

End-to-End Data Engineering Project: Building a Complete Analytics Pipeline with PySpark and DBT

Mohd Laraib Khan | khan.laraib.laraib77@gmail.com | www.linkedin.com/in/khanofdataengineering

Project Overview

This project demonstrates the construction of a complete data engineering pipeline using a pseudo-Uber dataset. The pipeline follows the Medallion Architecture (Bronze for raw data, Silver for cleaned and transformed data, and Gold for modeled analytics-ready data). It handles incremental data ingestion, data cleaning, deduplication, upserts, and dimensional modeling with slowly changing dimensions (SCD Type 2). The dataset includes CSV files for entities like trips, customers, drivers, vehicles, payments, and locations with completed Views for BI Analytics use at the end.

The project is built on Databricks (free tier) and emphasizes production-grade practices such as idempotency, cost optimization, and modular code. It showcases real-world data engineering skills applicable to cloud platforms like Azure, GCP, and AWS.

Tools and Technologies Used

- **PySpark:** For structured streaming, data ingestion, transformations, deduplication, and upserts.
- **Databricks:** Platform for compute, storage (Delta Lake), catalogs, schemas, volumes, and notebooks.
- **DBT (Data Build Tool):** For data modeling, incremental builds, snapshots, sources, and orchestration (using DBT Cloud free tier).
- **Jinja:** Templating engine within DBT for dynamic SQL generation (e.g., loops and conditionals).
- **Python:** For modular utilities, classes, and reusable transformation logic.
- **Delta Lake:** Storage format for reliable, ACID-compliant tables with merge capabilities.
- **YAML:** For DBT configurations, sources, and snapshots.
- **SQL:** Core language for DBT models and PySpark transformations.
- **Git:** Version control for DBT projects (integrated in DBT Cloud).
- **CSV Files:** Source data format for the pseudo-Uber dataset.

Relevant Skills Demonstrated

- PySpark Structured Streaming for incremental ingestion.
- Dynamic schema handling and notebook automation.
- Data cleaning techniques (e.g., regex, concatenation, conditional columns).
- Object-oriented programming in Python for reusable transformations.
- Deduplication using window functions and hashing.
- Upsert operations with Delta Tables and conditional merges.
- DBT modeling including incremental materializations, SCD Type 2 via snapshots, and lineage tracking.
- Jinja templating for dynamic SQL.
- Databricks environment management (catalogs, schemas, volumes, checkpoints).
- Data quality and governance (audit timestamps, idempotent processing).
- Modular and scalable pipeline design for production readiness.

Step-by-Step Project Build Process

Step 1: Environment Setup on Databricks

I started by setting up the environment in Databricks using the free community edition. This involved creating a serverless compute cluster with autoscaling to handle processing efficiently without manual VM management.

- Created a catalog named `pyspark_dbt`.
- Defined schemas: `source` for raw data, `bronze` for ingested data, `silver` for transformed data, and `gold` for modeled data.
- Set up a volume `source_data` under the `source` schema, with subdirectories for each entity (e.g., `trips`, `customers`, `drivers`, `vehicles`, `payments`, `locations`).
- Uploaded CSV files into the respective subdirectories.
- Created a checkpoint volume under the `bronze` schema to store metadata for incremental processing.

databricks

Free Edition

Search d

New

Home

Workspace

Recents

Catalog

Jobs & Pipelines

Compute

Marketplace

SQL

SQL Editor

Queries

Dashboards

Genie

Alerts

Query History

SQL Warehouses

Catalog

Serverless Starter Warehouse

Serverless

2XS

Type to search...

For you

All

My organization

>

workspace

>

system

▼

pysparkdbt

▼

bronze

▼

Tables (6)

customers

drivers

locations

payments

trips

vehicles




▼


Volumes (1)

checkpoint

Databricks catalog and schemas setup

Catalog Explorer > pysparkdbt > source >


source_data






Share ▾


Upload to this volume

Overview Files Details Permissions

Description


 AI generate

Add

/Volumes/pysparkdbt/source/source_data / **customers**


Refresh

Create directory

Name	Size	Last modified
 customers.csv	20.10 KB	12 days ago

Source data volume with uploaded CSV files

Step 2: Bronze Layer - Incremental Data Ingestion with PySpark

In this step, I used PySpark Structured Streaming to ingest raw CSV data incrementally into Delta tables in the bronze layer. This ensures only new files are processed, making the pipeline efficient.

- Defined a list of entities: ["customers", "trips", "vehicles", "drivers", "payments", "locations"].
- In a Databricks notebook, used a loop to dynamically process each entity:
 - Read streaming data with `spark.readStream.format("csv").option("header", "true").option("inferSchema", "true").`
 - Extracted schema dynamically from a batch DataFrame to avoid hardcoding.
 - Wrote to bronze tables using `format("delta").outputMode("append").option("checkpointLocation", "<path>").trigger(once=True).start().`
- This created 6 Delta tables in the bronze schema, preserving raw data with checkpoints for fault tolerance.

