# Data Partitioning

Divyansh Dahiya
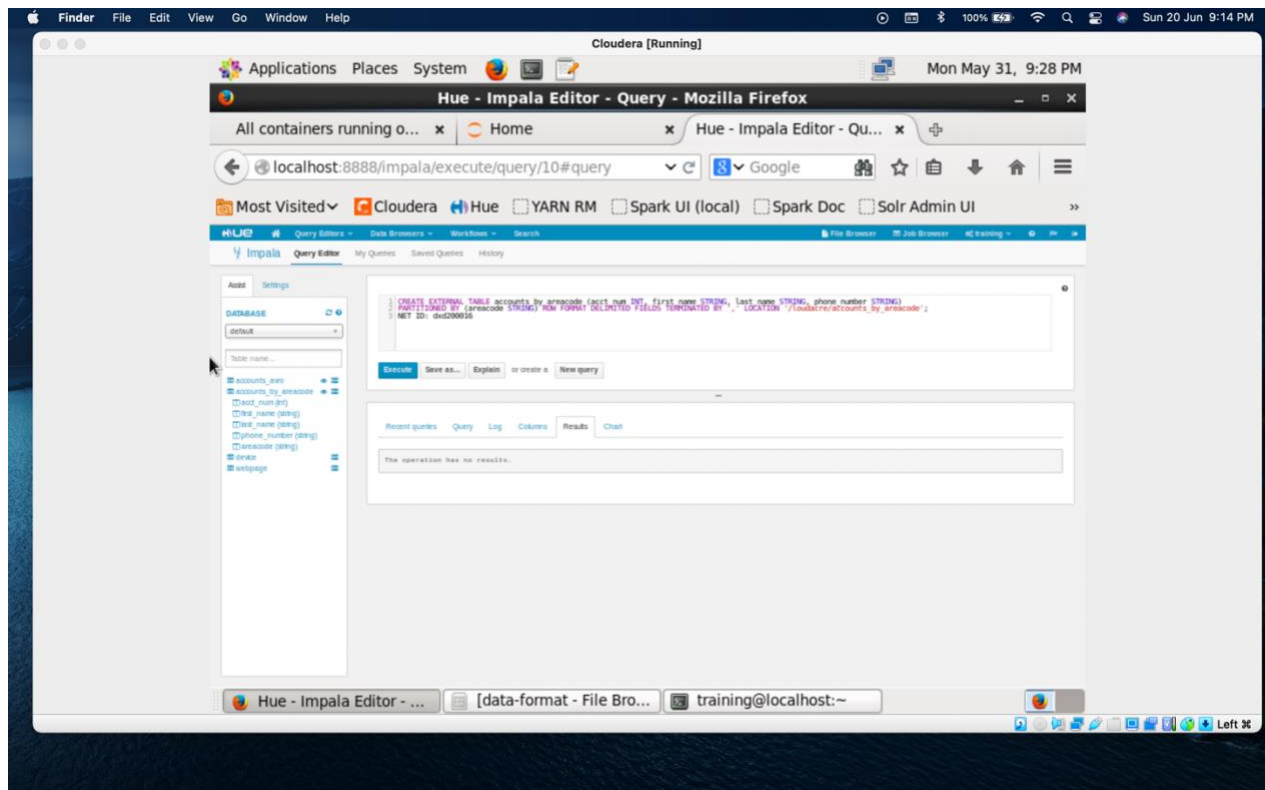
## LAB Chapter 08: PARTITIONS
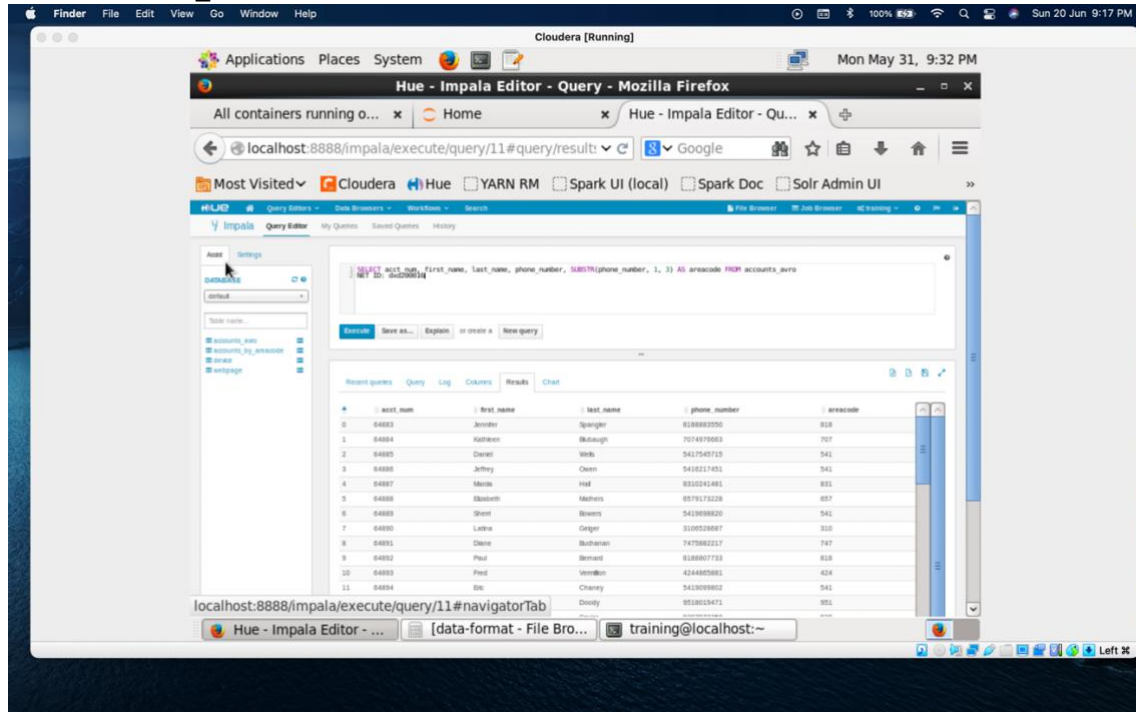
## Partition Data in Impala or Hive

