# Module 2 - Assignment Project

### Team 4, Data Science and AI (Full time, Batch 1)

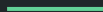## E-commerce Data pipeline and Analysis

# Members

Team 4

ian
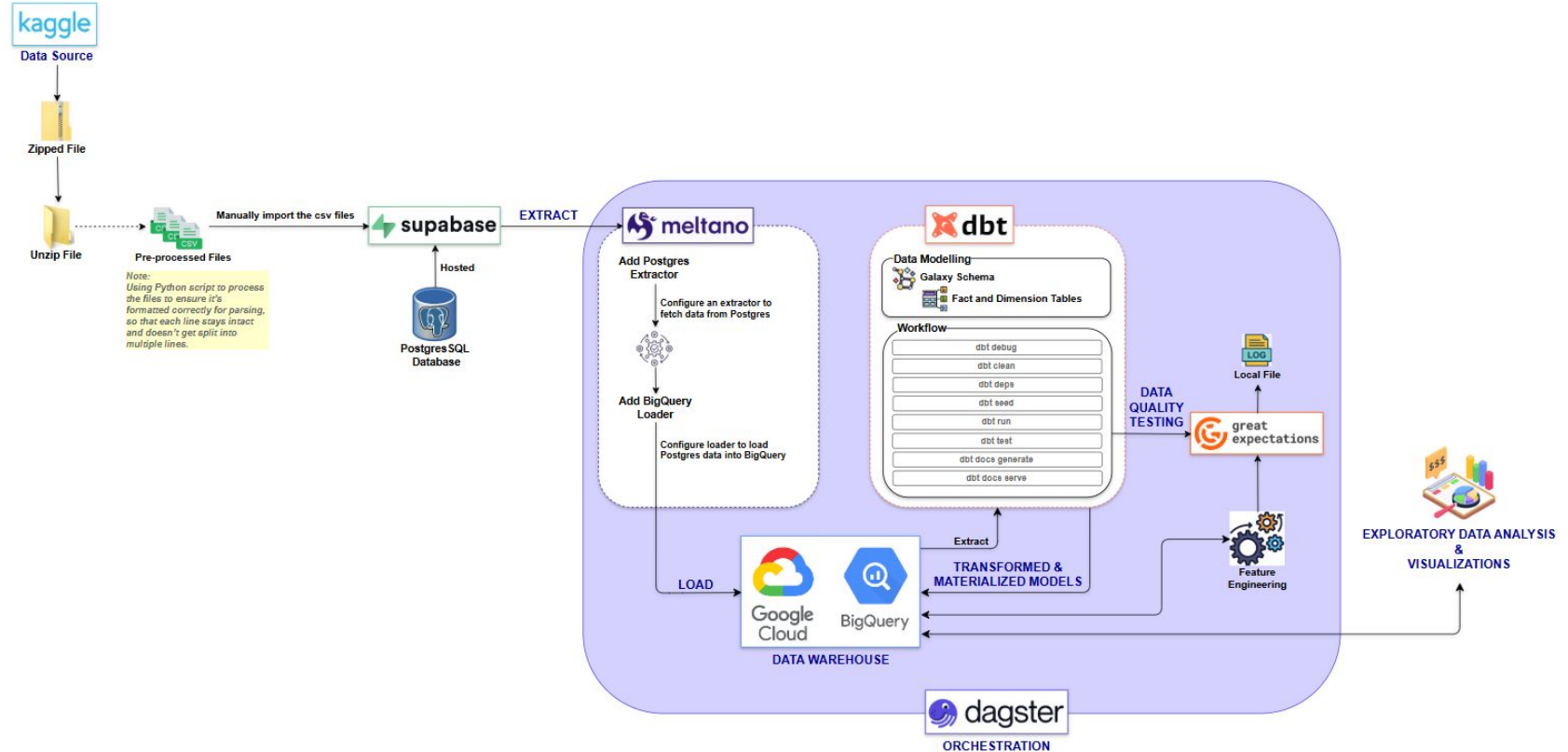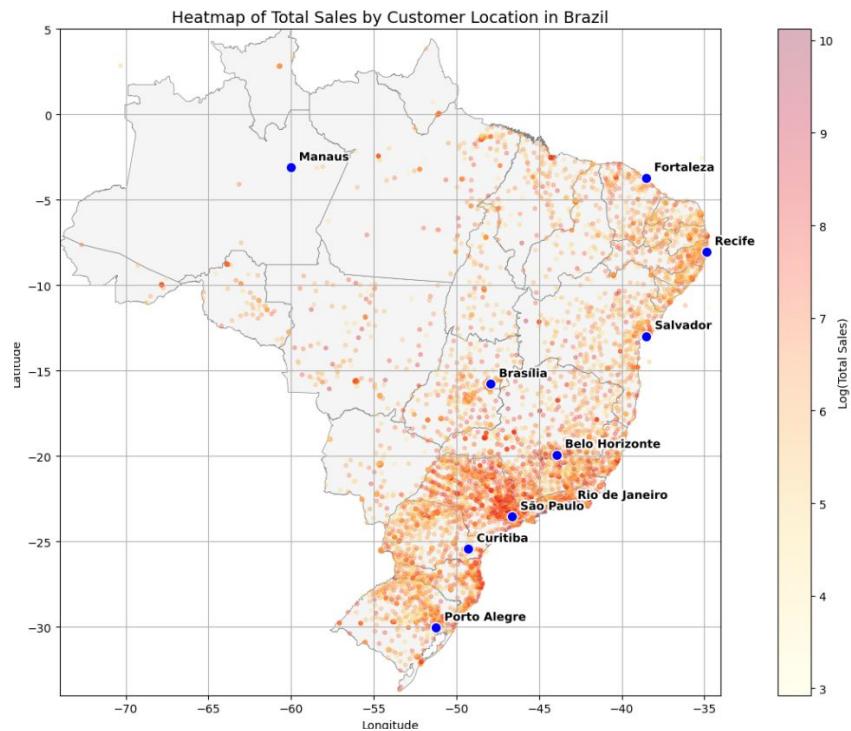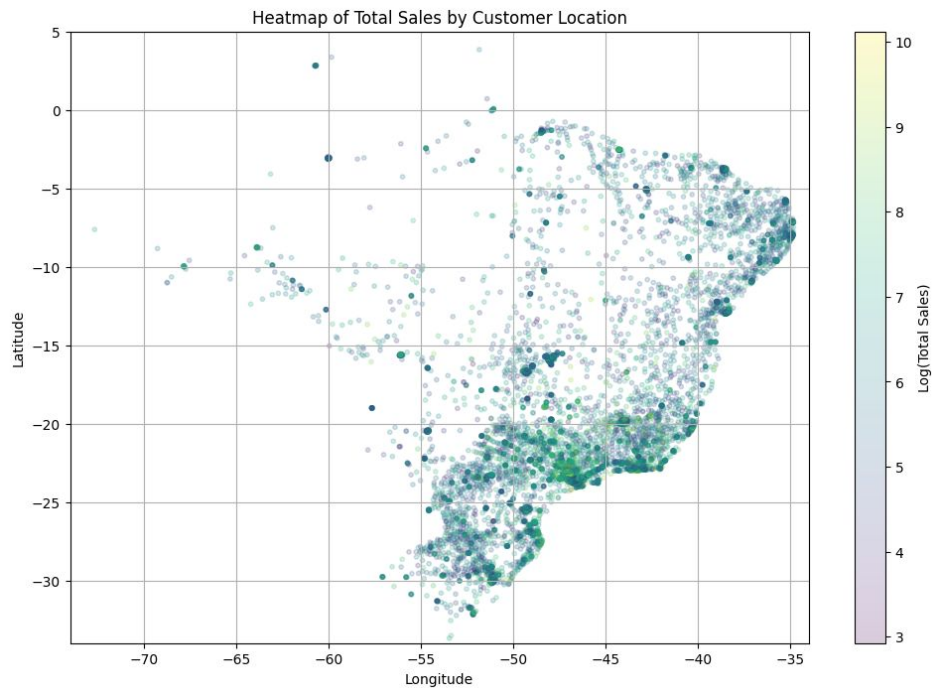
siew wen