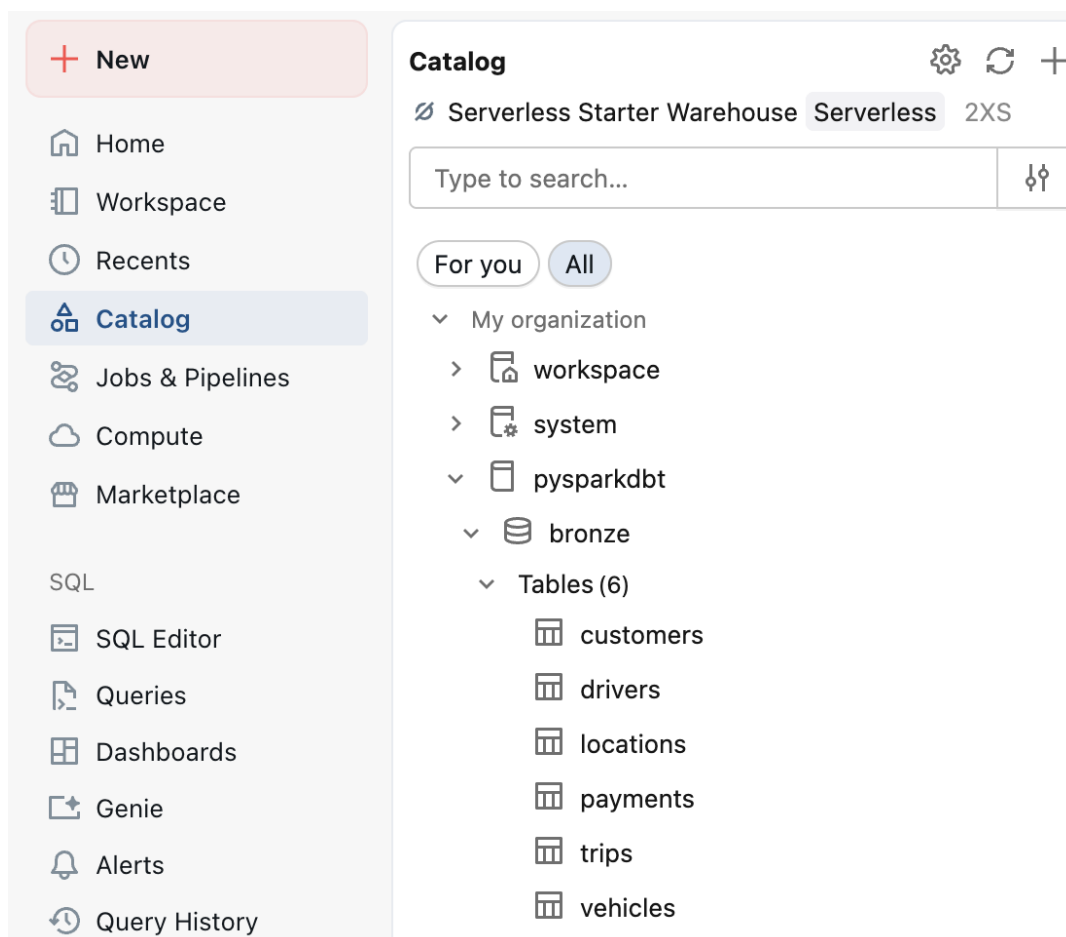
```
bronze_ingestion x +
File Edit View Run Help Python Tabs: ON ☆ Last edit was 11 days ago Run all Serverless Schedule Share
Dec 10, 2025 6
for entity in entities:
    df_batch = spark.read.format("csv")\
        .option("header", True)\
        .option("inferSchema", True)\
        .load(f"/Volumes/pysparkdbt/source/source_data/{entity}/")

    schema_entity = df_batch.schema

    df = spark.readStream.format("csv")\
        .option("header", True)\
        .schema(schema_entity)\
        .load(f"/Volumes/pysparkdbt/source/source_data/{entity}")

    df.writeStream.format("delta")\
        .outputMode("append")\
        .option("checkpointLocation", f"/Volumes/pysparkdbt/bronze/checkpoint/{entity}")\
        .option("mergeSchema", "true")\
        .trigger(once=True)\
        .toTable(f"pysparkdbt.bronze.{entity}")
> See performance (6)
```

PySpark code for dynamic ingestion loop

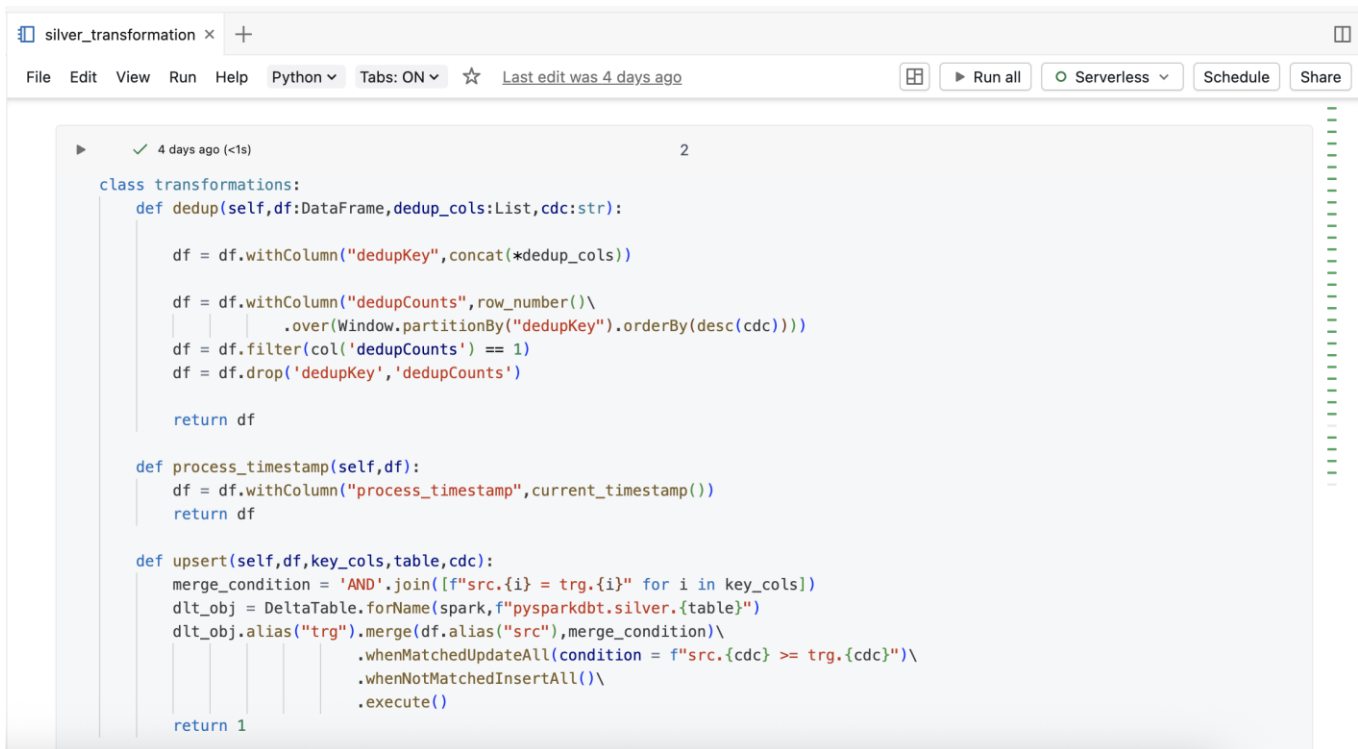


Bronze layer Delta tables in Databricks catalog

Step 3: Silver Layer - Data Transformation with PySpark

Here, I focused on cleaning and enriching the data using modular Python code. I created a utility file `custom_utils.py` with a Transformations class to encapsulate reusable logic.