In this lab, I will create and load an Impala table with account data, partitioned by area code.

1. First, I opened Impala Query Editor inside Hue Web Interface and ran the following command to create an empty table accounts_by_areacode
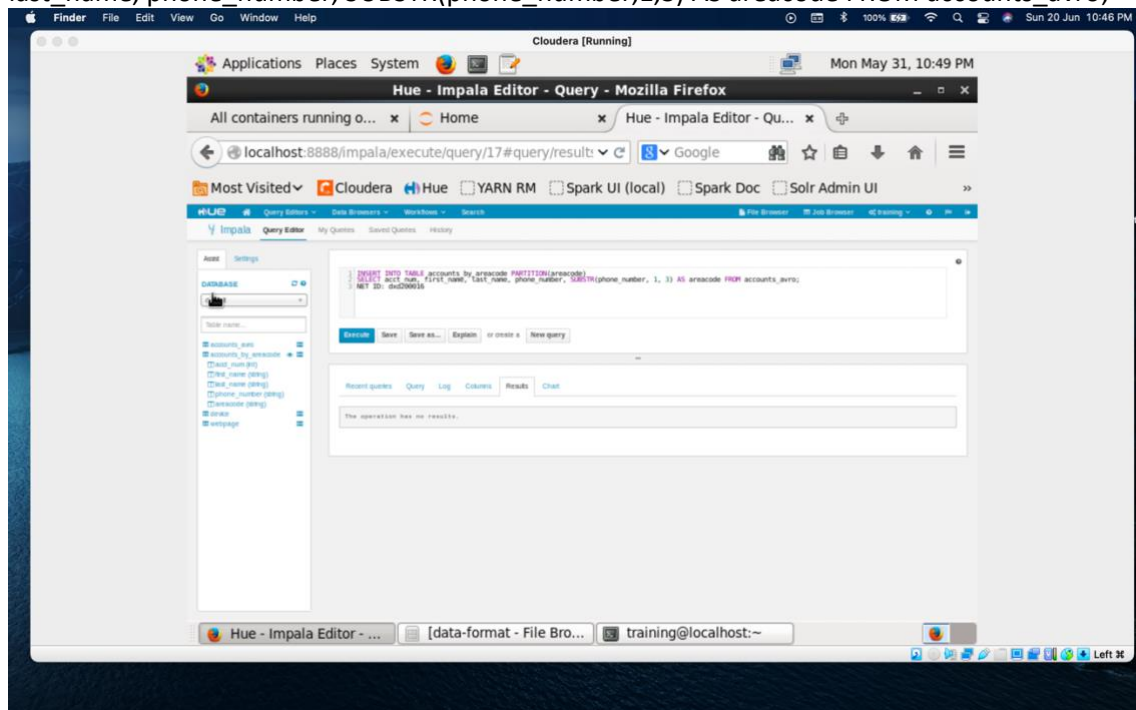   Command: CREATE EXTERNAL TABLE accounts_by_areacode ( acct_num INT, first_name STRING, last_name STRING, phone_number STRING) PARTITIONED BY (areacode STRING) ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
   LOCATION '/loudacre/accounts_by_areacode';

2. I ran the following command to select the specific columns and extracted the area code from the phone number.
Command: SELECT acct_num, first_name, last_name, phone_number, SUBSTR(phone_number,1,3) AS areacode FROM accounts_avro



3. Then, I ran the following command to insert the specific columns from accounts_avro table to accounts_by_areacode table, dynamically partitioned by area code.
Command: INSERT INTO TABLE accounts_by_areacode PARTITION(areacode) SELECT acct_num, first_name, last_name, phone_number, SUBSTR(phone_number,1,3) AS areacode FROM accounts_avro;

4. Then, to verify If the records are inserted successfully, I ran the following command to view the first 10 records from the accounts_by_areacode table.
Command: SELECT * FROM accounts_by_areacode LIMIT 10



5. Then, I ran the following command to view the current partitions in the accounts_by_areacode table.
Command: SHOW PARTITIONS accounts_by_areacode

6. Finally, I ran the following command to view all the columns and their data types of the accounts_by_areacode table.
   Command: DESCRIBE accounts_by_areacode