

CART CATALYST

GENERATION SINGAPORE

GROUP MEMBERS:

ALEX YONG

FAHMY M

HO CHOO GEOK

JANET KEH



CONTENT

- 1) Introduction
 - Executive Summary
 - Problem Statement
 - Project Objectives
- 2) Project Architecture Design
 - Project Tech Stack
 - Pipeline Architecture
 - Medallion Architecture
 - Bronze Lakehouse
 - Silver Lakehouse
 - Gold Lakehouse
 - Fabric Pipeline Design
- 3) Market Basket Analysis
- 4) Monitoring
- 5) Testing
- 6) Validation
 - Validation of Pipeline Activities
 - Validation of Pipeline Flow
- 7) Project Deliverables
 - Project Pipeline
 - Visualization and Reporting
- 8) Challenges and Future Roadmap

INTRODUCTION

EXECUTIVE SUMMARY

Olist, a leading Brazilian e-commerce platform, faced significant challenges due to its fragmented data architecture. With disparate datasets spread across multiple sources—including customer profiles, order details, product catalogs, seller information, and geolocation data—performance tracking and strategic decision-making were severely constrained. The lack of centralized governance and real-time visibility made it difficult for leadership to derive actionable insights and respond swiftly to market dynamics.

To address these limitations, the project will design and implement a modular, scalable data pipeline using Microsoft Fabric, a unified analytics platform that seamlessly integrates data engineering, data science, and business intelligence.

PROBLEM STATEMENT

Olist's leadership faces significant challenges in tracking business performance due to inconsistent, fragmented, and messy data spread across multiple sources. The absence of a centralized data architecture has made it difficult to generate reliable insights, hindering strategic planning and operational efficiency. Existing reporting mechanisms are limited in scope, lacking clarity, drill-down functionality, and real-time visibility, leaving decision-makers without the timely, actionable intelligence they need to steer the business effectively.

PROJECT OBJECTIVES

The primary objective of this project is to establish a unified data foundation by consolidating Olist's fragmented data sources into a Medallion architecture Lakehouse using Microsoft Fabric. This ensures consistent data access and governance across the organization.

Next, we aim to design a scalable and modular data pipeline by building reusable ETL components within Fabric Notebooks. Transformation logic, validation checks, and error handling will be abstracted into modular pipeline blocks to promote flexibility, maintainability, and ease of expansion.

Finally, the project seeks to enable drill-down reporting by delivering real-time dashboards equipped with drill-through capabilities and semantic modeling. These dashboards will empower business leaders with deeper insights and interactive visualizations for informed decision-making.

PROJECT TECH STACK

The project leverages a modern, cloud-native tech stack designed for scalability, modularity, and real-time analytics. Each component plays a distinct role across the data lifecycle—from ingestion to transformation, visualization, and monitoring.

1) Infrastructure & Storage

This foundation ensures scalable, ACID-compliant storage with schema flexibility and version control.

Azure Blob Storage: Initial landing zone for raw CSVs from Olist.

OneLake: Unified storage layer enabling seamless access across Microsoft Fabric workloads.

Lakehouse: Centralized repository for structured Delta tables across Bronze, Silver, and Gold layers.

2) Workflow Orchestration

Enables modular workflows with fabric pipeline design for flexible scheduling and execution.

Azure Fabric: Execution environment for notebooks and pipeline automation.

3) Ingestion & Transformation

Supports reusable pipeline blocks with embedded validation and error handling logic, enabling rapid iteration and debugging.

Fabric Notebooks: Modular ETL logic for cleaning, enrichment, and validation.

Apache PySpark: Distributed processing for scalable transformations.

Pandas: Lightweight data manipulation for targeted operations.

4) Visualization & Semantic Modeling

Bridges technical outputs with business insights, democratizing access to KPIs and enabling data-driven decisions.

Power BI: Interactive dashboards for performance monitoring and decision support.

Semantic Model: Star schema with DAX measures and hierarchies for drill-through analytics.

Fabric Experiment: ML experimentation environment for market basket analysis.

Power BI App: Centralized access to curated dashboards for stakeholders.

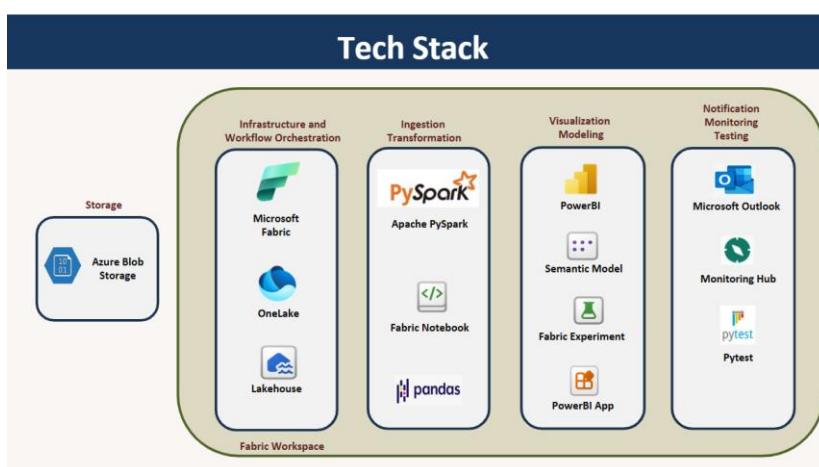
5) Notification, Monitoring & Testing

Ensures proactive observability, robust error handling, and traceable validation across the Bronze-Silver-Gold pipeline.

Microsoft Outlook: Email alerts triggered by pipeline failures or validation errors.

Monitoring Hub: Centralized log repository with Power BI integration for pipeline health tracking.

Pytest: Unit and integration testing of transformation logic and schema compliance.



PROJECT PIPELINE ARCHITECTURE

The pipeline is designed using a medallion architecture built on Microsoft Fabric. This modular structure ensures scalable ingestion, robust transformation, and ML-ready outputs, all governed by embedded validation and error handling.

1) Data Sources

- Azure Blob Storage serves as the initial landing zone for raw CSVs from Olist (orders, customers, products, payments, reviews, etc.).
- Files are ingested into Lakehouse using copy data, which load files from Azure Blob Storage into the Bronze Lakehouse Files folder.

2) Bronze Lakehouse: Ingestion Layer

Tasks:

- Copy raw CSVs from Blob Storage.
- Convert to Delta format and store in Lakehouse tables.

Validation & Error Handling:

- Schema-on-read validation.
- Overwrite mode ensures clean ingestion each run.
- Errors logged and surfaced via monitoring dashboard.

3) Silver Lakehouse: Cleaning Layer

Tasks:

- Standardized data types (e.g., timestamps, numerical fields).
- Normalize text fields and enrich with derived columns (e.g., purchase_year).
- Remove nulls from critical columns.

Validation & Error Handling:

- Type consistency checks.
- Null value enforcement.
- Schema validation between Bronze and Silver layers.

4) Gold Lakehouse: Transformation Layer

Tasks:

- Read curated tables from Silver Lakehouse.
- Merge into a unified, enriched dataset for ML and reporting.

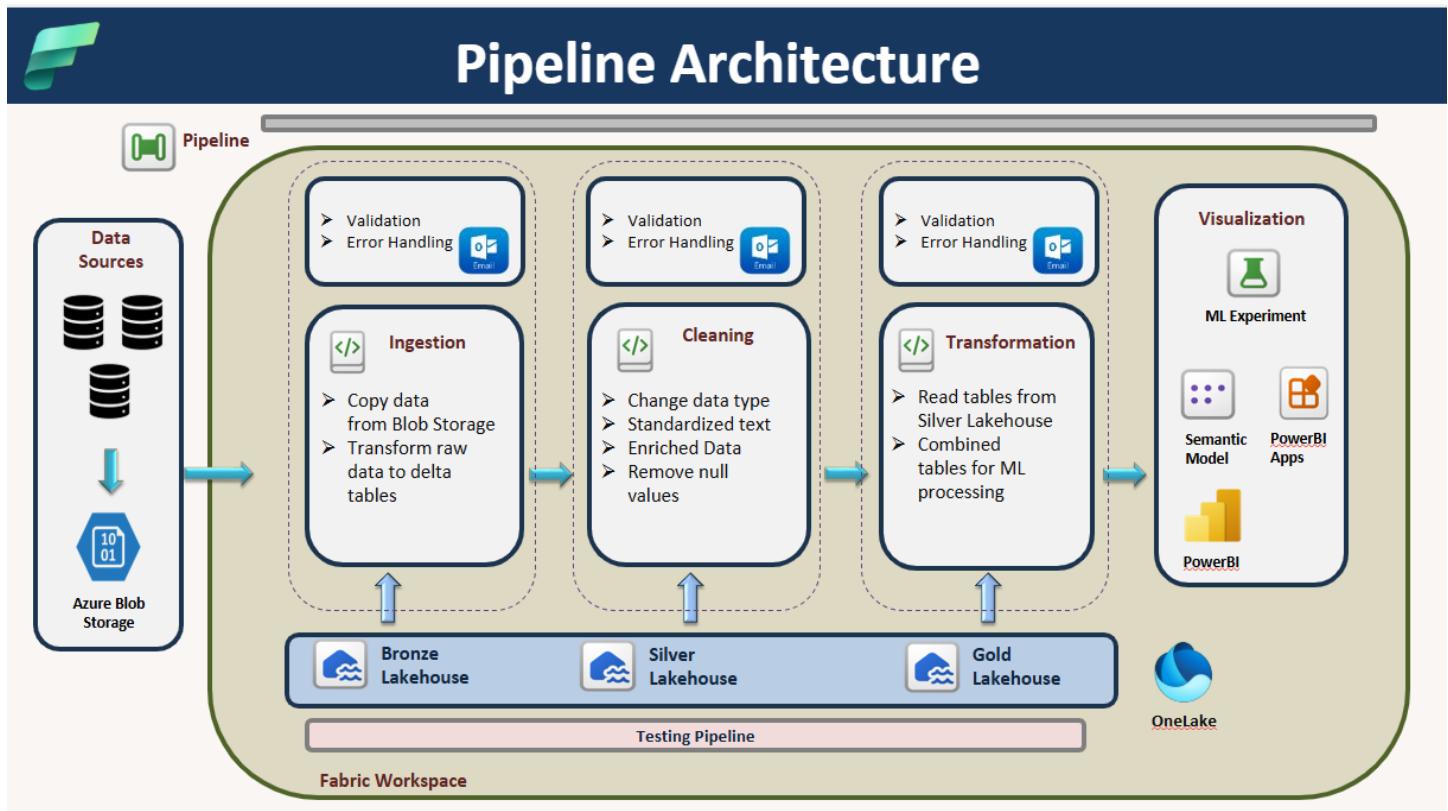
Validation & Error Handling:

- Row count reconciliation between Silver and Gold.
- Schema alignment checks.
- Error logging and alerting via Outlook integration.

5) Visualization & ML Integration

Outputs:

- ML Experiment: Market Basket Analysis using Apriori algorithm.
- Semantic Model: Star schema with DAX measures and hierarchies.
- Power BI Apps: Dashboards for ingestion health, sales performance, and product insights.



MEDALLION ARCHITECTURE

The project adopts Medallion Architecture (Bronze, Silver, and Gold layers) within Microsoft Fabric Lakehouse to ensure data quality, flexibility, and governance across the analytics lifecycle. Each layer serves as a progressive checkpoint, refining data and enabling modular, traceable workflows.

Key Reasons to Adopt Medallion Architecture

1) Improved Data Quality & Integrity

- Each layer progressively refines the data:
 - **Bronze:** Raw, unfiltered source data
 - **Silver:** Cleaned, deduplicated, and standardized
 - **Gold:** Curated, enriched, and business-ready
- This layered refinement ensures consistency, traceability, and trust in downstream analytics.

2) Scalability & Performance

- Supports distributed processing for large datasets.
- Enables optimized queries by isolating transformations and aggregations in the Gold layer.

3) Flexibility & Reusability

- Raw data in Bronze can be reused to regenerate Silver or Gold when business logic changes.
- Modular design makes it easier to plug in new sources or update logic without breaking downstream flows.

4) Clear Separation of Concerns

- Engineers focus on ingestion and transformation (Bronze/Silver).
- Analysts and business users consume curated outputs (Gold).
- This role alignment reduces friction and improves collaboration.

5) Enhanced Data Governance & Lineage

- Each layer provides a checkpoint for auditing, validation, and compliance.
- It is easier to trace errors or anomalies back to their origin.

6) Cost Efficiency

- Only refined data is promoted to higher layers, reducing storage and computing costs.
- Avoids redundant transformations by centralizing logic in Silver.

Lakehouse in Medallion Architecture

1) Bronze Lakehouse: Raw Ingestion Layer

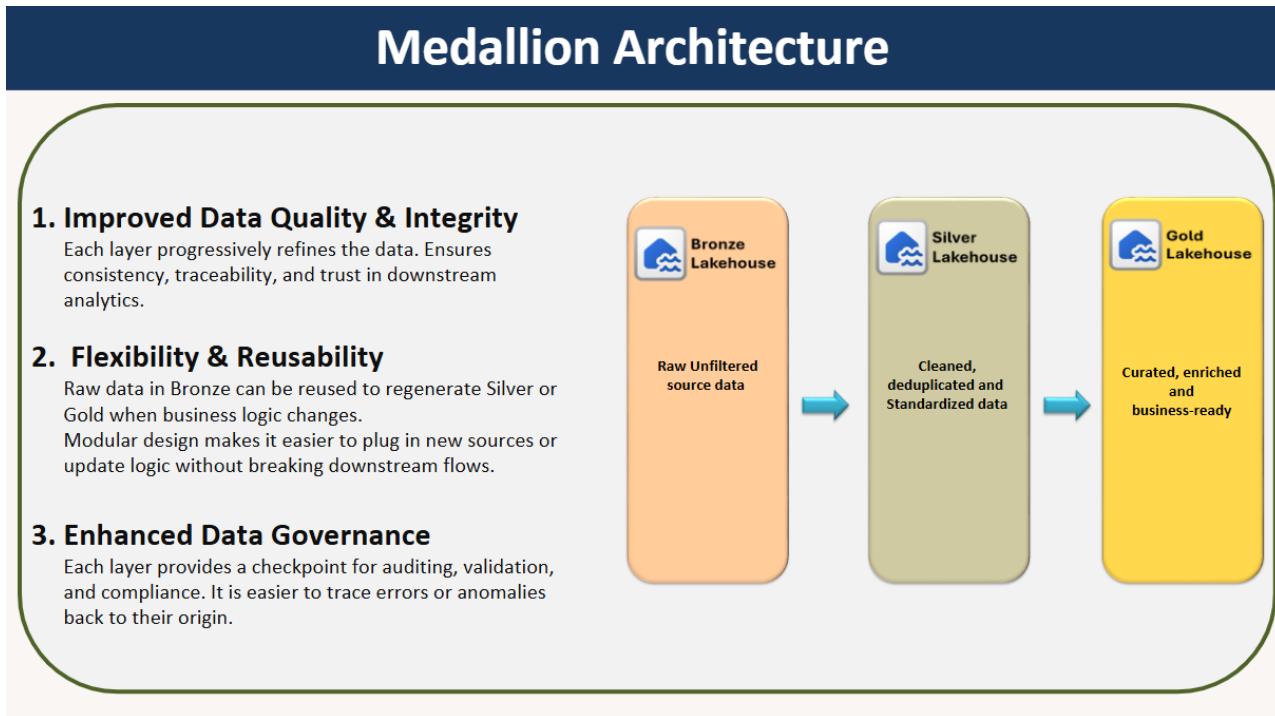
- **Purpose:** Capture raw, unfiltered source data from Azure Blob Storage.
- **Characteristics:**
 - Schema-on-read ingestion using inferred headers.
 - Delta format conversion for ACID compliance.
 - Overwrite mode ensures clean ingestion per run.
- **Role in Architecture:** Acts as the foundational layer for all downstream transformations, preserving original data fidelity.

2) Silver Lakehouse: Cleaning & Standardization Layer

- **Purpose:** Transform raw data into structured, validated datasets.
- **Characteristics:**
 - Type casting for timestamps, numeric, and strings.
 - Null value handling in critical fields.
 - Enrichment with derived columns (e.g., purchase_year).
- **Role in Architecture:** Establishes consistency and prepares data for analytical modeling and ML workflows

3) Gold Lakehouse: Enrichment & Business-Ready Layer

- **Purpose:** Consolidate and curate data for machine learning and reporting.
- **Characteristics:**
 - Merges multiple Silver Lakehouse tables into a unified schema.
 - Validates schema alignment and row count integrity.
 - Supports dimensional modeling for semantic layer integration.
- **Role in Architecture:** Delivers analytics-ready datasets for Power BI dashboards and ML experiments



BRONZE LAKEHOUSE

The Bronze Lakehouse serves as the foundational layer in the project pipeline, designed to ingest raw CSV files and convert them into Delta format for structured, traceable downstream processing. This layer prioritizes flexibility, simplicity, and fidelity to source data.

Process Flow

1. Ingestion

- Raw CSV files are loaded from Azure Blob Storage to the Files section of the Bronze Lakehouse using copy operation.
- Data originates from Olist's transactional systems (orders, customers, products, payments, reviews, etc.).

2. Conversion to Delta

- CSVs are transformed into Delta tables, enabling ACID compliance and version control.
- Tables are stored in the Bronze Lakehouse for downstream access.

3. Schema & Structure

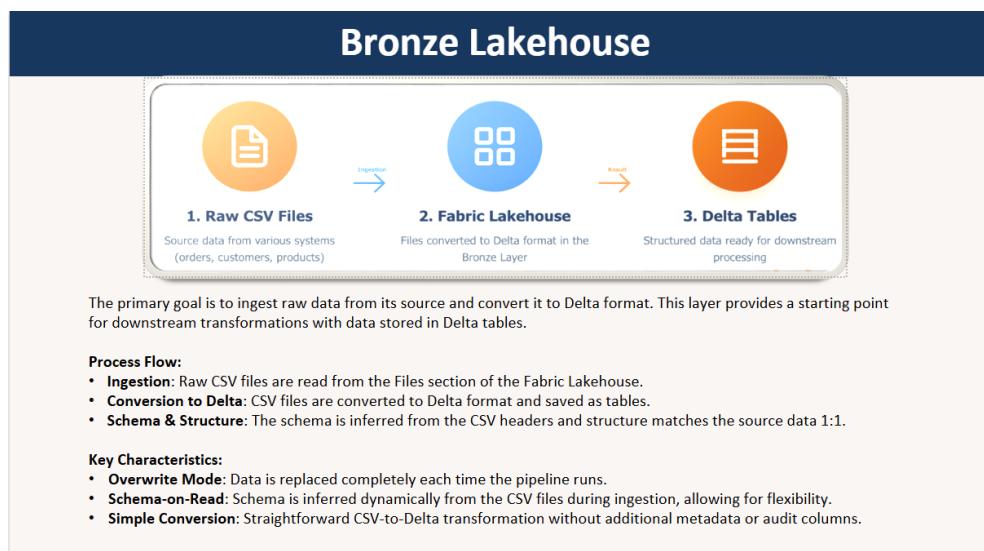
- Schema is inferred dynamically from CSV headers.
- Structure mirrors the source data 1:1, preserving original column names and formats.

Key Characteristics

- Overwrite Mode**
 - Each pipeline run replaces the existing Bronze tables entirely, ensuring clean ingestion and avoiding duplication.
- Schema-on-Read**
 - Schema is not predefined; instead, it is inferred at runtime from the CSV files, allowing flexibility in handling evolving source formats.
- Simple Conversion**
 - No enrichment, metadata tagging, or audit columns are added at this stage.
 - Focus remains on raw fidelity and straightforward transformation.

Role in Architecture

- Acts as the single source of truth for raw data.
- Enables reusability by allowing regeneration of Silver and Gold layers when business logic changes.
- Serves as a checkpoint for schema validation and ingestion monitoring.



NOTEBOOK SOURCE CODE FOR BRONZE LAKEHOUSE

Transform Raw CSV Files to Delta Tables for Bronze Lakehouse

1. Data has been loaded from the Azure Blob Storage to LH_Bronze files folder
2. The loaded csv files will be save as Delta tables in the LH_Bronze Tables folder

```
1 # To load all the raw csv files in the folder to delta tables
2
3 # Import of libraries
4 import logging
5 import re
6
7
8 # Configure of logging for monitoring
9 logging.basicConfig(level=logging.INFO)
10
11 # To list the file path
12 files = notebookutils.fs.ls("Files/")
13
14 # To loop through all CSV files in the Files folder
15 for f in files:
16     if f.name.endswith(".csv"):
17         try:
18             # Sanitize table name
19             table_name = re.sub(r"\.csv$", "", f.name)
20             table_name = re.sub(r"[^a-zA-Z0-9_]", "_", table_name)
21
22             # To read CSV file
23             df = spark.read.option("header", "true").csv(f.path)
24
25             # Write the loaded CSV file to Delta table
26             df.write.mode("overwrite").format("delta").saveAsTable(table_name)
27
28         except Exception as e:
29             print("Error {e}")
30             raise Exception("Error in notebook NB_Bronze")
```

✓ - Session ready in 12 sec 636 ms. Command executed in 1 min 31 sec 893 ms by ChooGeok H (JDE07) on 7:44:05 PM, 10/07/25

SILVER LAKEHOUSE

The Silver Lakehouse is the second stage in the project medallion architecture, designed to transform raw ingested data from the Bronze layer into clean, reliable Delta tables. This layer emphasizes type consistency, basic quality checks, and enrichment logic to prepare data for downstream analytics and machine learning.

