# BUAN 6346 Big Data Analytics PHASE 2

SHIVA KUMAR REDDY KOPPULA
03/27/2024
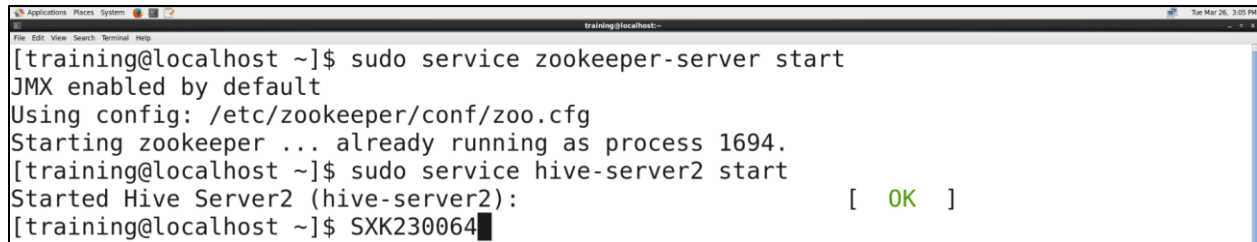
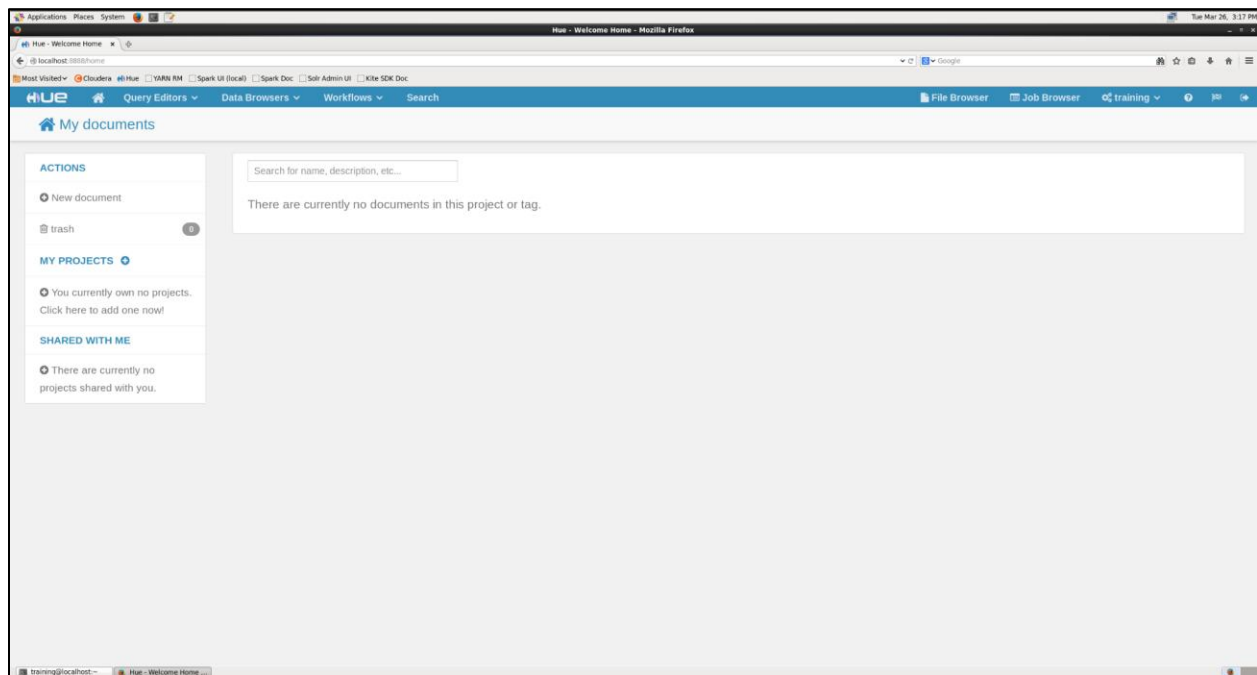# Table of Contents

# CHAPTER 6 - CREATE AND POPULATE TABLES IN IMPALA OR HIVE

## Create and Query a Table in Impala or Hive

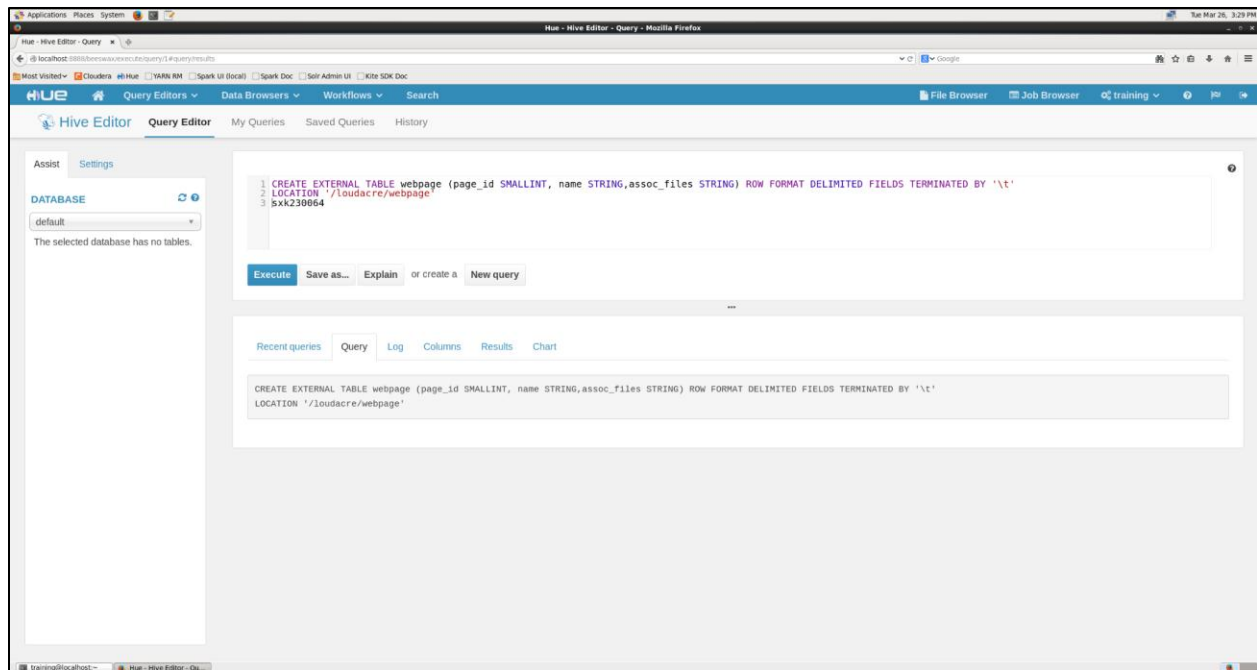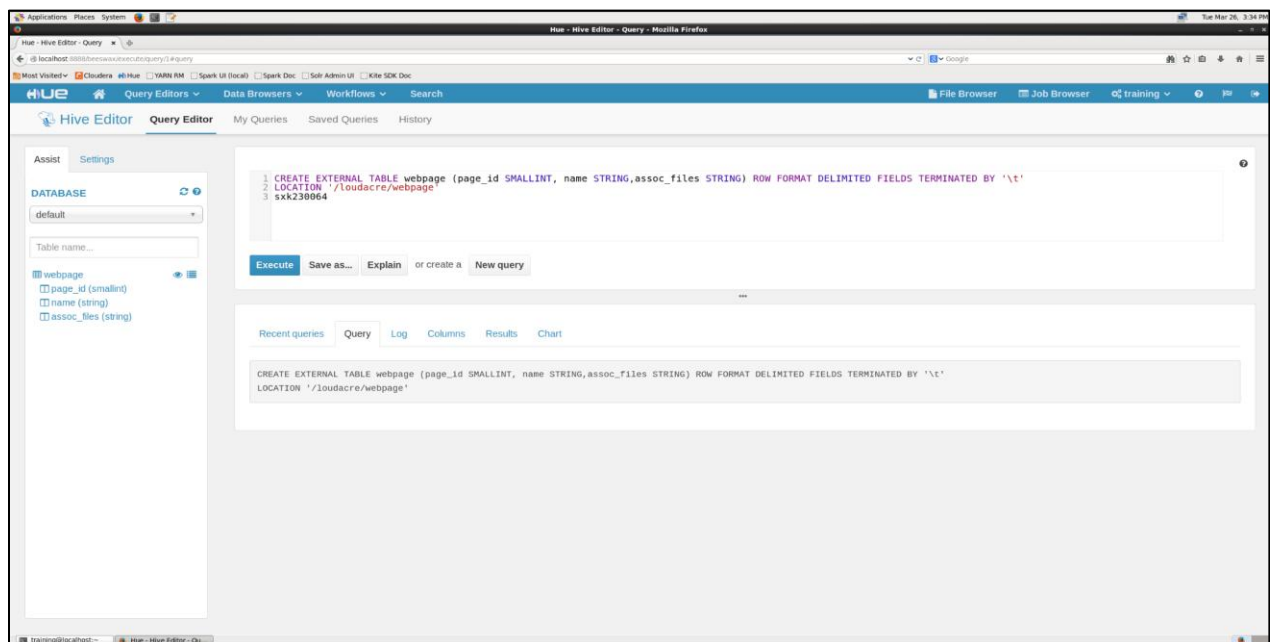1. Started the Hive server by executing the necessary commands in the terminal.