eve

selena

andrew

# Project Overview

an end to end approach

| | |
|---:|:---|
| SOURCE | kaggle |
| DATA | supabase |
| WAREHOUSE | bigquery |
| EXTRACT & LOAD | meltano |
| TRANSFORM | dbt |
| VALIDATE | great expectations |
| ANALYSIS | pandas |
| ORCHESTRATE | dagster |

# Data Pipeline Architecture

# Geo heatmap sales

red= higher sales



Heatmap of Total Sales by Customer Location



Heatmap of Total Sales by Customer Location in Brazil

# Monthly sales trend for top 5 cities
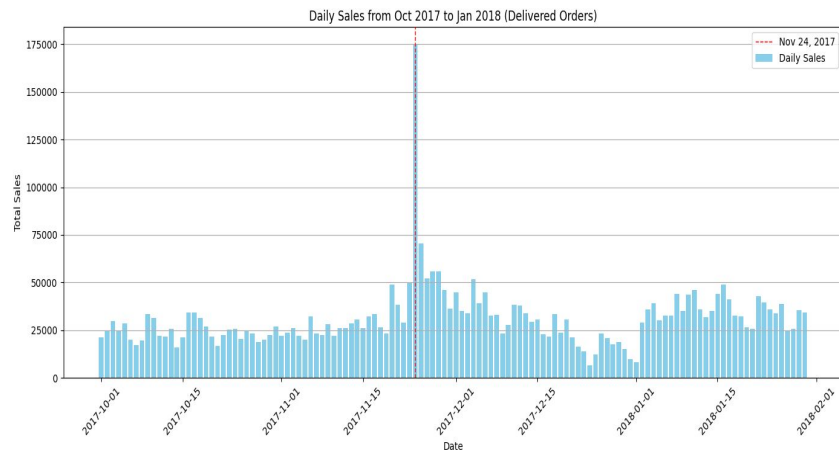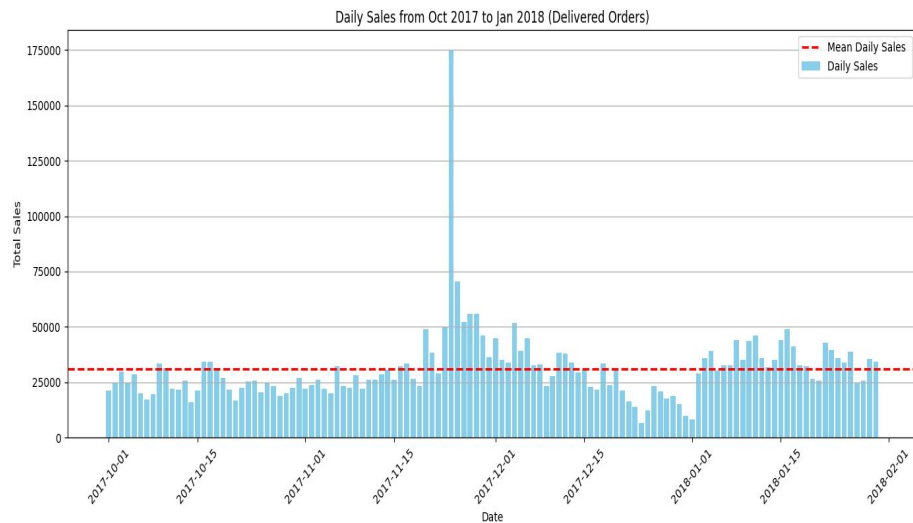


Monthly Sales Trend for Top 5 Buyer Cities

# Daily sales highlighting black friday sales



Nov 24
black friday

# Data Ingestion

# Grant User Access

## Logs Explorer

# Olist Brazilian Ecommerce - Galaxy Schema

# Galaxy Schema -  When One Star Isn't Enough



**Why It Matters:**

Share Dimension

Avoid Cartesian Joins

Granular Control

# Data Lineage in dbt

# Pre-Aggregated Metrics

## Why It Matters:

| Improve Report Performance | Consistency across reports | Simplified Storytelling |

| payment_status | total_payment | order_amt | freight_amt | amt_wf_freight | balance_amt |
|---|---|---|---|---|---|
| completed | 14296050.9 | 12268849.16 | 2024130.57 | 14292979.73 | -3071.17 |
| in progress | 1550229.27 | 1322659.57 | 227770.48 | 1550430.05 | 200.78 |
| order items not found | 162591.95 | null | null | null | null |
| canceled | 0.0 | null | null | null | null |
| payments not found | null | 134.97 | 8.49 | 143.46 | null |

Screenshot: Pre-aggregated columns from DIM_ORDER in BigQuery.

# Lookup Tables for Decoupling Logic

The LKP_STATUS_DESC table adds business logic without cluttering DIM_ORDERS.

| There columns are derived by query | | | These columns are user input | |
|---|---|---|---|---|
| order_status | payment_status | record_count | profit_lost | status_description |
| Invoiced | Completed | 273 | Profit | Invoice processed, payment received. |