Process Flow

1. Data Sourcing

- Reads raw Delta tables from the Bronze Lakehouse.
- Inherits schema and structure directly from ingested CSVs.

2. Cleaning & Conforming

- **Data Type Standardization:** Ensures consistent formatting for timestamps, numeric, and strings across all tables.
- **Null Value Handling:** Drops or flags of nulls in critical fields (e.g., orders and reviews) based on business rules.
- **Transformations:**
 - Adds derived columns like purchase_year and purchase_month to enhance temporal analysis.
 - Corrects column name typos (e.g., lenght → length) to ensure schema clarity and semantic alignment

Key Characteristics

1. Type Consistency

- a. Enforces uniform data types across all tables to support reliable joins, aggregations, and semantic modeling.

2. Basic Quality Checks

- a. Validates the presence of critical fields.
- b. Ensures schema integrity before passing data to the Gold layer.

3. Enrichment Logic

- Introduces business-friendly features to support reporting and ML workflows.

Role in Architecture

- Acts as the refinement checkpoint between raw ingestion and curated analytics.
- Enables modular validation and transformation logic that can be reused or extended.
- Serves as the foundation for semantic modeling, ensuring clean inputs for Power BI and ML experiments.

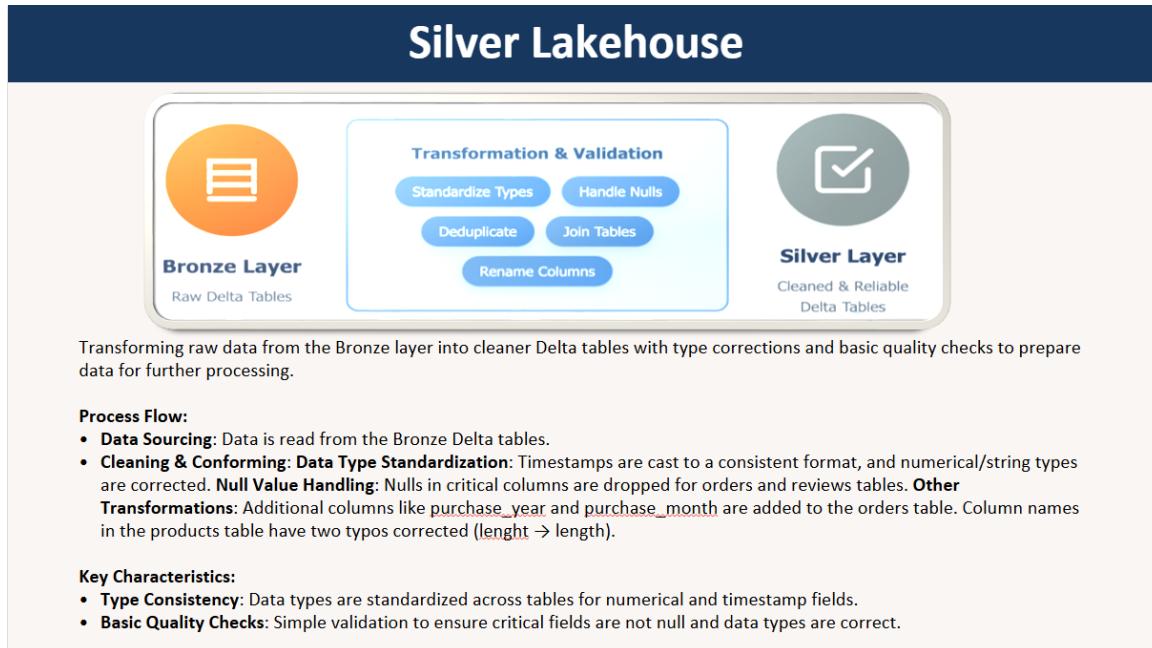
Cleaning of tables

Table Name	Column Name	Cleaning
silver_order_items	Shipping_limit_date	Change to date time
silver_order_items	Price, freight_value	Change to float, 2 decimal places
silver_order_items	Order_item_id	Change to integer
Silver_order_payments	Payment_sequential, payment_installments	Change to integer
Silver_order_payments	Payment_value	Change to float, 2 decimal places
Silver_order_reviews	Review_creation_date, review_answer_timestamp	Change to date time
Silver_order_reviews	Review_score	Change to integer, if not number remove rows
Silver_order_reviews	rows	Remove rows if order_id is null
Silver_orders	Order_approved_at, order_delivered_carrier_date, order_delivered_customer_date, order_estimated_delivery_date	Change to date time
Silver_products	Product_name_lenght, product_description_lenght	Change column name lenght to length

PROJECT ARCHITECTURE DESIGN

Silver Lakehouse Transformation

Table Name	Column(s)	Cleaning Instruction	Code Implementation
silver_order_items	Shipping_limit_date	Change to date time	.withColumn("Shipping_limit_date", to_timestamp(col("Shipping_limit_date"))))
silver_order_items	Price, freight_value	Change to float, 2 decimal places	.withColumn("Price", round(col("Price").cast("float"), 2)) .withColumn("freight_value", round(col("freight_value").cast("float"), 2)))
silver_order_items	Order_item_id	Change to integer	.withColumn("Order_item_id", col("Order_item_id").cast("integer"))
Silver_order_payments	Payment_sequential, payment_installments	Change to integer	.withColumn("Payment_sequential", col("Payment_sequential").cast("integer")) .withColumn("payment_installments", col("payment_installments").cast("integer"))
Silver_order_payments	Payment_value	Change to float, 2 decimal places	2.withColumn("Payment_value", round(col("Payment_value").cast("float"), 2))
Silver_order_reviews	Review_creation_date, review_answer_timestamp	Change to date time	.withColumn("Review_creation_date", to_timestamp(col("Review_creation_date")))) .withColumn("review_answer_timestamp", to_timestamp(col("review_answer_timestamp"))))
Silver_order_reviews	Review_score	Change to integer, if not number remove rows	.withColumn("Review_score", col("Review_score").cast("integer")) .na.drop(subset=["Review_score"]))
Silver_order_reviews	rows	Remove rows if order_id is null	.na.drop(subset=["order_id"]))
Silver_orders	Order_approved_at, order_delivered_carrier_date, order_delivered_customer_date, order_estimated_delivery_date	Change to date time	.withColumn("...", to_timestamp(col("...")))) for each specified date column.
Silver_products	Product_name_lenght, product_description_lenght	Change column name lenght to length	.withColumnRenamed("Product_name_lenght", "product_name_length") .withColumnRenamed("product_description_lenght", "product_description_length")



NOTEBOOK SOURCE CODE FOR SILVER LAKEHOUSE

Cleaning and Transformation Steps

For the orders Table

1. Corrected Data Types: The `order_purchase_timestamp` column was converted from text into a proper timestamp, making it usable for date-based calculations.
2. Standardized Text: The `order_status` column was converted to lowercase to ensure consistency (e.g., "Delivered" and "delivered" are treated as the same).
3. Enriched Data: Two new columns, `purchase_year` and `purchase_month`, were created by extracting the year and month from the timestamp. This makes filtering and grouping by date much easier.
4. Removed Bad Data: Rows with missing (null) values in critical columns like `order_id`, `customer_id`, or `order_purchase_timestamp` were deleted to improve data integrity.

For All Other Tables

Renamed for Clarity: The primary transformation for the other tables was renaming them from their raw source names (e.g., `olist_customers_dataset`) to clean, intuitive names for the Silver layer (e.g., `silver_customers`).

```

1  # Import of libraries
2  import logging
3  from pyspark.sql.functions import col, to_timestamp, lower, year, month, round
4
5  # Configure of logging for monitoring
6  logging.basicConfig(level=logging.ERROR)
7
8  # Configuration mapping Bronze names to clean Silver names
9  config = {
10     "orders": {"bronze": "olist_orders_dataset", "silver": "silver_orders"},
11     "customers": {"bronze": "olist_customers_dataset", "silver": "silver_customers"},
12     "order_items": {"bronze": "olist_order_items_dataset", "silver": "silver_order_items"},
13     "products": {"bronze": "olist_products_dataset", "silver": "silver_products"},
14     "sellers": {"bronze": "olist_sellers_dataset", "silver": "silver_sellers"},
15     "payments": {"bronze": "olist_order_payments_dataset", "silver": "silver_order_payments"},
16     "reviews": {"bronze": "olist_order_reviews_dataset", "silver": "silver_order_reviews"},
17     "geolocation": {"bronze": "olist_geolocation_dataset", "silver": "silver_geolocation"},
18     "translations": {"bronze": "product_category_name_translation", "silver": "silver_product_translations"}
19 }
20

21
22 for key, names in config.items():
23     bronze_table_name = names["bronze"]
24     silver_table_name = names["silver"]
25
26     try:
27         bronze_df = spark.read.table(f'LH_Bronze.{bronze_table_name}')
28
29         df_to_write = None
30
31         if key == "orders":
32             # This block now includes all transformations
33             cleaned_df = bronze_df \
34                 .withColumn("order_purchase_timestamp", to_timestamp(col("order_purchase_timestamp"))) \
35                 .withColumn("order_approved_at", to_timestamp(col("order_approved_at"))) \
36                 .withColumn("order_delivered_carrier_date", to_timestamp(col("order_delivered_carrier_date"))) \
37                 .withColumn("order_delivered_customer_date", to_timestamp(col("order_delivered_customer_date"))) \
38                 .withColumn("order_estimated_delivery_date", to_timestamp(col("order_estimated_delivery_date"))) \
39                 .withColumn("order_status", lower(col("order_status"))) \
40                 .withColumn("purchase_year", year(col("order_purchase_timestamp"))) \
41                 .withColumn("purchase_month", month(col("order_purchase_timestamp"))) \
42                 .na.drop(subset=["order_id", "customer_id", "order_purchase_timestamp"])
43             df_to_write = cleaned_df
44
45         elif key == "order_items":
46             cleaned_df = bronze_df \
47                 .withColumn("Shipping_limit_date", to_timestamp(col("Shipping_limit_date"))) \
48                 .withColumn("Price", round(col("Price").cast("float"), 2)) \
49                 .withColumn("freight_value", round(col("freight_value").cast("float"), 2)) \
50                 .withColumn("Order_item_id", col("Order_item_id").cast("integer"))
51             df_to_write = cleaned_df
52
53         elif key == "payments":
54             cleaned_df = bronze_df \
55                 .withColumn("Payment_sequential", col("Payment_sequential").cast("integer")) \
56                 .withColumn("payment_installments", col("payment_installments").cast("integer")) \
57                 .withColumn("Payment_value", round(col("Payment_value").cast("float"), 2))
58             df_to_write = cleaned_df
59

```

PROJECT ARCHITECTURE DESIGN

```
60     elif key == "reviews":
61         cleaned_df = bronze_df \
62             .withColumn("Review_creation_date", to_timestamp(col("Review_creation_date"))) \
63             .withColumn("review_answer_timestamp", to_timestamp(col("review_answer_timestamp"))) \
64             .withColumn("Review_score", col("Review_score").cast("integer")) \
65             .na.drop(subset=["order_id", "Review_score"])
66         df_to_write = cleaned_df
67
68     elif key == "products":
69         cleaned_df = bronze_df \
70             .withColumnRenamed("Product_name_lenght", "product_name_length") \
71             .withColumnRenamed("product_description_lenght", "product_description_length")
72         df_to_write = cleaned_df
73
74     else:
75         df_to_write = bronze_df
76
77
78     df_to_write.write \
79         .mode("overwrite") \
80         .format("delta") \
81         .option("overwriteSchema", "true") \
82         .saveAsTable(f'LH_Silver.{silver_table_name}')
83
84 except Exception as e:
85     print(f"ERROR processing {bronze_table_name}: {e}")
86     raise Exception(f"Silver Lakehouse files failed to load: {bronze_table_name}") from e
```

GOLD LAKEHOUSE

The Gold Lakehouse represents the final refinement stage in the project's pipeline. Its primary goal is to consolidate cleaned data from the Silver Lakehouse into a single, enriched dataset optimized for machine learning workflows and business reporting. This layer ensures schema consistency, data integrity, and analytical readiness.

Transformation Process

1. Data Sourcing

- Reads curated Delta tables from the Silver Lakehouse.
- Includes key entities such as products, orders, order items, and product translations.

2. Transformation Logic

- Merges multiple Silver tables into a unified schema.
- Applies enrichment logic to support downstream ML experiments and semantic modeling.

Key Characteristics

• Basic Quality Checks

- To validate that the schema of the Gold Lakehouse table mirrors its Silver source post-transformation.
- Ensures column alignment and data type consistency.

• Data Integrity Checks

- Reconciles row counts between Silver and Gold tables to confirm no data loss.
- Flags anomalies or mismatches for immediate resolution.

Role in Architecture

- Serves as the launchpad for ML experiments, including market basket analysis and recommendation engines.
- Provides a business-ready dataset for semantic modeling and Power BI dashboards.
- Act as a final checkpoint for validation before data enters the reporting and analytics layer

Gold Lakehouse



The primary goal is to prepare a unified, ML-ready dataset by transforming and consolidating multiple Silver Lakehouse tables.

- **Data Sourcing:** Data is read from the Silver Delta tables.
- **Transformation Logic:** Merge the above tables into a single, enriched table suitable for machine learning workflows.

Key Characteristics:

- **Basic Quality Checks:** Validate that the schema of the Gold Lakehouse table mirrors the Silver Lakehouse source post-transformation.
- **Data Integrity Checks:** Ensure row counts between Silver and Gold tables remain consistent to confirm no data loss during transformation.

GOLD LAKEHOUSE DIMENSIONAL SCHEMA

Foundation for ML & Analytics

The Gold Lakehouse schema is the culmination of the project pipeline, designed to organize enriched data into a dimensional model that supports advanced analytics, machine learning, and stakeholder reporting. This schema enables efficient querying, slicing, and aggregation across key business domains.

Purpose

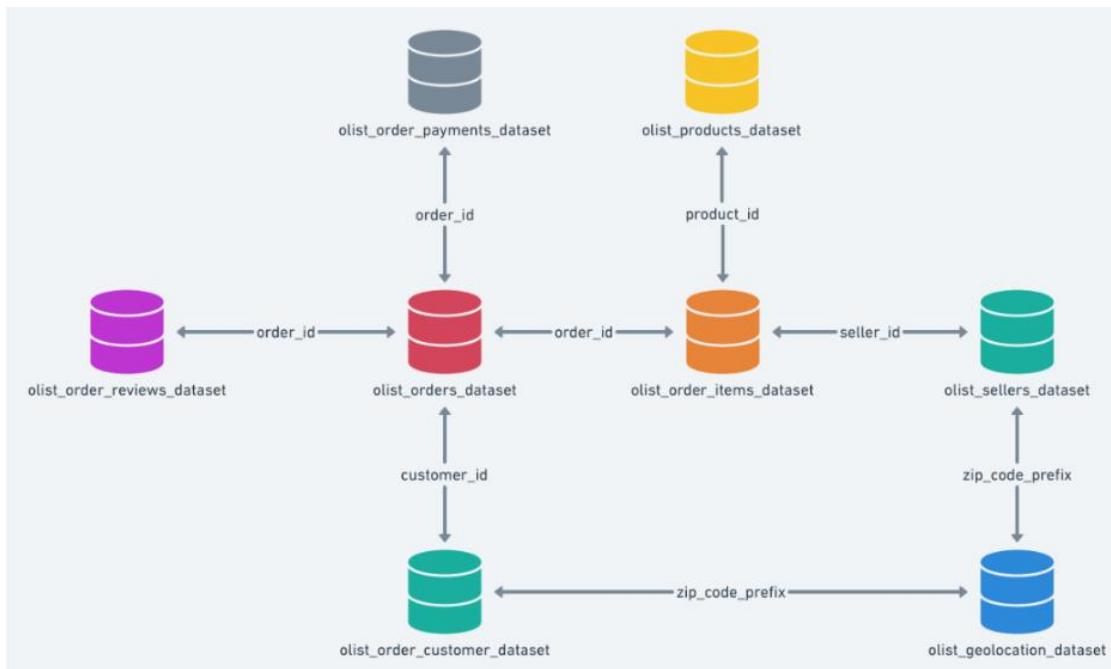
- To deliver a machine learning ready dataset by consolidating cleaned Silver Lakehouse tables into a structured, dimensional format.
- To support semantic modeling and Power BI dashboards with optimized relationships and hierarchies.

Schema Design Highlights

- **Star Schema Structure**
 - Central fact tables (e.g., orders, order_items) linked to dimension tables (e.g., customers, products).
 - Enables efficient joins and aggregations for reporting and ML workflows.
- **Semantic Alignment**
 - Cleaned column names and standardized data types ensure compatibility with Power BI and DAX modeling

Benefits

- **Query Optimization**
 - Dimensional design reduces complexity and improves performance for dashboard visuals and ML feature extraction.
- **Business Readiness**
 - Tables are curated for stakeholder consumption, enabling drill-through insights and KPI tracking.
- **ML Enablement**
 - Unified schema supports feature engineering for market basket analysis and recommendation engines.



PROJECT ARCHITECTURE DESIGN

Dataset Descriptions

Entity Name	Description	IsNull	Unique Values
Table Name: customers			
customer_id	key to the orders dataset. Each order has a unique customer_id.	No	99441
customer_unique_id	unique identifier of a customer.	No	96096
customer_zip_code_prefix	first five digits of customer zip code	No	
customer_city	customer city name	No	4119
customer_state	customer state	No	27
Table Name: geolocation			
geolocation_zip_code_prefix	first 5 digits of zip code	No	
geolocation_lat	latitude	No	96096
geolocation_lng	longitude	No	
geolocation_city	city name	No	8011
geolocation_state	state	No	27
Table Name: order_items			
order_id	order unique identifier	No	98666
order_item_id	sequential no. identifying no. of items included in the same order.	No	
product_id	product unique identifier	No	32951
seller_id	seller unique identifier	No	3095
shipping_limit_date	Shows the seller shipping limit date for handling the order over to the logistic partner.	No	
price	item price	No	
freight_value	item freight value item (if an order has more than one item the freight value is splitted between items)	No	
Table Name: order_payments			
order_id	unique identifier of an order.	No	99440
payment_sequential	customer may pay an order with more than one payment method. If he does so, a sequence will be created to accommodate all payments.	No	
payment_type	method of payment chosen by the customer.	No	5
payment_installments	No. of installments chosen by the customer.	No	
payment_value	transaction value.	No	
Table Name: order_reviews			
review_id	unique review identifier	No	98410
order_id	unique order identifier	No	98673
review_score	Note ranging from 1 to 5 given by the customer on a satisfaction survey.	No	

PROJECT ARCHITECTURE DESIGN

review_comment_title	Comment title from the review left by the customer, in Portuguese.	87.7k	4527
review_comment_message	Comment message from the review left by the customer, in Portuguese.	58.2k	36.2k
review_creation_date	Shows the date in which the satisfaction survey was sent to the customer.	No	
review_answer_timestamp	Shows satisfaction survey answer timestamp.	No	

Table Name: orders

order_id	unique identifier of the order.	No	99441
customer_id	key to the customer dataset. Each order has a unique customer_id.	No	99441
order_status	Reference to the order status (delivered, shipped, etc).	No	
order_purchase_timestamp	Shows the purchase timestamp.	No	
order_approved_at	Shows the payment approval timestamp	No	
order_delivered_carrier_date	Shows the order posting timestamp. When it was handled to the logistic partner.	1783	
order_delivered_customer_date	Shows the actual order delivery date to the customer.	2965	
order_estimated_delivery_date	Shows the estimated delivery date that was informed to customer at the purchase moment.	No	

Table Name: products

product_id	unique product identifier	No	32951
product_category_name	root category of product, in Portuguese.	610	73
product_name_length	number of characters extracted from the product name.	610	
product_description_length	number of characters extracted from the product description.	610	
product_photos_qty	number of product published photos	610	
product_weight_g	product weight measured in grams.	2	
product_length_cm	product length measured in centimeters.	2	
product_height_cm	product height measured in centimeters.	2	
product_width_cm	product width measured in centimeters.	2	

Table Name: sellers

seller_id	seller unique identifier	No	3095
seller_zip_code_prefix	first 5 digits of seller zip code	No	
seller_city	seller city name	No	611
seller_state	seller state	No	23

Table Name: product_category

product_category_name	category name in Portuguese	No	71
product_category_name_english	category name in English	No	71