2. Accessed the HUE interface in Chrome, clicking on the HUE icon to navigate to the homepage.



3. Opened the query editor menu, selecting "Hive" from the dropdown, which directed me to the Hive query editor. Then, I entered an SQL command in the query editor pane to create a table for the previously imported webpage data.
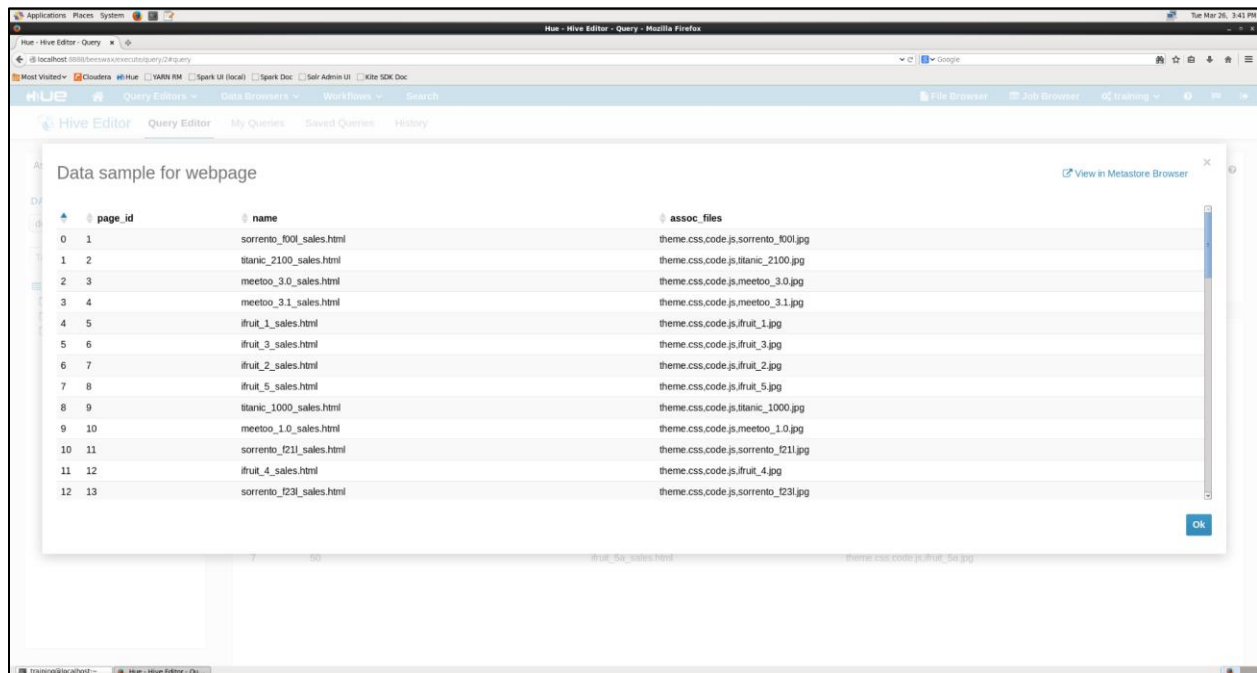
4. Executed the query by clicking the "Execute" button. To view the newly created table, I clicked on the 'refresh' button next to 'Database' on the left-hand side of the page, locating the webpage table with columns (page_id), name, and (assoc_files).I clicked on the webpage table to review the column definitions which reside below it.



5. Initiated a test query by clicking the 'New Query' button, observing the findings in the "Results" tab.

6. Previewed sample data by clicking on the preview sample data icon adjacent to the table name.



## Use Sqoop to Import Directly into Hive and Impala

7. Imported the device table directly into the Hive metastore using the terminal, employing the below command.

8. Utilized HUE to navigate to the specific data location (/user/hive/warehouse/device), to review the imported data files residing in the default hive warehouse.