- Read data from bronze tables: `spark.read.table("bronze.<entity>")`.
- Applied generic transformations:
 - Deduplication: Used window functions (`row_number` over partitioned hash keys) to keep the latest records based on a change data capture (CDC) timestamp.
 - Added `process_timestamp` using `current_timestamp()` for auditing.
 - Upserts: Checked if silver table exists; if yes, performed Delta merge with dynamic conditions (e.g., update if source timestamp is newer).
- Entity-specific cleaning (examples):
 - Customers: Extracted email domains with `split`, cleaned phone numbers with `regexp_replace`, concatenated full names with `concat_ws`.
 - Payments: Added conditional `online_payment_status` using `when().otherwise()`.
 - Vehicles: Converted `make` to uppercase.
- Looped over entities to apply transformations and upsert into `silver.<entity>` tables.



```
class transformations:
    def dedup(self, df: DataFrame, dedup_cols: List, cdc: str):
        df = df.withColumn("dedupKey", concat(*dedup_cols))

        df = df.withColumn("dedupCounts", row_number()\
            .over(Window.partitionBy("dedupKey").orderBy(desc(cdc))))
        df = df.filter(col('dedupCounts') == 1)
        df = df.drop('dedupKey', 'dedupCounts')

        return df

    def process_timestamp(self, df):
        df = df.withColumn("process_timestamp", current_timestamp())
        return df

    def upsert(self, df, key_cols, table, cdc):
        merge_condition = 'AND'.join([f"src.{i} = trg.{i}" for i in key_cols])
        dlt_obj = DeltaTable.forName(spark, f"pysparkdbt.silver.{table}")
        dlt_obj.alias("trg").merge(df.alias("src"), merge_condition)\
            .whenMatchedUpdateAll(condition = f"src.{cdc} >= trg.{cdc}")\
            .whenNotMatchedInsertAll()\
            .execute()

        return 1
```

Transformations class in custom_utils.py

The screenshot shows a Databricks notebook titled "silver_transformation" with a menu bar (File, Edit, View, Run, Help) and a language dropdown set to "Python". The notebook contains five code cells, each with a play button, a status indicator (green checkmark), a timestamp, and a cell number. The code in the cells is as follows:

- Cell 6: `df_cust = spark.read.table("pysparkdbt.bronze.customers")`
Output: `df_cust: pyspark.sql.connect.dataframe.DataFrame = [customer_id: integer, first_name: string ... 6 more fields]`
- Cell 7: `df_cust = df_cust.withColumn('domain', split(col('email'), '@')[1])`
Output: `df_cust: pyspark.sql.connect.dataframe.DataFrame = [customer_id: integer, first_name: string ... 7 more fields]`
- Cell 8: `df_cust = df_cust.withColumn("phone_number", regexp_replace("phone_number", r"^[0-9]", ""))`
Output: `df_cust: pyspark.sql.connect.dataframe.DataFrame = [customer_id: integer, first_name: string ... 7 more fields]`
- Cell 9: `df_cust = df_cust.withColumn("full_name", concat_ws(col('first_name'), col('last_name')))`
`df_cust = df_cust.drop('first_name', 'last_name')`
Output: `df_cust: pyspark.sql.connect.dataframe.DataFrame = [customer_id: integer, email: string ... 6 more fields]`
- Cell 10: `cust_obj = transformations()`

`cust_df_trns = cust_obj.dedup(df_cust[['customer_id', 'last_updated_timestamp']])`
`display(cust_df_trns)`

A "Microsoft OneNote" watermark is visible in the bottom right corner of the notebook interface.

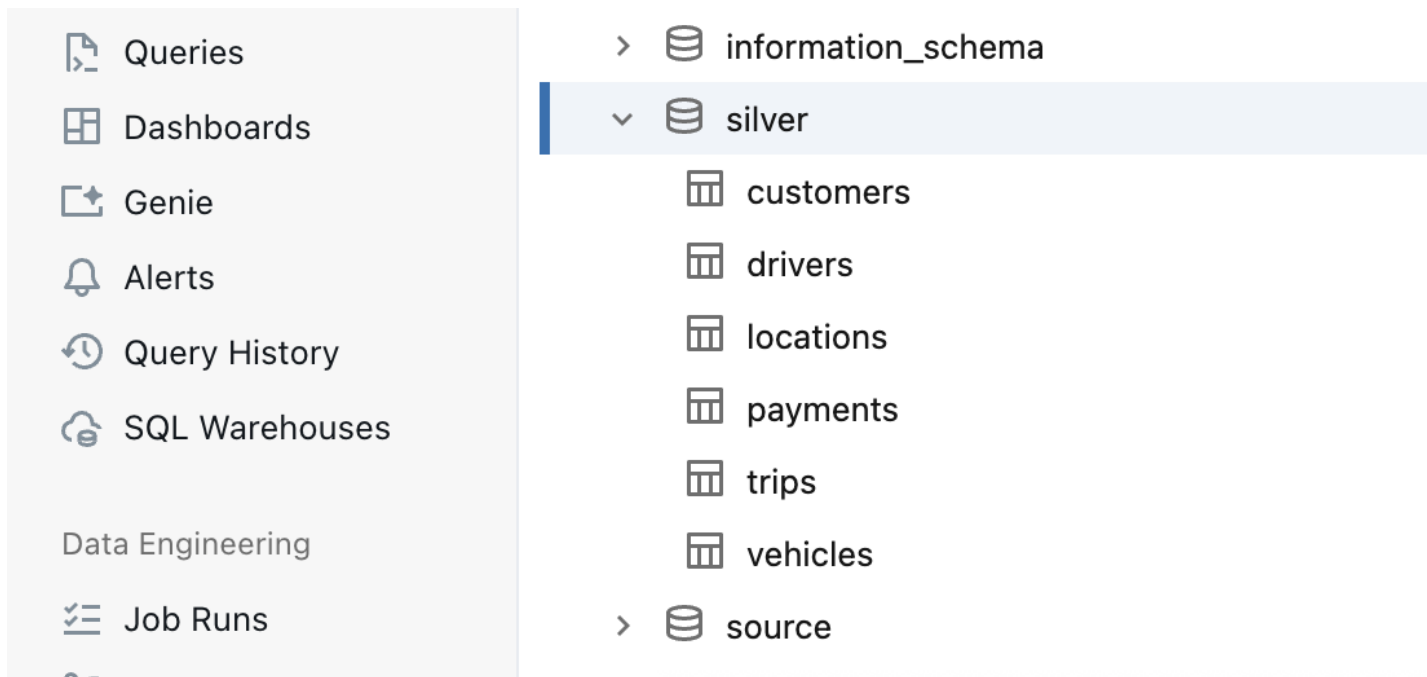
The screenshot shows the final code cell (Cell 12) of the "silver_transformation" notebook. The code implements an upsert function:

```
#this is the upsert function

if not spark.catalog.tableExists("pysparkdbt.silver.customers"):
    df_cust.write.format("delta")\
        .mode("append")\
        .saveAsTable("pysparkdbt.silver.customers")
else:
    cust_obj.upsert(df_cust, ['customer_id'], 'customers', 'last_updated_timestamp')
```

