# LAB 4
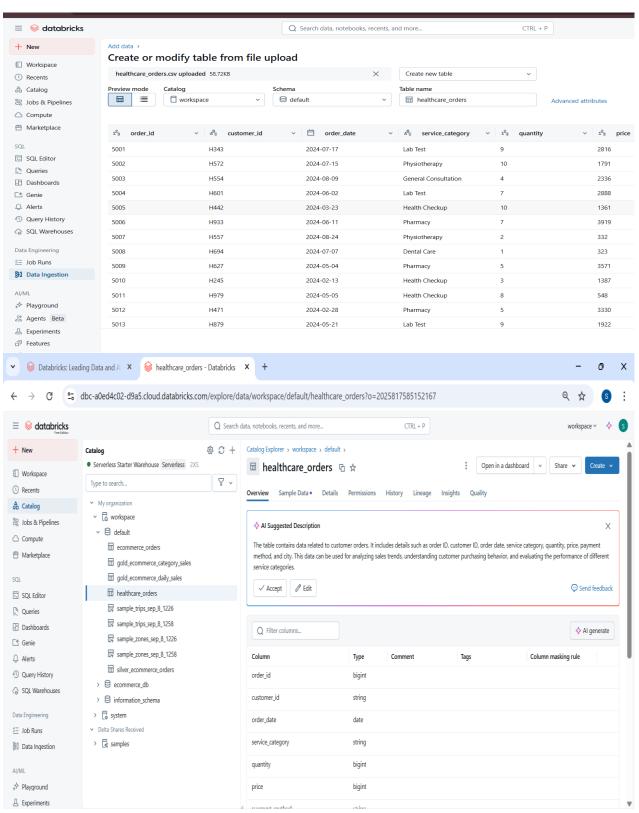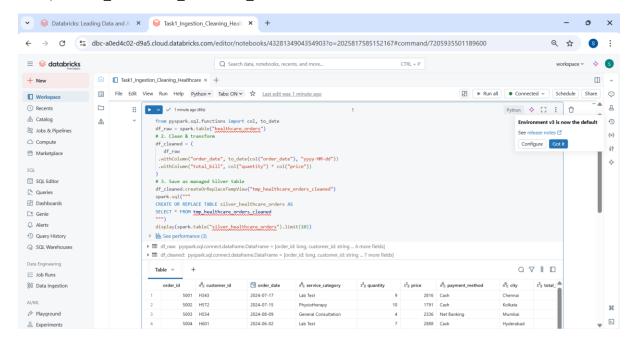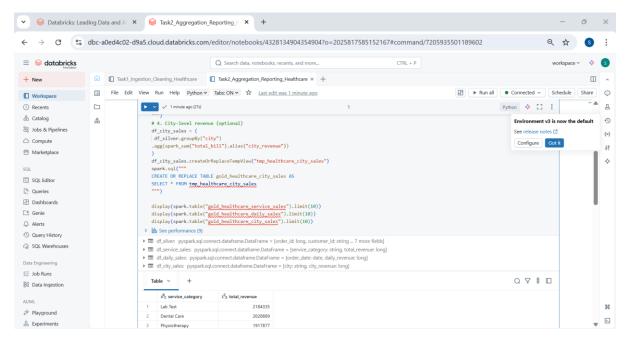
Name : Sahanaa M S

USN : 1RVU23CSE394

1) UPLOAD healthcare_orders.csv TO WORKSPACE:
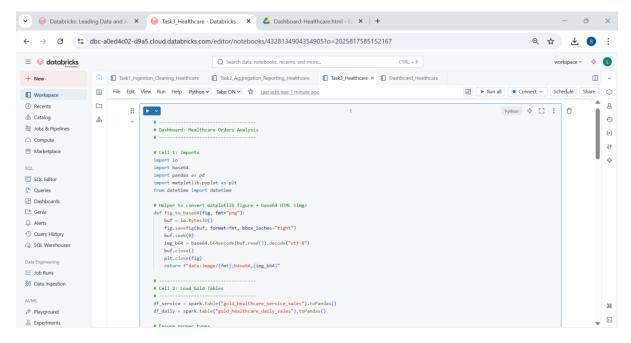
2) TASK1_INGESTION_CLEANING_HEALTHCARE:
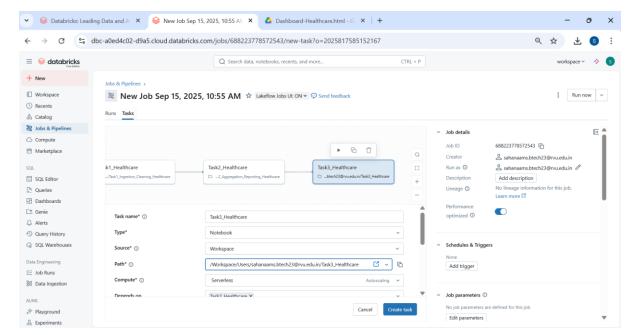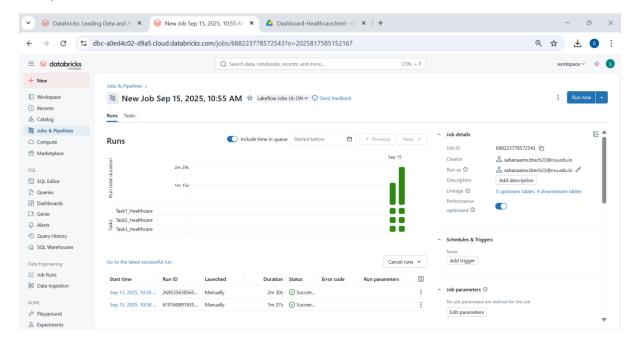


3) TASK2_AGGREGATION_REPORTING_HEALTHCARE:

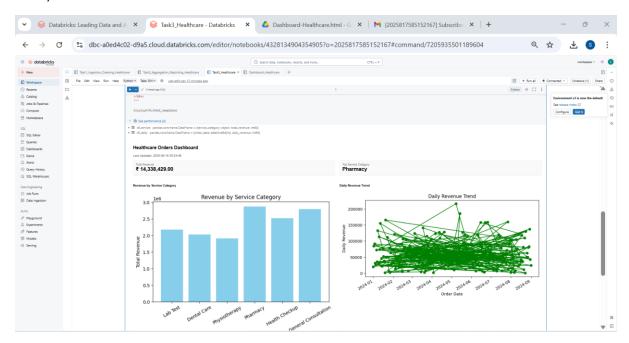## 4) TASK3_HEALTHCARE AND DASHBOARD_HEALTHCARE :



## 5) JOB AND PIPELINES – CREATE JOB:

## 6) JOB AND PIPELINES – SUCCESSFUL JOB CREATED:



## 7) DASHBOARD REPORT:

8) EMAIL FROM DATABRICKS FOR SUCCESSFUL JOB CREATION: