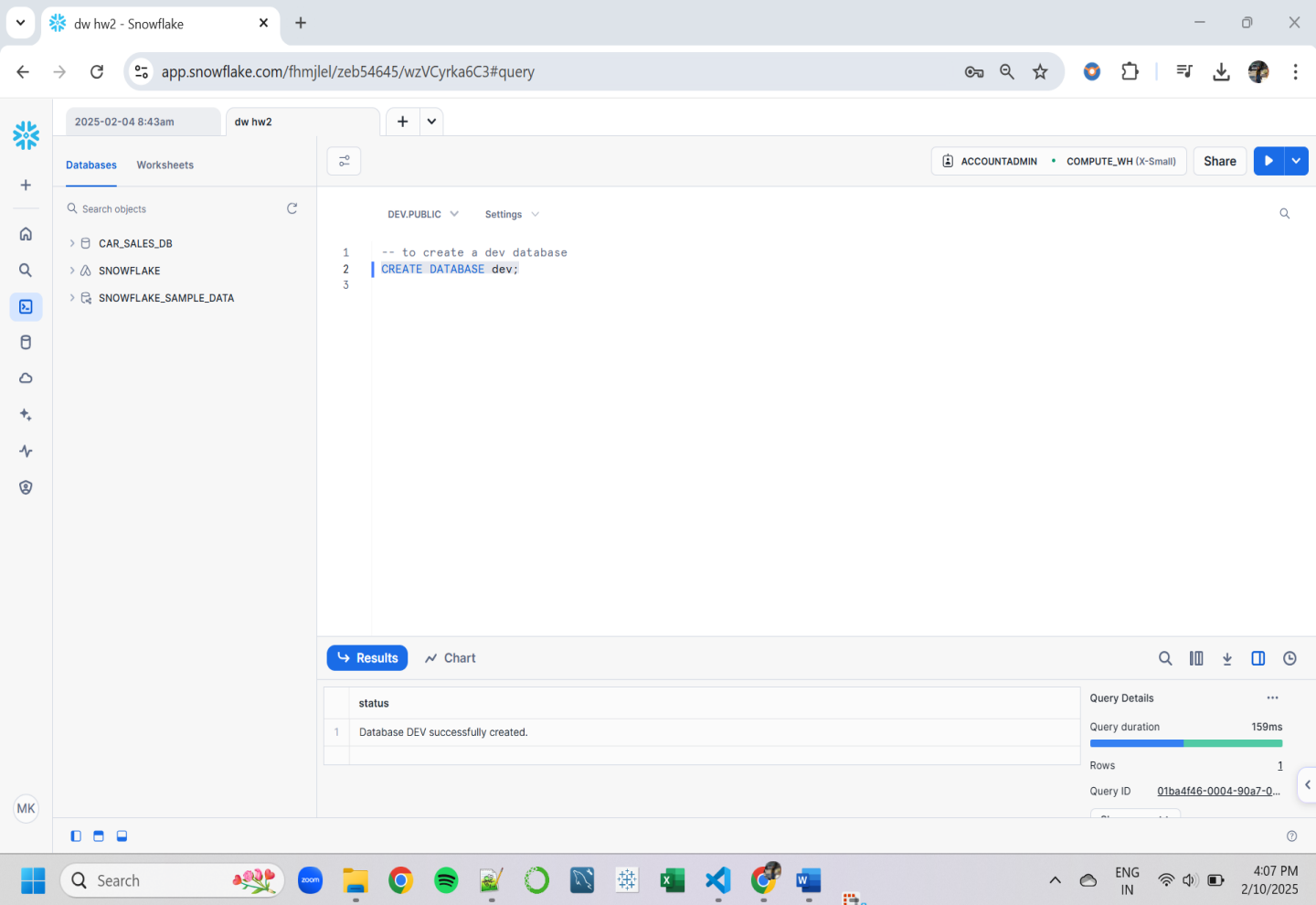Name: Mrunali Katta

ID: 017516785

# DATA 226 – Homework 2

1.(+1) Create database dev and