NOTEBOOK SOURCE CODE FOR GOLD LAKEHOUSE

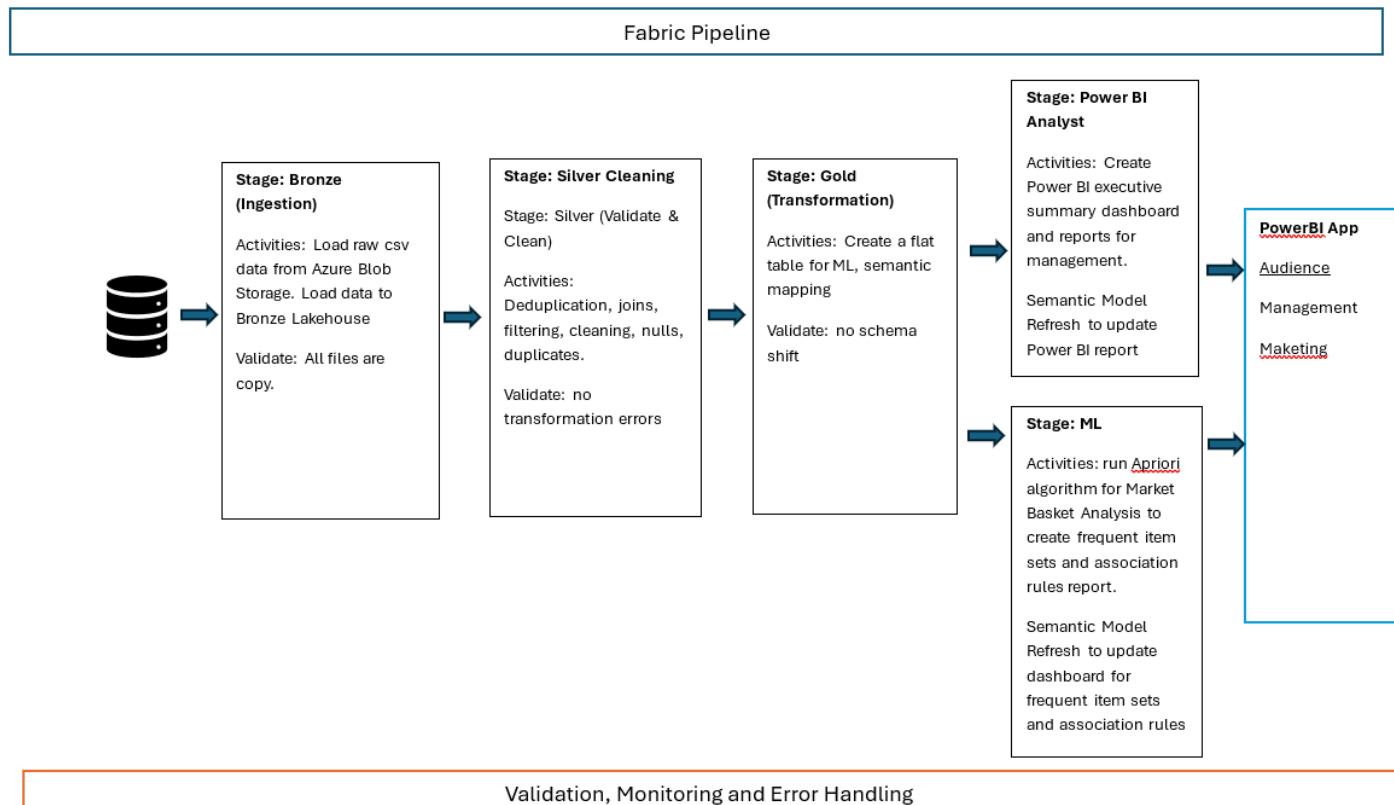
To read tables from Silver Lakehouse and save to Gold Lakehouse

```
1 # To read tables from LH_Silver lakehouse and Write to LH_Gold lakehouse
2 # Use of concurrent.futures to parallelize reads and writes across tables
3
4 from pyspark.sql import SparkSession
5 from concurrent.futures import ThreadPoolExecutor, as_completed
6 import logging
7
8 logging.basicConfig(level=logging.ERROR)
9
10 # Start Spark session
11 spark = SparkSession.builder.getOrCreate()
12
13 # Define silver_tables for the executor function
14 silver_tables = [t.name for t in spark.catalog.listTables('LH_Silver')]
15
16 # Function to iterate through the tables, read and rename, then save to the LH_Gold lakehouse
17 def migrate_table(table):
18     try:
19         logging.info(f"Processing {table}")
20         df = spark.read.table(f"LH_Silver.{table}")
21         gold_table = table.replace('silver_', 'gold_', 1) if table.startswith('silver_') else f'gold_{table}'
22         df.write.mode("overwrite").format("delta").saveAsTable(f"LH_Gold.{gold_table}")
23         logging.info(f"Finished saving {gold_table}")
24     except Exception as e:
25         # exc_info ensures the full traceback is logged, which is crucial for debugging
26         logging.error(f'Error processing table {table}:{e}', exc_info=True)
27         raise Exception(f"Gold Lakehouse files failed to load: {f.name}") from e
28
29
30 # To map through the tables in LH_Silver
31 with ThreadPoolExecutor() as executor:
32     futures = {executor.submit(migrate_table, table): table for table in silver_tables}
33     for future in as_completed(futures):
34         try:
35             # This will raise any exception from the thread
36             future.result()
37         except Exception as e:
38             logging.error(f"Notebook failed due to: {e}", exc_info=True)
39             raise Exception(f"Gold Lakehouse files failed to load: {f.name}") from e
40
```

FABRIC PIPELINE DESIGN

This structure of the end-to-end workflow for our project pipeline design is as shown below. The pipeline is segmented into Five key phases, Medallion stages, Validation, Market Basket Analysis, Error Handling, and Power BI Analysis, each representing a critical stage in transforming raw data into actionable insights.

Pipeline Architecture Design



1) Medallion Stages

- Data Input & Extraction: Gathering raw datasets from various sources.
- Validation & Cleaning: Identifying and correcting inconsistencies, missing values, and formatting issues.
- Transformation: Structuring the data into a usable format for downstream processes.

2) Validation

This initial phase validates the data integrity and schema shift of the various stages. It includes:

- **Bronze Lakehouse Validation**: Establish a clean ingestion baseline with schema-on-read flexibility and overwrite mode enforcement
- **Silver Lakehouse Validation**: Ensure clean and standardized data maintains structural integrity before enrichment and transformation
- **Gold Lakehouse Validation**: Guarantee that the final ML-ready dataset is complete, consistent, and trustworthy for semantic modeling and reporting.

3) Market Basket Analysis

This phase applies Apriori association rule mining techniques to uncover patterns in transactional data:

- Data Preparation: Merge datasets to create a table for data mining.
- Analysis Execution: Applying algorithms to identify product affinities.
- Insight Generation: Producing recommendations and strategic findings.

PROJECT ARCHITECTURE DESIGN

4) Power BI Analysis

The final stage focuses on data visualization and reporting:

- Dashboard Creation: Building interactive visualizations in Power BI.
- Insight Delivery: Presenting findings to stakeholders for informed decision-making.

5) Error Handling

Error handling is embedded throughout the pipeline to ensure robustness:

- Validation Phase: Early detection of anomalies.
- Dedicated Module: A separate error-handling loop captures and logs issues for review and resolution.

PROJECT ARCHITECTURE DESIGN

CREATION OF AZURE BLOB STORAGE

1. Create an azure blob storage to store the raw files

The screenshot shows the Azure Storage Account overview page for 'storagegegenjde07'. The left sidebar contains links for Overview, Activity log, Tags, Diagnose and solve problems, Access Control (IAM), Data migration, Events, Storage browser, Partner solutions, Resource visualizer, Data storage, Security + networking, Data management, Settings, Monitoring, Monitoring (classic), Automation, Help, and Add tags. The main content area displays the following details:

Setting	Value
Resource group (move)	RG_FinalProject
Location	southeastasia
Subscription (move)	Gen_sub
Subscription ID	25726918-4f63-40f5-803a-e6fb9b6a8c2
Disk state	Available
Tags (edit)	Add tags
Properties	
Hierarchical namespace	Enabled
Default access tier	Hot
Blob anonymous access	Disabled
Blob soft delete	Enabled (7 days)
Container soft delete	Enabled (7 days)
Versing	Disabled
Change feed	Disabled
NFS v3	Disabled
SFTP	Disabled
Storage tasks assignments	None

2. Microsoft Azure Storage Explorer

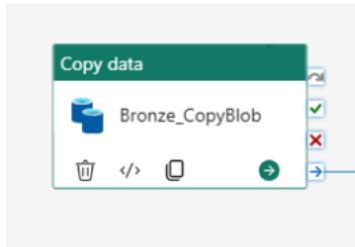
The raw files are loaded into the container

The screenshot shows the Microsoft Azure Storage Explorer interface. The left sidebar displays a tree view of storage accounts and containers. The 'storagegegenjde07' account is selected, showing its blob containers, file shares, queues, and tables. The 'containeregenjde07' container is selected, showing its blobs. The right pane lists the blobs in the container, each with a preview icon, name, and access tier.

Name	Access Tier
olist_customers_dataset.csv	Hot (inferred)
olist_geolocation_dataset.csv	Hot (inferred)
olist_order_items_dataset.csv	Hot (inferred)
olist_order_payments_dataset.csv	Hot (inferred)
olist_order_reviews_dataset.csv	Hot (inferred)
olist_orders_dataset.csv	Hot (inferred)
olist_products_dataset.csv	Hot (inferred)
olist_sellers_dataset.csv	Hot (inferred)
product_category_name_translation.csv	Hot (inferred)

PROJECT ARCHITECTURE DESIGN

ACTIVITIES IN THE PIPELINE



1) To copy data from Azure Blob Storage to Bronze Lakehouse

- Create pipeline and select copy data

Choose data source, select Azure Blob Storage and the container. Connect to the data source. Preview of the data when connection is established.

A screenshot of the 'Copy data' pipeline. On the left, a vertical progress bar shows steps: 'Choose data source' (done), 'Connect to data source' (in progress), 'Choose data destination' (not yet started), 'Connect to data destination' (not yet started), and 'Review + save' (not yet started). The main area is titled 'Connect to data source'. It shows a preview of data from 'containerengerjde07' in 'DelimitedText' format with a comma delimiter. The first few rows of data are visible:

- Choose data destination. Select the Lakehouse (LH_Bronze)

A screenshot of the 'Copy data' pipeline. The vertical progress bar shows: 'Choose data source' (done), 'Connect to data source' (done), 'Choose data destination' (in progress), 'Connect to data destination' (not yet started), and 'Review + save' (not yet started). The main area shows the 'New Fabric item' section with 'Lakehouse' selected. Below it, 'New destinations' include 'SQL Server database', 'Dataverse', 'Azure SQL database', 'Folder', 'Oracle database', and 'Odbc'. The 'OneLake catalog' section shows two entries: 'LH_Alex' and 'LH_olist'. A 'Back' button is at the bottom right.

PROJECT ARCHITECTURE DESIGN

c. Connect to data destination

Select the Root folder as Files to save the raw files into the Files folder

The screenshot shows the 'Copy data' wizard with the following steps completed:

- Choose data source
- Connect to data source
- Choose data destination
- Connect to data destination**
Select and map to folder path or table.
- Review + save

The right panel displays the 'Connect to data destination' configuration:

- Connection: LH_Bronze
- Root folder: Files Tables
- Folder path: (empty)
- File name: (empty)
- Copy behavior: Select...

The second screenshot shows the same process after selecting 'Files' as the root folder:

The right panel displays the 'Connect to data destination' configuration with the 'Root folder' setting changed to 'Files':

- File format: DelimitedText
- Column delimiter: Comma (,)
- Row delimiter: Default (\r,\n, or \n\r)
- Add header to file
- Compression type: No compression

d. Review and save

The screenshot shows the 'Copy data' wizard with the following steps completed:

- Choose data source
- Connect to data source
- Choose data destination
- Connect to data destination
- Review + save
Confirm Copy summary

The right panel displays the 'Review + save' summary:

Copy Summary

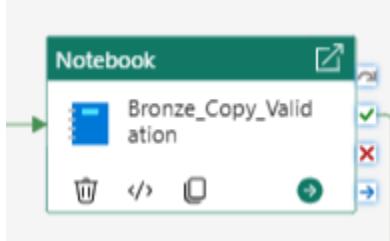
Source **Destination**

Source	Destination
Connection name: BlobStorage ChooGeokH (2)	Connection name: LH_Bronze
Container: containerengerjide07	

PROJECT ARCHITECTURE DESIGN

- 2) To validate the copy operation from Azure Blob Storage to Bronze Lakehouse

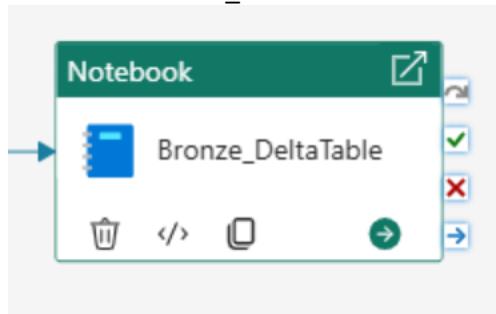
Run notebook: NB_Bronze_Copy_Validation



- 3) To run notebook to load raw files into Delta Tables

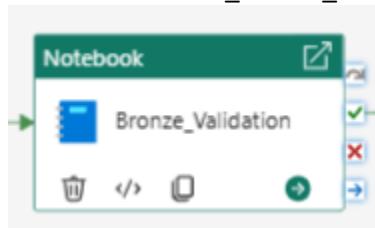
The notebook will load all csv files in the Files folder to Delta Tables.

Run notebook: NB_Bronze



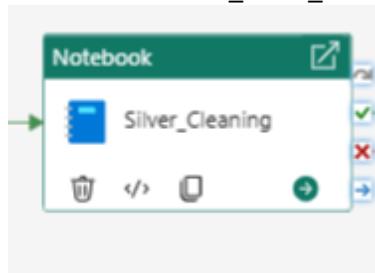
- 4) To validate the all csv files have been converted to delta tables in the Bronze Lakehouse

Run notebook: NB_Bronze_Validation



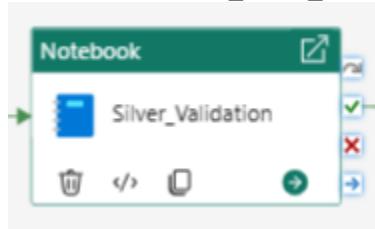
- 5) Load data from Bronze Lakehouse and save it to Silver Lakehouse. Perform Cleaning of dataset

Run Notebook: NB_Silver_Transform



- 6) To validation the tables in Silver Lakehouse

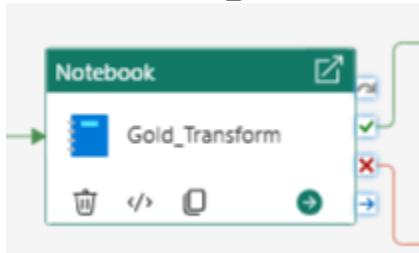
Run Notebook: NB_Silver_Validation



PROJECT ARCHITECTURE DESIGN

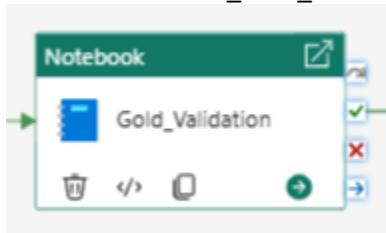
- 7) Load dataset from Silver Lakehouse and save into Gold Lakehouse. Perform transformation of dataset.

Run Notebook: NB_Gold



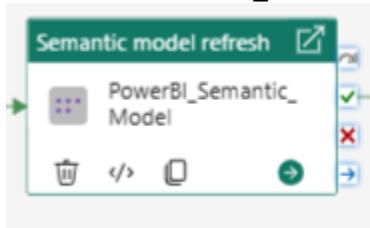
- 8) Validation for the tables in LH_Gold Lakehouse

Run Notebook: NB_Gold_Validation

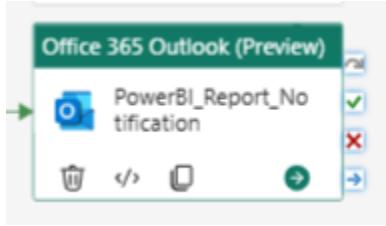


- 9) Semantic Model Refresh for Power BI report

Semantic Model: SM_PowerBI

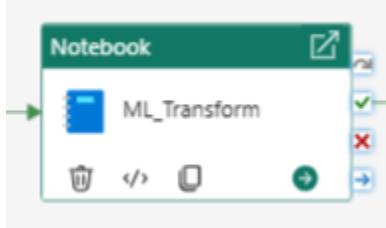


- 10) Email Notification for semantic model refresh and update of report



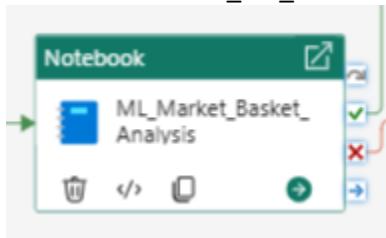
- 11) To load Delta tables from LH_Gold Lakehouse and transform into data suitable for ML model training

Run Notebook: NB_DL_Clean



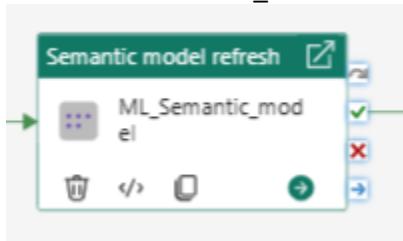
12) To use Apriori Model for Market Basket Analysis

Run Notebook: NB_ML_MarketBasket

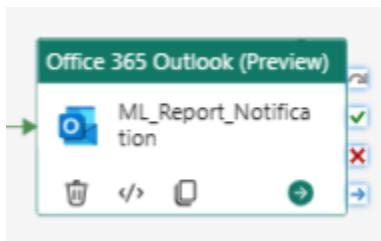


13) Semantic Model refresh for ML report creation

Semantic Model: SM_ML



14) Email Notification for semantic model refresh and update of ML report



MARKET BASKET ANALYSIS

MARKET BASKET ANALYSIS

Market Basket Analysis (MBA) is a powerful data mining technique used to uncover associations between items in transactional data. By identifying patterns in customer purchasing behavior, MBA enables businesses to optimize product placement, personalize recommendations, and drive revenue through targeted promotions.

What Is Market Basket Analysis?

- **Definition:** MBA analyzes transactional data to discover which products are frequently bought together.
- **Goal:** Reveal hidden relationships between items to inform bundling, upselling, and cross-selling strategies

Use Cases & Business Impact

- 1) **Recommendation Engine:** Suggest complementary products to increase basket size.
- 2) **Personalized Promotions:** Tailor discounts and bundles based on buying patterns.
- 3) **Cross-Selling & Upselling:** Identify high-lift product pairs to boost revenue.
- 4) **Product Strategy:** Inform inventory planning and category expansion.

Example: A 10% discount on frequently paired items can nudge customers toward higher-value purchases.

Quote: *"Igniting every cart with smarter, bigger buys."*

5) Strategic Value

- Enhance customer experience through personalization.
- Enables data-driven decisions for merchandising and marketing.
- Supports revenue growth via smarter promotions and product bundling.

Market Basket Analysis (MBA)

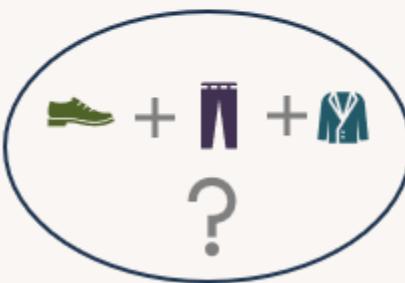


Ecommerce platform which connects sellers to multiple marketplaces and centralizes ecommerce operations

Business Goal: \$ Profit



+



MBA



"Igniting every cart with smarter, bigger buys."

MARKET BASKET ANALYSIS

Market Basket Analysis with Apriori Algorithm

Market Basket Analysis (MBA) is a core component of the project, designed to uncover product associations and drive smarter, personalized shopping experiences. Using the Apriori algorithm, the team mined transactional data to identify frequent item sets and generate actionable association rules.

Conceptual Framework

1) Frequent Item Sets

- **Definition:** Groups of products that appear together in transactions more often than a defined threshold.
- **Support:** Measure of how frequently an item set occurs in the dataset.
- Example: Support = 0.10 → Item set appears in 10% of transactions.

2) Association Rules

- **Antecedents → Consequents:** If a customer buys item A, they are likely to buy item B.
- **Confidence:** Probability that item B is purchased when item A is purchased.
 - Example: Confidence = 0.86
- **Lift:** Strength of association beyond random chance.
- Example: Lift = 3.15 → Strong positive correlation between items

Market Basket Analysis (MBA)

What is Market Basket Analysis?

How likely is a customer buying product B if the shopping cart has product A.



What are the use cases and impact to Olist ?

1. Cross-Selling and Upselling Recommendations

- Identify frequently co-purchased items.
- Recommend bundles or complementary products during checkout.



Frequent item sets: Which products are often bought together

2. Marketplace Insights and Seller Enablement

- Provide sellers with actionable insights on which is the best product combinations.
- Help optimize listings, bundle strategies, and pricing to increase visibility and conversion.



MARKET BASKET ANALYSIS

Market Basket Analysis (MBA)

How to implement Market Basket Analysis ?