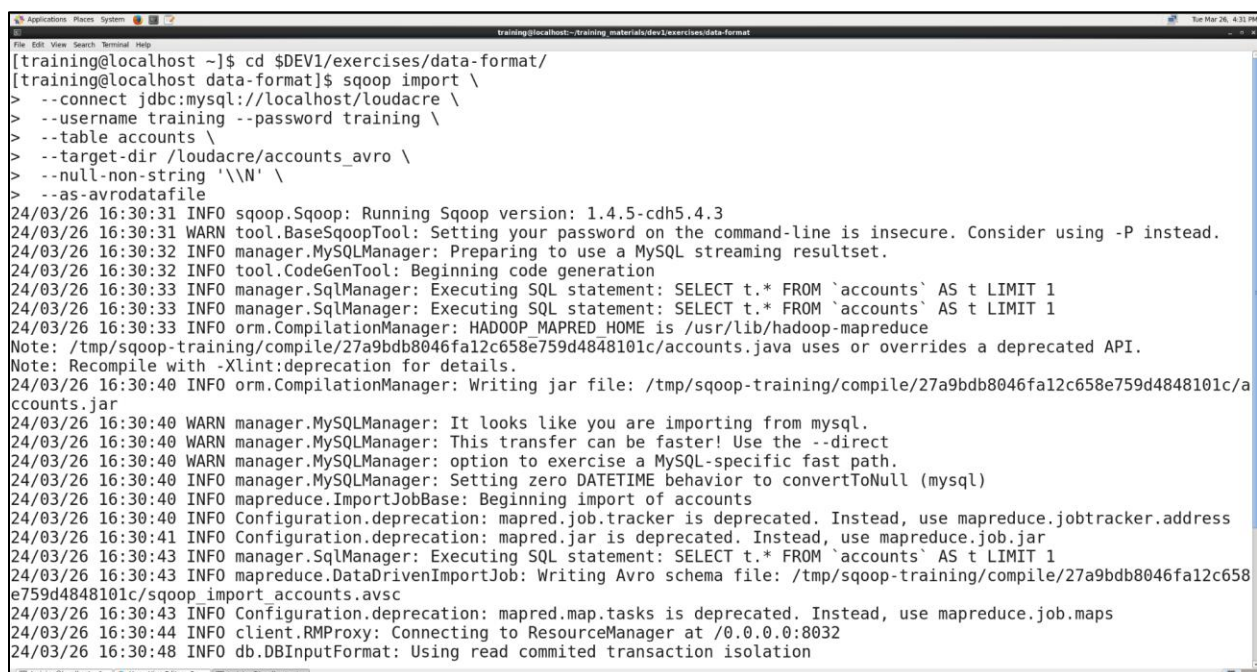9. Ran a test query to view all columns of the device table.

# CHAPTER 7 - SELECT A FORMAT FOR A DATA FILE

1.  I used the 'cd' command to navigate to the exercise directory.



2.  Utilizing the Sqoop import command, I imported the accounts table into an Avro data format.

```
               HDFS: Number of bytes read=470
               HDFS: Number of bytes written=12713125
               HDFS: Number of read operations=16
               HDFS: Number of large read operations=0
               HDFS: Number of write operations=8
       Job Counters
               Launched map tasks=4
               Other local map tasks=4
               Total time spent by all maps in occupied slots (ms)=0
               Total time spent by all reduces in occupied slots (ms)=0
               Total time spent by all map tasks (ms)=72307
               Total vcore-seconds taken by all map tasks=72307
               Total megabyte-seconds taken by all map tasks=18510592