schemas RAW, CURATION and ANALYTICS

<u>Created a dev database</u>
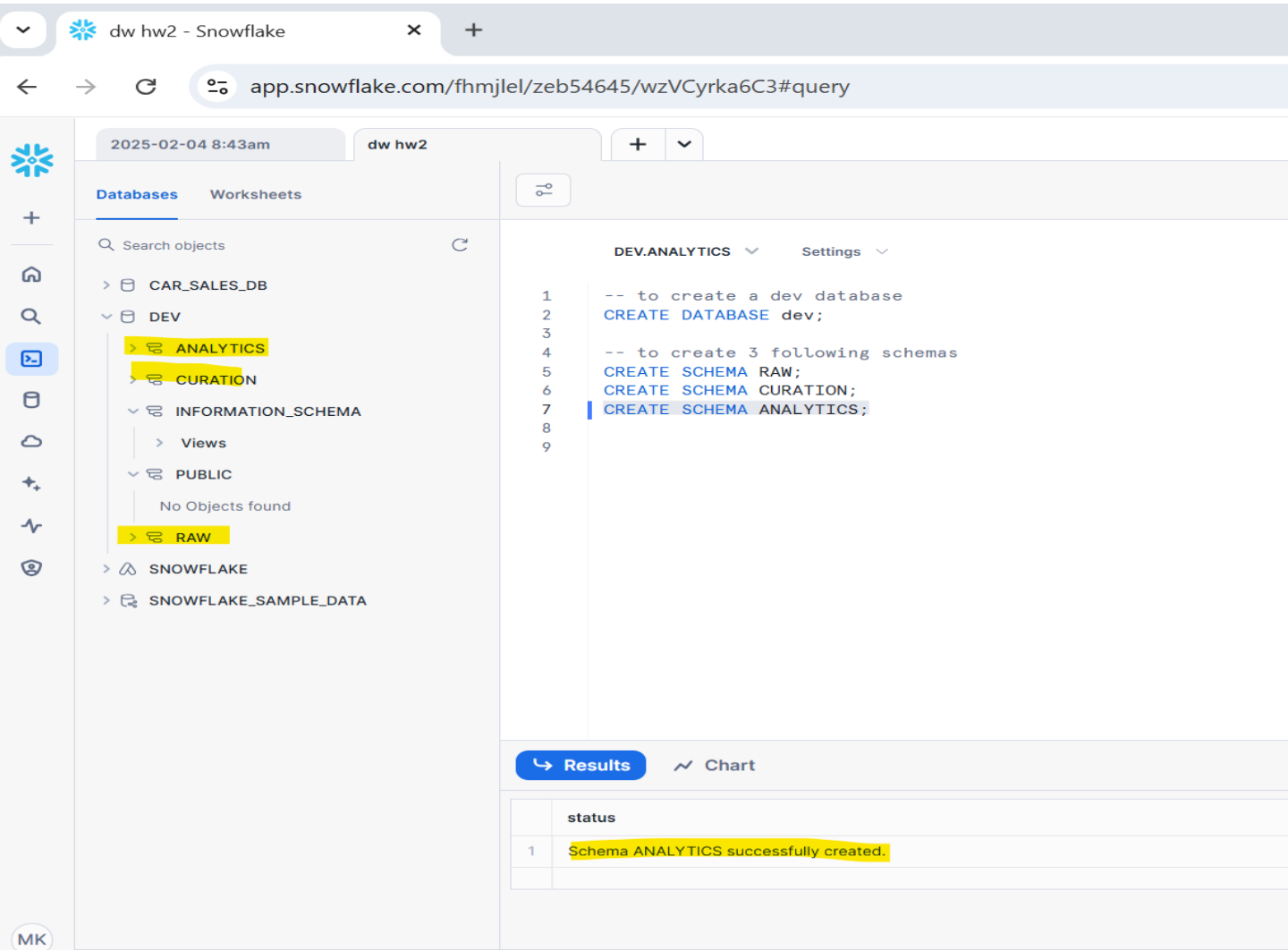
CREATE DATABASE dev;



<u>Also created the three schema's mentioned **RAW, CURATION and ANALYTICS**</u>

CREATE SCHEMA RAW;
CREATE SCHEMA CURATION;
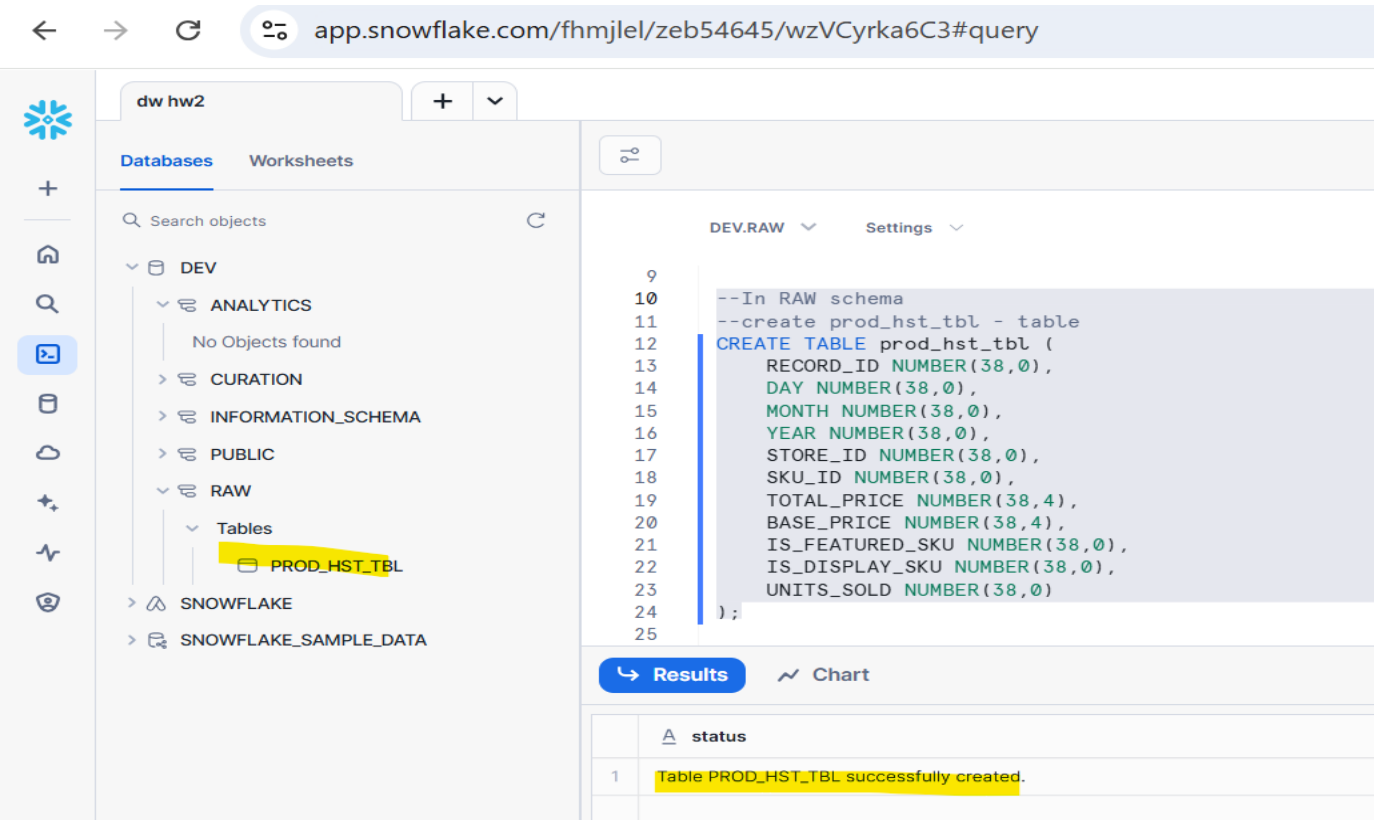CREATE SCHEMA ANALYTICS;