The cell output is a link: [See performance \(2\)](#).

Example PySpark code for customer transformations and upsert



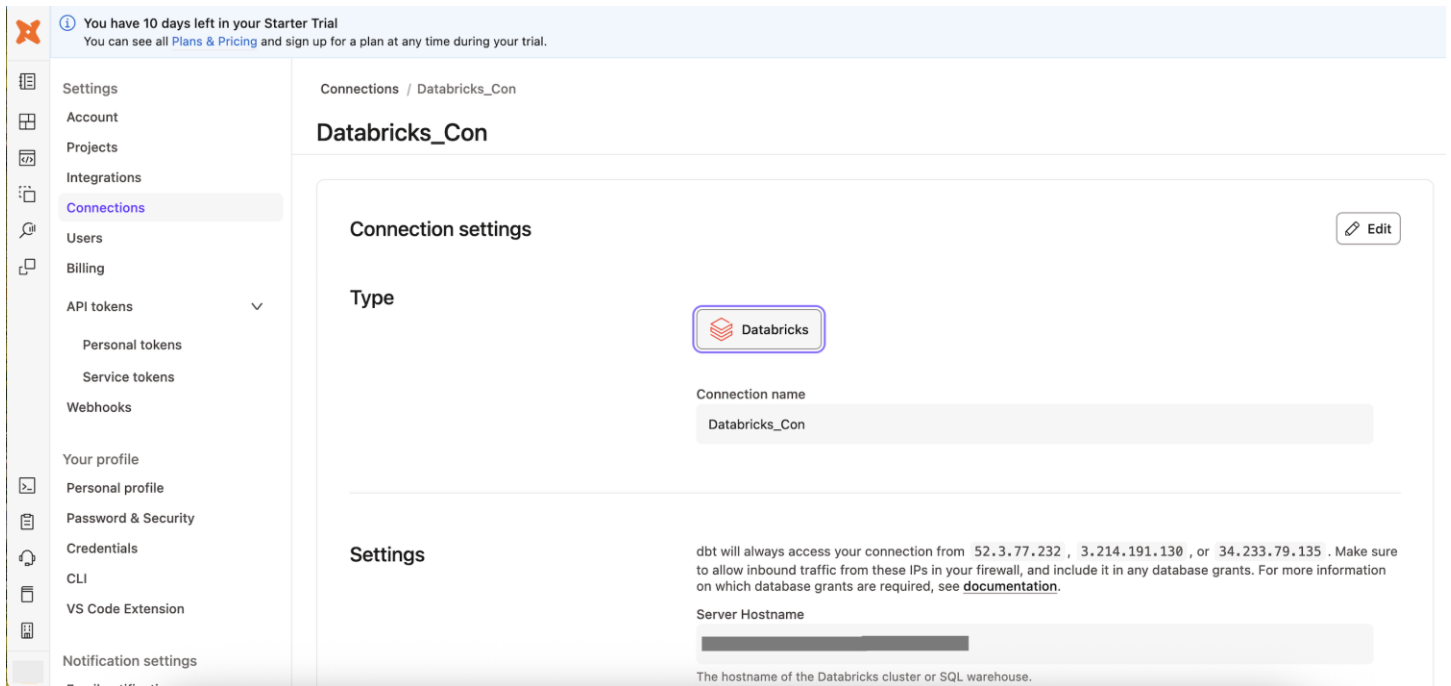
Silver layer tables after processing

Step 4: DBT Setup and Integration

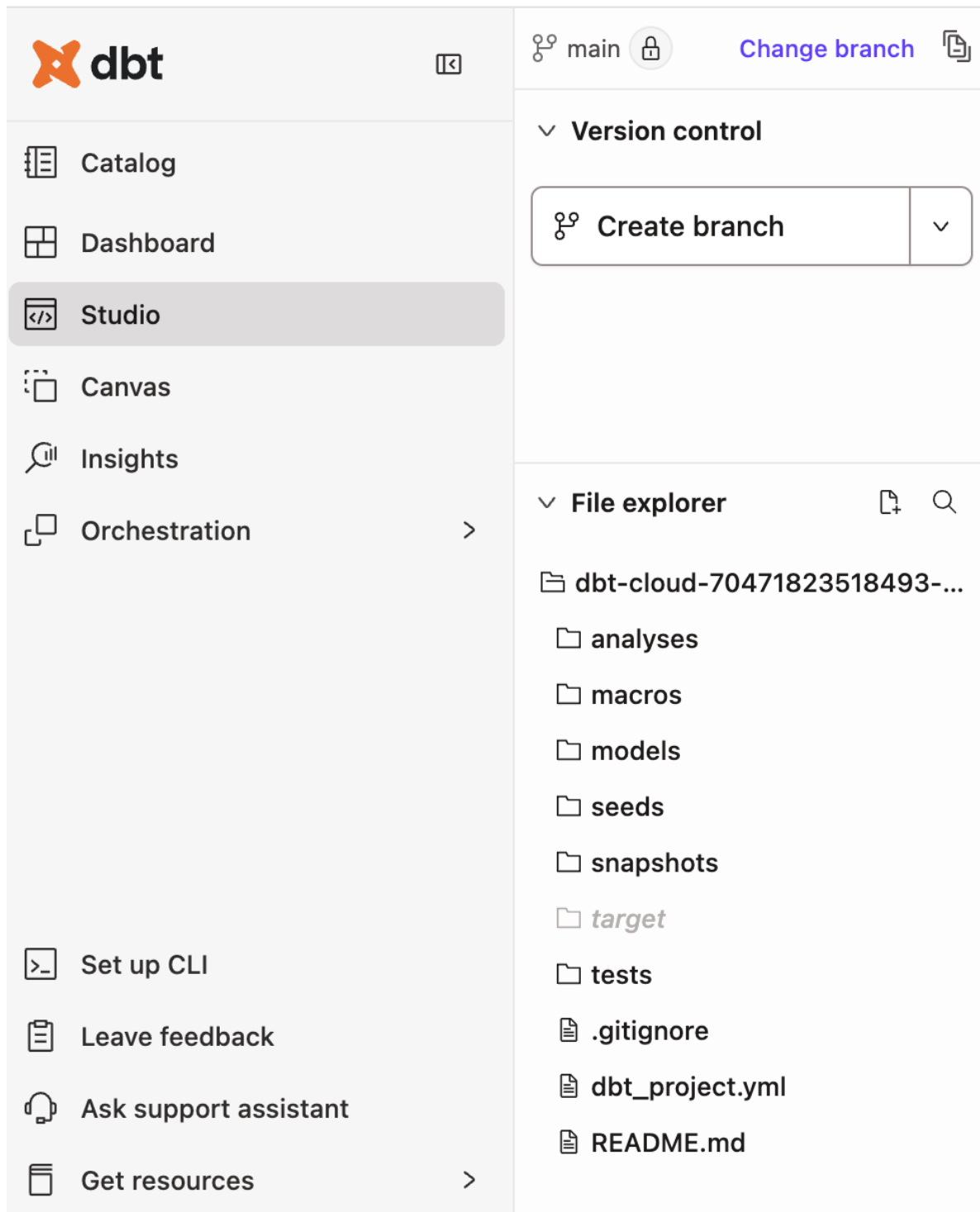
I shifted to DBT Cloud (free tier) for the modeling phase. Connected DBT to Databricks using hostname, HTTP path, access token, catalog, and schema details.

- Initialized a DBT project in DBT Cloud Studio.
- Created a development branch for safe experimentation.