Apriori Algorithm

- 1) Generate Frequent Item Sets in the database
- 2) Generate Association Rules from frequent item sets

LHS	RHS	FREQUENCY	SUPPORT	CONFIDENCE	LIFT
Shirts	Pants	4	0.4	0.5	1

What is Association Rules ?

Order Number				
001	✓			✓
002	✓	✓		
003	✓		✓	✓
004	✓			
005		✓	✓	
006	✓			
007	✓			✓
008	✓	✓		
009			✓	
010	✓			✓

1

2

3

4

Frequency: No of times shirts and pants occurred in the transaction.
= $\text{Freq}(\text{LHS}, \text{RHS}) = 4$

Support: Frequency of item sets (LHS + RHS) in total transactions.
= $\text{Freq}(\text{LHS}, \text{RHS}) / \text{Total Transactions}$
= $4/10 = 0.4$

Confidence: Likelihood of RHS being bought when item LHS is bought.
= $\text{Freq}(\text{LHS}, \text{RHS}) / \text{Freq}(\text{LHS})$
= $4/8 = 0.5$

Lift: Strength of association beyond random chance.
= $\text{Support} / (\text{Support}(\text{LHS}) \times \text{Support}(\text{RHS}))$
= $0.4 / (8/10 \times 4/10) = 1.25$

MARKET BASKET ANALYSIS

Implementation Workflow

1) Stage 1: Environment Setup

Purpose: Establish a reproducible and traceable analytics workspace.

- Created a new notebook environment in Microsoft Fabric.
- Imported external module mlxtend for Apriori and association rule mining.
- Configure a Fabric experiment to track execution and enable logging for traceability.

2) Stage 2: Load Tables

Purpose: Prepare high-quality input data for association rule mining.

- Pull relevant datasets from the Gold Lakehouse layer.
- Merge multiple tables into a unified transactional view.
- Clean the dataset by removing unnecessary columns.
- Construct a market basket matrix and filter out low-frequency items to reduce noise.

3) Stage 3: Apriori Algorithm

Purpose: Discover actionable product relationships and store them for semantic modeling.

- Apply the Apriori algorithm to extract frequent item sets and generate association rules.
- Persist results as **Delta tables** in the Gold Lakehouse for downstream use.

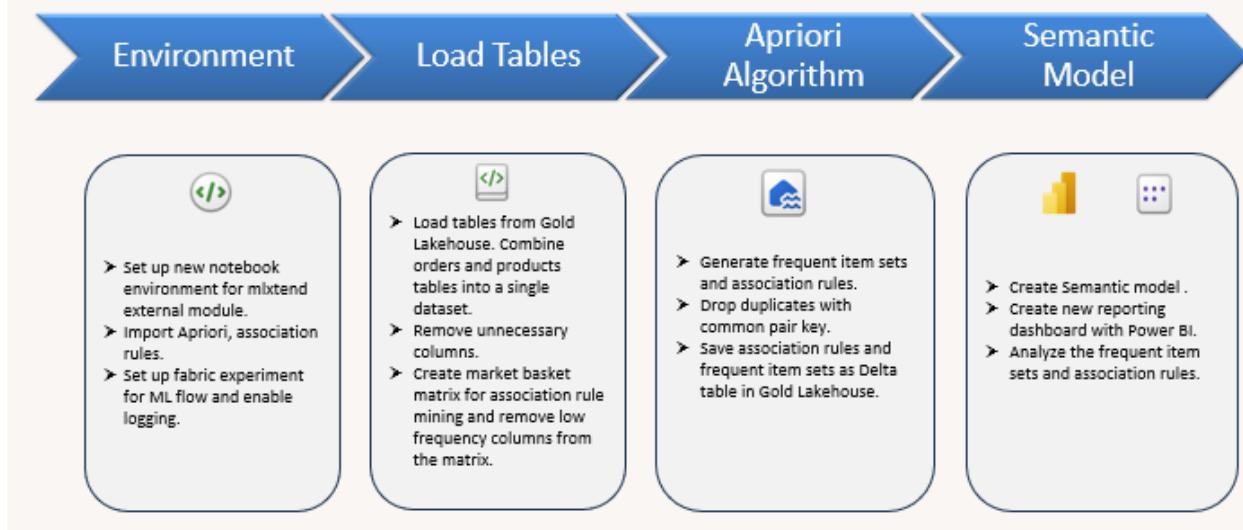
4) Stage 4: Semantic Model & Reporting

Purpose: Democratize insights and empower data-driven marketing, merchandising, and recommendation strategies.

- Build a semantic model to structure and tag the rules meaningfully.
- Create a Power BI dashboard to visualize frequent item sets and rule metrics (e.g., confidence, lift).
- Enable interactive analysis for business users to explore product affinities and bundling opportunities.

Market Basket Analysis (MBA)

What is the process flow in Microsoft Fabric ?



MARKET BASKET ANALYSIS

olist

Market Basket Analysis

Frequent Single Item

Top 10 Frequent Single Item

Product Category	Support
bed_bath_table	0.38
furniture_decor	0.28
housewares	0.20
baby	0.15
home_confort	0.15
garden_tools	0.14
cool_stuff	0.12
toys	0.10
health_beauty	0.10
sports_leisure	0.08

IMPLICATIONS

- Categories for cross-selling and bundle promotions.
- Prioritize in home page placement, recommendation engines or email campaigns.

Frequent Item Sets

item_1	item_2	support
bed_bath_table	furniture_decor	0.10
home_confort	bed_bath_table	0.06
housewares	furniture_decor	0.03
cool_stuff	baby	0.03
housewares	bed_bath_table	0.03
toys	baby	0.03
bed_bath_table	baby	0.02
garden_tools	furniture_decor	0.02

IMPLICATIONS

- Use these associations to design combo deals or discount bundles.
- Trigger real-time recommendations during check out.
- Inform home page design to group related items.

Confidence and Lift

High Confidence, High Lift represent strong, reliable rules

- Recommendation engines
- Cross-selling strategies
- Personalized promotions

Rules cluster in the confidence range of 0.2–0.4, with lift values above 1, indicating meaningful associations.

ASSESS THE QUALITY OF RULES

Confidence: Measures how often the rule has been found to be true.
Lift: Indicates how much more likely the consequents is to appear within the antecedent than by random chance.

Association Rules

Association Rules with Lift > 1			
antecedents	consequents	confidence	lift
home_confort	bed_bath_table	0.86	3.15
toys	baby	0.38	2.96
cool_stuff	baby	0.31	2.40
bed_bath_table	furniture_decor	0.35	1.26

Rules with **lift greater than 1**, identify item pairings that occur more frequently together than by random chance.

IMPLICATIONS

- Prioritize pair for bundling, homepage recommendation and inventory planning.
- Targeted promotions for baby products when toys are purchased.
- Lower confidence but still a strong lift—this might reflect seasonal or novelty-driven purchases. Consider event-driven campaigns.
- This rule has the lowest lift indicating a mild positive association. Useful for low-risk cross-selling.

MARKET BASKET ANALYSIS

NOTEBOOK SOURCE CODE FOR ML TRANSFORMATION

1) Merge gold_products and gold_product_translations table

1. Load gold_products and gold_product_translations table
2. Remove unnecessary columns from gold_products table
3. Merge both gold_products and gold_product_translations tables
4. The final table contains the product_id and product_category_name in English

```
1 # Load products and product_translations table
2 products_df = spark.read.table('gold_products')
3 translation_df = spark.read.table('gold_product_translations')
4
5 # To remove columns that are unnecessary before the merge
6 products_df = products_df.drop('product_photos_qty','product_weight_g','product_length_cm','product_height_cm','product_width_cm','product_name_length','product_description_length')
7
8 # To join both tables to create a new column for english product category name
9 joined_products_df = products_df.join(
10     translation_df,
11     products_df['product_category_name'] == translation_df['product_category_name'],
12     how='left'
13 ).drop(translation_df['product_category_name'])
14
15 # To drop the rows with na
16 joined_products_df = joined_products_df.na.drop()
```

2) Merge gold_orders and gold_order_items table

1. Load gold_orders and gold_order_items table
2. Drop unnecessary columns from both tables
3. Merge both gold_orders and gold_order_items table

```
1 # To load the orders, order_reviews and order_items table
2 orders_df = spark.read.table('gold_orders').alias('o')
3 order_items_df = spark.read.table("gold_order_items").alias('oi')
4
5 # To drop the unnecessary column from the orders table
6 # orders_df = orders_df.drop('order_estimated_delivery_date','order_delivered_customer_date','order_delivered_carrier_date','order_approved_at')
7
8 # To drop the unnecessary columns from the order_items table
9 # order_items_df = order_items_df.drop('seller_id','shipping_limit_date','price','freight_value')
10
11 # To rename the order_id column in order_items table as the column name is same as orders table
12 order_items_df_renamed = order_items_df.withColumnRenamed('order_id','order_items_order_id')
13
14 # To merge tables orders and order_items, drop the column for order_items_order_id as it is duplicate with order_id
15 joined_orders_df = order_items_df_renamed.join(
16     orders_df,
17     order_items_df_renamed['order_items_order_id'] == orders_df['order_id'],
18     how='left'
19 ).drop(order_items_df_renamed['order_items_order_id'])
```

3) Merge orders and products table

1. Rename the product_id of the products table as the column name is duplicated
2. Merge both orders and products table and remove the product_id from the products table
3. Filter rows that have the order status of canceled and unavailable
4. Save the final cleaned table

```
1 from pyspark.sql.functions import col, to_timestamp
2
3 timestamp_cols = ['order_purchase_timestamp','order_approved_at','order_delivered_carrier_date','order_delivered_customer_date','order_estimated_delivery_date','Shipping_limit_date']
4
5 # To rename the column name before merge due to duplicated column name
6 product_cleaned_df_rename = joined_products_df.withColumnRenamed('product_id', 'product_id_renamed')
7
8 # To merge both tables
9 joined_product_order_df = joined_orders_df.join(
10     product_cleaned_df_rename,
11     product_cleaned_df_rename['product_id_renamed'] == joined_orders_df['product_id'],
12     how='left'
13 ).drop(product_cleaned_df_rename['product_id_renamed'])
14
15 # To filter the rows that have the order status of canceled and unavailable
16 joined_product_order_df = joined_product_order_df.filter(~col('order_status').isin('canceled','unavailable'))
17
18 # To convert the data type for order_purchase_timestamp column
19 for col in timestamp_cols:
20     joined_product_order_df = joined_product_order_df.withColumn(col,to_timestamp(col))
21
22 # To save the final cleaned table
23 joined_product_order_df.write.mode('overwrite').saveAsTable('products_orders_cleaned')
```

MARKET BASKET ANALYSIS

NOTEBOOK SOURCE CODE FOR ML MARKET BASKET ANALYSIS

Notebook for Olist Market Basket Analysis

About Olist

Founded on July 15, 2015 in Curitiba, Brazil by Tiago Dalvi, Olist was built with a clear mission: to empower small and medium-sized businesses by simplifying their access to major e-commerce marketplaces. Through a unified platform, Olist enables sellers to seamlessly list products, manage orders, and scale their online presence across multiple channels.

Market Basket Analysis

This dataset covers the period from September 2016 to September 2018, offering a snapshot of marketplace activity during Olist's early growth phase. While the timeframe is relatively short, it provides valuable insights into product performance, customer behavior, and operational dynamics.

1) Import required libraries

```
1 # To import the library from the environment
2 import mlflow
3
4 ✓ - Command executed in 297 ms by ChooGeek H (JDE07) on 11:23:37 AM, 10/05/25
```

```
1 import mlflow
2 import numpy as np
3 import pandas as pd
4 import seaborn as sns
5 import statsmodels.api as sm
6
7 # Import ML Libraries for Market Basket Analysis
8 from mlxtend.frequent_patterns import apriori
9 from mlxtend.frequent_patterns import association_rules
10
11 import warnings
12 warnings.filterwarnings("ignore")
1
2 ✓ - Command executed in 266 ms by ChooGeek H (JDE07) on 11:23:37 AM, 10/05/25
```

Set up MLflow autolog

Microsoft Fabric extends the MLflow autologging capabilities by automatically capturing the values of input parameters and output metrics of a machine learning model as it's being trained. This information is then logged to the workspace, where it can be accessed and visualized using the MLflow APIs or the corresponding experiment in the workspace.

Experiments in Fabric let you track the results of your AutoML process, providing a comprehensive view of all the metrics and parameters from your trials.

2) Set up experiment tracking for ML flow

```
1 # Declare experiment name as a variable
2 # Enable MLflow autologging
3
4 EXPERIMENT_NAME = "ML_MarketBasket"
5 mlflow.set_tracking_uri("file:///lakehouse/default/files/ML_flow_logs")
6 mlflow.autolog()
7
8
9 ✓ - Command executed in 245 ms by ChooGeek H (JDE07) on 11:23:37 AM, 10/05/25
10
11 2025/10/05 03:23:37 INFO mlflow.tracking.fluent: Autologging successfully enabled for statsmodels.
12 PySpark (PyTf)
```

3) To load the table from LH_Gold Lakehouse and drop columns that are not necessary for Market Basket Analysis

```
1 # Convert the delta table to pandas dataframe
2
3 orders = spark.read.table('products_orders_cleaned')
4 orders_df = orders.toPandas()
5
6 drop_cols = [
7     'seller_id',
8     'shipping_limit_date',
9     'price',
10    'customer_id',
11    'order_status',
12    'freight_value',
13    'order_approved_at',
14    'order_purchase_timestamp',
15    'order_delivered_carrier_date',
16    'order_delivered_customer_date',
17    'order_estimated_delivery_date',
18    'product_category_name',
19    'purchase_year',
20    'purchase_month',
21    'product_id'
22]
23
24 orders_df = orders_df.drop(drop_cols, axis=1)
```

MARKET BASKET ANALYSIS

4) To creates a market basket matrix for association rule mining and product affinity analysis

```
1 # Enables association rule mining and helps identify co-purchase patterns
2 # For basket analysis encoding the data to 1 and 0, makes the data compatible with mlxtend model
3 # Enables frequent itemset mining and supports association rule generation
4
5 market_basket = pd.pivot_table(data=orders_df,index='order_id',columns='product_category_name_english',values='Order_item_id',fill_value=0)
6
7 market_basket_encoded = market_basket.applymap(lambda x: 1 if x >= 1 else 0)
```

[✓] - Command executed in 2 sec 297 ms by ChooGeok H (JDE07) on 11:23:43 AM, 10/05/25

5) To remove low-frequency columns from the matrix

```
1 # To drop columns that have frequency of less than or equal 5 occurrences
2 # Eliminates single-item orders, which cannot form meaningful associations
3 # Ensures that association rules (like "If A, then B") are based on co-occurrence
4 # Improves support and confidence metrics in algorithms like Apriori
5
6 low_freq_cols = [col for col in market_basket_encoded.columns if market_basket_encoded[col].sum(skipna=True) <= 5]
7 market_basket_encoded = market_basket_encoded.drop(low_freq_cols, axis=1)
8
9 market_basket_encoded = market_basket_encoded[(market_basket_encoded > 0).sum(axis=1)>=2]
```

[✓] - Command executed in 288 ms by ChooGeok H (JDE07) on 11:23:43 AM, 10/05/25

6) Generate frequent itemsets and association rules from your frequent itemsets using the confidence metric

```
1 # Filters itemsets that appear in at least 2% of all transactions
2
3 with mlflow.start_run():
4     # Log parameters
5     mlflow.log_param("min_support", 0.02)
6     mlflow.log_param("metric", "'confidence'")
7     mlflow.log_param("min_threshold", 0.1)
8
9     # Run Apriori
10    freq_itemset = apriori(market_basket_encoded,min_support=0.02,use_colnames=True)
11    # Log metric: number of itemsets
12    mlflow.log_metric("num_itemsets", len(freq_itemset))
13
14    # min_threshold of 0.1 to includes rules with confidence greater or equal 10%
15    rules_conf = association_rules(freq_itemset,metric='confidence',min_threshold=0.1)
16    mlflow.log_metric("num_rules", len(rules_conf))
```

MARKET BASKET ANALYSIS

7) To drop duplicates with common pair_key and convert python sets to string

```
1 # To concatenated both antecedents and consequent column as a new column 'pair_key', to list all the possible combinations to filter the cross pairs
2 # To drop duplicates rows with common pair_key and to drop the pair_key column
3
4 rules_conf_cleaned = rules_conf.copy()
5
6 # To convert the python sets to string before saving the rules as spark doesn't natively support python sets
7 rules_conf_cleaned['antecedents'] = rules_conf_cleaned['antecedents'].apply(lambda x: ', '.join(list(x)))
8 rules_conf_cleaned['consequents'] = rules_conf_cleaned['consequents'].apply(lambda x: ', '.join(list(x)))
9
10 # To rename the columns to remove the space between words
11 rules_conf_cleaned = rules_conf_cleaned.rename(columns={'antecedent support':'antecedent_support','consequent support':'consequent_support'})
12
13 # To save frequent itemset into separate columns and drop the itemsets column which is a python sets
14 freq_itemset["item_1"] = freq_itemset["itemsets"].apply(lambda x: list(x)[0] if len(x) > 0 else None)
15 freq_itemset["item_2"] = freq_itemset["itemsets"].apply(lambda x: list(x)[1] if len(x) > 1 else None)
16 freq_itemset = freq_itemset.drop('itemsets',axis=1)
```

✓ - Command executed in 260 ms by ChooGeok H (JDE07) on 11:23:47 AM, 10/05/25

8) To save association rules and frequent itemsets as Delta Table into the LH_Gold Lakehouse

```
1 spark.createDataFrame(rules_conf_cleaned).write.format('delta').mode('overwrite').saveAsTable('market_basket_rules')
2 spark.createDataFrame(freq_itemset).write.format('delta').mode('overwrite').saveAsTable('market_basket_freqitem')
```

MONITORING

MONITORING

The screenshot shows the Microsoft Fabric Monitor interface. On the left is a sidebar with icons for Home, Workspaces, Copilot, OneLake catalog, Monitor, Real-Time, Workloads, Olist_Works pace, My workspace, and Fabric. The main area has a title 'Monitor' and a subtitle 'View and track the status of the activities across all the workspaces for which you have permissions within Microsoft Fabric.' It includes a 'Refresh' button, a search bar, and filter options. A note says 'Clear all To apply filters, select the values from the Filter dropdown menu.' Below is a table with the following data:

Activity name	Status	Item type	Start time	Submitted by	Location
Dashboard Usage Metrics Model	Succeeded	Semantic model	11/10/2025, 6:27 pm	Fahmy Mahamud	Olist_Workspace
NB_Unit_Tests_38084b57-eb28-4ae9-bde7-884e...	Succeeded	Notebook	11/10/2025, 6:10 pm	Alex Yong	Olist_Workspace
NB_Gold_Tests_99187a03-4cb5-4372-a0ac-0bcfa...	Succeeded	Notebook	11/10/2025, 6:07 pm	Alex Yong	Olist_Workspace
NB_Silver_Tests_1ca40c41-ee83-4611-b952-d17...	Succeeded	Notebook	11/10/2025, 6:05 pm	Alex Yong	Olist_Workspace
NB_Bronze_Tests_7d58b9aa-2eef-41c0-a32f-ea0...	Succeeded	Notebook	11/10/2025, 6:03 pm	Alex Yong	Olist_Workspace
Testing_Pipeline	Succeeded	Pipeline	11/10/2025, 6:01 pm	Alex Yong	Olist_Workspace
SM_ML	Succeeded	Semantic model	11/10/2025, 5:10 pm	ChooGeok Ho	Olist_Workspace
NB_ML_MarketBasket_3a41b9f2-c5c3-4b8d-af7c...	Succeeded	Notebook	11/10/2025, 5:01 pm	ChooGeok Ho	Olist_Workspace
SM_PowerBI	... Succeeded	Semantic model	11/10/2025, 4:59 pm	ChooGeok Ho	Olist_Workspace
NB_ML_Clean_00afcbb3-da2e-45ca-bb17-47334...	Succeeded	Notebook	11/10/2025, 4:59 pm	ChooGeok Ho	Olist_Workspace
NB_Gold_Validation_d373a4d1-6111-45f1-a2d2-...	Succeeded	Notebook	11/10/2025, 4:57 pm	ChooGeok Ho	Olist_Workspace
NB_Gold_d9df4e3f-ab9f-46f6-911c-79c1bbe84f88	Succeeded	Notebook	11/10/2025, 4:55 pm	ChooGeok Ho	Olist_Workspace

Key Capabilities and Information

- The hub allows you to view and track the status of activities across the workspaces you have permission for.
- You can filter the activities by a keyword or use the **Filter** dropdown for advanced options.
- The data can be refreshed using the **Refresh** button.
- It is possible to customize the displayed columns using the **Column Options** menu.