2.(+4) In RAW schema, create

- prod_hst_tbl – table

created table **"prod_hst_tbl"** in RAW schema
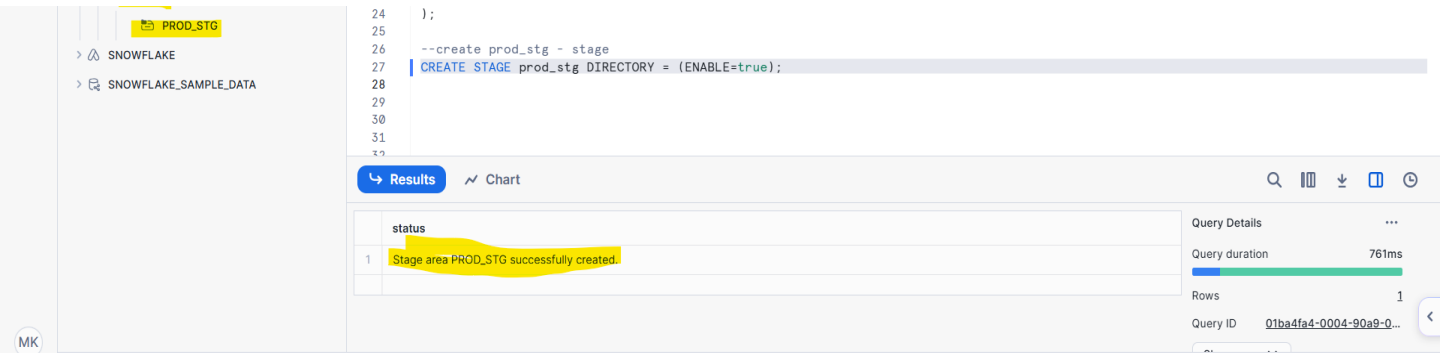
```
CREATE TABLE prod_hst_tbl (
    RECORD_ID NUMBER(38,0),
    DAY NUMBER(38,0),
    MONTH NUMBER(38,0),
    YEAR NUMBER(38,0),
    STORE_ID NUMBER(38,0),
    SKU_ID NUMBER(38,0),
    TOTAL_PRICE NUMBER(38,4),
    BASE_PRICE NUMBER(38,4),
    IS_FEATURED_SKU NUMBER(38,0),
    IS_DISPLAY_SKU NUMBER(38,0),
    UNITS_SOLD NUMBER(38,0)
);
```



- prod_stg – stage
created stage **"prod_stg"** in RAW schema

```
CREATE STAGE prod_stg DIRECTORY = (ENABLE=true);
```



- prod_raw_task - task (alter the task)

created task **"prod_raw_task"** in RAW schema

```
CREATE TASK prod_raw_task
    WAREHOUSE = COMPUTE_WH
    SCHEDULE = 'USING CRON * * * * * UTC'  --scheduling using CRON
AS
    COPY INTO prod_hst_tbl
    FROM @prod_stg
    FILE_FORMAT = (TYPE = 'CSV' SKIP_HEADER = 1 FIELD_OPTIONALLY_ENCLOSED_BY = '"');
```

```
29
30
31
32    -- create prod_raw_task - task (alter the task)
33    CREATE TASK prod_raw_task
34        WAREHOUSE = COMPUTE_WH
35        SCHEDULE = 'USING CRON * * * * * UTC'   --scheduling using CRON
36    AS
37        COPY INTO prod_hst_tbl
38        FROM @prod_stg
39        FILE_FORMAT = (TYPE = 'CSV' SKIP_HEADER = 1 FIELD_OPTIONALLY_ENCLOSED_BY = '"');
40
```

| | status |
| --- | --- |
| 1 | Task PROD_RAW_TASK successfully created. |



```
39
40    show tasks;
41
42    -- create prod_stream - stream
43
```

| | created_on | name | id | database_name | schema_name | owner |
| --- | --- | --- | --- | --- | --- | --- |
| 1 | 2025-02-10 17:54:08.638 -0800 | PROD_RAW_TASK | 01ba4fb2-37fa-13bf-0000-000000000034 | DEV | RAW | ACCOUN |