       Map-Reduce Framework
               Map input records=129764
               Map output records=129764
               Input split bytes=470
               Spilled Records=0
               Failed Shuffles=0
               Merged Map outputs=0
               GC time elapsed (ms)=1770
               CPU time spent (ms)=32010
               Physical memory (bytes) snapshot=542732288
               Virtual memory (bytes) snapshot=3383013376
               Total committed heap usage (bytes)=191889408
       File Input Format Counters
               Bytes Read=0
       File Output Format Counters
               Bytes Written=12713125
24/03/26 16:32:28 INFO mapreduce.ImportJobBase: Transferred 12.1242 MB in 104.4246 seconds (118.8911 KB/sec)
24/03/26 16:32:28 INFO mapreduce.ImportJobBase: Retrieved 129764 records.
[training@localhost data-format]$ sxk230064
```
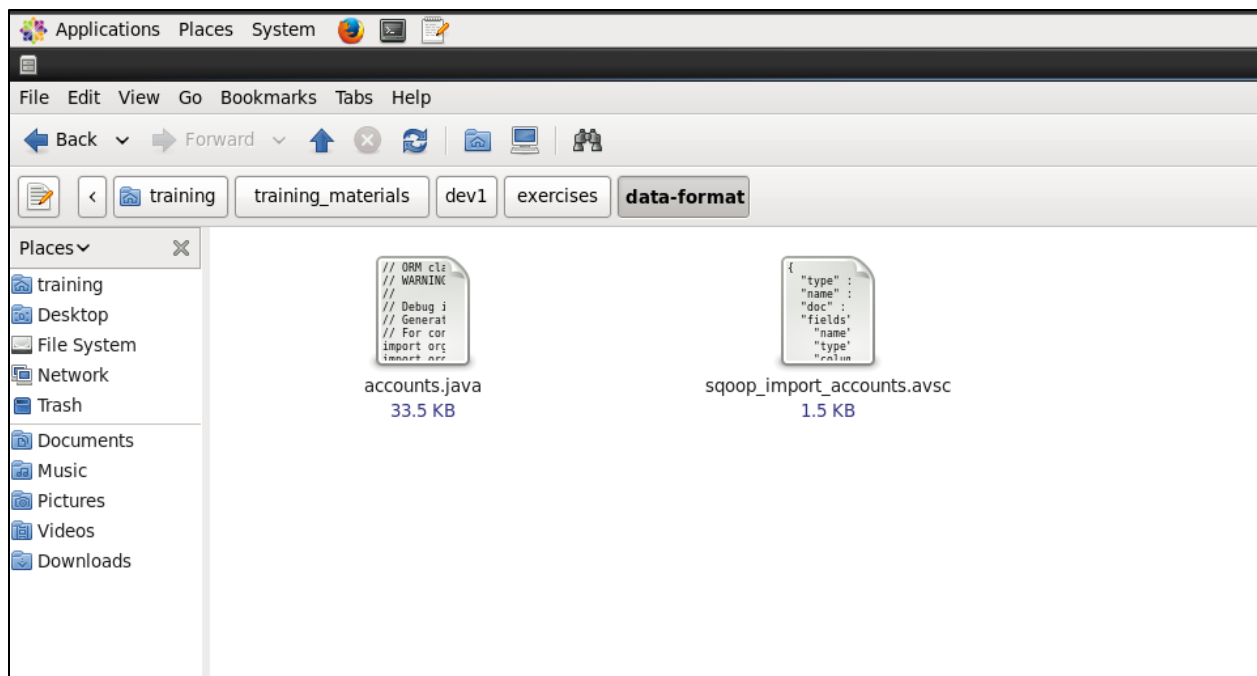
3.  Then, I accessed the file browser to view the imported files by Sqoop into HDFS.



