# BIG DATA HADOOP AND SPARK DEVLOPMENT

# ASSIGNMENT 8

Table of Contents:

# BIG DATA HADOOPAND SPARK DEVELOPMENT

## 1. Introduction

In this assignment, the given tasks are performed and Output of the tasks are recorded in the form of Screenshots.

## 2. Objective

This Assignment consolidates the deeper understanding of the Session – 8 Introduction to the HIVE

## 3. Problem Statement

- ## 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.

4. Expected Output
   - Task 1
   - **Create a database named 'custom'**.
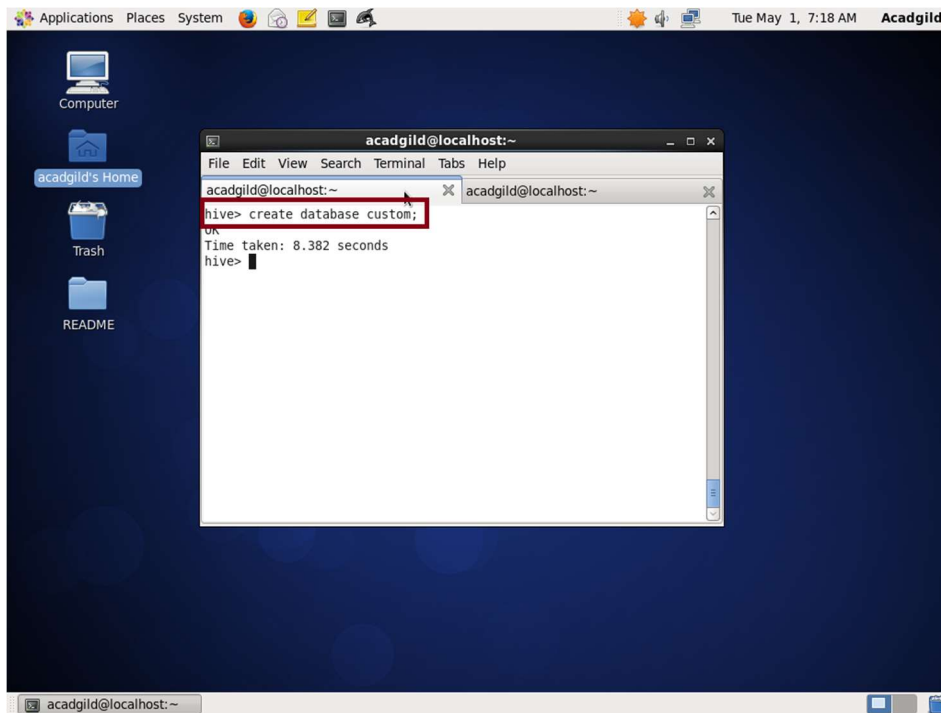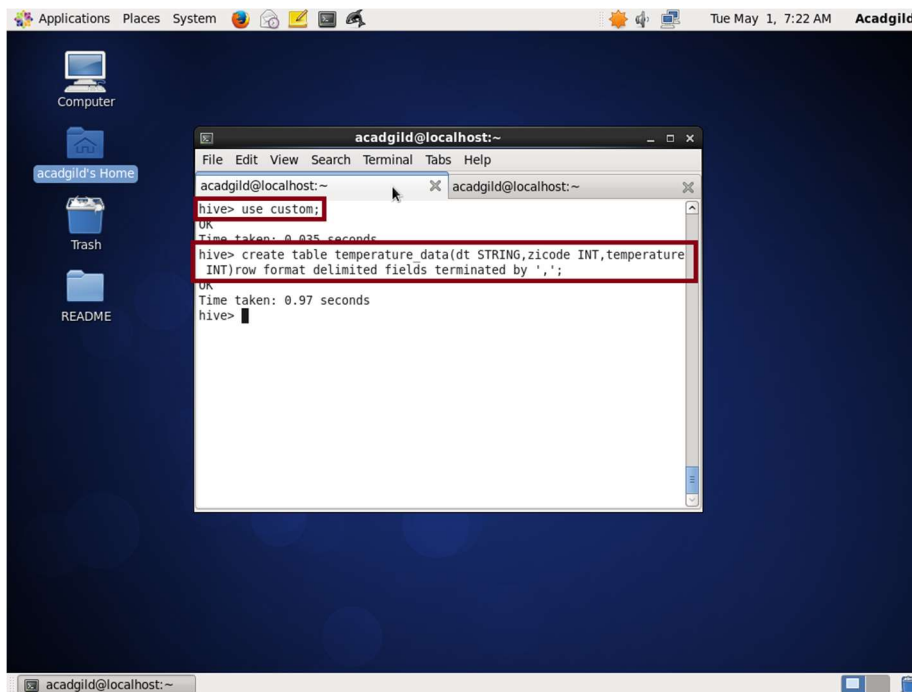
# Create Database Statement

Create Database is a statement used to create a database in Hive. A database in Hive is a namespace or a collection of tables.