- Defined project structure: models/silver/, models/gold/, sources/, snapshots/, macros/.

- Customized schema generation with a macro generate_schema_name.sql to use clean names like silver without prefixes.



DBT Cloud connection setup to Databricks

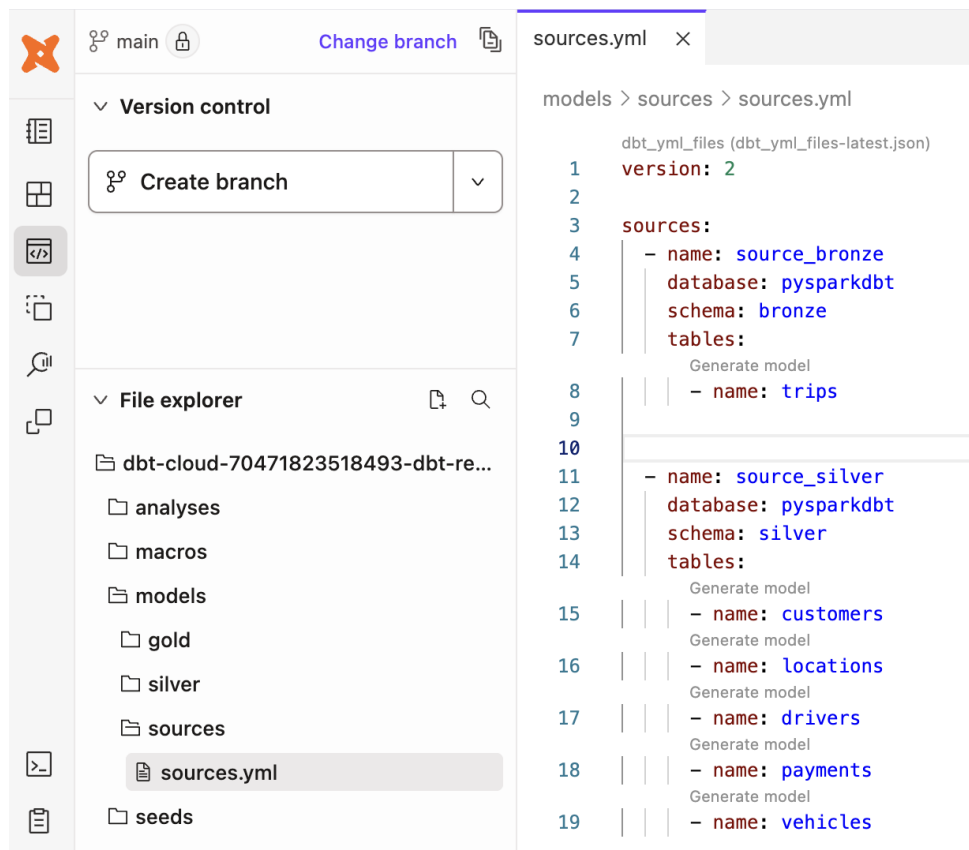


DBT project structure in Studio

Step 5: Gold Layer - Dimensional Modeling with DBT

This step involved building a star schema with dimension and fact tables using DBT models, sources, and snapshots.