4.  Upon locating the file part-m-00000.avro, I proceeded to view its contents, experiencing a lengthy loading time which may be due to large number of records.

5. Sqoop automatically generated a schema named "sqoop_import_accounts.avsc" in the current directory.



6. To examine the contents of the generated schema file, I accessed and viewed it accordingly.
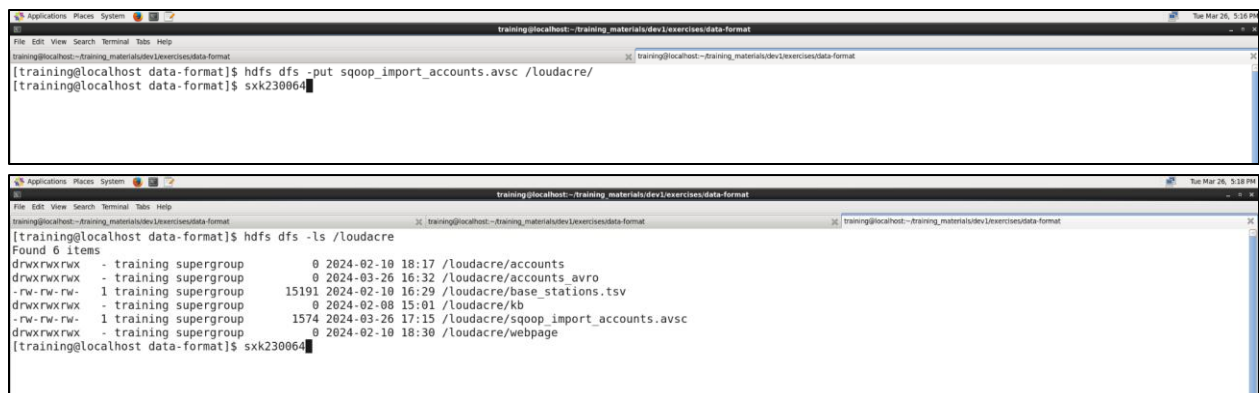
7. Employing the '-put' command, I copied the schema file to the /loudacre directory in HDFS. Utilizing the '-ls' command, I verified the successful transfer.
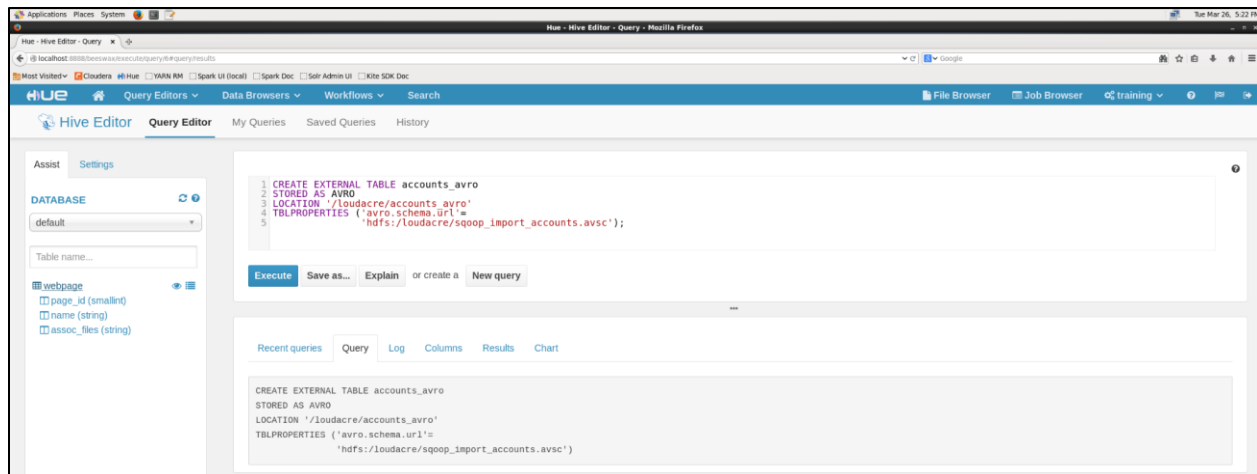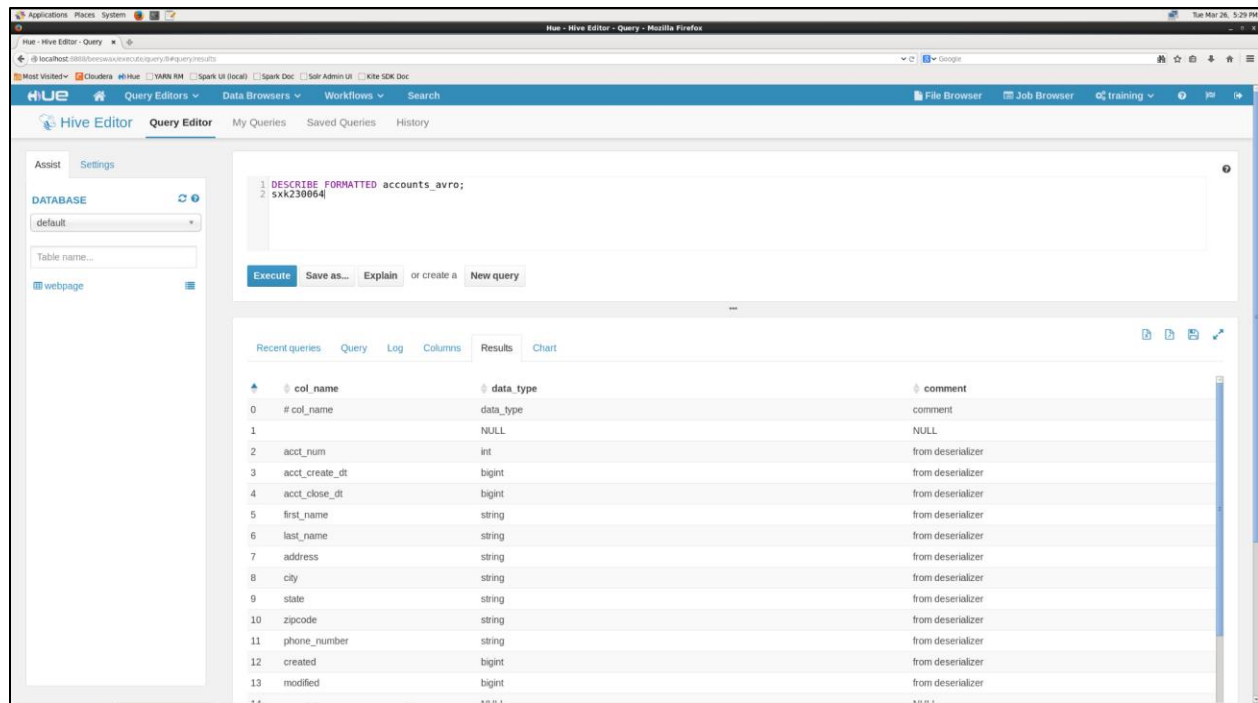


8. Within Hive, I proceeded to create the table "accounts_avro" utilizing the provided Avro schema.

9. After creating the table, I performed a test query using a select statement to ensure its successful created.
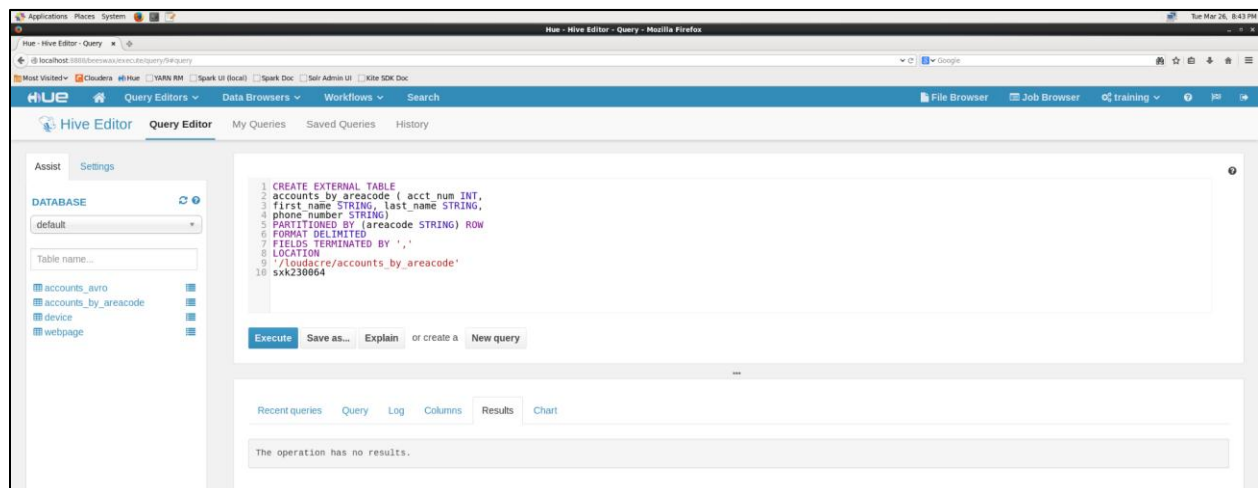


10. I utilized the 'DESCRIBE FORMATTED' command to list the columns and data types of the accounts_avro table, derived from the Avro schema.

## CHAPTER 8 - PARTITION DATA IN IMPALA OR HIVE

1. Initially, I created a new empty table in Hive using the CREATE EXTERNAL TABLE statement.



2. To extract the area code from phone numbers, I executed a query as per the below snippet.

3. Then, I employed the SELECT statement within an INSERT INTO TABLE command to transfer the specified columns into the newly created table. Notably, the process involved dynamic partitioning by area code, as depicted in the following screenshot.



4. I ran a test query to ensure that the table was populated correctly.

5. Using Hue, I confirm that the index structure of the accounts_by_areacode table encompassed partition directories. Additionally, I reviewed the data within the directories to affirm that the partitioning is correct.
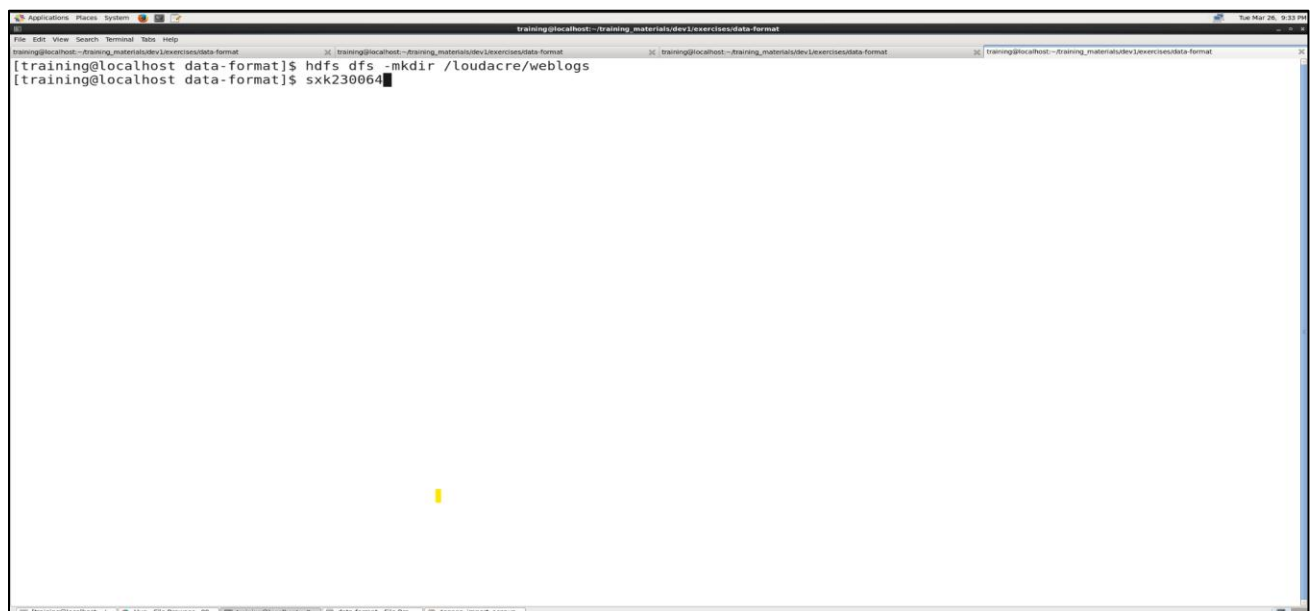
# CHAPTER 9 - COLLECT WEB SERVER LOGS WITH FLUME