Breakdown of Monitored Activity Fields

The activity log provides the following detailed information for each executed job:

- Activity name:** This identifies the specific job or item being executed
- Status:** This indicates the execution result, which is uniformly **Succeeded** for all items visible in the provided log.
- Item type:** This defines the type of Fabric item, such as a **Notebook**, **Pipeline**, or **Semantic model**.
- Start time:** This timestamp shows when the activity began
- Submitted by:** This is the user who started the activity
- Location:** This specifies the workspace where the item resides, which is Olist_Workspace in all shown examples

Examples of Monitored Project Items

The screenshot displays the monitoring of various data lifecycle stages:

- Notebooks** are monitored, including unit tests across the data layers: NB_Unit_Tests, NB_Silver_Tests, and NB_Bronze_Tests
- A specific **Pipeline** named Testing_Pipeline is tracked
- Semantic models** are monitored for refresh status, such as Dashboard Usage Metrics Model, SM_ML, and SM_PowerBI

TESTING

TESTING

Imports all necessary Python and PySpark modules required for the transformations and testing environment.

```
1 # Imports
2 import pytest
3 from pyspark.sql import SparkSession
4 from pyspark.sql.types import StructType, StructField, StringType, IntegerType, DoubleType, TimestampType,
5 from pyspark.sql.functions import col, to_timestamp, lower, year, month, round as spark_round
6 from datetime import datetime
7 import re
8
9
```

Loads the actual transformation functions, which were presumably extracted from the main pipeline notebooks (NB_Bronze.ipynb, NB_Silver_Transform.ipynb, and NB_Gold.ipynb).

```
2 # Cell 3: Transformation Functions from Your Notebooks
3 """
4 These are the actual transformation functions extracted from:
5 - NB_Bronze.ipynb
6 - NB_Silver_Transform.ipynb
7 - NB_Gold.ipynb
8 """
9
10 # =====
11 # BRONZE TRANSFORMATIONS
12 # =====
13
14 def sanitize_table_name(filename: str) -> str:
15     """
16     From NB_Bronze.ipynb
17     Sanitize CSV filename to valid table name
18     """
19     table_name = re.sub(r"\.csv$", "", filename)
20     table_name = re.sub(r"^[^a-zA-Z0-9_]", "_", table_name)
21     return table_name
22
```

TESTING

```
24 # SILVER TRANSFORMATIONS (extracted from NB_Silver_Transform.ipynb)
25 # =====
26
27 def transform_orders_timestamps(df):
28     """Convert order timestamp columns to timestamp type"""
29     return df \
30         .withColumn("order_purchase_timestamp", to_timestamp(col("order_purchase_timestamp"))) \
31         .withColumn("order_approved_at", to_timestamp(col("order_approved_at"))) \
32         .withColumn("order_delivered_carrier_date", to_timestamp(col("order_delivered_carrier_date"))) \
33         .withColumn("order_delivered_customer_date", to_timestamp(col("order_delivered_customer_date"))) \
34         .withColumn("order_estimated_delivery_date", to_timestamp(col("order_estimated_delivery_date")))
35
36 def transform_order_status_lowercase(df):
37     """Convert order_status to lowercase"""
38     return df.withColumn("order_status", lower(col("order_status")))
39
40 def add_purchase_year_month(df):
41     """Add purchase_year and purchase_month columns"""
42     return df \
43         .withColumn("purchase_year", year(col("order_purchase_timestamp"))) \
44         .withColumn("purchase_month", month(col("order_purchase_timestamp")))
45
46 def remove_order_nulls(df):
47     """Remove rows with null critical columns"""
48     return df.na.drop(subset=["order_id", "customer_id", "order_purchase_timestamp"])
49
50 def transform_order_items(df):
51     """Transform order items - timestamps, rounding, casting"""
52     return df \
53         .withColumn("Shipping_limit_date", to_timestamp(col("Shipping_limit_date"))) \
54         .withColumn("Price", spark_round(col("Price").cast("float"), 2)) \
55         .withColumn("freight_value", spark_round(col("freight_value").cast("float"), 2)) \
56         .withColumn("Order_item_id", col("Order_item_id").cast("integer"))
57
58 def transform_payments(df):
59     """Transform payments - casting and rounding"""
60     return df \
61         .withColumn("Payment_sequential", col("Payment_sequential").cast("integer")) \
62         .withColumn("payment_installments", col("payment_installments").cast("integer")) \
63         .withColumn("Payment_value", spark_round(col("Payment_value").cast("float"), 2))
64
65 def transform_reviews(df):
66     """Transform reviews - timestamps and casting"""
67     return df \
68         .withColumn("Review_creation_date", to_timestamp(col("Review_creation_date"))) \
69         .withColumn("review_answer_timestamp", to_timestamp(col("review_answer_timestamp"))) \
70         .withColumn("Review_score", col("Review_score").cast("integer")) \
71         .na.drop(subset=["order_id", "Review_score"])
72
73 def transform_products(df):
74     """Transform products - column renaming"""
75     return df \
76         .withColumnRenamed("Product_name_lenght", "product_name_length") \
77         .withColumnRenamed("product_description_lenght", "product_description_length")
```

TESTING

```
79 # =====
80 # GOLD TRANSFORMATIONS
81 # =====
82
83 def gold_table_name(silver_table_name: str) -> str:
84     """
85     From NB_Gold.ipynb
86     Convert silver table name to gold table name
87     """
88     if silver_table_name.startswith('silver_'):
89         return silver_table_name.replace('silver_', 'gold_', 1)
90     else:
91         return f'gold_{silver_table_name}'
92
93 print("/ Transformation Functions Loaded")
94 print("=" * 80)
95
96 # Cell 4: Unit Tests - Bronze Layer Functions
97 print("=" * 80)
98 print("UNIT TESTS - BRONZE LAYER TRANSFORMATIONS")
99 print("=" * 80)
100
101 failed_tests = []
102 passed_tests = []
103
104 # Test 1: sanitize_table_name - basic CSV removal
105 try:
106     result = sanitize_table_name("olist_customers_dataset.csv")
107     assert result == "olist_customers_dataset", f"Expected 'olist_customers_dataset', got '{result}'"
108     print("/ PASSED: sanitize_table_name - basic CSV removal")
109     passed_tests.append("sanitize_table_name - basic CSV")
110 except AssertionError as e:
111     print(f"\x1B[31m FAILED: {e}\x1B[0m")
112     failed_tests.append(str(e))
113
114 # Test 2: sanitize_table_name - hyphen replacement
115 try:
116     result = sanitize_table_name("product-category.csv")
117     assert result == "product_category", f"Expected 'product_category', got '{result}'"
118     print("/ PASSED: sanitize_table_name - hyphen replacement")
119     passed_tests.append("sanitize_table_name - hyphen")
120 except AssertionError as e:
121     print(f"\x1B[31m FAILED: {e}\x1B[0m")
122     failed_tests.append(str(e))
123
124 # Test 3: sanitize_table_name - space replacement
125 try:
126     result = sanitize_table_name("file with spaces.csv")
127     assert result == "file_with_spaces", f"Expected 'file_with_spaces', got '{result}'"
128     print("/ PASSED: sanitize_table_name - space replacement")
129     passed_tests.append("sanitize_table_name - spaces")
130 except AssertionError as e:
131     print(f"\x1B[31m FAILED: {e}\x1B[0m")
132     failed_tests.append(str(e))
133
134 # Test 4: sanitize_table_name - special characters
135 try:
136     result = sanitize_table_name("special@chars#here.csv")
137     assert result == "special_chars_here", f"Expected 'special_chars_here', got '{result}'"
138     print("/ PASSED: sanitize_table_name - special characters")
139     passed_tests.append("sanitize_table_name - special chars")
140 except AssertionError as e:
141     print(f"\x1B[31m FAILED: {e}\x1B[0m")
142     failed_tests.append(str(e))
143
```

TESTING

```
144 # Test 5: sanitize_table_name - no extension
145 try:
146     result = sanitize_table_name("table_name")
147     assert result == "table_name", f"Expected 'table_name', got '{result}'"
148     print("\u2713 PASSED: sanitize_table_name - no extension")
149     passed_tests.append("sanitize_table_name - no extension")
150 except AssertionError as e:
151     print(f"\u2718 FAILED: {e}")
152     failed_tests.append(str(e))
153
154 print("\n" + "=" * 80)
155
156 # Cell 5: Unit Tests - Silver Order Transformations
157 print("=" * 80)
158 print("UNIT TESTS - SILVER ORDERS TRANSFORMATIONS")
159 print("=" * 80)
160
161 spark = SparkSession.builder.getOrCreate()
162
163 # Test 6: Order status lowercase transformation
164 try:
165     test_data = [("ORDER123", "DELIVERED"), ("ORDER124", "PENDING"), ("ORDER125", "CANCELLED")]
166     schema = StructType([
167         StructField("order_id", StringType(), True),
168         StructField("order_status", StringType(), True)
169     ])
170     test_df = spark.createDataFrame(test_data, schema)
171
172     result_df = transform_order_status_lowercase(test_df)
173     results = [row.order_status for row in result_df.collect()]
174
175     assert results == ["delivered", "pending", "cancelled"], f"Expected lowercase, got {results}"
176     print("\u2713 PASSED: transform_order_status_lowercase")
177     passed_tests.append("transform_order_status_lowercase")
178 except AssertionError as e:
179     print(f"\u2718 FAILED: {e}")
180     failed_tests.append(str(e))
181
182 # Test 7: Timestamp conversion for orders
183 try:
184     test_data = [("2018-09-04 21:15:22",)]
185     schema = StructType([StructField("order_purchase_timestamp", StringType(), True)])
186     test_df = spark.createDataFrame(test_data, schema)
187
188     result_df = test_df.withColumn("order_purchase_timestamp", to_timestamp(col("order_purchase_timestamp")))
189
190     result_type = str(result_df.schema["order_purchase_timestamp"].dataType)
191     assert "timestamp" in result_type.lower(), f"Expected timestamp type, got {result_type}"
192     print("\u2713 PASSED: Timestamp conversion for orders")
193     passed_tests.append("Order timestamp conversion")
194 except AssertionError as e:
195     print(f"\u2718 FAILED: {e}")
196     failed_tests.append(str(e))
197
198 # Test 8: Add purchase year and month
199 try:
200     test_data = [(datetime(2018, 5, 15, 10, 30, 0),)]
201     schema = StructType([StructField("order_purchase_timestamp", TimestampType(), True)])
202     test_df = spark.createDataFrame(test_data, schema)
203
204     result_df = add_purchase_year_month(test_df)
205
206     assert "purchase_year" in result_df.columns, "Missing purchase_year column"
207     assert "purchase_month" in result_df.columns, "Missing purchase_month column"
208
209     row = result_df.collect()[0]
210     assert row.purchase_year == 2018, f"Expected year 2018, got {row.purchase_year}"
211     assert row.purchase_month == 5, f"Expected month 5, got {row.purchase_month}"
212
213     print("\u2713 PASSED: add_purchase_year_month")
214     passed_tests.append("add_purchase_year_month")
215 except AssertionError as e:
216     print(f"\u2718 FAILED: {e}")
217     failed_tests.append(str(e))
```

TESTING

```
219 # Test 9: Remove order nulls
220 try:
221     test_data = [
222         ("ORDER1", "CUST1", datetime(2018, 5, 15)),
223         ("ORDER2", None, datetime(2018, 5, 16)),
224         (None, "CUST3", datetime(2018, 5, 17)),
225         ("ORDER4", "CUST4", datetime(2018, 5, 18))
226     ]
227     schema = StructType([
228         StructField("order_id", StringType(), True),
229         StructField("customer_id", StringType(), True),
230         StructField("order_purchase_timestamp", TimestampType(), True)
231     ])
232     test_df = spark.createDataFrame(test_data, schema)
233
234     result_df = remove_order_nulls(test_df)
235     count = result_df.count()
236
237     assert count == 2, f"Expected 2 rows after null removal, got {count}"
238     print("\u2713 PASSED: remove_order_nulls")
239     passed_tests.append("remove_order_nulls")
240 except AssertionError as e:
241     print(f"\u2717 FAILED: {e}")
242     failed_tests.append(str(e))
243
244 print("\n" + "=" * 80)
```

```
246 # Cell 6: Unit Tests - Silver Order Items Transformations
247 print("=" * 80)
248 print("UNIT TESTS - SILVER ORDER ITEMS TRANSFORMATIONS")
249 print("=" * 80)
250
251 # Test 10: Order items price rounding
252 try:
253     test_data = [("ITEM1", "12.3456", "5.6789", "1")]
254     schema = StructType([
255         StructField("order_id", StringType(), True),
256         StructField("Price", StringType(), True),
257         StructField("freight_value", StringType(), True),
258         StructField("Order_item_id", StringType(), True)
259     ])
260     test_df = spark.createDataFrame(test_data, schema)
261
262     result_df = test_df \
263         .withColumn("Price", spark_round(col("Price").cast("float"), 2)) \
264         .withColumn("freight_value", spark_round(col("freight_value").cast("float"), 2))
265
266     row = result_df.collect()[0]
267     assert abs(row.Price - 12.35) < 0.01, f"Expected 12.35, got {row.Price}"
268     assert abs(row.freight_value - 5.68) < 0.01, f"Expected 5.68, got {row.freight_value}"
269
270     print("\u2713 PASSED: Order items price rounding")
271     passed_tests.append("Order items price rounding")
272 except AssertionError as e:
273     print(f"\u2717 FAILED: {e}")
274     failed_tests.append(str(e))
```

TESTING

```
276 # Test 11: Order item ID casting to integer
277 try:
278     test_data = [("1",), ("2",), ("3",)]
279     schema = StructType([StructField("Order_item_id", StringType(), True)])
280     test_df = spark.createDataFrame(test_data, schema)
281
282     result_df = test_df.withColumn("Order_item_id", col("Order_item_id").cast("integer"))
283
284     result_type = str(result_df.schema["Order_item_id"].dataType)
285     assert "int" in result_type.lower(), f"Expected integer type, got {result_type}"
286
287     print("\u2713 PASSED: Order item ID casting to integer")
288     passed_tests.append("Order item ID integer cast")
289 except AssertionError as e:
290     print(f"\u2717 FAILED: {e}")
291     failed_tests.append(str(e))
292
293 # Test 12: Shipping limit date timestamp conversion
294 try:
295     test_data = [("2018-09-04 21:15:22")]
296     schema = StructType([StructField("Shipping_limit_date", StringType(), True)])
297     test_df = spark.createDataFrame(test_data, schema)
298
299     result_df = test_df.withColumn("Shipping_limit_date", to_timestamp(col("Shipping_limit_date")))
300
301     result_type = str(result_df.schema["Shipping_limit_date"].dataType)
302     assert "timestamp" in result_type.lower(), f"Expected timestamp, got {result_type}"
303
304     print("\u2713 PASSED: Shipping limit date timestamp conversion")
305     passed_tests.append("Shipping limit timestamp")
306 except AssertionError as e:
307     print(f"\u2717 FAILED: {e}")
308     failed_tests.append(str(e))
309
310 print("\n" + "=" * 80)
--=
```

```
312 # Cell 7: Unit Tests - Silver Payments Transformations
313 print("=" * 80)
314 print("UNIT TESTS - SILVER PAYMENTS TRANSFORMATIONS")
315 print("=" * 80)
316
317 # Test 13: Payment sequential integer casting
318 try:
319     test_data = [("1", "2", "150.50")]
320     schema = StructType([
321         StructField("Payment_sequential", StringType(), True),
322         StructField("payment_installments", StringType(), True),
323         StructField("Payment_value", StringType(), True)
324     ])
325     test_df = spark.createDataFrame(test_data, schema)
326
327     result_df = test_df \
328         .withColumn("Payment_sequential", col("Payment_sequential").cast("integer")) \
329         .withColumn("payment_installments", col("payment_installments").cast("integer"))
330
331     seq_type = str(result_df.schema["Payment_sequential"].dataType)
332     inst_type = str(result_df.schema["payment_installments"].dataType)
333
334     assert "int" in seq_type.lower(), f"Expected integer for Payment_sequential, got {seq_type}"
335     assert "int" in inst_type.lower(), f"Expected integer for payment_installments, got {inst_type}"
336
337     print("\u2713 PASSED: Payment integer casting")
338     passed_tests.append("Payment integer casting")
339 except AssertionError as e:
340     print(f"\u2717 FAILED: {e}")
341     failed_tests.append(str(e))
342
```

TESTING

```
343 # Test 14: Payment value rounding
344 try:
345     test_data = [("150.567",)]
346     schema = StructType([StructField("Payment_value", StringType(), True)])
347     test_df = spark.createDataFrame(test_data, schema)
348
349     result_df = test_df.withColumn("Payment_value", spark_round(col("Payment_value").cast("float"), 2))
350
351     row = result_df.collect()[0]
352     assert abs(row.Payment_value - 150.57) < 0.01, f"Expected 150.57, got {row.Payment_value}"
353
354     print("\u2713 PASSED: Payment value rounding")
355     passed_tests.append("Payment value rounding")
356 except AssertionError as e:
357     print(f"\u2718 FAILED: {e}")
358     failed_tests.append(str(e))
359
360 print("\n" + "=" * 80)
361
362 # Cell 8: Unit Tests – Silver Reviews Transformations
363 print("=" * 80)
364 print("UNIT TESTS – SILVER REVIEWS TRANSFORMATIONS")
365 print("=" * 80)

367 # Test 15: Review timestamp conversion
368 try:
369     test_data = [("2018-01-18 00:00:00", "2018-01-18 21:46:59")]
370     schema = StructType([
371         StructField("Review_creation_date", StringType(), True),
372         StructField("review_answer_timestamp", StringType(), True)
373     ])
374     test_df = spark.createDataFrame(test_data, schema)
375
376     result_df = test_df \
377         .withColumn("Review_creation_date", to_timestamp(col("Review_creation_date")))\ \
378         .withColumn("review_answer_timestamp", to_timestamp(col("review_answer_timestamp")))
379
380     create_type = str(result_df.schema["Review_creation_date"].dataType)
381     answer_type = str(result_df.schema["review_answer_timestamp"].dataType)
382
383     assert "timestamp" in create_type.lower(), f"Expected timestamp, got {create_type}"
384     assert "timestamp" in answer_type.lower(), f"Expected timestamp, got {answer_type}"
385
386     print("\u2713 PASSED: Review timestamp conversion")
387     passed_tests.append("Review timestamp conversion")
388 except AssertionError as e:
389     print(f"\u2718 FAILED: {e}")
390     failed_tests.append(str(e))

392 # Test 16: Review score integer casting
393 try:
394     test_data = [(4,), (5,), (3,)]
395     schema = StructType([StructField("Review_score", StringType(), True)])
396     test_df = spark.createDataFrame(test_data, schema)
397
398     result_df = test_df.withColumn("Review_score", col("Review_score").cast("integer"))
399
400     result_type = str(result_df.schema["Review_score"].dataType)
401     assert "int" in result_type.lower(), f"Expected integer, got {result_type}"
402
403     print("\u2713 PASSED: Review score integer casting")
404     passed_tests.append("Review score integer cast")
405 except AssertionError as e:
406     print(f"\u2718 FAILED: {e}")
407     failed_tests.append(str(e))
408
```

TESTING

```
409 # Test 17: Review nulls removal
410 try:
411     test_data = [
412         ("ORDER1", 4),
413         ("ORDER2", None),
414         (None, 5),
415         ("ORDER4", 3)
416     ]
417     schema = StructType([
418         StructField("order_id", StringType(), True),
419         StructField("Review_score", IntegerType(), True)
420     ])
421     test_df = spark.createDataFrame(test_data, schema)
422
423     result_df = test_df.na.drop(subset=["order_id", "Review_score"])
424     count = result_df.count()
425
426     assert count == 2, f"Expected 2 rows after null removal, got {count}"
427     print("✓ PASSED: Review nulls removal")
428     passed_tests.append("Review nulls removal")
429 except AssertionError as e:
430     print(f"\x1b[31m FAILED: {e}\x1b[0m")
431     failed_tests.append(str(e))
432
433 print("\n" + "=" * 80)
434
435 # Cell 9: Unit Tests - Silver Products Transformations
436 print("=" * 80)
437 print("UNIT TESTS - SILVER PRODUCTS TRANSFORMATIONS")
438 print("=" * 80)
```

```
440 # Test 18: Product column renaming
441 try:
442     test_data = [("PROD1", 50, 200)]
443     schema = StructType([
444         StructField("product_id", StringType(), True),
445         StructField("Product_name_lenght", IntegerType(), True),
446         StructField("product_description_lenght", IntegerType(), True)
447     ])
448     test_df = spark.createDataFrame(test_data, schema)
449
450     result_df = transform_products(test_df)
451
452     assert "product_name_length" in result_df.columns, "Missing product_name_length column"
453     assert "product_description_length" in result_df.columns, "Missing product_description_length column"
454     assert "Product_name_lenght" not in result_df.columns, "Old column name still exists"
455
456     print("✓ PASSED: Product column renaming")
457     passed_tests.append("Product column renaming")
458 except AssertionError as e:
459     print(f"\x1b[31m FAILED: {e}\x1b[0m")
460     failed_tests.append(str(e))
461
462 print("\n" + "=" * 80)
463
464 # Cell 10: Unit Tests - Gold Layer Transformations
465 print("=" * 80)
466 print("UNIT TESTS - GOLD LAYER TRANSFORMATIONS")
467 print("=" * 80)
```