Syntax:

```
CREATE DATABASE|SCHEMA [IF NOT EXISTS] <database name>
```

Here the following command is used to create the database custom

**create database 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.**

## Create Table Statement

Create Table is a statement used to create a table in Hive.

Syntax:

```
CREATE [TEMPORARY] [EXTERNAL] TABLE [IF NOT EXISTS] [db_name.] table_name
```

The following command is used to create the table temperature_data with date,zipcode,temperature.

**create table temperature_data(dt STRING, zipcode INT, temperature INT) row format delimited fields terminated by ',';**

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

## Load Data Statement

Generally, after creating a table in SQL, we can insert data using the Insert statement. But in Hive, we can insert data using the LOAD DATA statement.

While inserting data into Hive, it is better to use LOAD DATA to store bulk records. There are two ways to load data: one is from local file system and second is from Hadoop file system.

**Syntax:**

```
LOAD DATA [LOCAL] INPATH 'filepath' [OVERWRITE] INTO TABLE tablename
[PARTITION (partcol1=val1, partcol2=val2 ...)]
```

- LOCAL is identifier to specify the local path. It is optional.
- OVERWRITE is optional to overwrite the data in the table.
- PARTITION is optional.

By using the following command, the data for the table temperature_data is loaded.

**Load data local inpath '/home/acadgild/text files/dataset.txt' into table temperature_data;**

SELECT statement is used to retrieve the data from a table.

By the below command the data in the table are shown in the console by select command.

**select * from temperature_data;**

- Task 2
- **Fetch date and temperature from temperature_data where zip code is greater than 300000 and less than 399999.**

By using the following command, we can fetch the date and temperature from the table temperature_data where zip code is greater than 300000 and less than 399999.

## SELECT statement with WHERE clause

**SELECT** statement is used to retrieve the data from a table. **WHERE** clause works similar to a condition. It filters the data using the condition and gives you a finite result.

Syntax:

```
SELECT [ALL | DISTINCT] select_expr, select_expr, ...

FROM table_reference

[WHERE where_condition]

[GROUP BY col_list]

[HAVING having_condition]

[CLUSTER BY col_list | [DISTRIBUTE BY col_list] [SORT BY col_list]]

[LIMIT number];
```

**select dt,temperature from temperature_data where zipcode < 399999 and zipcode > 300000;**

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

By using the following command, Maximum temperature is displayed by using Max function and year is displayed by using Substring function as dt is saved as string format.

**Select SUBSTRING(dt,7,4), MAX(temperature) FROM custom.temperature_data GROUP BY SUBSTRING(dt,7,4);**

- **Calculate maximum temperature from temperature_data table corresponding to those years which have at least 2 entries in the table.**

By using following command, the maximum temperature is calculated from temperature_data table corresponding to those years which have at least 2 entries in the table.

**select dt, MAX(t.temperature)as temperature from (select SUBSTRING(dt,7,4)dt,temperature from temperature_data)t GROUP BY HAVING count (t.dt)>=2;**

