| Ex No:  4

Date: 11-09-2025 | **Building and automating a pipeline in Databricks for both E-commerce and Healthcare datasets** |
|---|---|

## Objective

This lab provides hands-on experience in implementing the **data engineering lifecycle** using **Databricks**. Participants will simulate responsibilities of data engineers, data scientists, and business analysts by ingesting raw datasets (E-commerce & Healthcare), cleaning and transforming them, aggregating for analytics, and finally creating dashboards with automation and alerts.

## Outcomes

- Identify and describe each stage of the **data engineering lifecycle** (Ingestion, Transformation, Aggregation, Visualization, Automation).

- Implement the **Architecture** (Bronze → Silver → Gold) in Databricks using PySpark and SQL.

- Collaborate to define a **business problem** using raw datasets:
  *E-commerce*: Revenue by product category, daily revenue trends.

  *Healthcare*: Service category performance, daily hospital revenue trends.

## Materials

- **Raw datasets**:
  ecommerce_orders.csv → E-commerce sales transactions

  healthcare_orders.csv → Healthcare services transactions

- **Tools & Environment**:
  Databricks Workspace
  PySpark & SQL

- **Artifacts**:
  Bronze Layer → Raw ingestion tables
  Silver Layer → Cleaned, transformed tables (silver_ecommerce_orders, silver_healthcare_orders)
  Gold Layer → Aggregated analytics tables (gold_ecommerce_category_sales, gold_healthcare_service_sales, etc.)

- **Dashboard & Visualization**:
  Bar charts (Revenue by category/service)
  Line charts (Daily revenue trends)
  Filters (City, Payment Method)

- **Automation Setup**:
  Databricks Jobs for scheduled refresh
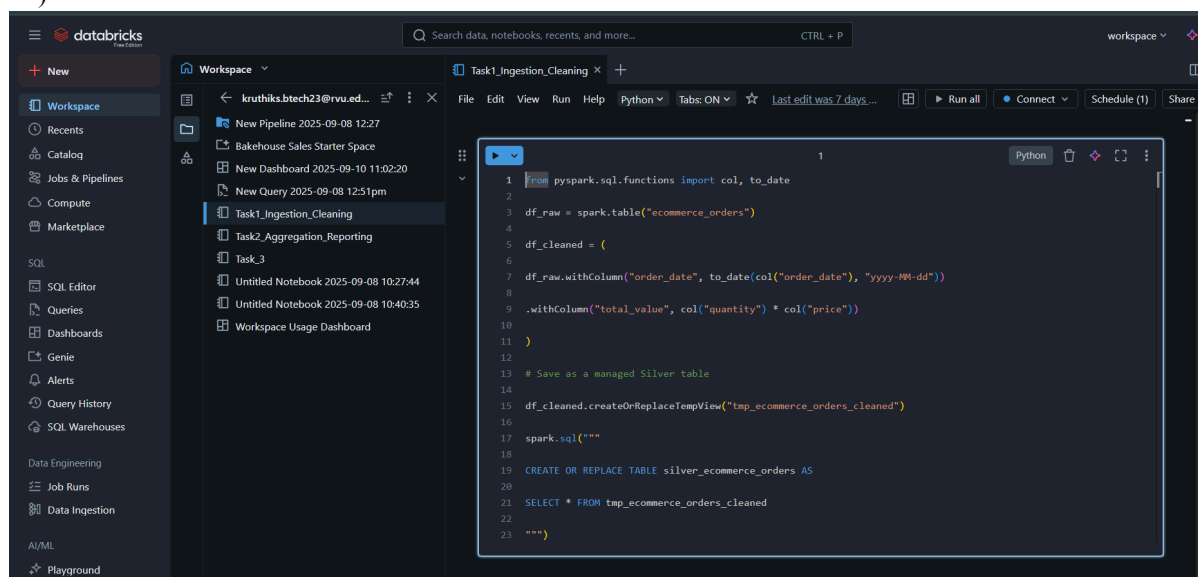  Email notifications for success/failure alerts