| Processing | Completed | 269 | Profit | Order finalized with successful payment. |
| Unavailable | Completed | 6 | Profit | Payment received despite item unavailability. |
| Canceled | Completed | 411 | Profit | Payment was completed despite order cancellation. |
| Shipped | Completed | 1004 | Profit | Shipment delivered with payment finalized. |
| Delivered | Completed | 86973 | Profit | Successfully fulfilled with payment received. |
| Approved | Completed | 1 | Profit | Successfully processed and finalized. |
| Processing | In Progress | 32 | In Progress | Awaiting completion of payment processing. |
| Approved | In Progress | 1 | In Progress | Awaiting payment completion. |
| Shipped | In Progress | 102 | In Progress | Goods dispatched, awaiting payment. |
| Canceled | In Progress | 50 | In Progress | Payment processing despite order cancellation. |
| Delivered | In Progress | 9504 | In Progress | Payment processing post-delivery. |
| Invoiced | In Progress | 39 | In Progress | Payment still being finalized. |
| Invoiced | Order Items Not Found | 2 | Unknown | Invoiced but items missing from order. |
| Canceled | Order Items Not Found | 161 | Unknown | Items missing, causing uncertainty in status. |
| Unavailable | Order Items Not Found | 603 | Unknown | Missing items prevent proper status evaluation. |
| Created | Order Items Not Found | 5 | Unknown | Newly created order but items not located. |
| Shipped | Order Items Not Found | 1 | Unknown | Shipped but item data missing. |
| Delivered | Payments Not Found | 1 | Lost | Payment missing despite delivery completion. |
| Canceled | Canceled | 3 | Lost | Transaction voided; no revenue gained. |

# Ensuring Primary Key Integrity in dbt

## **DBT Test & DBT Expectations**

Implemented 13 dbt tests to validate that all primary keys are

🔑 **unique**

🚫 **not null**

❌ **Failure Handling:**
Failing rows are automatically logged in
📁 *target/run_results.json*