6. I created a directory named /loudacre/weblogs in HDFS, intending to store the data files ingested by Flume.
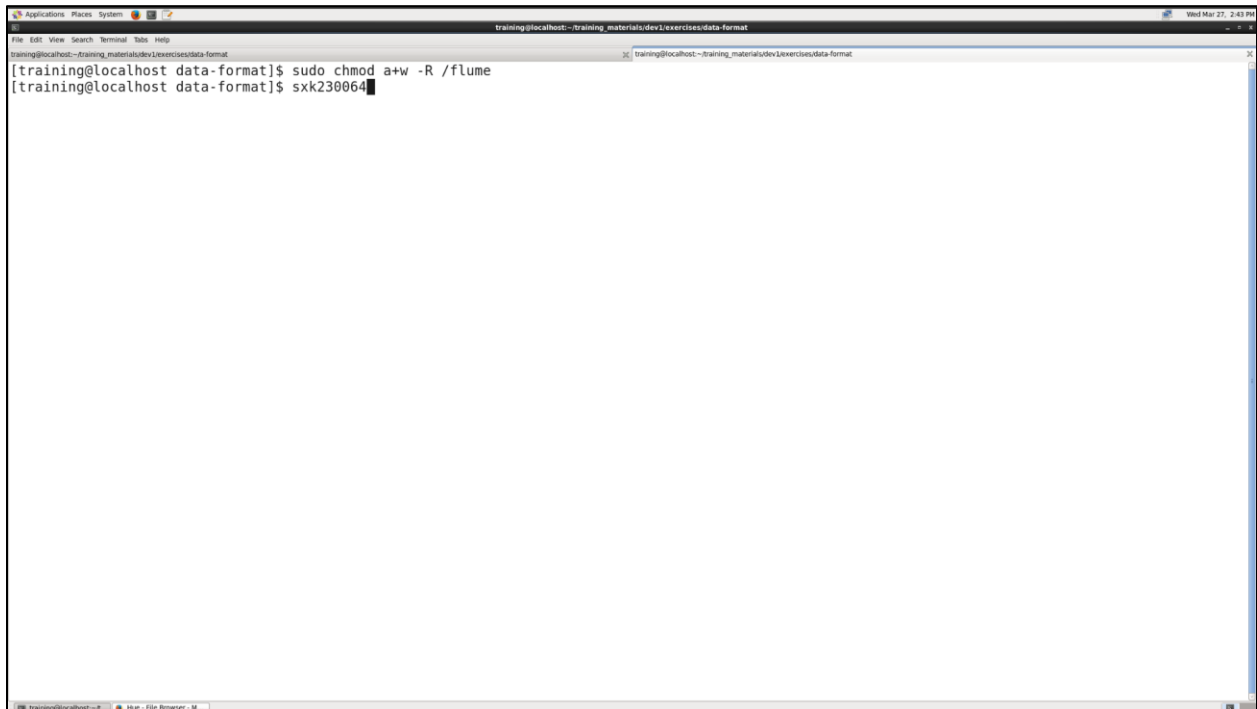
7. Furthermore, I created a spool directory to accommodate the data files that the weblog simulator will generate for Flume ingestion.



8. To facilitate seamless data ingestion, I granted all users permissions to write to the /flume/weblogs_spooldir directory, executing the command $ sudo chmod a+w -R /flume in the command line.

## Configure Flume

9. I created a flume configuration file with the specified properties for sink, source, and channel.



## Run the Agent

10. Following configuration setup, I navigated to the /training_materials/dev1/exercises/flume directory.



11. I launched the Flume agent utilizing the previously created configuration.

12. Towards the end of the lines, it's confirmed that the agent is successfully running with the configuration named src1.



**Simulate Apache web server output**

13. I opened a new terminal window, changed to the exercise directory, and executed the script to place the web log files in the /flume/weblogs_spooldir directory.



14. Upon completion, I terminated the process by clicking CTRL + C. I used the HUE file browser to validate that Flume successfully copied the weblogs into the weblogs directory. Each imported file is tagged with a Unix timestamp corresponding to the time it was imported.