## Lab Procedure:

# Step 1: Ingestion & Cleaning (Bronze → Silver)

```python
from pyspark.sql.functions import col, to_date

# Bronze layer: raw table (already uploaded as healthcare_orders.csv)
df_raw = spark.read.csv("/FileStore/tables/healthcare_orders.csv", header=True, inferSchema=True)

# Clean & transform (Silver layer)
df_cleaned = (
    df_raw.withColumn("order_date", to_date(col("order_date"), "yyyy-MM-dd"))
        .withColumn("total_value", col("quantity") * col("price"))
)

# Save as Silver table
df_cleaned.createOrReplaceTempView("tmp_healthcare_orders_cleaned")
spark.sql("""
CREATE OR REPLACE TABLE silver_healthcare_orders AS
SELECT * FROM tmp_healthcare_orders_cleaned
""")
```



# Step 2: Aggregation & Enrichment (Silver → Gold)

USN NUMBER: 1RVU23CSE227
NAME: KRUTHIK S

```python
from pyspark.sql.functions import sum as spark_sum

# Load Silver table
df_cleaned = spark.table("silver_healthcare_orders")

# Gold 1: Revenue by service category
df_service_sales = (
    df_cleaned.groupBy("service_category")
            .agg(spark_sum("total_value").alias("total_revenue"))
)
df_service_sales.createOrReplaceTempView("tmp_healthcare_service_sales")
spark.sql("""
CREATE OR REPLACE TABLE gold_healthcare_service_sales AS
SELECT * FROM tmp_healthcare_service_sales
""")

# Gold 2: Daily revenue trends
df_daily_sales = (
    df_cleaned.groupBy("order_date")
            .agg(spark_sum("total_value").alias("daily_revenue"))
)
df_daily_sales.createOrReplaceTempView("tmp_healthcare_daily_sales")
spark.sql("""
CREATE OR REPLACE TABLE gold_healthcare_daily_sales AS
SELECT * FROM tmp_healthcare_daily_sales
""")
```
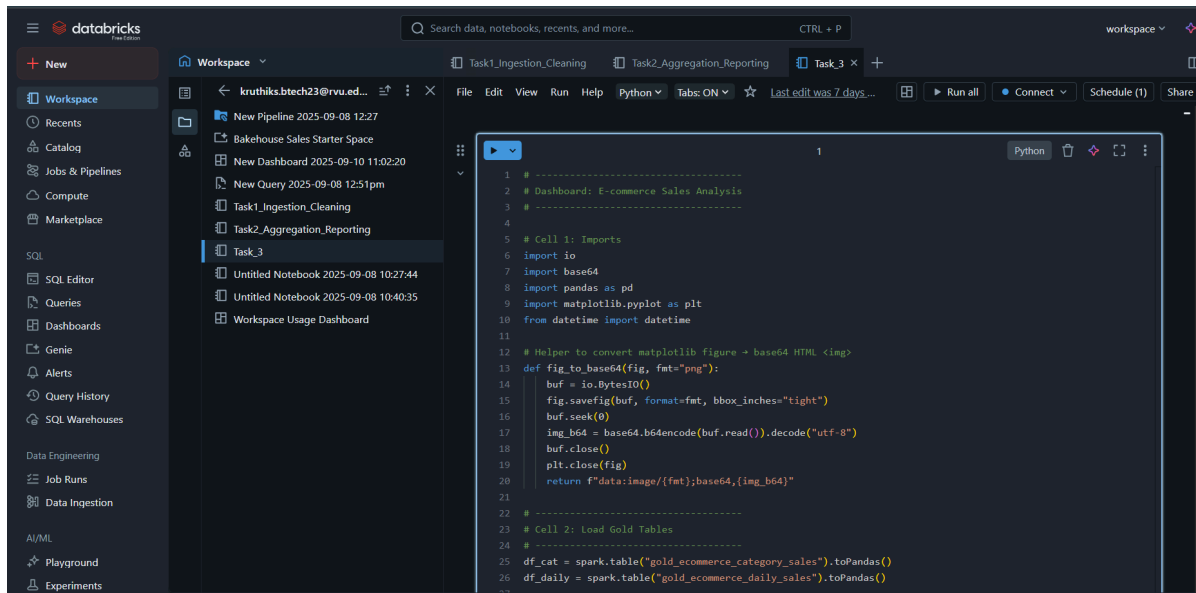


# Step 3: Dashboard & Visualization

In Databricks Notebook:

- **Bar Chart**: gold_healthcare_service_sales → service_category vs total_revenue