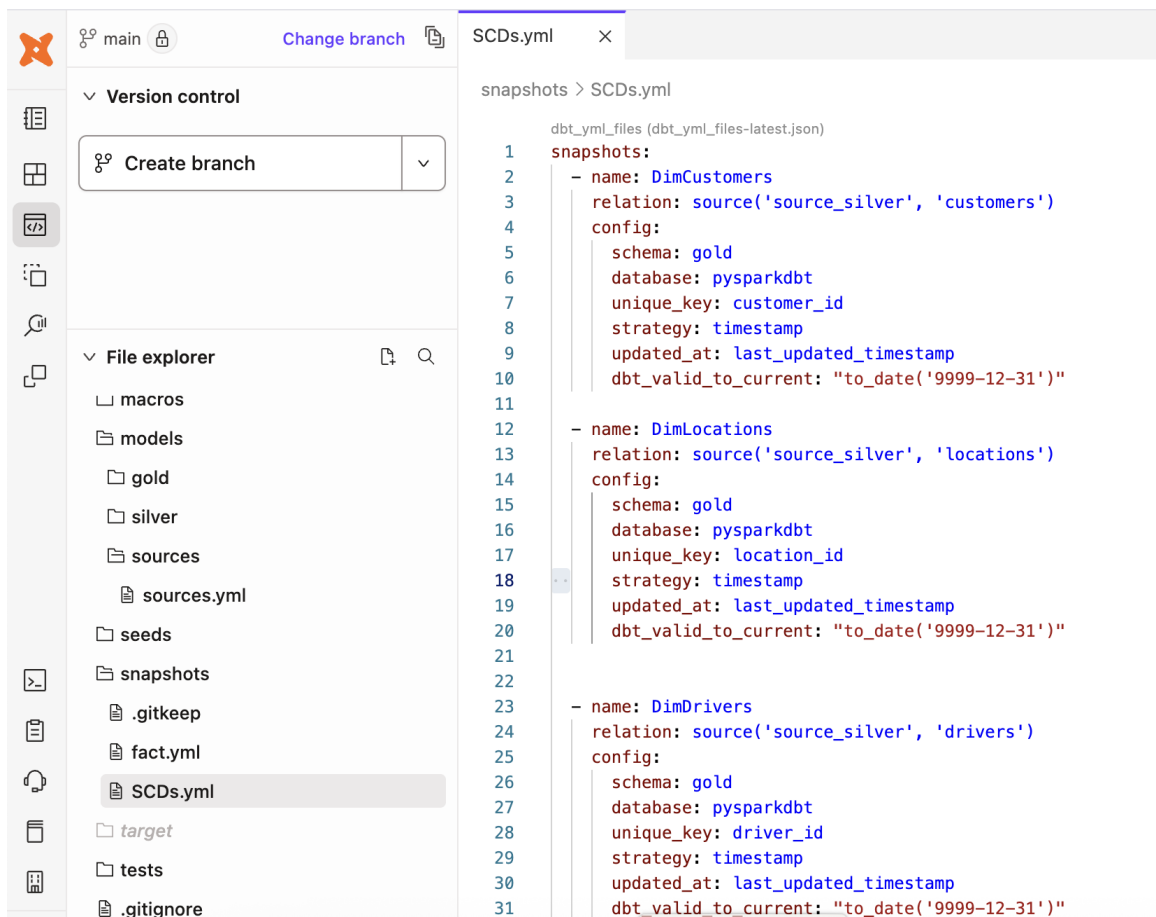
- Defined sources in sources.yml to register silver tables for lineage.
- Created incremental models (e.g., silver/trips.sql):
 - Used config(materialized='incremental', unique_key='trip_id').
 - Applied Jinja for dynamic SELECT: Looped over column lists with {% for %}...{% if not loop.last %},{% endif %}.
 - Filtered incremental data with is_incremental() macro and last_updated_timestamp.
- Implemented SCD Type 2 for dimensions via snapshots (e.g., dim_customers.yml):
 - Strategy: timestamp with unique_key and updated_at.
 - Added dbt_valid_from and dbt_valid_to for versioning.
- Created fact table fact_trips referencing silver models with ref('trips') and joining dimensions.
- Ran dbt run for models and dbt snapshot for SCDs.
- Viewed lineage and documentation in DBT.