TESTING

```
469 # Test 19: Gold table name - with silver_ prefix
470 try:
471     result = gold_table_name("silver_orders")
472     assert result == "gold_orders", f"Expected 'gold_orders', got '{result}'"
473     print("\u2713 PASSED: gold_table_name - with silver_ prefix")
474     passed_tests.append("gold_table_name - with prefix")
475 except AssertionError as e:
476     print(f"\u2717 FAILED: {e}")
477     failed_tests.append(str(e))
478
479 # Test 20: Gold table name - without silver_ prefix
480 try:
481     result = gold_table_name("custom_table")
482     assert result == "gold_custom_table", f"Expected 'gold_custom_table', got '{result}'"
483     print("\u2713 PASSED: gold_table_name - without silver_ prefix")
484     passed_tests.append("gold_table_name - without prefix")
485 except AssertionError as e:
486     print(f"\u2717 FAILED: {e}")
487     failed_tests.append(str(e))
488
489 # Test 21: Gold table name - already has gold prefix
490 try:
491     result = gold_table_name("gold_orders")
492     assert result == "gold_gold_orders", f"Expected 'gold_gold_orders', got '{result}'"
493     print("\u2713 PASSED: gold_table_name - already has gold prefix")
494     passed_tests.append("gold_table_name - already gold")
495 except AssertionError as e:
496     print(f"\u2717 FAILED: {e}")
497     failed_tests.append(str(e))
498
499 print("\n" + "=" * 80)
```



```
501 # Cell 11: Summary Report
502 print("=" * 80)
503 print("UNIT TEST SUMMARY - ACTUAL PIPELINE TRANSFORMATIONS")
504 print("=" * 80)
505 print(f"Total Tests Run: {len(passed_tests) + len(failed_tests)}")
506 print(f"Tests Passed: {len(passed_tests)}")
507 print(f"Tests Failed: {len(failed_tests)}")
508 print("\nCoverage:")
509 print(f"  Bronze Layer: 5 tests")
510 print(f"  Silver Orders: 4 tests")
511 print(f"  Silver Order Items: 3 tests")
512 print(f"  Silver Payments: 2 tests")
513 print(f"  Silver Reviews: 3 tests")
514 print(f"  Silver Products: 1 test")
515 print(f"  Gold Layer: 3 tests")
516 print("=" * 80)
517
518 if len(failed_tests) > 0:
519     print("\nFailed Tests:")
520     for test in failed_tests:
521         print(f"  - {test}")
522
523 # Cell 12: Handle Results
524 if len(failed_tests) > 0:
525     print("\n" + "!" * 80)
526     print("UNIT TESTS FAILED")
527     print("!" * 80)
528     raise Exception(f"Unit tests failed: {len(failed_tests)} test(s) failed")
529 else:
530     print("\n" + "=" * 80)
531     print("\u2713 ALL UNIT TESTS PASSED")
532     print(f"All {len(passed_tests)} transformation functions validated")
533     print(f"Completed: {datetime.now().strftime('%Y-%m-%d %H:%M:%S')}")
534     print("=" * 80)
```

TESTING

```
=====  
✓ PASSED: Payment integer casting  
✓ PASSED: Payment value rounding
```

```
=====  
=====  
UNIT TESTS - SILVER REVIEWS TRANSFORMATIONS  
=====
```

```
✓ PASSED: Review timestamp conversion  
✓ PASSED: Review score integer casting  
✓ PASSED: Review nulls removal
```

```
=====  
=====  
UNIT TESTS - SILVER PRODUCTS TRANSFORMATIONS  
=====
```

```
✓ PASSED: Product column renaming
```

```
=====  
=====  
UNIT TESTS - GOLD LAYER TRANSFORMATIONS  
=====
```

```
✓ PASSED: gold_table_name - with silver_ prefix  
✓ PASSED: gold_table_name - without silver_ prefix  
✓ PASSED: gold_table_name - already has gold prefix
```

```
=====  
=====  
UNIT TEST SUMMARY - ACTUAL PIPELINE TRANSFORMATIONS  
=====
```

```
Total Tests Run: 21  
Tests Passed: 21  
Tests Failed: 0
```

```
Coverage:
```

```
Bronze Layer: 5 tests  
Silver Orders: 4 tests  
Silver Order Items: 3 tests  
Silver Payments: 2 tests  
Silver Reviews: 3 tests  
Silver Products: 1 test  
Gold Layer: 3 tests
```

```
=====  
=====  
✓ ALL UNIT TESTS PASSED  
All 21 transformation functions validated  
Completed: 2025-10-07 05:23:27  
=====
```

VALIDATION

VALIDATION

VALIDATION OF PIPELINE ACTIVITIES

To ensure data quality, schema integrity, and reliable downstream analytics, the project pipeline incorporates structured validation checkpoints across all three Lakehouse layers—Bronze, Silver, and Gold. Each stage includes notebook-driven validation logic that confirms successful data transfer, schema alignment, and row count consistency.

1) Bronze Lakehouse Validation

- Notebook: NB_Bronze_Validation
- Checks Performed:
 - Confirm raw CSV files are successfully converted to Delta format.
 - Validate that no schema shift occurred during ingestion.
- Purpose: Establish a clean ingestion baseline with schema-on-read flexibility and overwrite mode enforcement.

2) Silver Lakehouse Validation

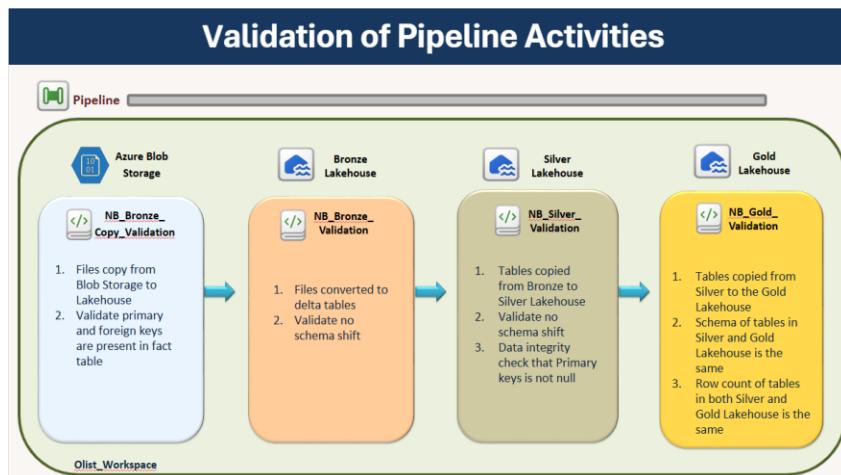
- Notebook: NB_Silver_Validation
- Checks Performed:
 - Verify that tables are correctly copied from Bronze to Silver.
 - Validate schema consistency between Bronze and Silver layers.
- Purpose: Ensure clean and standardized data maintains structural integrity before enrichment and transformation.

3) Gold Lakehouse Validation

- Notebook: NB_Gold_Validation
- Checks Performed:
 - Confirm successful transfer of tables from Silver to Gold.
 - Validate that the schema in Gold matches the Silver source post-transformation.
 - Reconcile row counts to ensure no data loss during consolidation.
- Purpose: Guarantee that the final ML-ready dataset is complete, consistent, and trustworthy for semantic modeling and reporting.

End-to-End Assurance

- Each notebook is integrated into the pipeline flow.
- Failures trigger automated alerts via Outlook, enabling rapid response and root cause analysis.



NOTEBOOK SOURCE CODE FOR COPY DATA FROM AZURE BLOB STORAGE TO BRONZE LAKEHOUSE

Validation for copy of data from Azure Blob Storage to LH_Bronze Lakehouse

```

1 # Import libraries
2 from pyspark.sql.functions import input_file_name
3 import os
4
5 - Session ready in 13 sec 660 ms. Command executed in 391 ms by ChooGeok H (JDE07) on 9:00:19 AM, 10/07/25

```

1) Validate that the expected files have been successfully copied to the LH_Bronze Lakehouse

```

1 # Define the expected filenames to be read from the blob storage
2 file_list = [
3     'olist_customers_dataset.csv',
4     'olist_geolocation_dataset.csv',
5     'olist_order_items_dataset.csv',
6     'olist_order_payments_dataset.csv',
7     'olist_order_reviews_dataset.csv',
8     'olist_orders_dataset.csv',
9     'olist_products_dataset.csv',
10    'olist_sellers_dataset.csv',
11    'product_category_name_translation.csv'
12 ]
13
14
15 # Declare the variable for Lakehouse Files folder path
16 lakehouse_path = "abfss://Olist_Worksapce@onelake.dfs.fabric.microsoft.com/LH_Bronze.Lakehouse/Files/"
17
18 # Read all CSV files from the Bronze Lakehouse
19 df = spark.read.option("header", "true").csv(lakehouse_path + "*.csv")
20
21 # Extract distinct file names from the Lakehouse Files folder, remove the version with the delimiter of ?
22 actual_files = [
23     os.path.basename(path).split("?")[0]
24     for path in df.select(input_file_name()).distinct().rdd.flatMap(lambda x: x).collect()
25 ]
26
27 # To compare the distinct set of file names from the actual and declared list
28 missing_files = list(set(file_list) - set(actual_files))
29 extra_files = list(set(actual_files) - set(file_list))
30
31 if missing_files:
32     print(f"Missing files: {missing_files}")
33 else:
34     print("All expected files are present.")
35

```

VALIDATION

```
35 assert not missing_files, f"Missing files in LH_Bronze during the copy operation: {missing_files}"
36 assert not extra_files, f"Extra files in LH_Bronze during the copy operation: {extra_files}"
37
✓ - Command executed in 6 sec 55 ms by ChooGeok H (JDE07) on 9:00:26 AM, 10/07/25
```

All expected files are present.

2) Validate that the primary and foreign keys are present in the fact table (orders dataset)

```
1 # Define the primary key and foreign keys in the fact table
2 # Validate the both keys are present in the dataset
3
4 keys = ['order_id','customer_id']
5 df_orders = spark.read.option("header", "true").csv(lakehouse_path + 'olist_orders_dataset.csv')
6 missing_keys = [key for key in keys if key not in df_orders.columns]
7 assert not missing_keys, f"Missing primary and foreign keys in orders table: {missing_keys}"
```

NOTEBOOK SOURCE CODE FOR BRONZE LAKEHOUSE VALIDATION

Validation for Bronze Lakehouse

```
1 # Import of libraries
2 import logging
3 import os
```

1) Validate that the files loaded from the Azure Blob Storage is converted to delta tables and saved in the Bronze Lakehouse

1. Ensures that raw csv files were actually ingested and materialized as Delta tables
2. Prevents silent failures where files are read but not persisted
3. Avoids broken lineage or missing inputs in transformation steps

```
1 # To compare the filenames of the csv files with the filenames of the delta tables
2 # To validate that all csv files have been saved as delta tables
3
4 # Logging of error messages for monitoring
5 logging.basicConfig(level=logging.ERROR)
6
7 # To list the file path
8 files = notebookutils.fs.ls("Files/")
9 expected_tables = []
10
11 # To validate and list the csv files are present in the folder
12 expected_files = [f.name for f in files if f.name.endswith(".csv")]
13 print("Expected CSV files:", expected_files)
14
15 # List actual tables created
16 actual_tables = [t.name for t in spark.catalog.listTables()]
17
18 # To extract the filename with extension
19 for file in expected_files:
20     new_name = os.path.splitext(file)[0]
21     expected_tables.append(new_name)
22
23 # Compare the list of table name with the list of csv files
24 missing_tables = [tbl for tbl in expected_tables if tbl not in actual_tables]
25 if missing_tables:
26     print("Missing tables:", missing_tables)
27 else:
28     print("All expected tables loaded successfully.")
29
```

2) Validate no schema shift and that the schema of the delta tables are the same as the schema declared

1. Guarantees that the structure of ingested data aligns with expectations
2. Prevents mismatches between declared schema and actual data
3. Early detection of schema shift avoids cascading failures in Silver and Gold layers

```
1 # Read the file using Spark
2 df = spark.read.text("abfss://Olist_Worksplace@onelake.dfs.fabric.microsoft.com/LH_Bronze.Lakehouse/Files/bronze_schemas.py")
3
4 # Convert to string and execute
5 schema_code = "\n".join(row['value'] for row in df.collect())
6 exec(schema_code)
7
8 for table_name, expected_schema in schemas.items():
9     try:
10         actual_schema = spark.table(f'LH_Bronze.{table_name}').schema
11         assert actual_schema == expected_schema, f"Schema mismatch for {table_name}\n"
12     except Exception as e:
13         print(f"Error validating schema for {table_name}: {e}")
14         raise Exception("Schema shifted for {table_name} in Bronze Lakehouse:{e}")
```

SCHEMA FOR BRONZE LAKEHOUSE

```

from pyspark.sql.types import StructType, StructField, StringType, IntegerType, DoubleType, DateType, TimestampType

schemas = {
    'silver_customers' : StructType([
        StructField('customer_id', StringType(), True),
        StructField('customer_unique_id', StringType(), True),
        StructField('customer_zip_code_prefix', StringType(), True),
        StructField('customer_city', StringType(), True),
        StructField('customer_state', StringType(), True)
    ]),
    'silver_geolocation' : StructType([
        StructField('geolocation_zip_code_prefix', StringType(), True),
        StructField('geolocation_lat', StringType(), True),
        StructField('geolocation_lng', StringType(), True),
        StructField('geolocation_city', StringType(), True),
        StructField('geolocation_state', StringType(), True)
    ]),
    'silver_order_items' : StructType([
        StructField('order_id', StringType(), True),
        StructField('Order_item_id', IntegerType(), True),
        StructField('product_id', StringType(), True),
        StructField('seller_id', StringType(), True),
        StructField('Shipping_limit_date', TimestampType(), True),
        StructField('Price', FloatType(), True),
        StructField('freight_value', FloatType(), True)
    ]),
    'silver_order_payments' : StructType([
        StructField('order_id', StringType(), True),
        StructField('Payment_sequential', IntegerType(), True),
        StructField('payment_type', StringType(), True),
        StructField('payment_installments', IntegerType(), True),
        StructField('Payment_value', FloatType(), True)
    ]),
    'silver_order_reviews' : StructType([
        StructField('order_id', StringType(), True),
        StructField('Payment_sequential', IntegerType(), True),
        StructField('payment_type', StringType(), True),
        StructField('payment_installments', IntegerType(), True),
        StructField('Payment_value', FloatType(), True)
    ])
}

```

VALIDATION

```
'silver_orders' : StructType([StructField('order_id', StringType(), True),
                                StructField('customer_id', StringType(), True),
                                StructField('order_status', StringType(), True),
                                StructField('order_purchase_timestamp', TimestampType(), True),
                                StructField('order_approved_at', TimestampType(), True),
                                StructField('order_delivered_carrier_date', TimestampType(), True),
                                StructField('order_delivered_customer_date', TimestampType(), True),
                                StructField('order_estimated_delivery_date', TimestampType(), True),
                                StructField('purchase_year', IntegerType(), True),
                                StructField('purchase_month', IntegerType(), True)
                               ]),
'silver_products' : StructType([
                                StructField('product_id', StringType(), True),
                                StructField('product_category_name', StringType(), True),
                                StructField('product_name_length', StringType(), True),
                                StructField('product_description_length', StringType(), True),
                                StructField('product_photos_qty', StringType(), True),
                                StructField('product_weight_g', StringType(), True),
                                StructField('product_length_cm', StringType(), True),
                                StructField('product_height_cm', StringType(), True),
                                StructField('product_width_cm', StringType(), True)
                               ]),
'silver_sellers' : StructType([
                                StructField('seller_id', StringType(), True),
                                StructField('seller_zip_code_prefix', StringType(), True),
                                StructField('seller_city', StringType(), True),
                                StructField('seller_state', StringType(), True)
                               ]),
'silver_product_translations' : StructType([
                                StructField('product_category_name', StringType(), True),
                                StructField('product_category_name_english', StringType(), True)
                               ])
}
```

NOTEBOOK SOURCE CODE FOR SILVER LAKEHOUSE VALIDATION

Validation for Silver Lakehouse

1) Validate that the tables in LH_Bronze Lakehouse have been successfully copied to the LH_Silver Lakehouse

1. Ensures no tables were lost, duplicated, or corrupted during the copy
2. Detects silent failures or partial writes due to timeouts, memory issues, or schema drift
3. Validates that the pipeline executed all steps as expected

```

1 # To load the table names in both LH_Silver and LH_Bronze
2 # Compare the table names present in LH_Bronze with the table names in LH_Silver
3 # Assert that there is no missing tables between both Lakehouse after running of the notebook
4
5 tables_name_bronze = []
6 tables_name_silver = []
7
8 bronze_tables = [t.name for t in spark.catalog.listTables('LH_Bronze')]
9 silver_tables = [t.name for t in spark.catalog.listTables('LH_Silver')]
10 tables_name_bronze = [t.removeprefix('olist_').removesuffix('_dataset') for t in bronze_tables]
11 tables_name_silver = [t.removeprefix('silver_') for t in silver_tables]
12 tables_name_silver = ['product_category_name_translation' if t == 'product_translations' else t for t in tables_name_silver]
13 missing_tables = [t for t in tables_name_silver if t not in tables_name_bronze]
14 assert not missing_tables, f"Missing tables in LH_Silver: {missing_tables}"

```

| ✓ - Session ready in 13 sec 430 ms. Command executed in 47 sec 45 ms by ChooGeok H (JDE07) on 9:43:21 AM, 10/06/25

2) Validate no schema shift and that the schema of the delta tables are the same as the schema declared

1. Guarantees that the structure of ingested data aligns with expectations
2. Prevents mismatches between declared schema and actual data
3. Early detection of schema shift avoids cascading failures in Gold layers

```

1 # To read the schema file from the LH_Silver file folders
2 df = spark.read.text("abfss://olist_Workspace@onelake.dfs.fabric.microsoft.com/LH_Silver.Lakehouse/Files/silver_schemas.py")
3
4 # Convert to string and execute to load the schema
5 schema_code = "\n".join(row['value'] for row in df.collect())
6 exec(schema_code)
7

```

```
8 # To extract the table name and schema
9 for table_name, expected_schema in schemas.items():
10     try:
11         actual_schema = spark.table(f'LH_Silver.{table_name}').schema
12         assert actual_schema == expected_schema, f"Schema mismatch for {table_name}\n"
13
14     except Exception as e:
15         print(f"Error validating schema for {table_name}: {e}")
16         raise Exception(f"Schema shifted for {table_name} in Silver Lakehouse:{e}")
✓ - Command executed in 16 sec 96 ms by ChooGeok H (JDE07) on 1:48:14 PM, 10/01/25
```

3) Row-Level Integrity Checks

1. Ensure critical fields are populated and Null checks on primary keys in the fact table

```
1 # To validate that the primary key of the fact table is not null
2 from pyspark.sql.functions import col
3 orders = spark.read.table('LH_Silver.silver_orders')
4 assert orders.filter(col("order_id").isNull()).count() == 0, "Missing order_id"
```

SCHEMA FOR SILVER LAKEHOUSE

```
from pyspark.sql.types import StructType, StructField, StringType, IntegerType, TimestampType, FloatType

schemas = {
    'silver_customers' : StructType([
        StructField('customer_id', StringType(), True),
        StructField('customer_unique_id', StringType(), True),
        StructField('customer_zip_code_prefix', StringType(), True),
        StructField('customer_city', StringType(), True),
        StructField('customer_state', StringType(), True)
    ]),
    'silver_geolocation' : StructType([
        StructField('geolocation_zip_code_prefix', StringType(), True),
        StructField('geolocation_lat', StringType(), True),
        StructField('geolocation_lng', StringType(), True),
        StructField('geolocation_city', StringType(), True),
        StructField('geolocation_state', StringType(), True)
    ]),
    'silver_order_items' : StructType([
        StructField('order_id', StringType(), True),
        StructField('Order_item_id', IntegerType(), True),
        StructField('product_id', StringType(), True),
        StructField('seller_id', StringType(), True),
        StructField('Shipping_limit_date', TimestampType(), True),
        StructField('Price', FloatType(), True),
        StructField('freight_value', FloatType(), True)
    ]),
    'silver_order_payments' : StructType([
        StructField('order_id', StringType(), True),
        StructField('Payment_sequential', IntegerType(), True),
        StructField('payment_type', StringType(), True),
        StructField('payment_installments', IntegerType(), True),
        StructField('Payment_value', FloatType(), True)
    ]),
    'silver_order_reviews' : StructType([
        StructField('review_id', StringType(), True),
        StructField('order_id', StringType(), True),
        StructField('Review_score', IntegerType(), True),
        StructField('review_comment_title', StringType(), True),
        StructField('review_comment_message', StringType(), True),
        StructField('Review_creation_date', TimestampType(), True),
        StructField('review_answer_timestamp', TimestampType(), True)
    ])
}
```