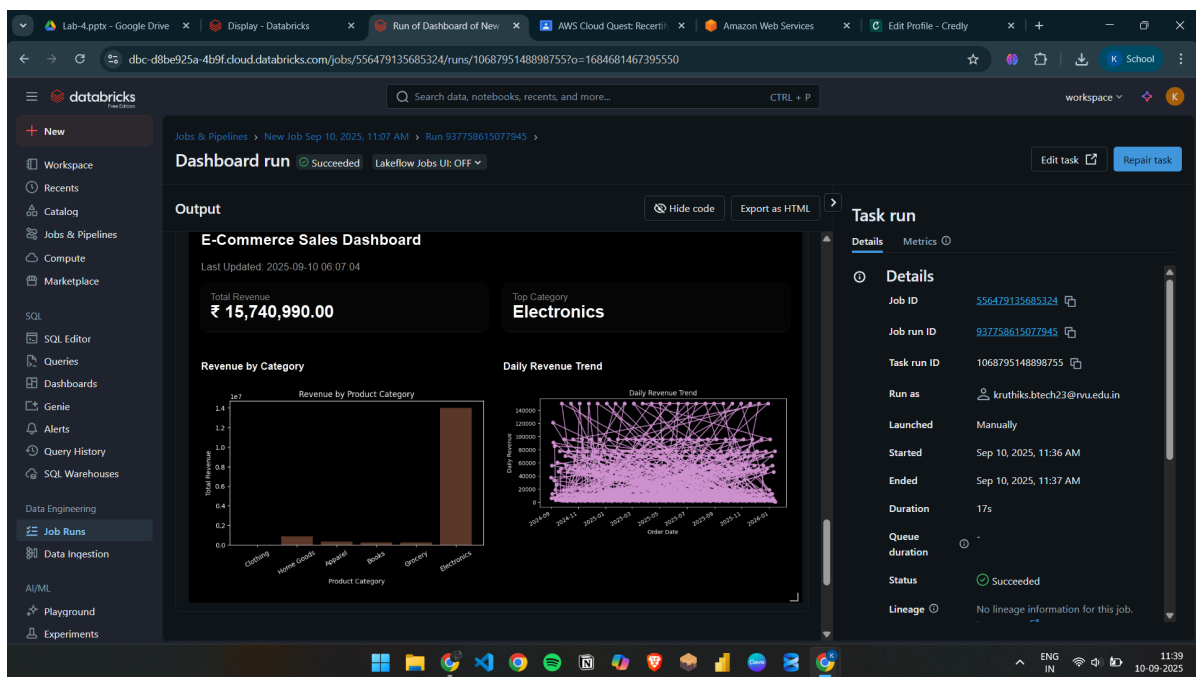- **Line Chart**: gold_healthcare_daily_sales → order_date vs daily_revenue
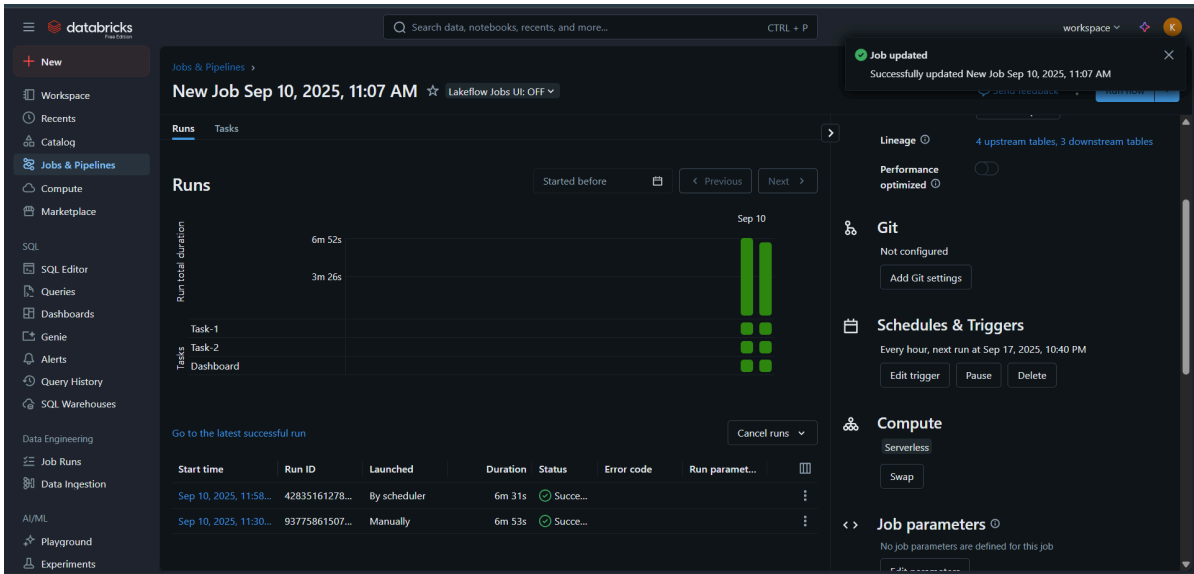- Add **filters**: city, payment_method





