Queries:

Average Fulfillment Time by Product Category

From the plot we can see that the artes_e_artesanato category have the lowest fulfillment within 5 days and then movies_escritorio have the highest fulfillment may take up to 20 days. With average fulfillment in 10 days range for all other categories.

# DATA WAREHOUSE IMPLEMENTATION PROJECT

Year-Over-Year Revenue:

From the below graph we can see that there is huge gap between 2016 and 2017 this might be due to the fact that we are missing data form 2016. From 2017 to 2018 the revenue has increased from around 6500000 to about 8750000 that's equal to 30% increase.

# DATA WAREHOUSE IMPLEMENTATION PROJECT

Average Order Value By Month To Get Seasonal Trends:

We can see that our data starts from 9/2016 this might've been contributed to the low revenue in previous analysis. We can see that from 1-2017 every month we have average orders of around 120/month.

# DATA WAREHOUSE IMPLEMENTATION PROJECT

Summary:

The project implements a data warehouse by converting  a raw, normalized e-commerce dataset into a star schema dimensional model. The central fact table (fact_sales) stores key transactional data, while surrounding dimension tables (e.g., customer_dim, products_dim, orders_dim) provide descriptive context for each dimension. Using PySpark, I've efficiently processed complex queries such as Year-over-Year revenue, top-selling categories, and fulfillment times. This approach streamlined the data model, improved performance, and enabled flexible business intelligence reporting.

Challenges included data quality issues, ambiguous column names, and data type mismatches, which were resolved by renaming columns and casting types. The star schema reduced query complexity and improved performance. Future improvements include implementing surrogate keys, enhancing the ETL pipeline for better data quality and optimizing query performance and report automation for better scalability and maintenance.