VALIDATION

```
'silver_orders' : StructType([
    StructField('order_id', StringType(), True),
    StructField('customer_id', StringType(), True),
    StructField('order_status', StringType(), True),
    StructField('order_purchase_timestamp', TimestampType(), True),
    StructField('order_approved_at', TimestampType(), True),
    StructField('order_delivered_carrier_date', TimestampType(), True),
    StructField('order_delivered_customer_date', TimestampType(), True),
    StructField('order_estimated_delivery_date', TimestampType(), True),
    StructField('purchase_year', IntegerType(), True),
    StructField('purchase_month', IntegerType(), True)
]),
'silver_products' : StructType([
    StructField('product_id', StringType(), True),
    StructField('product_category_name', StringType(), True),
    StructField('product_name_length', StringType(), True),
    StructField('product_description_length', StringType(), True),
    StructField('product_photos_qty', StringType(), True),
    StructField('product_weight_g', StringType(), True),
    StructField('product_length_cm', StringType(), True),
    StructField('product_height_cm', StringType(), True),
    StructField('product_width_cm', StringType(), True)
]),
'silver_sellers' : StructType([
    StructField('seller_id', StringType(), True),
    StructField('seller_zip_code_prefix', StringType(), True),
    StructField('seller_city', StringType(), True),
    StructField('seller_state', StringType(), True)
]),
'silver_product_translations' : StructType([
    StructField('product_category_name', StringType(), True),
    StructField('product_category_name_english', StringType(), True)
])
}
```

NOTEBOOK SOURCE CODE FOR GOLD LAKEHOUSE VALIDATION

Validation for Gold Lakehouse

1) Validate that the tables in LH_Silver Lakehouse have been successfully copied to the LH_Gold Lakehouse

```

1 # To load the table names in both LH_Silver and LH_Gold
2 # Compare the table names present in LH_Silver with the table names in LH_Gold
3 # Assert that there is no missing tables between both Lakehouse after running of the notebook
4
5 tables_name_silver = []
6 tables_name_gold = []
7
8 silver_tables = [t.name for t in spark.catalog.listTables('LH_Silver')]
9 gold_tables_all = [t.name for t in spark.catalog.listTables('LH_Gold')]
10 gold_tables = [name for name in gold_tables_all if name.startswith('gold_')]
11
12 tables_name_silver = [t.removeprefix('silver_') for t in silver_tables]
13 tables_name_gold = [t.removeprefix('gold_') for t in gold_tables]
14
15 missing_tables = [t for t in tables_name_silver if t not in tables_name_gold]
16 assert not missing_tables, f"Missing tables in LH_Gold: {missing_tables}"
!6] ✓ - Command executed in 9 sec 783 ms by ChooGeok H (JDE07) on 3:32:34 PM, 9/30/25

```

[+](#) Code [+](#) Markdown

2) Validate that the schema of tables in both LH_Silver and LH_Gold Lakehouse is the same

```

1 # To extract the schema (column names, types) from both LH_Silver and LH_Gold Lakehouse
2 # Assert that the schemas in both Lakehouse is the same
3
4 for table_name in tables_name_silver:
5     df_silver_schema = spark.table(f'LH_Silver.silver_{table_name}').schema
6     df_gold_schema = spark.table(f'LH_Gold.gold_{table_name}').schema
7     assert df_silver_schema == df_gold_schema

```

3) Validate that the row count of tables in both LH_Silver and LH_Gold Lakehouse is the same

```

1 # To ensure row-level fidelity with row count comparison
2 # Assert that the number of rows in the same table name of both Lakehouse is the same
3
4 for table_name in tables_name_silver:
5     df_silver = spark.table(f'LH_Silver.silver_{table_name}')
6     df_gold = spark.table(f'LH_Gold.gold_{table_name}')
7     assert df_silver.count() == df_gold.count(), f"Row count mismatch in {table_name}"
!5] ✓ - Command executed in 7 sec 933 ms by ChooGeok H (JDE07) on 3:32:45 PM, 9/30/25

```

VALIDATION OF PIPELINE FLOW

To ensure transparency, reliability, and rapid response across the project pipeline, a structured validation framework was implemented alongside automated email notifications. Each Lakehouse layer—Bronze, Silver, and Gold—undergoes rigorous validation, and successful execution triggers targeted alerts to the data engineering team.

Layer-by-Layer Validation Summary

1) Bronze Lakehouse

Validation Steps:

- Bronze_DeltaTable: Failed
- Bronze_CopyBlob: Succeeded

Activity:

Email_Notification_Bronze is triggered when the Bronze_DeltaTable fails.

2) Silver Lakehouse

Validation Steps:

- Silver_Cleaning: Failed
- Bronze_Validation: Succeeded
- Bronze_DeltaTable: Succeeded
- Bronze_CopyBlob: Succeeded

Activity:

Email_Notification_Silver is triggered when the Silver_Cleaning fails.

3) Gold Lakehouse

Validation Steps:

- Gold_Transform: Failed
- Silver_Validation: Succeeded
- Silver_Cleaning: Succeeded
- Bronze_Validation: Succeeded
- Bronze_DeltaTable: Succeeded
- Bronze_CopyBlob: Succeeded

Activity:

Email_Notification_Gold is triggered when the Gold_Transform fails.

4) Semantic Model Refresh

Validation Steps:

- ML_Semantic_model: Succeeded
- Gold_Validation: Succeeded
- Gold_Transform: Succeeded
- Silver_Validation: Succeeded
- Silver_Cleaning: Succeeded
- Bronze_Validation: Succeeded
- Bronze_DeltaTable: Succeeded
- Bronze_CopyBlob: Succeeded

Activity:

New_Report_Notification is triggered when the whole pipeline runs is completed.

VALIDATION

Email Notification Logic

- Each validation notebook is linked to a corresponding notification trigger.
- Upon failed execution of the notebook, an email is sent via Outlook to the designated recipients.
- The email includes:
 - Pipeline stage (Bronze, Silver, Gold, Semantic Model)
 - Execution status
 - Timestamp and activity name
- This ensures the teams involved are informed in real time and can act immediately if any stage fails.

Benefits

- **Operational Transparency:** Clear visibility into pipeline health across all layers.
- **Rapid Response:** Immediate alerts enable faster troubleshooting and resolution.
- **Auditability:** Logged validation outcomes support governance and historical analysis.

Validation of Pipeline Flow

Trigger of Email Notification

Lakehouse Stage	Activity Name	Activity Status
Bronze Lakehouse	Error_Email_Notification_Bronze	Succeeded
Bronze Lakehouse	Bronze_DeltaTable	Failed
Bronze Lakehouse	Bronze_CopyBlob	Succeeded
Silver Lakehouse	Error_Email_Notification_Silver	Succeeded
Silver Lakehouse	Silver_Cleaning	Failed
Silver Lakehouse	Bronze_Validation	Succeeded
Silver Lakehouse	Bronze_DeltaTable	Succeeded
Silver Lakehouse	Bronze_CopyBlob	Succeeded
Gold Lakehouse	Error_Email_Notification_Gold	Succeeded
Gold Lakehouse	Gold_Transform	Failed
Gold Lakehouse	Silver_Validation	Succeeded
Gold Lakehouse	Silver_Cleaning	Succeeded
Gold Lakehouse	Bronze_Validation	Succeeded
Gold Lakehouse	Bronze_DeltaTable	Succeeded
Gold Lakehouse	Bronze_CopyBlob	Succeeded
Semantic Model Refresh	New_Report_Notification	Succeeded
Semantic Model Refresh	ML_Semantic_model	Succeeded
Semantic Model Refresh	Gold_Validation	Succeeded
Semantic Model Refresh	Gold_Transform	Succeeded
Semantic Model Refresh	Silver_Validation	Succeeded
Semantic Model Refresh	Silver_Cleaning	Succeeded
Semantic Model Refresh	Bronze_Validation	Succeeded
Semantic Model Refresh	Bronze_DeltaTable	Succeeded
Semantic Model Refresh	Bronze_CopyBlob	Succeeded

Validation of Pipeline Flow

Error Email Notification

Failed Notebook Run for Bronze Lakehouse

Ignore Block Delete Archive Report Reply Reply Forward Meeting Chat Share in Teams Zoom

Deleted Failed Summarize

Email Address

Error Notification for the Bronze_DeltaTable Activity

Hi Team,
The scheduled Fabric pipeline encountered an error during notebook execution.
Date: 2025-10-02T08:24:18.702177Z
Workspace: Olist_Worksplace
Workspace ID: 3a670224-6cd8-472d-b8ed-f3df1feef537e
Notebook Activity Name: Bronze_DeltaTable
Notebook Name: NB_Bronze
Pipeline Name: Project_Pipeline
Please review the notebook logs and validate the input dataset.

Regards,
IT Support Department

New Report Email Notification

New Report Created

Ignore Block Delete Archive Report Reply Reply Forward Meeting Chat Share in Teams Zoom Move Sweep Rules

Deleted New Report Created Summarize

Email Address

New Analysis Report Update

Hi Team,
A new update of the Business Analysis Power BI report is now available for your analysis and review.
Date: 2025-10-02T08:38:22.366673Z
Workspace: Olist_Worksplace
Workspace ID: 3a670224-6cd8-472d-b8ed-f3df1feef537e
Pipeline Name: Project_Pipeline
Please access the report via the Power BI Apps and share any insights or feedback.

Regards,
IT Support Department

Project Deliverable

PROJECT DELIVERABLE

PROJECT PIPELINE OVERVIEW

The project pipeline is a modular, orchestrated system that spans ingestion, validation, analytics, and alerting. Built on Microsoft Fabric, it integrates notebooks, semantic modeling, and Power BI dashboards with automated email notifications for robust observability and stakeholder communication.

Medallion Layer

The Medallion Layer structures the data lifecycle into Bronze, Silver, and Gold stages, enabling progressive refinement from raw ingestion to enrich, analytics-ready datasets.

- **Bronze_DeltaTable:** Convert CSV files into Delta Table and save them in Bronze Lakehouse.
- **Silver_Cleaning:** Applies type casting, null handling, and enrichment logic.
- **Gold_Transform:** Copy delta tables from Silver Lakehouse into Gold Lakehouse and perform transformation.

Validation Layer

This foundational layer ensures data integrity and schema consistency across the Bronze, Silver, and Gold Lakehouse stages.

- **Bronze_Copy_Validation:** Initiates data ingestion from source files.
- **Bronze_Validation:** Confirms successful conversion to Delta format and schema alignment.
- **Silver_Validation:** Validates schema and row count consistency between Bronze and Silver.
- **Gold_Validation:** Ensures schema and data integrity between Silver and Gold.

Market Basket Analysis Layer

This layer applies Apriori-based association rule mining to uncover product affinities and drive personalized recommendations.

- **ML_Transform:** Consolidates cleaned data into a unified schema for ML and reporting.
- **ML_Market_Basket_Analysis:** Executes the Apriori algorithm using mlxtend within a Fabric notebook. Generates frequent item sets and association rules.
- **ML_Semantic_Model:** Refresh of data for dashboard.
- **ML_Report_Notification:** Sends email notifications with insights and rule summaries to stakeholders.

Power BI Analysis Layer

This layer delivers curated dashboards for leadership decision support.

- **ML_Semantic_Model:** Refresh of data for Power BI analysis report.
- **PowerBI_Report_Notification:** Triggers email alerts upon successful semantic model refresh and dashboard updates.

Error Handling Layer

This layer ensures proactive alerting and rapid response to pipeline failures.

- **Outlook:** Automatically send error notifications when validation or transformation notebooks fail. Emails include pipeline stage, error context, and timestamp for root cause analysis.

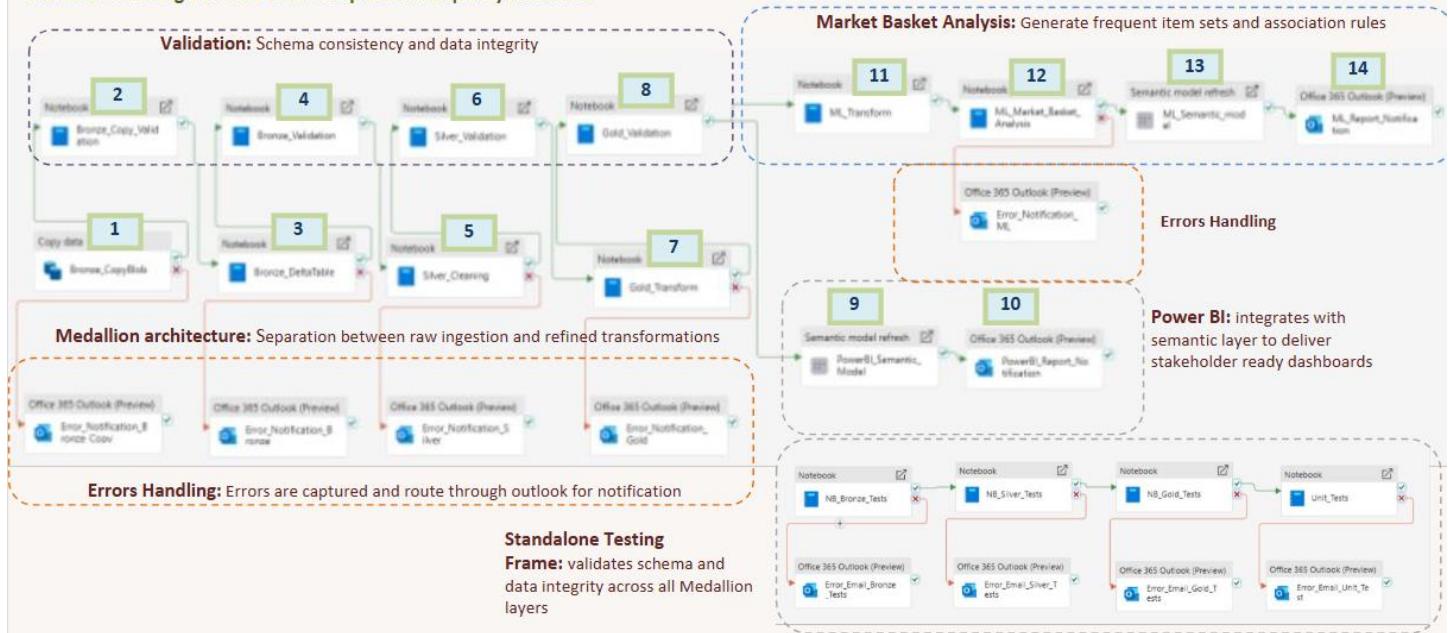
Architectural Highlights

- **Modular Notebook Design:** Each task is encapsulated in reusable notebooks for validation, transformation, and ML logic.
- **Integrated Alerting:** Outlook notifications are triggered at key checkpoints to ensure transparency and accountability.
- **Semantic Modeling:** Power BI semantic layer supports drill-through analytics and stakeholder-friendly reporting.
- **End-to-End Traceability:** Logs and validation results are stored in Delta tables and visualized via monitoring dashboards.

Screen shot of Project Pipeline

PROJECT PIPELINE

End-to-end pipeline from validation to reporting, engineered with a robust, modular, and scalable architecture.
It spans the full lifecycle from Medallion-based transformations.
Standalone testing framework for comprehensive quality assurance.

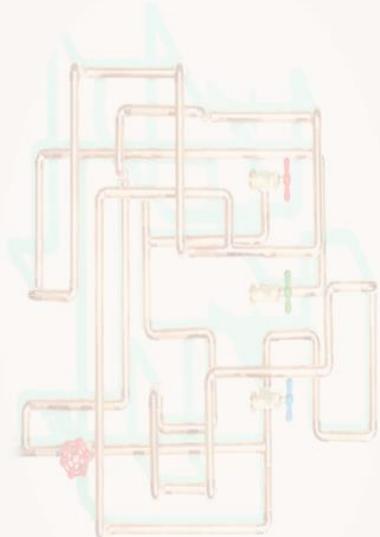


Screen shot of Project Pipeline Status

PROJECT PIPELINE

Screenshot showing the pipeline status with all activities marked as succeeded

Parameters	Variables	Settings	Output	Library variables
Pipeline run ID	764453da-711d-4212-b2d5-1a6d112f8abf			View run detail
	<input type="text"/> Filter by keyword			Showing 1 - 14 items
Activity name	Activity status	Run start		
ML_Report_Notification	✓ Succeeded	10/8/2025, 11:43:09 AM		
ML_Semantic_model	✓ Succeeded	10/8/2025, 11:42:25 AM		
ML_Market_Basket_Analysis	✓ Succeeded	10/8/2025, 11:33:53 AM		
PowerBI_Report_Notification	✓ Succeeded	10/8/2025, 11:33:04 AM		
ML_Transform	✓ Succeeded	10/8/2025, 11:32:23 AM		
PowerBI_Semantic_Model	✓ Succeeded	10/8/2025, 11:32:23 AM		
Gold_Transform	✓ Succeeded	10/8/2025, 11:30:37 AM		
Silver_Validation	✓ Succeeded	10/8/2025, 11:28:35 AM		
Silver_Cleaning	✓ Succeeded	10/8/2025, 11:27:23 AM		
Bronze_Validation	✓ Succeeded	10/8/2025, 11:25:09 AM		
Bronze_DeltaTable	✓ Succeeded	10/8/2025, 11:24:13 AM		
Bronze_Copy_Validation	✓ Succeeded	10/8/2025, 11:21:12 AM		
Bronze_CopyBlob	✓ Succeeded	10/8/2025, 11:20:49 AM		



VISUALIZATION AND REPORTING

Setting the context right...

 Persona 4: Matthew Taylor – Strategy Director

As Strategy Director, Matthew's role is to **drive marketplace growth, profitability, and customer trust**. He needs a single, consolidated view of **business performance across sales, customer retention, delivery efficiency, and seller reliability**.

Problem Statement

- Leadership **lacks a consolidated view** of sales, customer retention, delivery, and seller performance on one page.

Project Objective

- To deliver an **executive-level summary** that brings together all marketplace KPIs into a single, intuitive dashboard.

1. Are we growing in the right places?

"As Strategy Director, I need to quickly see whether our marketplace is expanding in the right categories and regions — not just topline revenue."

2. Are customers staying with us?

"Our growth is hollow if customers don't return. I need clarity on churn vs. loyalty."

3. Are our goods delivered timely as promise?

"Delays erode trust fast. I want an early-warning system on delivery performance."

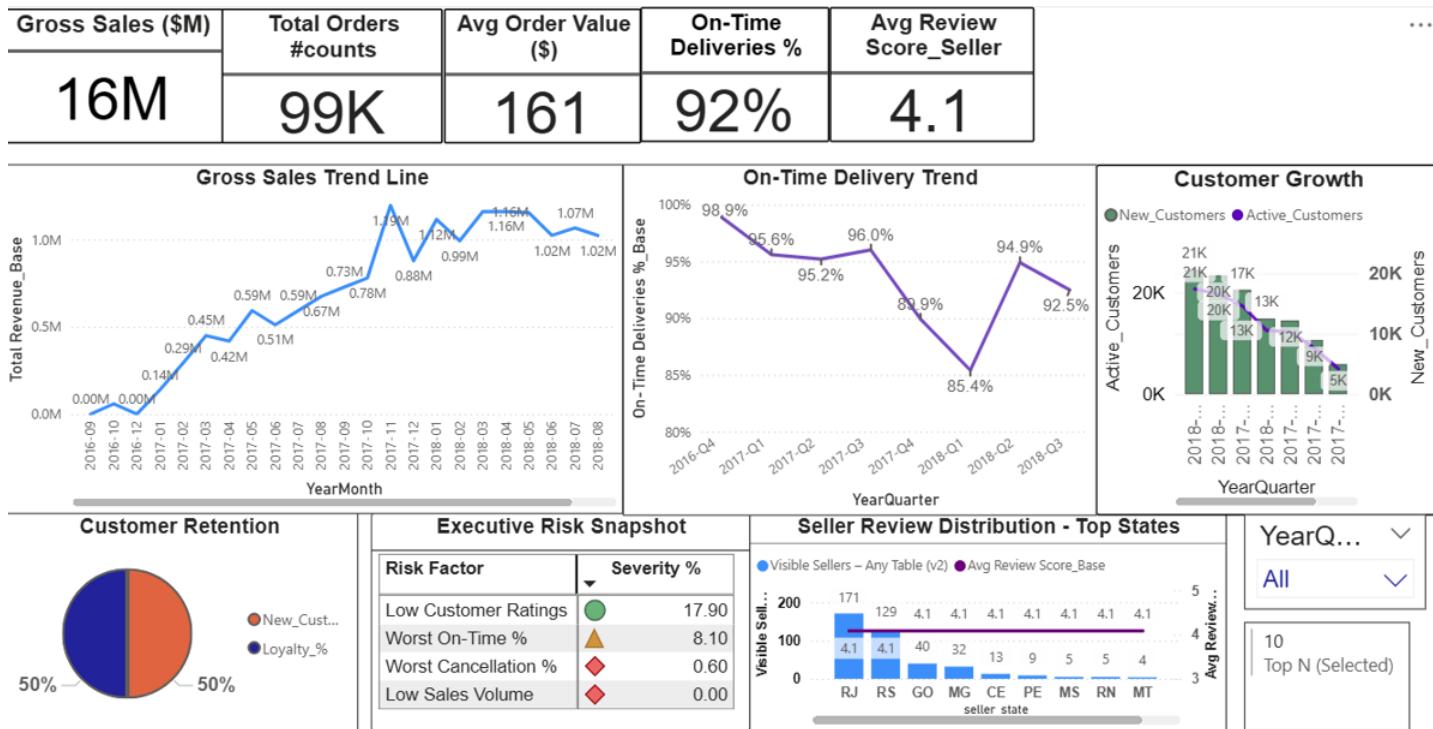
4. Are sellers strengthening or harming trust?

