## Big Data Hadoop 'Session 8 HIVE BASICS'

DATASET

| | | |
|---|---|---|
| 10-01-1990 | 123112 | 10 |
| 14-02-1991 | 283901 | 11 |
| 10-03-1990 | 381920 | 15 |
| 10-01-1991 | 302918 | 22 |
| 12-02-1990 | 384902 | 9 |
| 10-01-1991 | 123112 | 11 |
| 14-02-1990 | 283901 | 12 |
| 10-03-1991 | 381920 | 16 |
| 10-01-1990 | 302918 | 23 |
| 12-02-1991 | 384902 | 10 |
| 10-01-1993 | 123112 | 11 |
| 14-02-1994 | 283901 | 12 |
| 10-03-1993 | 381920 | 16 |
| 10-01-1994 | 302918 | 23 |
| 12-02-1991 | 384902 | 10 |
| 10-01-1991 | 123112 | 11 |
| 14-02-1990 | 283901 | 12 |
| 10-03-1991 | 381920 | 16 |
| 10-01-1990 | 302918 | 23 |
| 12-02-1991 | 384902 | 10 |

**Task 1**

- Create a database named 'custom'



- Create a table named temperature_data inside custom having below fields:

    1. date (mm-dd-yyyy) format

    2. zip code

    3. temperature
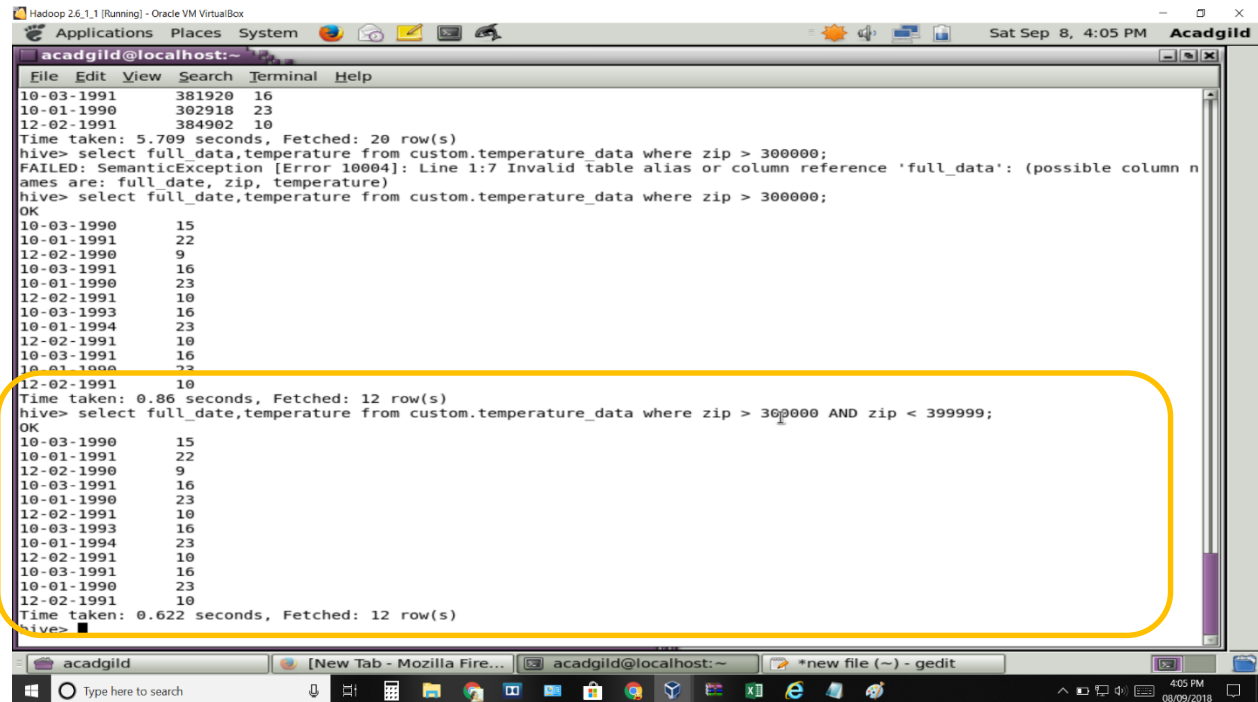
    The table will be loaded from comma-delimited file.

- Load the dataset.txt (which is ',' delimited) in the table.

**Task 2**

● Fetch date and temperature from temperature_data where zip code is greater than

300000 and less than 399999.

- Calculate maximum temperature corresponding to every year from temperature_data table.



● Calculate maximum temperature from temperature_data table corresponding to those

years which have at least 2 entries in the table.

● Create a view on the top of last query, name it temperature_data_vw.



● Export contents from temperature_data_vw to a file in local file system, such that each

file is '|' delimited.

Screenshot 1 — Terminal (acadgild@localhost:~)

```
OK
Time taken: 1.062 seconds
hive> CREATE VIEW temperature_data_vw AS
    > select count(year) as cntof_year_greaterthan2,MAX(temperature) from (select substring(full_date,7) as year,temperat
ure from temperature_data)t2 group by year having cntof_year_greaterthan2>2 ;
OK
Time taken: 0.475 seconds
hive> INSERT OVERWRITE LOCAL DIRECTORY '/home/acadgild/hive_output_view.txt'
    > ROW FORMAT DELIMITED
    > FIELDS TERMINATED BY '|'
    > SELECT * FROM temperature_data_vw;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different e
xecution engine (i.e. spark, tez) or using Hive 1.X releases.
Query ID = acadgild_20180911062723_970a936d-02bb-4721-b37e-a5a8e2f6cb6a
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1536626025229_0003, Tracking URL = http://localhost:8088/proxy/application_1536626025229_0003/
Kill Command = /home/acadgild/install/hadoop/hadoop-2.6.5/bin/hadoop job  -kill job_1536626025229_0003
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2018-09-11 06:27:39,757 Stage-1 map = 0%,  reduce = 0%
2018-09-11 06:27:52,687 Stage-1 map = 100%,  reduce = 0%, Cumulative CPU 3.81 sec
2018-09-11 06:28:07,906 Stage-1 map = 100%,  reduce = 100%, Cumulative CPU 8.05 sec
MapReduce Total cumulative CPU time: 8 seconds 50 msec
Ended Job = job_1536626025229_0003
Moving data to local directory /home/acadgild/hive_output_view.txt
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1  Reduce: 1   Cumulative CPU: 8.05 sec   HDFS Read: 9912 HDFS Write: 10 SUCCESS
Total MapReduce CPU Time Spent: 8 seconds 50 msec
OK
Time taken: 45.442 seconds
hive>
```



Screenshot 2 — File browser and gedit (000000_0)

```
7|23
9|22
```