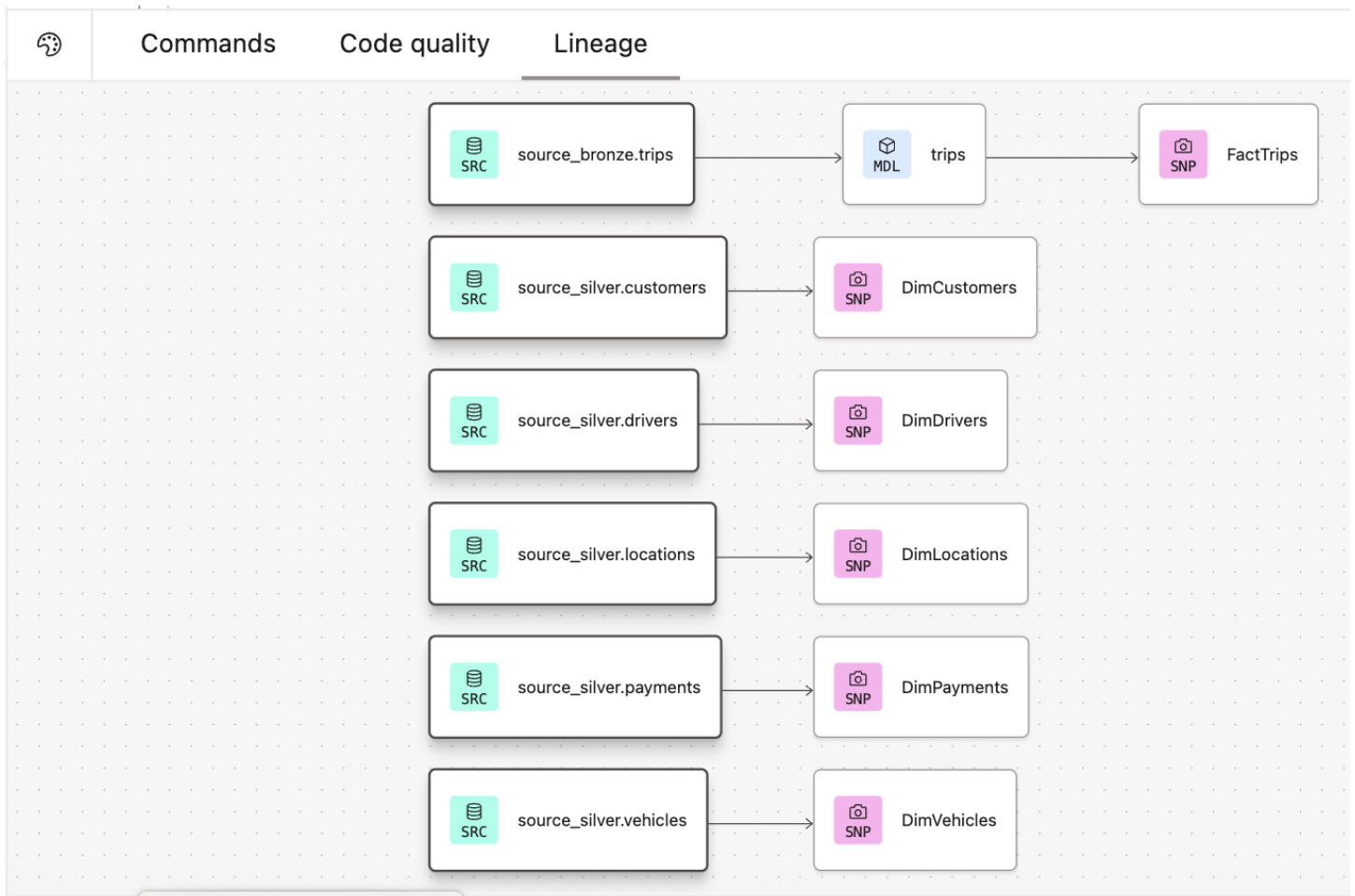
sources.yml file



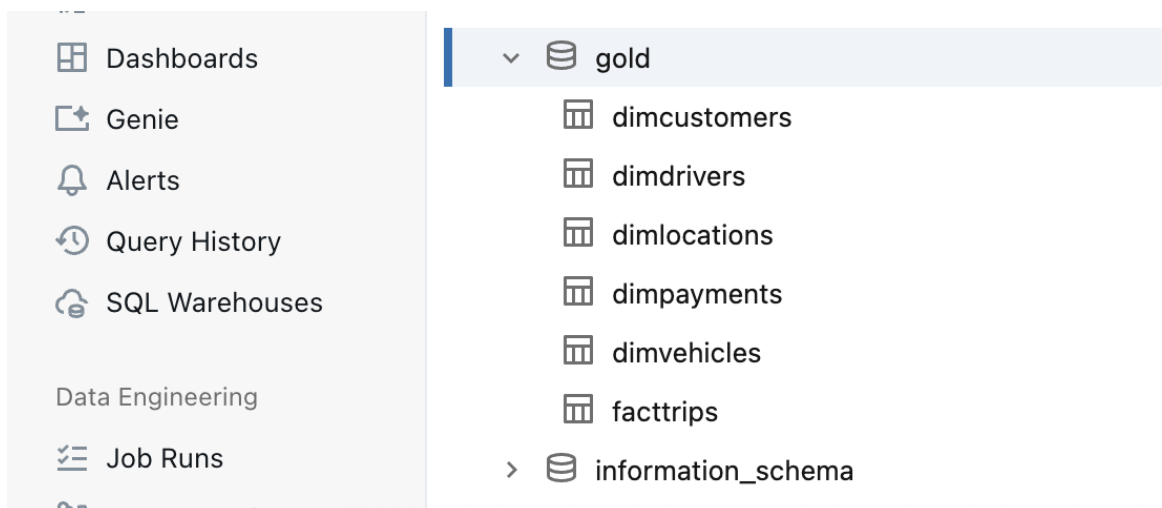
Example DBT model with Jinja templating



DBT snapshot configuration for SCD



DBT lineage graph



Gold layer tables in Databricks (dimensions and facts)

Step 6: Final Deployment and Testing

- Committed changes to Git in DBT Cloud and merged to main.
 - Ran full pipeline: Ingested new data in PySpark, then executed DBT jobs.
 - Verified data in Databricks catalog and shared connection details for BI tools (e.g., Power BI).
 - Ensured idempotency by re-running with no duplicates or errors.
- > System logs

All 6 Pass 6 Warn 0 Error 0 Skip 0 Running 0

>	✓ DimCustomers	17.26s
>	✓ DimDrivers	16.03s
>	✓ DimLocations	15.90s
>	✓ DimPayments	16.45s
>	✓ DimVehicles	8.98s
>	✓ FactTrips	6.92s

DBT run output and success logs

Step 7: Created Analytical Views for Business Intelligence

- Developed 3 SQL views in Databricks gold schema to support BI and analytics.
 - Utilized **CTEs** for modular query structure and **window functions** (RANK(), ROW_NUMBER(), AVG() OVER) for rankings and rolling averages.
- Views answer key business questions:

1. **vw_city_payment_performance**: Regional payment success rates and 3-month rolling fare averages for operational insights.

```
Business View 3 x +
Run 09:19 PM (1s) workspace: default New SQL editor: ON Last edit was 32 minutes ago
--Q3:What is the average fare amount and payment success rate per city, including a rolling average of fares over the last 3 months (assuming transaction_time in payments), to monitor regional performance?
CREATE VIEW pysparkdbt.gold.vw_city_payment_performance AS
WITH trip_payments AS (
  SELECT
    t.trip_id,
    d.city AS driver_city, -- Using driver city as proxy for trip city
    p.amount,
    p.payment_status,
    p.transaction_time,
    CASE WHEN p.payment_status = 'success' THEN 1 ELSE 0 END AS success_flag
  FROM pysparkdbt.gold.facttrips t
  JOIN pysparkdbt.gold.dimpayments p ON t.trip_id = p.trip_id
  JOIN pysparkdbt.gold.dimdrivers d ON t.driver_id = d.driver_id
),
monthly_aggregates AS (
  SELECT
    driver_city AS city,
    DATE_TRUNC('month', transaction_time) AS month,
    ROUND(AVG(amount),2) AS avg_fare,
    AVG(CAST(success_flag AS DOUBLE)) AS success_rate
  FROM trip_payments
  GROUP BY driver_city, DATE_TRUNC('month', transaction_time)
)
SELECT
  city,
  month,
  avg_fare,
  success_rate,
  ROUND(AVG(avg_fare) OVER (PARTITION BY city ORDER BY month ROWS BETWEEN 2 PRECEDING AND CURRENT ROW), 2) AS rolling_3m_avg_fare
FROM monthly_aggregates
ORDER BY city, month DESC;
```

View 1 Query Results

View 3 Query x +

Run 09:32 PM (11s) workspace: default New SQL editor: ON Edit

1 | SELECT * FROM pysparkdbt.gold.vw_city_payment_performance;