"Our sellers are the marketplace. Good ones build trust; weak ones destroy it."

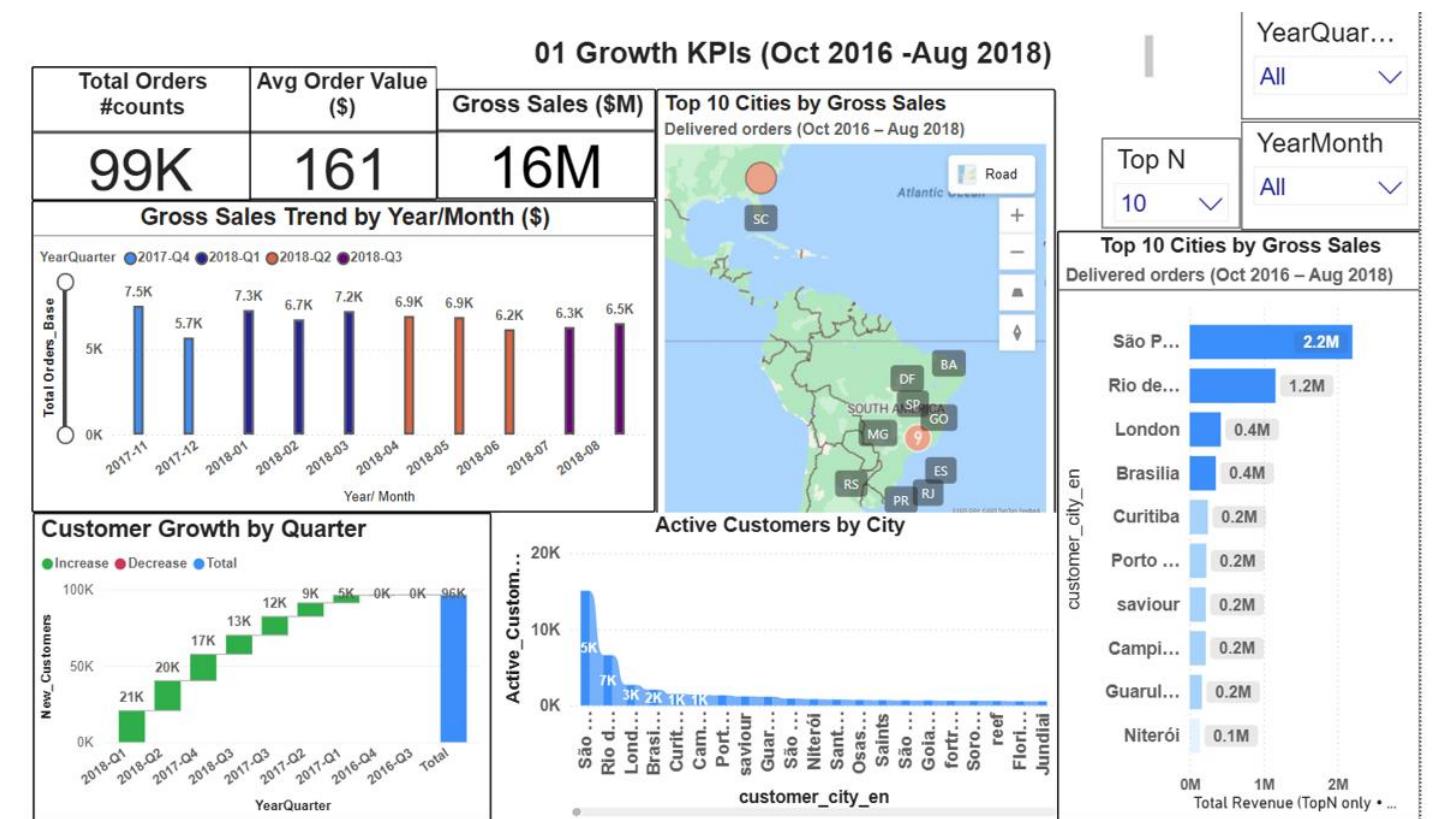
5. Where are the biggest risks this month?

"I don't want 20 pages of KPIs — I need a red flag list of risks now."

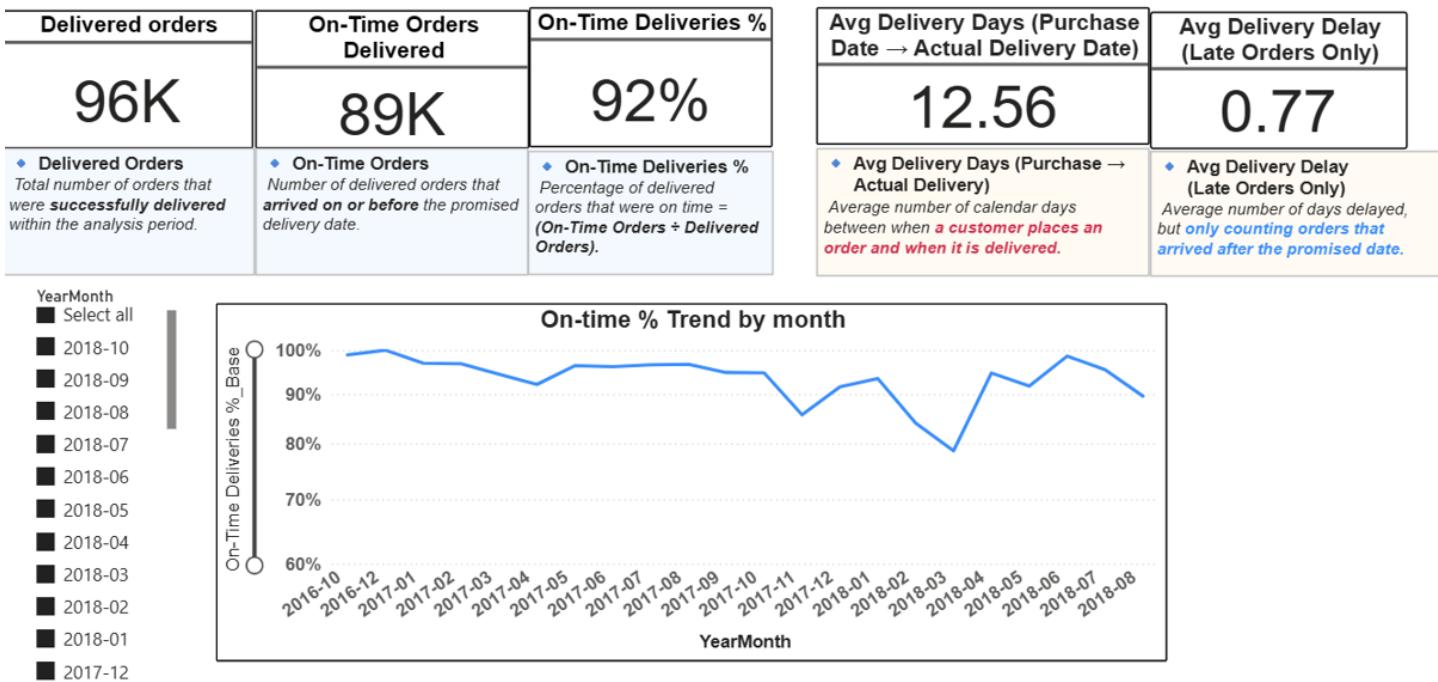
Executive Summary (Oct 2016 -Aug 2018)



PROJECT DELIVERABLE



02 Delivery Performance (Oct 2016 -Aug 2018)



03 Customer Retention (Oct 2016 -Aug 2018)

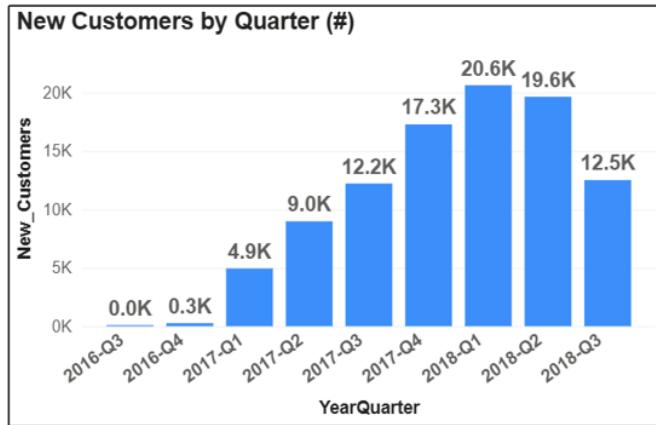
Quarterly new-customer acquisition is slowing.

At the same time, the real story in the data is that active customers average 6–7K per month, but are trending downward into 2018.

Repeat purchases are rare, so overall lift depends on **winning back inactive customers** and securing second purchases quickly.

Possible Reasons For Decline in New Customers

- **Market Saturation** – Rapid early growth left fewer untapped customers.
- **Seasonality** – Q3/Q4 shifts in demand (holidays, school, budget cycles).
- **Operational Issues** – Bottlenecks in logistics or onboarding slowed growth.
- **Marketing Shift** – Less ad spend, more focus on retention/profitability.
- **Competition** – New rivals and promotions diverted potential customers.



04 Seller Trust & Ranking (Oct 2016 -Aug 2018)

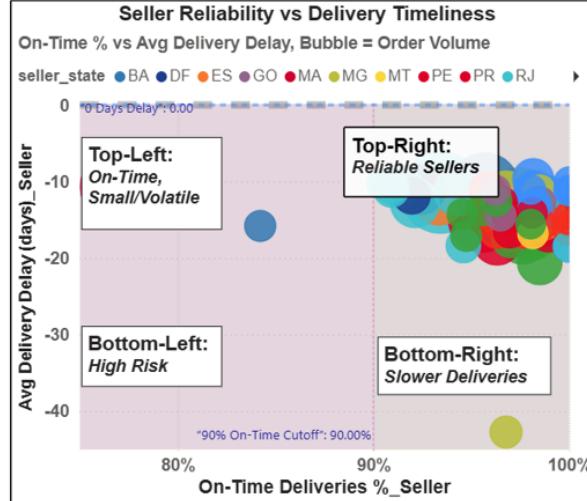
On-Time Deliveries %	Delivered Orders
92%	96K

% Orders Cancelled	Delivery Delay (days)
0.5%	-11.20

YearMonth
Multiple s... ▾

Seller Reliability & Order Performance

seller_id	Delivered Orders	On-Time Deliveries %
6560211a19b47992c3666cc44a7e94c0	1819	0.94
4a3ca9315b744ce9f8e9374361493884	1772	0.89
cc419e0650a3c5ba77189a1882b7556a	1651	0.94
1f50f920176fa81dab994f9023523100	1399	0.89
da8622b14eb17a	1311	0.92
Total	32788	0.92



Top N
10 ▾

seller_id
All ▾

05 TOP 5 Risks

✖ Worst On-Time % — Top 10 Sellers			
seller_id	Visible Sellers – Any Table (v2)	On-Time Deliveries %_Seller	seller_state
001cca7ae9ae17fb1caed9dfb1094831	1	0.93	ES
001e6ad469a905060d959994f1b41e4f	1		RJ
02a2272692e13558373c66db98f05e2e	1		RJ
02d35243ea2e497335cd0f076b45675d	1	0.64	RN
Total	424	0.83	

⌚ Worst Avg Delay (days) — Top 10 Sellers			
seller_id	Visible Sellers – Any Table (v2)	Avg Delivery Delay (days)_Seller	seller_state
001e6ad469a905060d959994f1b41e4f	1		RJ
02a2272692e13558373c66db98f05e2e	1		RJ
02d35243ea2e497335cd0f076b45675d	1	-5.45	RN
0417b067eeab773d2f7061a726dc477f	1	-3.21	SC
04843805947f0fc584fc1969hb6e50fe7	1	0.63	RS
Total	183	-10.47	

🚫 Worst Cancellation Rate % — Top 10 Sellers			
seller_id	Visible Sellers – Any Table (v2)	Cancellation Rate %_Seller	seller_state
001cca7ae9ae17fb1caed9dfb1094831	1		ES
001e6ad469a905060d959994f1b41e4f	1	1.00	RJ
003554e2dce176b5555353e4f3555ac8	1		GO
01b600d251a011280000701b6060b2a47	1		CO
Total	581	0.02	

Total No. of Sellers

3095

Show Worst N

Top N ▾

10

Top N (Selected)

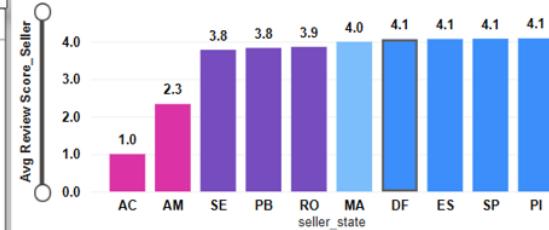
10 ▾

🔴 Low Sales Volume — Showing Bottom 10 Sellers

● Delivered Orders_Seller ● Top N (Selected)



Low Customer Ratings — Showing Bottom 10 Sellers



Risks Identified & Action Steps

The Risk Snapshot shows three hot spots: **low ratings, on-time delivery gaps, and cancellations.**

Action: roll out an **Earned Trust program** — coaching and incentives tied to OTD, cancellation, and review lifts.

Seller reviews are concentrated in top states but flat at ~4.1.

Focus needs to be on lifting underperforming high-volume states.

Net: growth is real, but we protect it by pushing OTD ≥95% and reviews ≥4.3.”

Monthly sales confirm seasonality spikes — we should replicate peak-month playbooks two months early. Revenue is highly concentrated in a handful of states.

Action:

Strategy: **land-and-expand** in the top five, and pilot two adjacent states with similar demographics.

Quarterly new customers are slowing, so overall lift hinges on retention. We must win back inactive customers and secure second purchases within 30 days.

By city, we see a long tail of small but meaningful customer bases.

Our next \$1M comes from targeting the **next 10 cities** with lightweight demand gen and faster delivery promises.

Decision: Double down on top-5 states and surgically activate the next 10 cities.”

Growth is strong, retention is steady, but seller reliability is our weak spot To scale, we must **protect trust** with OTD ≥95%, reviews ≥4.3, and a ruthless focus on reliable sellers.”

Risks Identified & Action Steps

"Here's execution reliability.

We delivered **96K orders, 89K on time** → **92% OTD**. Average delivery is **12.6 days** from purchase, and late orders are only **0.77 days delayed**.

Good news: delays are shallow; bad news: OTD is inconsistent.

Trend line shows months in the mid-80s OTD, then recovery to ~95%. This volatility is the biggest threat to trust.

Action: enforce stricter SLAs with lagging sellers/carriers, and fund predictive monitoring to flag risks before they damage NPS."

"Here's highlights our seller-level reliability.

Overall: **92% OTD, 96K orders, cancellation only 0.5%**. But within that, large sellers dip to **89% OTD**, below our 90% cutoff.

The scatterplot shows clusters of sellers consistently underperforming while still driving volume — structural risk.

Action: implement a **Seller Reliability Program** with a 95% OTD minimum threshold, tiered performance recognition, and penalties for laggards."

"Now to loyalty.

Headline metrics show **100% retention, new, and repeat purchase** (due to setup quirk). The real story is in the detail: **6–7K active customers per month, trending downward** into 2018.

New vs returning looks balanced, but churn is high — we're replacing lost customers with new ones instead of compounding loyalty.

Action: approve retention programs. Specifically:

- Loyalty reward for **second purchase within 30 days**
- Win-back campaigns for customers inactive beyond 90 days."

"Finally, here's where exposure lives.

- **Worst OTD %:** Some sellers as low as **83%**, concentrated in RJ, ES, RN.
- **Worst Avg Delays:** >10 days in RJ and RN.
- **Worst Cancellation Rates:** up to **100%**, abandoned sellers.
- **Low Sales Volume:** 3K+ sellers contribute negligible orders.
- **Low Ratings:** states AC and AM average **1.0–2.3 ratings**, damaging trust.

Action: approve a '**Cut or Coach**' program. We offboard bottom-decile sellers, retrain salvageable ones, and protect expansion states by removing sub-3.5 rating sellers.

Challenges And Future Roadmap

CHALLENGES AND FUTURE ROADMAP

As part of the project initiative, the team identified key technical challenges that could compromise pipeline reliability and reporting accuracy. To address these risks, robust solutions were implemented across the notebook execution and semantic modeling layers, ensuring proactive error handling and seamless dashboard refresh.

Challenge 1: Ensuring Robust Error Handling in Notebooks

Problem:

Notebook failures such as schema mismatches or missing files can silently disrupt pipeline execution, leading to incomplete data processing and delayed reporting.

Solution:

- Implemented try-except blocks with exception chaining to capture and propagate errors.
- Integrated Outlook alerting to notify the data engineering team with actionable error context, including pipeline name, timestamp, and failure details.
- Ensured that all validation notebooks log execution status to a centralized Delta table for traceability and trend analysis.

Challenge 2: Power BI Semantic Model Refresh Failures

Problem:

Power BI dashboards may fail to update due to schema changes or upstream data issues, resulting in outdated or missing insights for stakeholders.

Solution:

- Automated semantic model refresh using Fabric pipeline triggers.
- Logged refresh status and validation results to the monitoring dashboard.
- Configured Outlook notifications to alert stakeholders upon successful or failed refresh, enabling rapid response and resolution

Challenge 3: Inconsistent data validation across pipeline stages

Problem:

Data anomalies (e.g., nulls, outliers, schema drift) may pass undetected from Bronze to Gold layers, leading to misleading insights or report failures.

Solution:

Design modular validation checkpoints at each layer (Bronze, Silver, Gold) using reusable notebook functions. Trigger alerts via Outlook when anomalies exceed thresholds, and log semantic tags for downstream traceability.

Architectural Impact

- Strengthened pipeline resilience through modular error handling and alerting.
- Improved observability with centralized logging and dashboard integration.
- Enhanced stakeholder confidence by ensuring timely, accurate, and transparent reporting

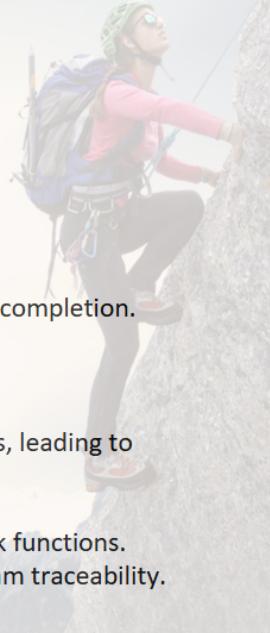
Challenges and Solutions

1. Ensuring robust errors handling in notebooks

Notebook failures (e.g., schema mismatch, missing files) can silently break pipelines.

Solution:

Implement try-except blocks and trigger Outlook alert with actionable context.



2. Power BI report fail to update with new data

Power BI semantic models may fail to refresh due to schema changes or upstream data issues.

Solution:

Automate refresh using Fabric pipeline triggers semantic model fresh and notify via Outlook upon completion.

3. Inconsistent data validation across pipeline stages

Data anomalies (e.g., nulls, outliers, schema drift) may pass undetected from Bronze to Gold layers, leading to misleading insights or report failures.

Solution:

Design modular validation checkpoints at each layer (Bronze, Silver, Gold) using reusable notebook functions. Trigger alerts via Outlook when anomalies exceed thresholds, and log semantic tags for downstream traceability.

IMPACT AND FUTURE ROADMAP

The project has delivered measurable improvements in operational efficiency and data reliability through modular design, automated workflows, and multi-layer validation. Looking ahead, the team is focused on expanding observability and alerting capabilities to further enhance responsiveness and stakeholder confidence.

Impact

1. Operational Efficiency

- Modular Notebooks: Reusable validation and transformation logic reduces manual intervention and accelerates development cycles.
- Automated Execution: Ingestion, transformation, and alerting workflows are fully automated, ensuring consistent and timely data delivery.

2. Data Trust & Quality

- Multi-Step Validation: Embedded schema checks, row count reconciliations, and semantic consistency validations ensure high data integrity across Bronze, Silver, and Gold layers.
- Semantic Alignment: Cleaned and enriched data supports reliable reporting and machine learning workflows.

Future Roadmap

1. Alerting & Observability

- Power Automate Integration: Extend alerting channels beyond Outlook to include Slack, Microsoft Teams, and SMS.
- Real-Time Monitoring: Enhance dashboards with live pipeline status, error trends, and semantic model refresh metrics.
- Proactive Diagnostics: Introduce anomaly detection and predictive alerts to identify issues before they impact reporting.

Impact and Future Roadmap

IMPACT

1. Operational Efficiency

- Modular notebooks and reusable validation logic reduce manual intervention.
- Automated ingestion, transformation, and alerting pipeline execution.

2. Data Trust & Quality

- Multi-step validation ensures schema and data integrity.

FUTURE ROADMAP

1. Alerting & Observability

- Use of Power Automate integration to include Slack, Teams, or SMS alerts.