# Step 4: Automation & Scheduling

- Go to **Databricks Jobs** → create a Job for this Notebook.
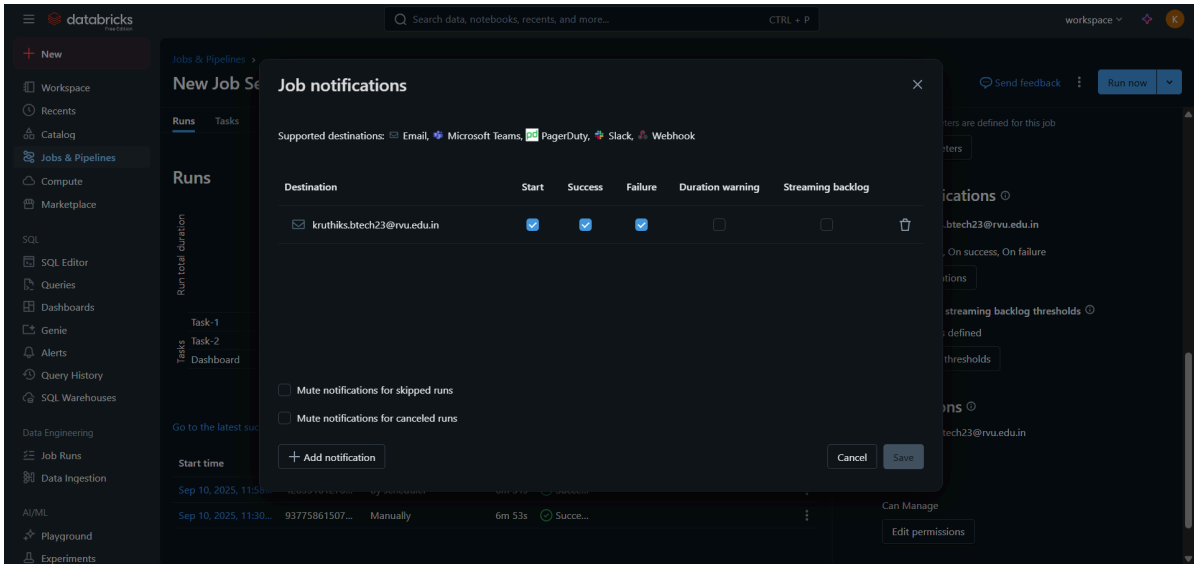- Set **schedule** = daily/weekly refresh.

