

## Quick notes

---

Sunday, April 20, 2025 10:36 PM

### Create data sources :

1. hospital-a-mysql-source
  - patients (incr)
  - providers (full)
  - encounters (incr)
  - transactions (incr)
  - departments (full)
2. hospital-b-mysql-source
  - patients (incr)
  - providers (full)
  - encounters (incr)
  - transactions (incr)
  - departments (full)
3. claims
  - monthly basis for each hospital
  - in one month
    - 1 claim.csv for hospital-a
    - 1 claim.csv for hospital-b
4. cpt\_codes
  - constant or standard codes for medical procedures
  - csv file

### Data Ingestion :

- create dataproc cluster
- ingestion-1 => configs => hospital-a-mysql (tables) ==> gcs\_bucket/landing/hospital-a/
- ingestion-2 => configs => hospital-b-mysql (tables) ==> gcs\_bucket/landing/hospital-b/
- claims on monthly basis for two hospitals => csv files ==> gcs\_bucket/landing/claims/\*
- cpt\_codes are constant => csv files ==> gcs\_bucket/landing/cptcodes/\*

### Bronze:

- gcs\_bucket/landing/hospital-a/\* ==> external tables ==> bigquery/bronze\_dataset/\*
- gcs\_bucket/landing/hospital-b/\* ==> external tables ==> bigquery/bronze\_dataset/\*
- gcs\_bucket/landing/claims/\*.csv ==> dataproc ==> bigquery/bronze\_dataset/claims
- gcs\_bucket/landing/cptcodes/\*.csv ==> dataproc ==> bigquery/bronze\_dataset/cpt\_codes

**Silver:**

- bronze\_dataset/{full/incr} ==> cdm, scd2, nulls, distinct, merge ==> silver\_dataset/

**Gold:**

- as per business usecases (silver tables ==> query ==> gold tables) ==> reports & dashboards
- ex : patients\_summary (4 tables ==> joins, aggr, operations ==> patients\_summary\_gold)

**Orchestration:**

- create composer environment (workflow Orchestration tool)
- pyspark\_dag ==> total 4 pyspark jobs
- bq\_dag ==> total 3 bq jobs

**CICD:**

- sett up github
- set up composer env
- cloud build trigger
- test

**gcs\_bucket**

- landing
- hospital-a

- archive
  - 2025
    - 03
      - 24
        - patients\_24032025.json
        - transactions\_24032025.json
      - 25
        - patients\_25032025.json
        - transactions\_25032025.json
  - patients
    - patients\_26032025.json
  - transactions
    - transactions\_26032025.json

hospital-a-mysql > patients ----> full data ----> gcs\_bucket/landing/hospital-a/patients

audit table -- in bigquery --