**Why It Matters:**
🛢️ Prevents duplicate records
🔗 Clean join across model
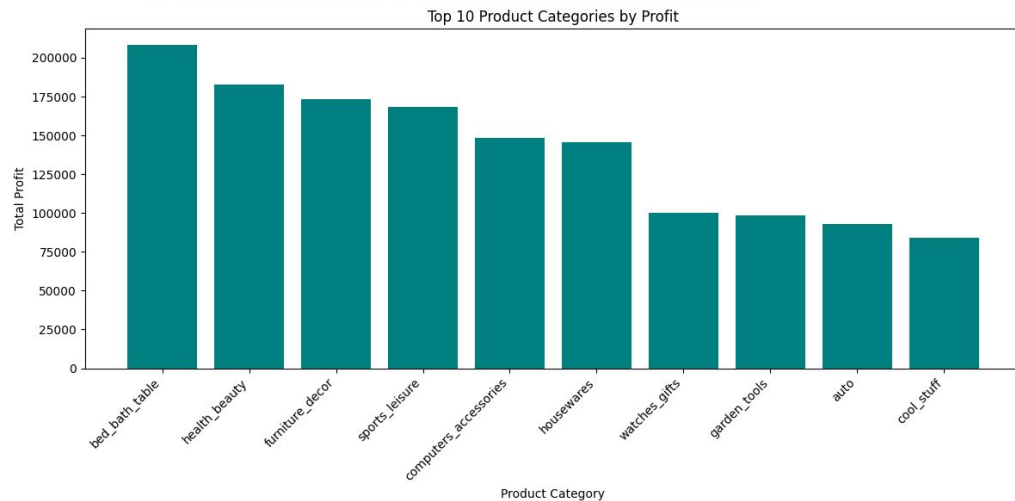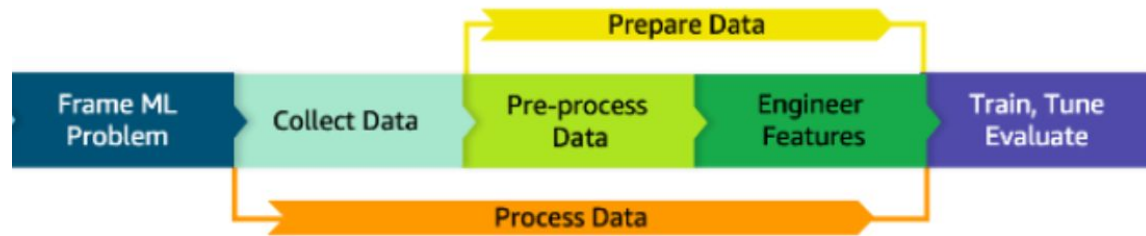📊 Trustworthy downstream reporting

# Role of feature engineering

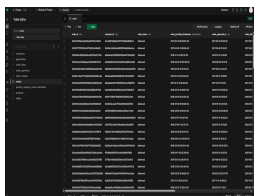Data reusability

Feature creation, transform

One to many mapping

Solution?

Aggregation





Top 10 Product Categories by Profit

# Data Quality

**great expectations**

Search...

Your data assets: database tables, flat files, dataframes...

Data validation with Great Expectations

High quality data in your data products

Data documentation & data quality reports

Logging & alerting

DATA QUALITY

1 DEFINITION
2 ASSESSMENT
3 ANALYSIS
4 IMPROVEMENT
5 IMPLEMENTATION
6 CONTROL

Data Governance
Ownership
Accessibility
Security
Quality
Knowledge
Technology
People
Process

great-expectations 1.5.2

Data logging

Cloud Logging
Machines & Equipment
Insights
Maintenance & Process optimisation
Customer

```
suite = gx.ExpectationSuite(name=suite_name)
suite = context.suites.add(suite)

# Add ExpectColumnValuesToBeOfType expectations for each expected column
for column, column_type in expected_columns.items():
    expectation = gx.expectations.ExpectColumnValuesToBeOfType(
        column=column, type_=column_type
    )
    suite.add_expectation(expectation)

# Create validation definition
definition_name = f"{table_name}_validation_definition"
validation_definition = gx.ValidationDefinition(
    data=batch_definition, suite=suite, name=definition_name
)

# Run validation
validation_results = validation_definition.run(batch_parameters=batch_parameters)
print(f"Validation results for {table_name}:")
print(validation_results)
```

```
# Save full results to file
output_folder = "gx_output"
os.makedirs(output_folder, exist_ok=True)
result_path = os.path.join(output_folder, "gx_results_geo.txt")

with open(result_path, "w") as f:
    f.write(pprint.pformat(validation_results))

print(f" Full GX test results saved to {result_path}")
```