- **Create a view on the top of last query, name it as temperature_data_vw.**

## Creating a View

You can create a view at the time of executing a SELECT statement. The syntax is as follows:

```
CREATE VIEW [IF NOT EXISTS] view_name [(column_name [COMMENT column_comment], ...) ]

[COMMENT table_comment]

AS SELECT ...
```
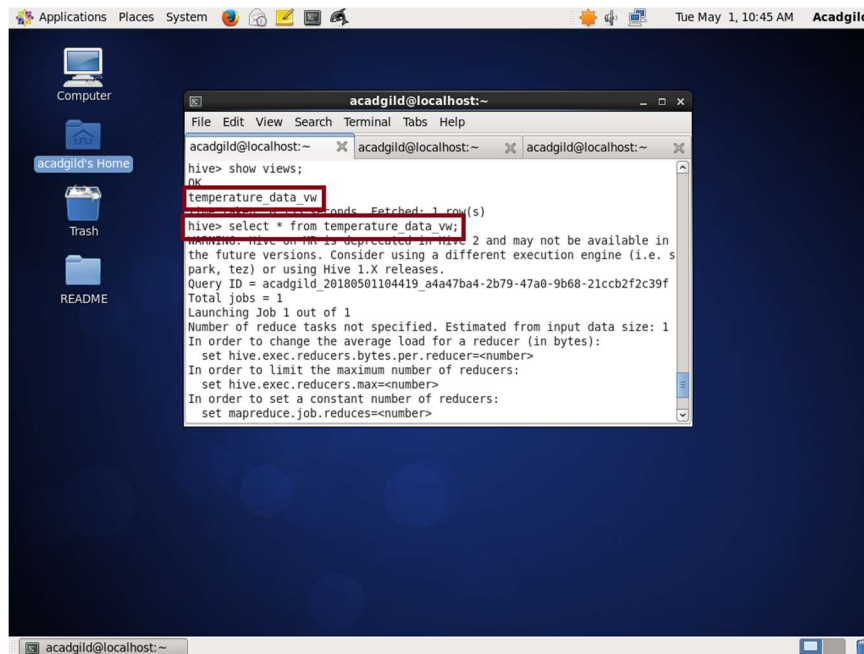
**By using the following command, we can create a view as temperature_data_vw on last query**

**CREATE VIEW temperature_data_vw as select dt,MAX(t.temperature)as temperature from (select SUBSTRING(dt,7,4),temperature from temperature_data)t GROUP BY dt HAVING count(t.dt)>=2;**

**By using the show views following command, the views are listed which are created under the custom database**

**show views;**



**To check the saved view, we are using select command to list the data in the view.**

**select * from temperature_data_vw;**

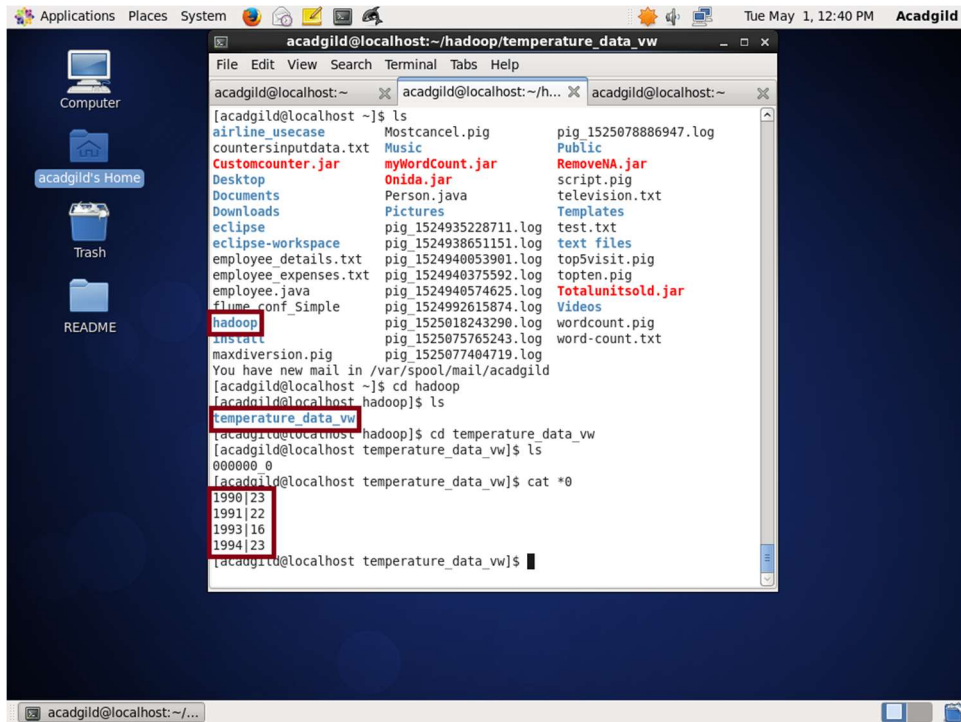- **Export contents from temperature_data_vw to a file in local file system, such that each file is '|' delimited.**

By using the following command, we can create a new directory temperature_data_vw and we can create a text file and save the data with delimited fields terminated by '|'

**insert overwrite local directory '/home/acadgild/Hadoop/temperature_data_vw' row format delimited fields terminated by '|' select * from temperature_data_vw;**

**Contents are exported from temperature_data_vw to a file in local file system, such that each file is '|' delimited.**

**The location of the file is /home/acadgild/Hadoop/temperature_data_vw**



By using,

 Ls command, the files are listed.

Cat command, the contents in the file are displayed in the console.