- prod_stream – stream

created stream **"prod_stream"** in RAW schema

CREATE STREAM prod_stream ON TABLE prod_hst_tbl;



```
35        SCHEDULE = 'USING CRON * * * * * UTC'   --scheduling using CRON
36    AS
37        COPY INTO prod_hst_tbl
38        FROM @prod_stg
39        FILE_FORMAT = (TYPE = 'CSV' SKIP_HEADER = 1 FIELD_OPTIONALLY_ENCLOSED_BY = '"');
40
41
42
43    -- create prod_stream - stream
44    CREATE STREAM prod_stream
45    ON TABLE prod_hst_tbl;
46
47
```

| | comme | warehouse | schedule | [] predecesso | state | definition |
| --- | --- | --- | --- | --- | --- | --- |
| 1 | MIN | COMPUTE_WH | USING CRON * * * * * UTC | [] | started | COPY INTO prod_hst_tbl  FROM @prod_stg  FILE_FO |

3. In CURATION schema, create

- prod_hst_tbl – table

created stream **"prod_hst_tbl"** in CURATION schema

```
CREATE CURATION.PROD_HST_TBL(
    RECORD_ID NUMBER(38,0),
    DAY NUMBER(38,0),
    MONTH NUMBER(38,0),
    YEAR NUMBER(38,0),
    STORE_ID NUMBER(38,0),
    SKU_ID NUMBER(38,0),
    TOTAL_PRICE NUMBER(38,4),
    BASE_PRICE NUMBER(38,4),
    UNITS_SOLD NUMBER(38,0)
);
```

- prod_curation_task using MERGE (alter the task)

```
-- to create task using merge
CREATE TASK RAW.prod_curation_task
    WAREHOUSE = COMPUTE_WH
    WHEN SYSTEM$STREAM_HAS_DATA('prod_stream')
AS MERGE INTO CURATION.prod_hst_tbl AS target
    USING ( SELECT RECORD_ID,DAY, MONTH, YEAR,STORE_ID,SKU_ID,TOTAL_PRICE,BASE_PRICE,UNITS_SOLD
        FROM RAW.prod_stream
        WHERE RECORD_ID IS NOT NULL
            OR DAY IS NOT NULL
            OR MONTH IS NOT NULL
            OR STORE_ID IS NOT NULL
            OR SKU_ID IS NOT NULL
            OR TOTAL_PRICE IS NOT NULL
            OR BASE_PRICE IS NOT NULL
            OR UNITS_SOLD IS NOT NULL
    ) AS source
    ON target.RECORD_ID = source.RECORD_ID
    WHEN MATCHED THEN
        UPDATE SET
            target.DAY = source.DAY,
            target.MONTH = source.MONTH,
            target.YEAR = source.YEAR,
            target.STORE_ID = source.STORE_ID,
            target.SKU_ID = source.SKU_ID,
            target.TOTAL_PRICE = source.TOTAL_PRICE,
            target.BASE_PRICE = source.BASE_PRICE,
            target.UNITS_SOLD = source.UNITS_SOLD
    WHEN NOT MATCHED THEN
        INSERT (RECORD_ID, DAY, MONTH, YEAR, STORE_ID, SKU_ID, TOTAL_PRICE, BASE_PRICE, UNITS_SOLD)
        VALUES (source.RECORD_ID, source.DAY, source.MONTH, source.YEAR, source.STORE_ID, source.SKU_ID,
source.TOTAL_PRICE, source.BASE_PRICE, source.UNITS_SOLD);
```