Add parameter

Table +

	city	month	1.2 avg_fare	1.2 success_rate	1.2 rolling_3m_avg_fare
1	Brownburgh	2025-09-01T00:00:00.000+00:00	65.54	0	57.99
2	Brownburgh	2025-08-01T00:00:00.000+00:00	54.7	0	58.31
3	Brownburgh	2025-07-01T00:00:00.000+00:00	53.72	0	60.11
4	Brownburgh	2025-06-01T00:00:00.000+00:00	66.5	0	66.5
5	Charlesborough	2025-09-01T00:00:00.000+00:00	44.42	0	54.94
6	Charlesborough	2025-08-01T00:00:00.000+00:00	47.8	0	56.12
7	Charlesborough	2025-07-01T00:00:00.000+00:00	72.6	0	60.28
8	Charlesborough	2025-06-01T00:00:00.000+00:00	47.96	0	47.96
9	Christianshire	2025-09-01T00:00:00.000+00:00	63.49	0	56.32
10	Christianshire	2025-08-01T00:00:00.000+00:00	59.18	0	52.73

2. vw_top_drivers_by_revenue: Top 10 drivers by total fare, trip counts, and ratings for performance analysis.

```

Business View 1 x +
Run 04:33 PM (1s) workspace default New SQL editor: ON Last edit was 6 hours ago Edit Serverless Star... 2XS Schedule Share Save*
1 --Q1: What are the top 10 drivers ranked by their total fare revenue, including their average rating and number of trips completed, to identify
2 high-performing drivers?
3 CREATE VIEW pysparkdbt.gold.vw_top_drivers_by_revenue AS
4 WITH driver_trips AS (
5     SELECT
6         d.driver_id,
7         d.full_name,
8         d.driver_rating,
9         d.city,
10        COUNT(t.trip_id) AS num_trips,
11        ROUND(SUM(t.fare_amount),2) AS total_fare,
12        ROUND(AVG(t.fare_amount),2) AS avg_fare_per_trip
13    FROM pysparkdbt.gold.facttrips t
14    JOIN pysparkdbt.gold.dimdrivers d ON t.driver_id = d.driver_id
15    GROUP BY d.driver_id, d.full_name, d.driver_rating, d.city
16 )
17 SELECT
18     driver_id,
19     full_name,
20     driver_rating,
21     city,
22     num_trips,
23     total_fare,
24     avg_fare_per_trip,
25     RANK() OVER (ORDER BY total_fare DESC) AS revenue_rank
26 FROM driver_trips
27 ORDER BY revenue_rank
28 LIMIT 10;

```

View 2 Query Results

View 1 Query x

+

Run

1 minute ago (16s)

workspace.default

New SQL editor: ON

☆

⋮

Edit

```
1 | select * from `pysparkdbt`.`gold`.`vw_top_drivers_by_revenue` limit 100;
```

Add parameter

Table

+

🔍

🏠

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

📄

3. **vw_customer_lifetime_value**: Customer LTV, trip frequency, and recency for segmentation and retention strategies.

```
Business View 2 x +
Run 04:33 PM (<1s) workspace.default New SQL editor: ON Last edit was 6 hours ago Edit Serverless Star... 2XS Schedule Share Save*
1  --Q2: For each customer, what is their lifetime value (total fare paid), average trip distance, and recency rank based on their last signup or trip date, to
2  analyze customer engagement?
3  CREATE VIEW pysparkdbt.gold.vw_customer_lifetime_value AS
4  WITH customer_stats AS (
5      SELECT
6          c.customer_id,
7          c.full_name,
8          c.city,
9          c.signup_date,
10         ROUND(SUM(t.fare_amount),2) AS lifetime_value,
11         ROUND(AVG(t.distance_km),0) AS avg_distance_km,
12         COUNT(t.trip_id) AS num_trips,
13         MAX(t.trip_id) AS last_trip_id -- Assuming trip_id increases over time as a proxy for recency
14     FROM pysparkdbt.gold.dimcustomers c
15     LEFT JOIN pysparkdbt.gold.facttrips t ON c.customer_id = t.customer_id
16     GROUP BY c.customer_id, c.full_name, c.city, c.signup_date
17 )
18 SELECT
19     customer_id,
20     full_name,
21     city,
22     signup_date,
23     lifetime_value,
24     avg_distance_km,
25     num_trips,
26     ROW_NUMBER() OVER (ORDER BY last_trip_id DESC) AS recency_rank
27 FROM customer_stats
28 ORDER BY lifetime_value DESC;
```

View 3 Query Results

View 2 Query x +

Run 04:34 PM (<1s) workspace.default New SQL editor: ON Edit Schedule Share Save*

1 | SELECT * FROM pysparkdbt.gold.vw_customer_lifetime_value;

Add parameter

	customer_id	full_name	city	signup_date	lifetime_value	avg_distance_km	num_trips	recency_rank
1	6	Blair	East Pamela	2025-09-17	2029.83	23	34	
2	47	Hayes	Smithfort	2022-08-13	1845.02	22	33	
3	4	Sanchez	Stephanieton	2024-08-31	1820.76	22	34	
4	44	Brown	Julieton	2021-11-21	1721.9	25	26	
5	39	Zimmerman	Shariborough	2021-11-20	1678.13	21	29	
6	22	Mayer	New Michaelbury	2024-04-05	1644.84	22	30	
7	48	Meyer	Donfurt	2024-01-25	1629.99	24	27	
8	8	Arnold	North Sarah	2023-09-21	1629.56	23	30	
9	18	Wade	Roseland	2024-09-28	1611.07	19	37	
10	23	Fowler	Katrinatown	2022-08-20	1585.91	21	31	

- Views query the star schema models, providing analysts clean, aggregated data without complex joins.

Project Summary

This end-to-end data engineering project transforms raw pseudo-Uber CSV data into a fully modeled analytics pipeline using PySpark for ingestion and transformation, and DBT for advanced modeling and orchestration on Databricks. It demonstrates proficiency in handling incremental loads, data quality, SCD Type 2, and dynamic SQL with Jinja, resulting in a scalable, production-ready data warehouse. Key outcomes include 6 bronze tables, cleaned silver layers with upserts, and a gold star schema with historical tracking. This portfolio project highlights my ability to build efficient, modular data pipelines, making it ideal for data engineering roles requiring expertise in big data tools, cloud platforms, and modern data stacks. The entire codebase is available in the repository for review.