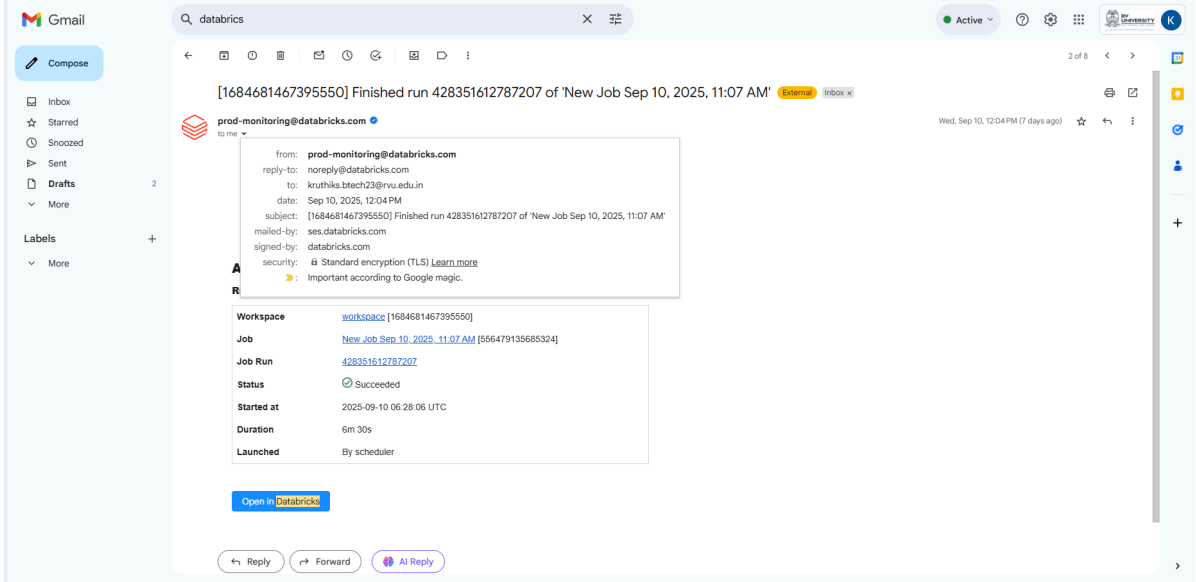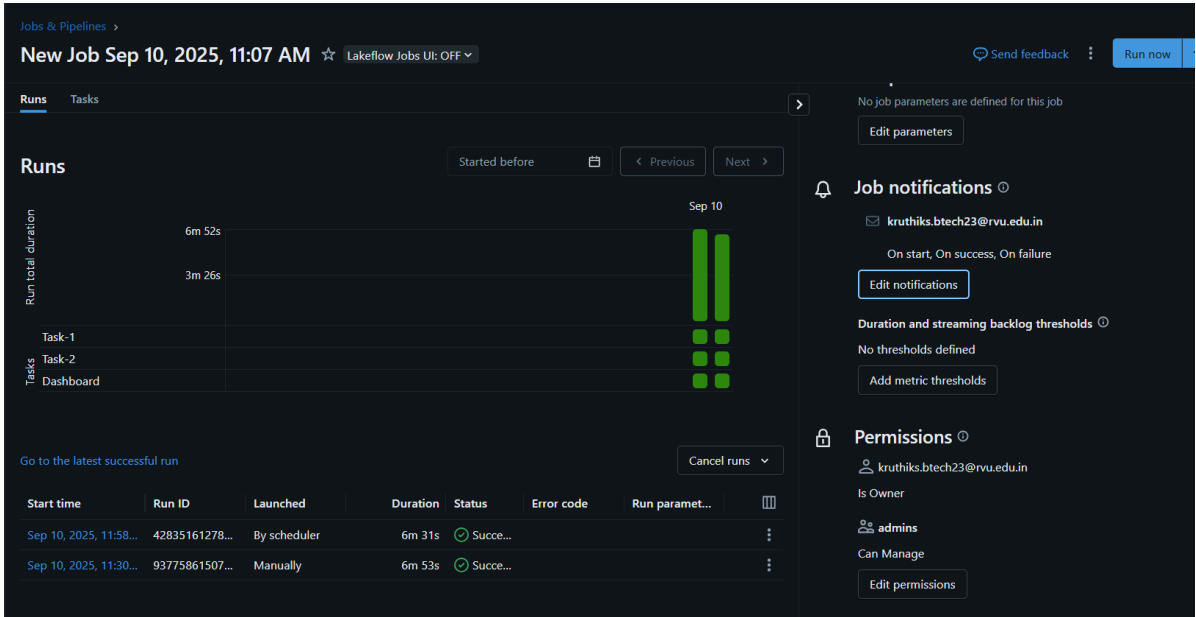- Configure **task order**: Bronze → Silver → Gold → Dashboard.



# Step 5: Notifications & Alerts

- In Job settings → **Notifications**.
- Add emails for *success* and *failure events*.
- Ensures admins are alerted if ingestion/transform fails.

USN NUMBER: 1RVU23CSE227
NAME: KRUTHIK S





**GitHub Link:** https://github.com/kruth-s/Data-Engg-Lab