**4.(+2) In ANALYTICS schema, create**

- book_dy_tbl – Dynamic Table

Created a table called **'book_dy_tbl'**. Later after loading file 'prod_old_data.csv' file the row count is 7.1k
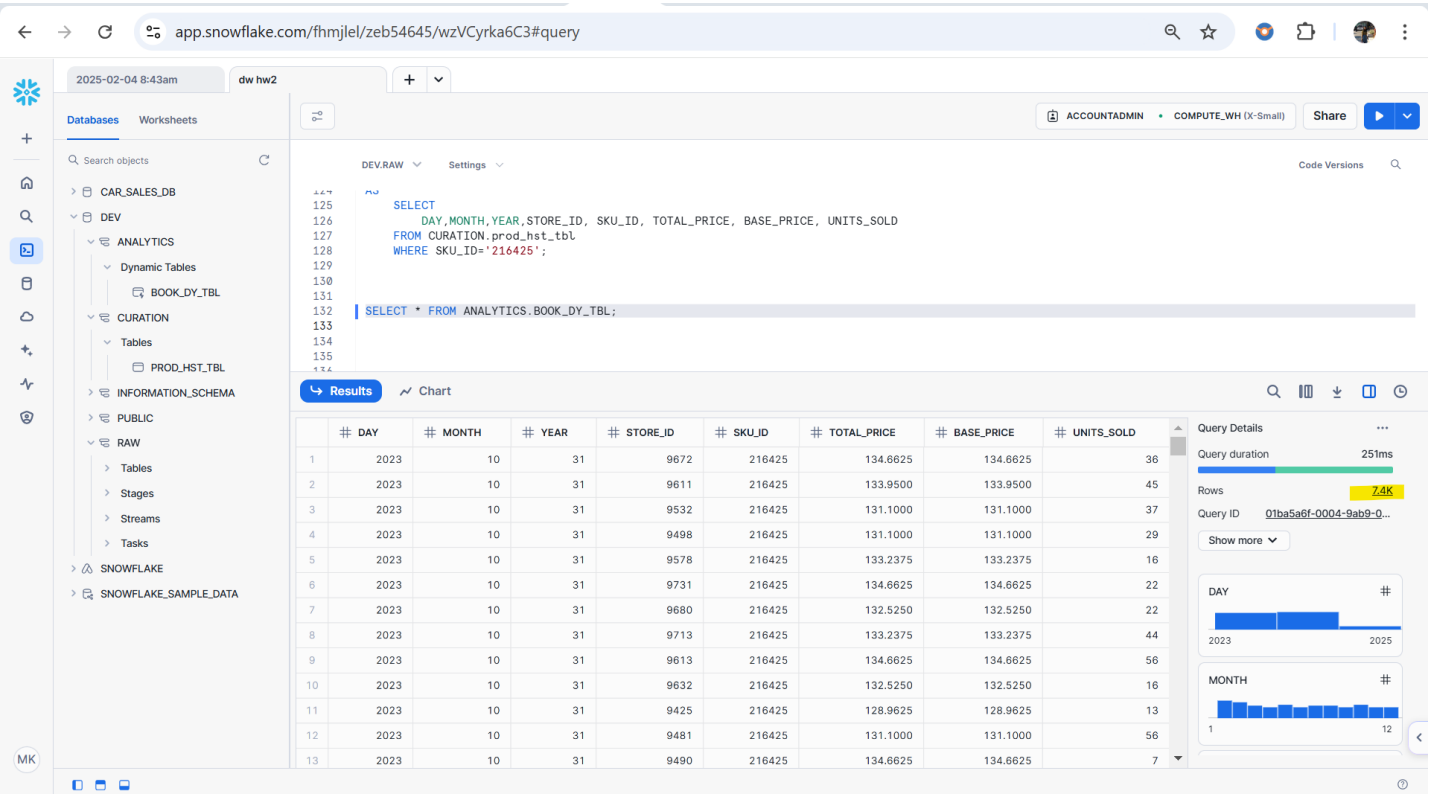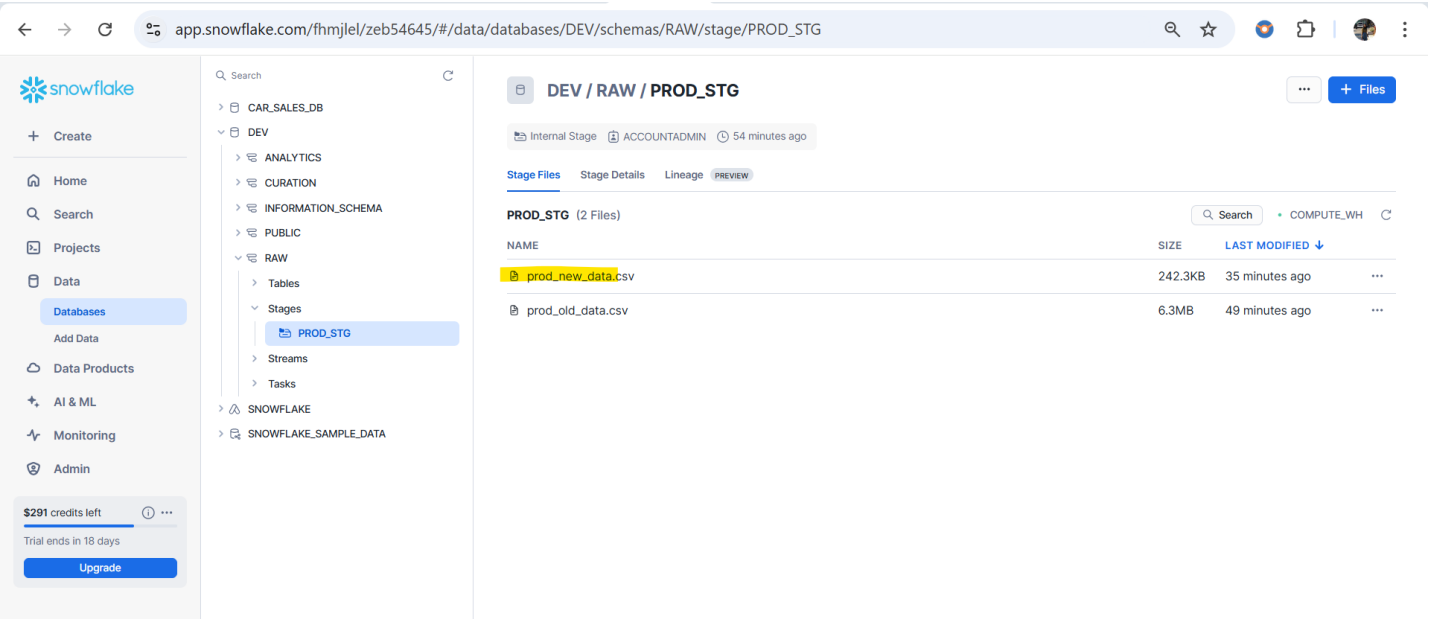
After uploading prod_new_data.csv file in 'prod_stage' stage the rows count is 7.4k.