Full GX test results saved to gx_output/gx_full_results.txt

# Orchestration

5 Stages:

Meltano

dbt run

dbt test

Feature engineering

Great expectation

# Q&A

How about some questions...

# Benefits of Galaxy Schema on Olist Dataset

| Shared Dimensions | Avoids Cartesian Joins | Granular Control |
|---|---|---|
| All dimension tables are shared across multiple fact tables, reducing redundancy and promoting consistency. | Clear one-to-many relationships between DIM_ORDERS and each fact table eliminate unintentional row multiplication. | Supports fine-grained analysis and improves query efficiency by separating transactional components. |
| | **Fact Tables Differ in Granular Detail:**<br>- DIM_ORDERS: One record per unique order<br>- FCT_ORDER_ITEMS: Captures multiple items per order<br>- FCT_PAYMENTS: Tracks multiple payments or installments per order<br>- FCT_REVIEWS: Stores multiple customer reviews per order | |

# Benefits of Lookup Tables

| Business Logic Separation | Code Translation | Explanation Standardization | Debugging |
|---|---|---|---|
| Keeps DIM_ORDERS clean by offloading business logic to a dedicated table. | Converts technical codes into business-friendly descriptions. | Ensures consistent language across dashboards, reports, and analytical outputs. | Improves traceability. |

# Pre-Aggregated & Derived Metrics

Listed below are the derived columns in DIM_ORDERS and the SQL logic used to generate them.

| Derived Column Name | Transformation |
|---|---|
| dim_order.total_payment | SELECT order_id, SUM(payment_value) AS total_payment<br>FROM {{ source('olist_brazilian_ecommerce', 'public_order_payments') }}<br>GROUP BY order_id |
| dim_order.order_amt | SELECT order_id, SUM(price) order_amt,<br>FROM {{ source('olist_brazilian_ecommerce', 'public_order_items') }}<br>GROUP BY order_id |
| dim_order.freight_amt | SELECT order_id, SUM(freight_value) as freight_amt,<br>FROM {{ source('olist_brazilian_ecommerce', 'public_order_items') }}<br>GROUP BY order_id |
| dim_order.<br>total_order_amt_wf_freight | SELECT order_id, SUM(price) + sum(freight_value) as total_order_amt_wf_freight<br>FROM {{ source('olist_brazilian_ecommerce', 'public_order_items') }}<br>GROUP BY order_id |
| dim_order.balance_amt | i.total_order_amt_wf_freight - p.total_payment |
| dim_order.payment_status | CASE<br> WHEN p.total_payment>0 and i.order_amt>0 and i.total_order_amt_wf_freight - p.total_payment = 0 THEN 'completed'<br> WHEN p.total_payment>0 and i.order_amt>0 and i.total_order_amt_wf_freight - p.total_payment < 0 THEN 'completed'<br> WHEN p.total_payment>0 and i.order_amt>0 and i.total_order_amt_wf_freight - p.total_payment > 0 THEN 'in progress'<br> WHEN p.total_payment > 0 and i.order_amt IS NULL THEN 'order items not found'<br> WHEN o.order_status in ('delivered','shipped') and (p.total_payment IS NULL OR p.total_payment <= 0) THEN 'payments not found'<br> WHEN o.order_status in ('canceled','unavailable') and (p.total_payment IS NULL OR p.total_payment <= 0) THEN 'canceled'<br> ELSE NULL<br>END AS payment_status, |