5.(+4) Explain end-to-end process based on your understanding.

→ Following is the end-to-end process explanation of the **Snowflake Basic Data Pipeline:**

1. **Database Setup**:
- First created a database called **"dev".**
2. **3 Schema creation**:
- Then created three schemas as follows: **"RAW", "CURATION", and "ANALYTICS".**
3. **RAW schema:**
- Then created table called **"prod_hst_tbl"** to store initial raw data.
- Also created a stage called **"prod_stg"** to stage the two .csv files which are **"prod_old_data"** and **"prod_new_data".**
- Later created task **"prod_raw_task"** to transfer data from the stage to the table every minute using CRON.
- Also created a stream called as **"prod_stream"** to monitor any changes in **"prod_hst_tbl"**.
4. **CURATION schema**:
- Here in this schema replicated **"prod_hst_tbl"** for data processing.
- And then created a task called **"prod_curation_task"** for merging updates from **RAW** schema to **CURATION** schema, to ensure data is consistent.
5. **ANALYTICS Schema**:
- In this schema, created a table called **"book_dy_tbl"** which is a a dynamic table for aggregating the earlier curated data.
- Then using a query extracted insights such as total sales, average prices, and transaction counts by SKU and store ID.
- Overall here the data moves from initial ingestion in RAW, to refinement in CURATION, and finally to aggregation for analytics in ANALYTICS.